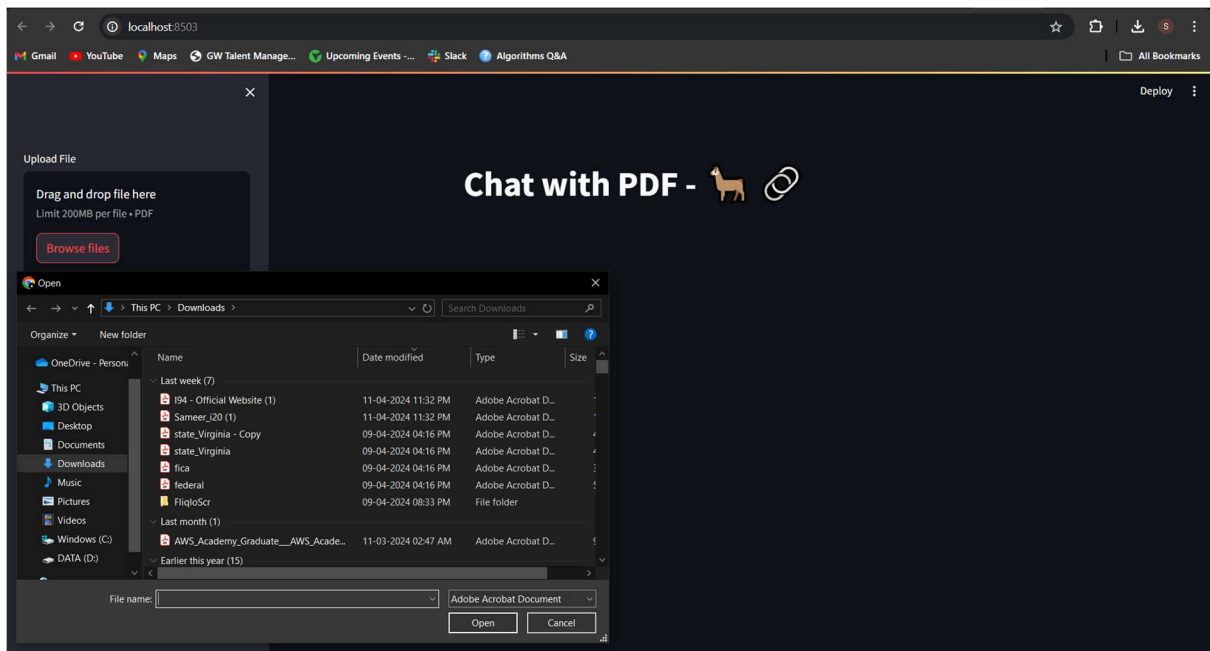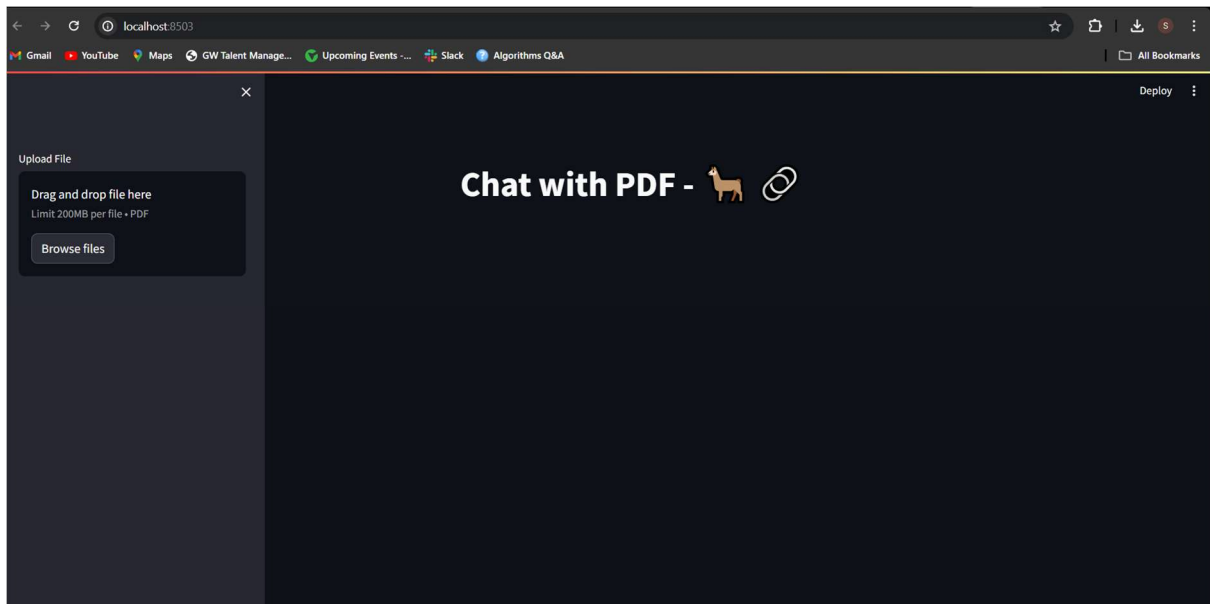**Name: Sameer Shaik**
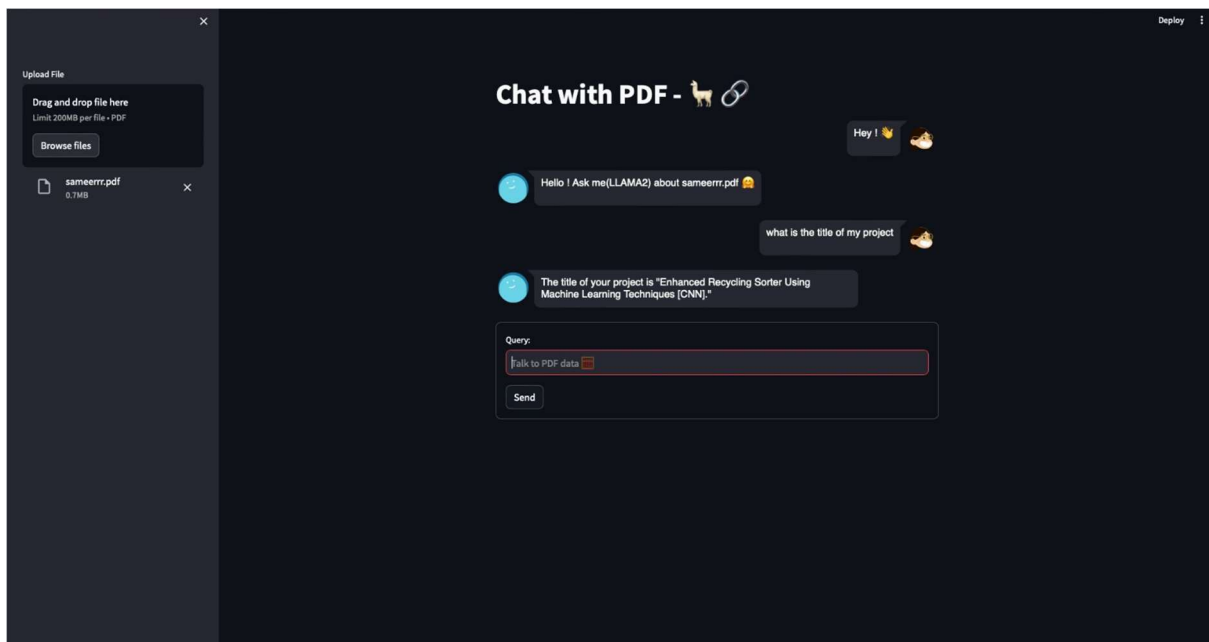
**GWID: G49843839**

**Sameer.shaik@gwmail.gwu.edu**

**Problem 2:**

```python
from langchain.llms import LlamaCpp
from langchain.callbacks.manager import CallbackManager
from langchain.callbacks.streaming_stdout import StreamingStdOutCallbackHandler

model_path = 'C:/Users/Sameer/Downloads/llama-2-7b-chat.Q4_K_M.gguf'


class LoadLlm:
    @staticmethod
    def load_llm():
        callback_manager = CallbackManager([StreamingStdOutCallbackHandler()])
        # Prepare the LLM

        llm = LlamaCpp(
            model_path=model_path,
            n_gpu_layers=40,
            n_batch=512,
            n_ctx=2048,
            f16_kv=True,  # MUST set to True, otherwise you will run into problem after a couple of calls
            callback_manager=callback_manager,
            verbose=True,
        )

        return llm
```

```
PS C:\Users\Sameer\OneDrive\Desktop\ASS> pip install llama-cpp-python
Collecting llama-cpp-python
  Downloading llama_cpp_python-0.2.63.tar.gz (37.5 MB)
                37.5/37.5 MB 15.6 MB/s eta 0:00:00
  Installing build dependencies ... done
  Getting requirements to build wheel ... done
  Installing backend dependencies ... done
  Preparing metadata (pyproject.toml) ... done
```

Deploy

Upload File

Drag and drop file here
Limit 200MB per file • PDF

Browse files

# Chat with PDF - 🦙 🔗

Deploy

Upload File

Drag and drop file here
Limit 200MB per file • PDF

Browse files

# Chat with PDF - 🦙 🔗

Open

This PC > Downloads

Search Downloads

Organize    New folder

OneDrive - Person

This PC
  3D Objects
  Desktop
  Documents
  Downloads
  Music
  Pictures
  Videos
  Windows (C:)
  DATA (D:)

| Name | Date modified | Type | Size |
|---|---|---|---|
| Last week (7) | | | |
| I94 - Official Website (1) | 11-04-2024 11:32 PM | Adobe Acrobat D... | |
| Sameer_i20 (1) | 11-04-2024 11:32 PM | Adobe Acrobat D... | |
| state_Virginia - Copy | 09-04-2024 04:16 PM | Adobe Acrobat D... | |
| state_Virginia | 09-04-2024 04:16 PM | Adobe Acrobat D... | |
| fica | 09-04-2024 04:16 PM | Adobe Acrobat D... | |
| federal | 09-04-2024 04:16 PM | Adobe Acrobat D... | |
| FliqloScr | 09-04-2024 08:33 PM | File folder | |
| Last month (1) | | | |
| AWS_Academy_Graduate__AWS_Acade... | 11-03-2024 02:47 AM | Adobe Acrobat D... | |
| Earlier this year (15) | | | |

File name: |                    | Adobe Acrobat Document

Open    Cancel

Prompt

What is the title of my project :

Response

The title of your project is " Enhanced Recycling Sorter Using Machine Learning Techniques [CNN]"

Metrics :



Prompt

What is the title of my project :

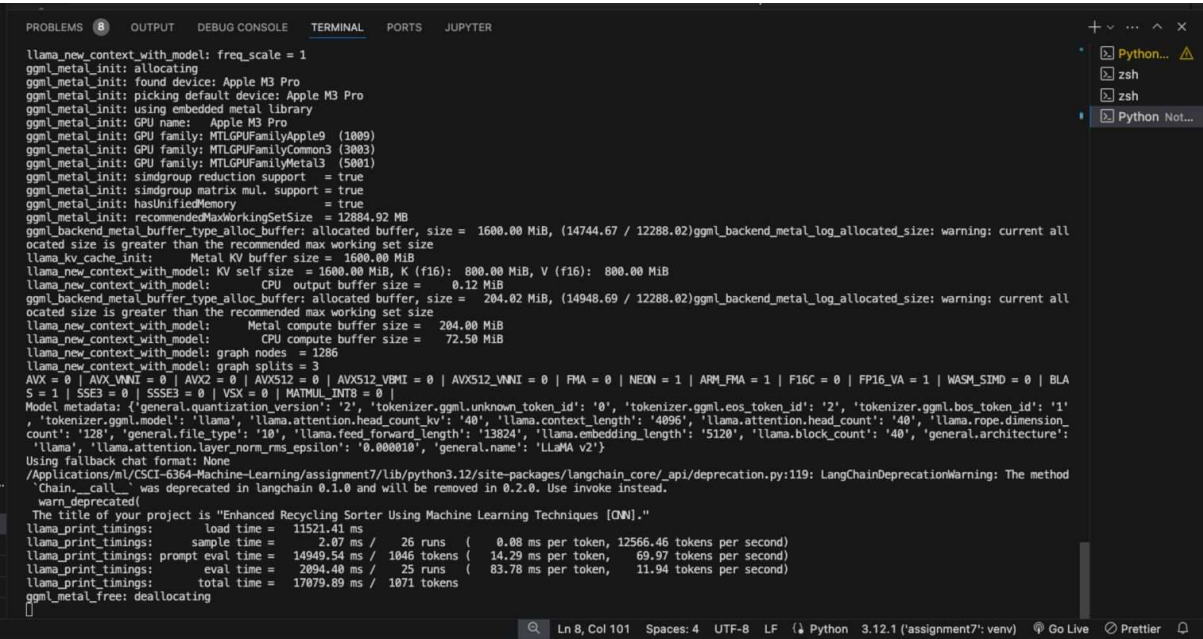Response

The title of your project is " Enhanced Recycling Sorter Using Machine Learning Techniques [CNN]"

Metrics :



Metrics



----------------------------------------------------------------------------------------------------------------------------

----------------------------------------------------------------------------------------------------------------------------

# Problem 3:

Prompt - Write about Neural Networks

**##### top_k=40, top_p=0.95, temperature=0.8, max_tokens=25**

Neural Networks and Their Applications

Introduction

Artificial neural networks have been a rapidly evolving field of research in recent years, with applications in various industries such as computer vision, natural language processing, speech recognition, and predictive analytics. This article will provide an overview of neural networks, their architecture, types, and applications.

What are Neural Networks?

Neural networks are a computational model inspired by the structure and function of the human brain. They consist of interconnected nodes or neurons that process information and learn from data. Each neuron receives input from other neurons or external sources, performs computations on that input, and then sends the output to other neurons or to the final output layer.

The architecture of a neural network typically includes three layers: input, hidden, and output. The input layer takes in the raw data, the hidden layer performs complex computations and feature extraction, and the output layer generates the final predictions or decisions. The connections between neurons are adjusted during training to minimize the error between the predicted output and the actual output.

----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

Prompt - Write about Neural Networks

**#### top_k=70, top_p=0.55, temperature=0.6, max_tokens=500**

Neural Networks and Their Applications

Introduction

Neural networks have revolutionized the field of artificial intelligence and machine learning. These networks are modeled after the human brain and are capable of learning and making decisions based on data. In this article, we will explore the basics of neural networks and their applications in various fields.

Basics of Neural Networks

A neural network consists of interconnected nodes or neurons that process information. Each neuron receives input from other neurons or external sources, performs a computation on that input, and then sends the output to other neurons or to the output layer. The network learns by adjusting the weights and biases of the connections between neurons based on the data it is trained on.

There are several types of neural networks, including:

1. Feedforward Networks: These networks process information in only one direction, from input to output, without any feedback loops.

2. Recurrent Networks: These networks have feedback connections, allowing them to process information in a loop. This allows the network to keep track of information over time.

3. Convolutional Networks: These networks are designed to process data with grid-like topology, such as images.

Applications of Neural Networks

Neural networks have numerous applications in various fields, including:

1. Image Recognition: Neural networks can be trained to recognize objects in images, such as faces, animals, or vehicles.

2. Natural Language Processing: Neural networks can be used to process and analyze natural language data, such as text or speech.

3. Predictive Modeling: Neural networks can be used to predict future events or outcomes based on past data.

4. Robotics: Neural networks can be used to control robots and enable them to learn from experience.

5. Healthcare: Neural networks can be used in medical imaging, disease diagnosis, and drug discovery.

6. Finance: Neural networks can be used for fraud detection, credit risk assessment, and portfolio optimization.

7. Recommendation Systems: Neural networks can be used to recommend products or services based on user behavior and preferences.

8. Autonomous Vehicles: Neural networks can be used to control autonomous vehicles and