# Application for Question Answering on Job Postings Dataset

## 1. Overview
This project aims to develop a Question Answering (QA) platform by implementing a Decoder-Based Transformer using the T5 model. The objective is to create an interactive application using Streamlit, leveraging advanced AI techniques such as fine-tuning and deployment within a Colab environment using ngrok.

---

## 2. Project Goals and Scope
- **Use Case**: The application is designed as a QA system that provides contextual answers to user queries about job postings.
- **Main Objectives**:
  - Develop a decoder-based transformer model using T5.
  - Fine-tune the model on a custom dataset (job_postings.csv).
  - Deploy the model as an interactive web application using Streamlit and ngrok for public access.

---

## 3. Process of Implementation
### Step 1: Preparing the Dataset
- **Dataset Used**: job_postings.csv, containing job titles, companies, locations, job levels, and job types.
- **Data Preprocessing**:
  - Combined relevant columns into a single context for each job posting:

```python
# Combine relevant columns into a single context for each job posting
df['context'] = (
    "Job Title: " + df['job_title'].fillna('N/A') + "\n" +
    "Company: " + df['company'].fillna('N/A') + "\n" +
    "Location: " + df['job_location'].fillna('N/A') + "\n" +
    "Job Level: " + df['job_level'].fillna('N/A') + "\n" +
    "Job Type: " + df['job_type'].fillna('N/A')
)
```

  - Generated Question-Answer pairs for training:
    - What is the job title?
    - Which company is hiring?
    - Where is the job located?
    - What is the job level?
    - Is it a remote or onsite job?

### Step 2: Model Training and Selection
- **Base Model**: T5-small (Text-to-Text Transfer Transformer).
- **Model Setup**:
  - Loaded pre-trained T5 tokenizer and model:

```
# Load the T5 tokenizer and model
tokenizer = T5Tokenizer.from_pretrained("t5-small")
model = T5ForConditionalGeneration.from_pretrained("t5-small").to(device)
```

```
/usr/local/lib/python3.11/dist-packages/huggingface_hub/utils/_auth.py:94: UserWarning:
The secret `HF_TOKEN` does not exist in your Colab secrets.
To authenticate with the Hugging Face Hub, create a token in your settings tab (https://huggingface.co/settings/tokens), set it as secret in your Google Colab and restart your session.
You will be able to reuse this secret in all of your notebooks.
Please note that authentication is recommended but still optional to access public models or datasets.
  warnings.warn(
tokenizer_config.json: 100%          2.32k/2.32k [00:00<00:00, 68.4kB/s]
spiece.model: 100%                   792k/792k [00:00<00:00, 3.22MB/s]
tokenizer.json: 100%                 1.39M/1.39M [00:00<00:00, 16.0MB/s]
You are using the default legacy behaviour of the <class 'transformers.models.t5.tokenization_t5.T5Tokenizer'>. This is expected, and simply means that the `legacy` (previous) behavior wil
config.json: 100%                    1.21k/1.21k [00:00<00:00, 44.5kB/s]
model.safetensors: 100%              242M/242M [00:01<00:00, 163MB/s]
generation_config.json: 100%         147/147 [00:00<00:00, 3.75kB/s]
```

- o Prepared input and output sequences using the context and questions.
- o Custom Dataset Class was created for tokenized inputs:

```
# Create a custom dataset class
class QADataset(torch.utils.data.Dataset):
    def __init__(self, encodings, targets):
        self.encodings = encodings
        self.targets = targets

    def __len__(self):
        return len(self.encodings.input_ids)

    def __getitem__(self, idx):
        item = {key: torch.tensor(val[idx]) for key, val in self.encodings.items()}
        item['labels'] = torch.tensor(self.targets.input_ids[idx])
        return item
```

- o Fine-tuned the model using Hugging Face Trainer API:



## Step 3: Development of the Application
- **Web Framework Used**: Streamlit
- **Key Features**:
  - o Interactive UI to select a job posting and ask questions.

- o Generated answers displayed in real-time using T5.
- o Streamlit components used:
  - ▪ selectbox: Select a job posting.
  - ▪ text_input: Input for user questions.
  - ▪ success: Display generated answers.

```python
# Inference: Example questions
example_context = df['context'].iloc[0]
example_question = "What is the job title?"

input_text = f"question: {example_question} context: {example_context}"
input_ids = tokenizer(input_text, return_tensors="pt").input_ids.to(device)
output_ids = model.generate(input_ids)
answer = tokenizer.decode(output_ids[0], skip_special_tokens=True)

print("\nExample Question:", example_question)
print("Generated Answer:", answer)
```

## Step 4: Testing and Deployment
- • **Colab Deployment**:
  - o Streamlit was used within Google Colab.
  - o **ngrok** was configured for external access:
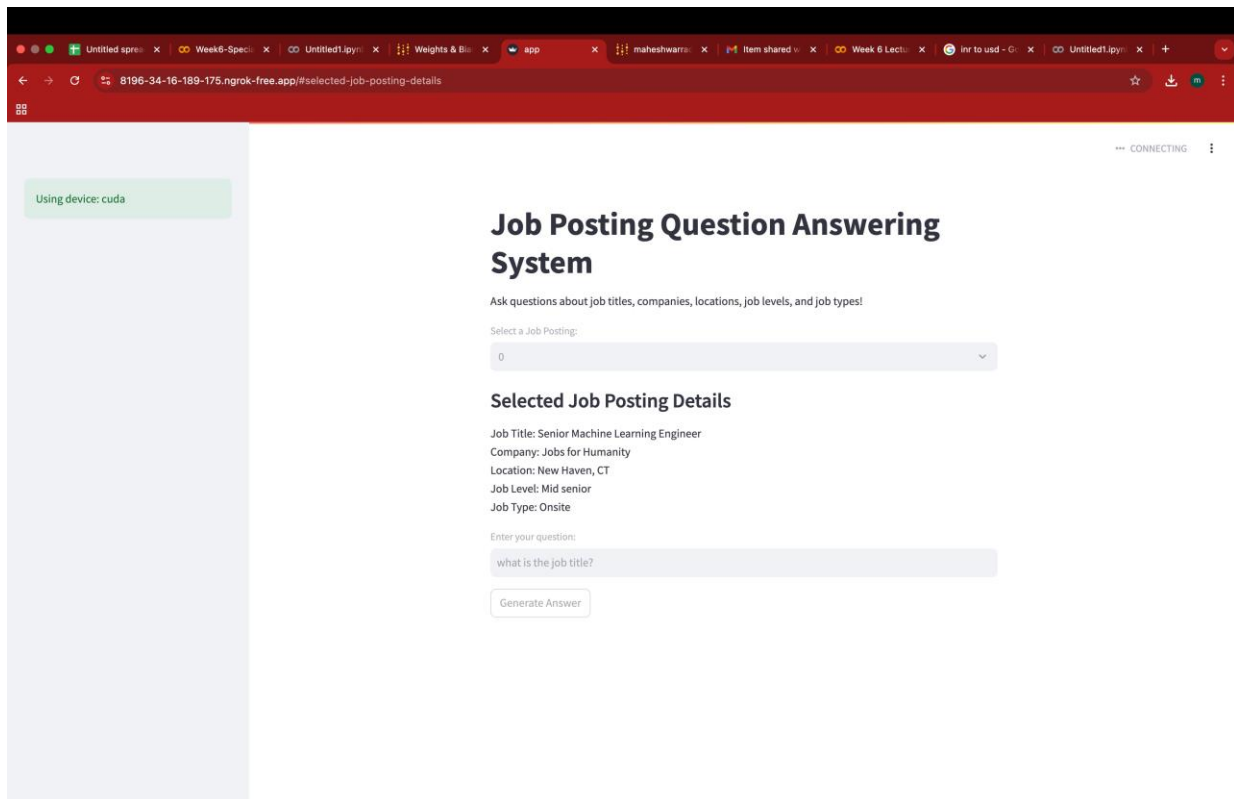


## 4. Challenges and Resolutions
## Key Challenges:
- • **KeyError during Dataset Processing**:
  - o Issue: Column mismatch in the dataset.
  - o Resolution: Adjusted column names to match the expected format.
- • **CUDA Out of Memory Error**:

- o Issue: GPU memory limitations.
- o Resolution: Switched to CPU-based inference when GPU was unavailable.
- **Streamlit Installation Problem**:
  - o Issue: ModuleNotFoundError for Streamlit.
  - o Resolution: Installed Streamlit manually using:

---

## 5. Conclusion and Results

- The project successfully implemented a Question Answering System using a Decoder-Based Transformer (T5).
- Achievements:
  - o Fine-tuned the T5 model for contextual QA on job postings.
  - o Developed an interactive Streamlit application for user interaction.
  - o Successfully deployed within Colab using ngrok.
- Outcome:
  - o Demonstrated the power of transformer-based models in contextual question answering.
  - o Showcased the deployment potential using Streamlit and ngrok.

**7. Recommendations for Future Work**
1. **Model Enhancement**:
   - o Experiment with larger T5 models (e.g., t5-base, t5-large) for improved accuracy.
   - o Fine-tune with more diverse datasets to enhance generalization.
2. **Feature Improvements**:
   - o Add filters for job type, location, and company in the Streamlit UI.
   - o Include a feedback system for user satisfaction with generated answers.
3. **Deployment Enhancements**:
   - o Deploy on Streamlit Cloud or Heroku for a more stable public URL.
   - o Implement CI/CD pipelines for continuous integration and deployment.

**8. Conclusion**

This project effectively demonstrates the application of Decoder-Based Transformers for building an intelligent QA system. It showcases the integration of T5 models with Streamlit and ngrok for real-time deployment. The system can be further expanded for other domains, including customer support, educational tutoring, and career guidance systems.