# Lecture 5:
# LLM Application Design (Cont'd)

**SPRING 2025**

**MOHAMED FARAG**
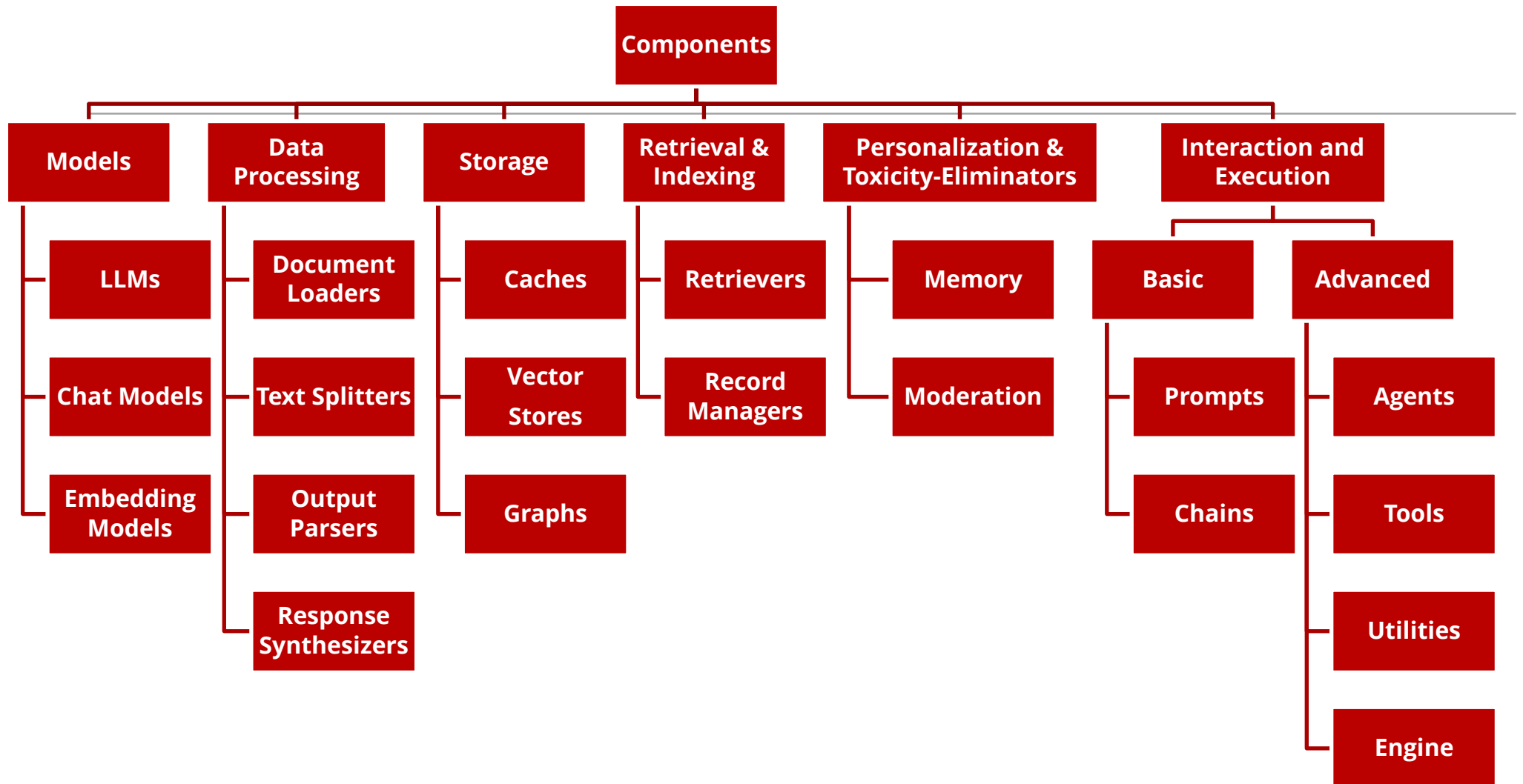
**FARAG@CMU.EDU**

# Agenda

- Enable Vertex AI for Document Processing
- Back to Basics: LLM Chain Design

  - Moderation

  - Output Parsers

  - If-Else Logic

- Multi-agent LLM Applications

- Connect Flowise AI to Local Applications

- Deploy Flowise AI Application

- FlowiseAI Limitations

# Reminder: FlowiseAI Components

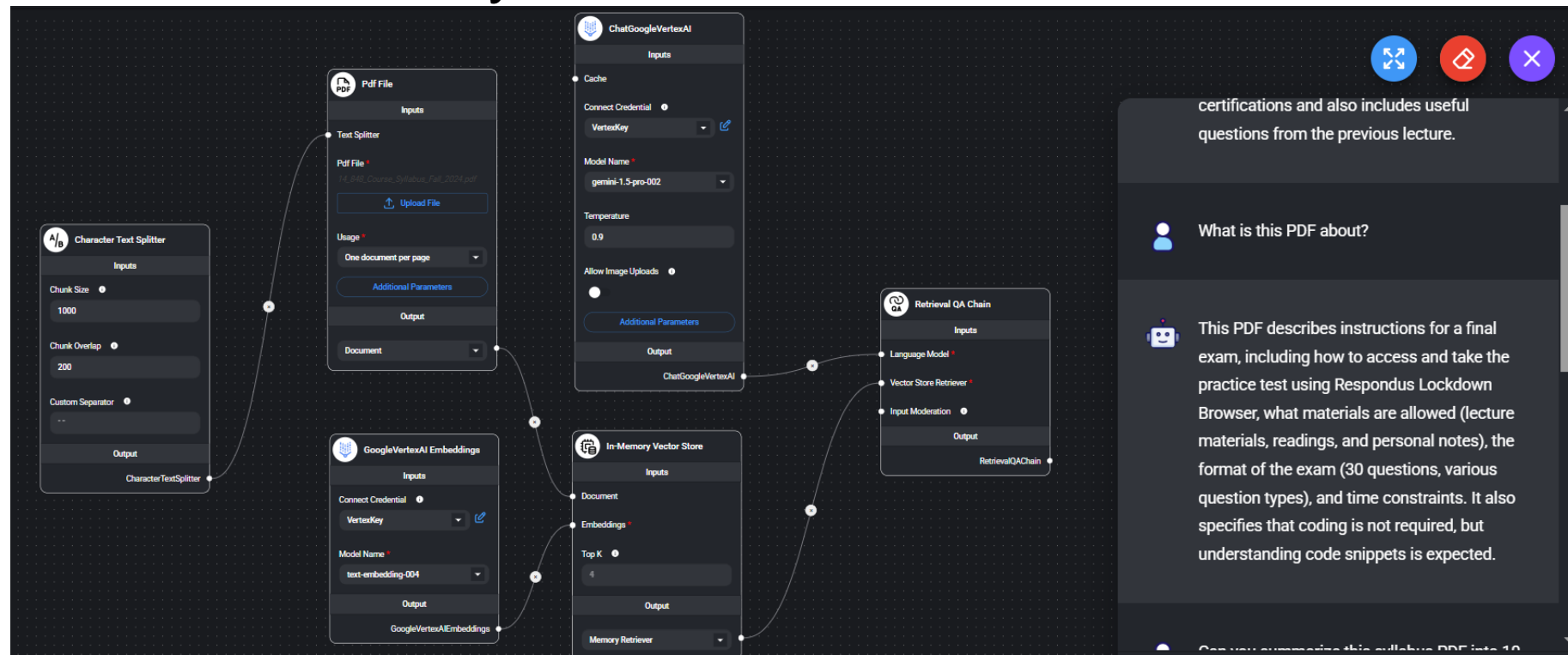# Enable Vertex AI for Document Processing

- Enable Cloud Documents AI API

- Use Gemini-models only.

# FlowiseAI Components: Prompts

# FlowiseAI Components: Chains

Combine LLMs and prompt template to enhance multistep workflows.

# Steps to Create LLM Chain

- Identify your LLM Model.

- Create a detailed prompt for your application.

  - A good prompt should provide information on the output format, structure, and any other useful details. LLM models have different output expectations.

- Assess your need for output parsers

- Add input moderation

- Debug FlowiseAI issues as needed.
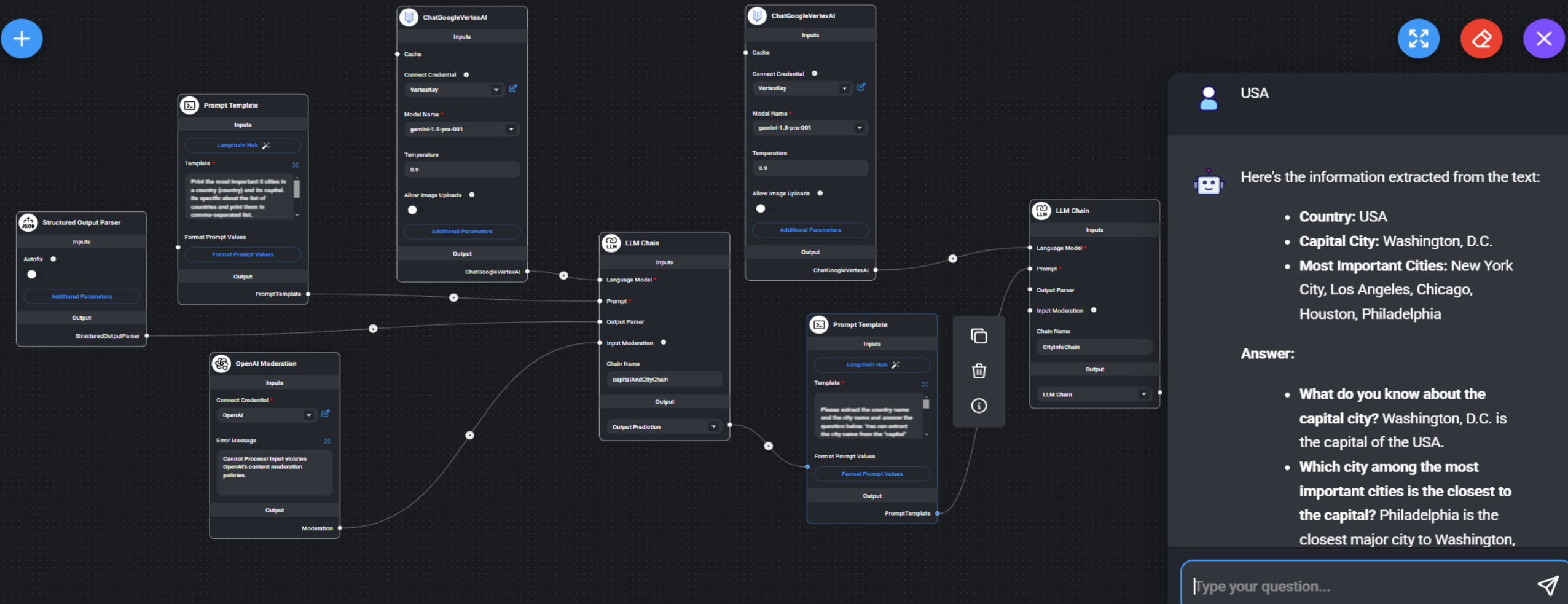
Carnegie Mellon University

# Chain Exercise

- Create an LLM Chain that identifies the capital and the most important 5 cities in a specific country.

- Then, list the closest city among the important cities to the capital
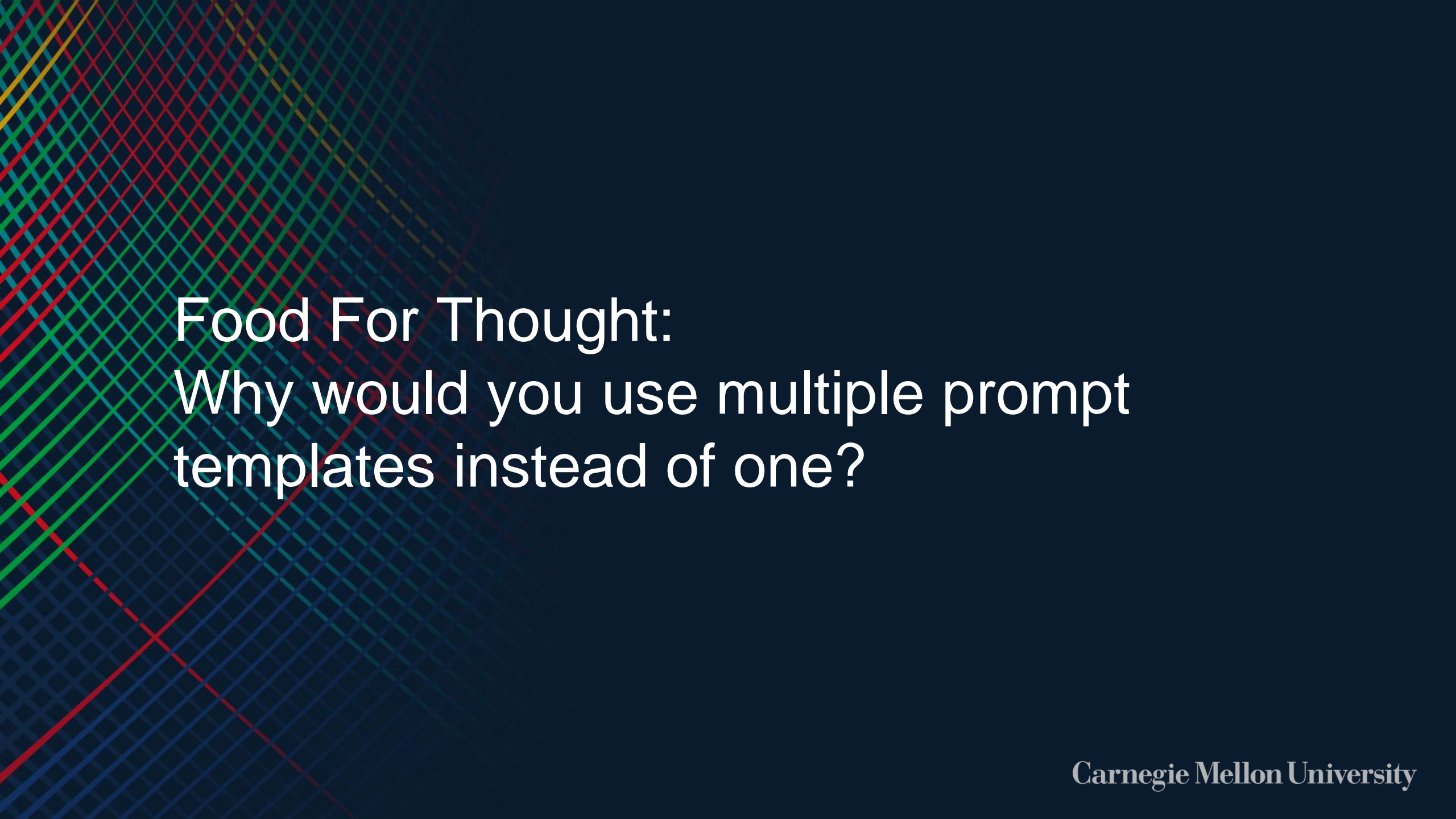
# Chain Exercise – Sample Solution
**(With Output Parsers)**

# Chain Exercise – Sample Solution – Cont'd

- **Selected LLM Model**: gemini-1.5-pro-001

- **LLM Chain-1 Prompt Template:**
  Print the most important 5 cities in a country {country} and its capital. Be specific about the list of countries and print them in comma-separated list. Produce the output in the format of two paragraphs. First paragraph starts with "capital:" and second paragraph starts with "cities:"

- **LLM Chain-2 Prompt Template:**
  Please extract the country name and the city name and answer the question below. You can extract the city name from the "capital" attribute in the string. Also, please extract the most important cities in this country from the "cities" attribute in the string.
  Country: {country}
  City: {cityChain}
  Most Important Cities: {cityChain}

  What do you know about the capital city? Which city among the most important cities is the closest to the capital?

Food For Thought:
Why would you use multiple prompt templates instead of one?

# LLM Application Debugging in FlowiseAI

# LLM Application Design using IfElse Utility
Note: Define Input Variables on ifElse Utility and Access Variables in JSON format

# Reminder: Custom Tool Exercise

Build Custom Tool for retrieving weather by Lat/Long

# Custom Tools Creation

- Create a custom tool with the API URL and proper arguments.
  API URL: http://api.open-meteo.com/v1/forecast?latitude=40.44&longitude=79.99&current_weather=true

# Custom Tools Creation – Cont'd

```
 9   const fetch = require('node-fetch');
10
11   //https://api.open-meteo.com/v1/forecast?latitude=40.44&longitude=79.99&current_weather=true
12   const url1 = 'https://api.open-meteo.com/v1/forecast?latitude=';
13   const url2 = '&longitude=';
14   const url3 = '&current_weather=true';
15   const options = {
16       method: 'GET',
17       headers: {
18           'Content-Type': 'application/json'
19       }
20   };
21   try {
22       const url = url1 + $latitude + url2 + $longitude+ url3;
23       console.log(url);
24       const response = await fetch(url, options);
25       const text = await response.text();
26       return text;
```

# Weather Application Design
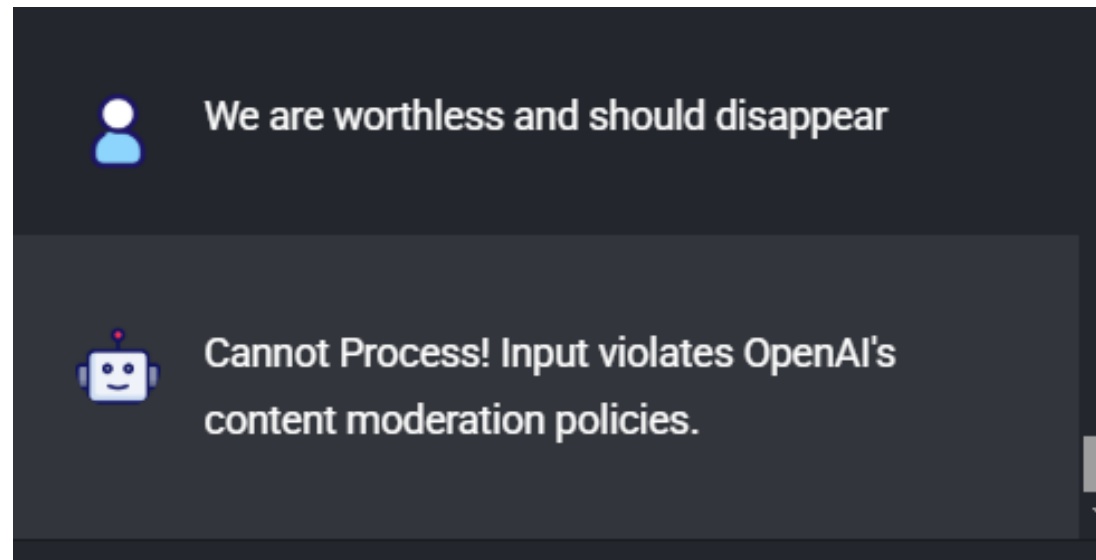
# Weather Application Design – Cont'd Execution Flow



**User Question** → **Moderation** → **Memory** → **LLM for Tool Selection** → **Tool Execution**

Carnegie Mellon University

# Weather Application Design – Cont'd

- **Input Moderation**

"We are worthless and should disappear"

Food For Thought:
How would you connect a chain to an agent?

# What about recursive Prompt Requests?

It's counter intuitive to manually require your LLM to generate better quality responses. To handle complex tasks and build recursive systems that can be composed of subsystems:

- Multi-agent System: A scalable framework that breaks down complex projects into manageable tasks, controlled by an agent.

- Sequential Agent System: similar to multi-agent systems but offers higher control and larger capabilities to build more complex systems.

# Multi-agent Systems

- User: Initiates requests.

- Supervisor AI: Orchestrates the workflow, breaking down requests into subtasks and assigning them to specialized workers.

- Worker AI Team: Specialized agents that independently execute specific tasks, using necessary tools and returning results to the Supervisor.
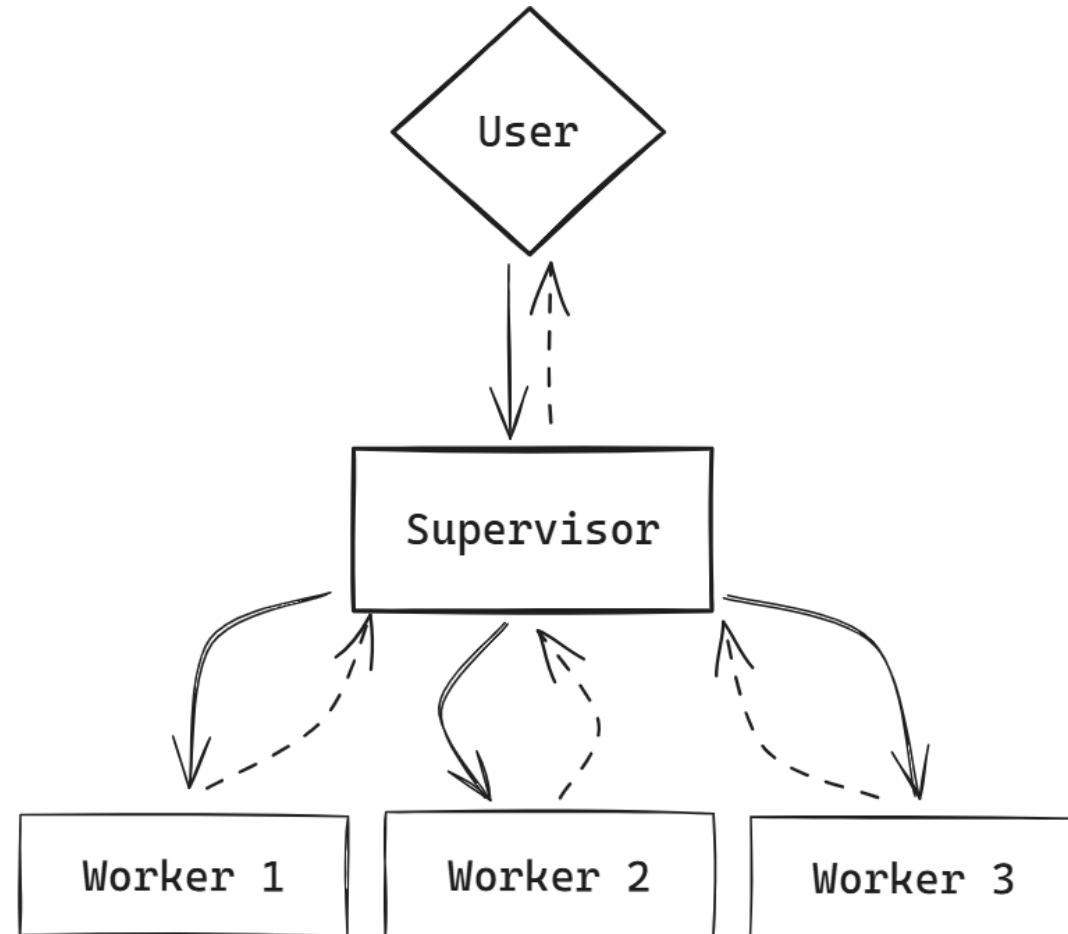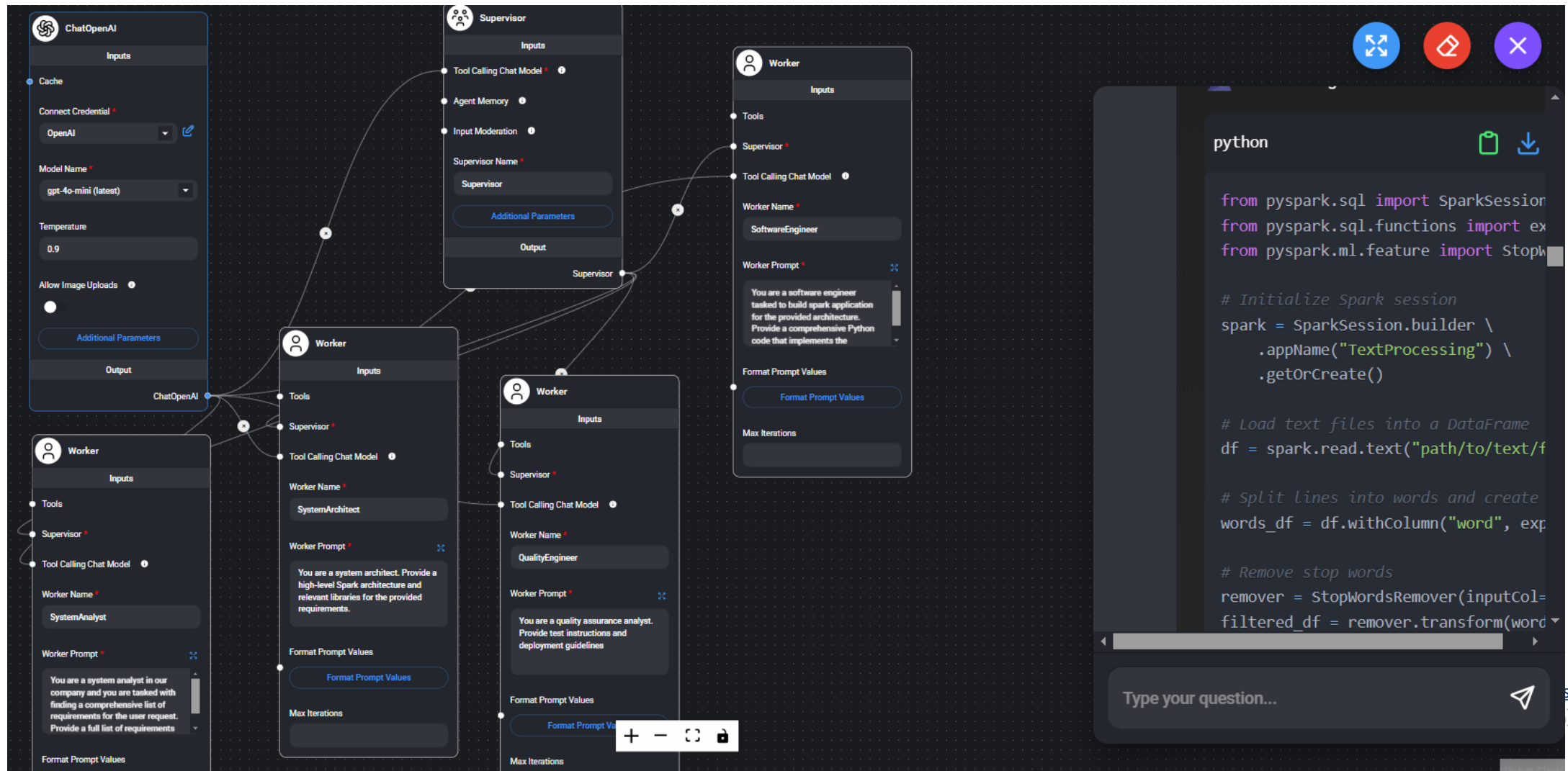


Image source: https://docs.flowiseai.com/using-flowise/agentflows/multi-agents

Carnegie Mellon University

# Multi-agent Systems Example

# Take Home Exercise:
# Deploy FlowiseAI to Google Kubernetes Engine

Follow the instructions on this page:

https://docs.flowiseai.com/configuration/deployment/gcp

# Connect FlowiseAI Apps to Local Applications

# Flowise AI Limitations

- Lack of output moderation

- Limited UI capabilities

- Limited deployment and scalability options

To overcome these limitations, **building LLM applications using coding frameworks is needed**

# Readings

- Multi-Agent Systems: https://docs.flowiseai.com/using-flowise/agentflows/multi-agents

- Sequential Agents: https://docs.flowiseai.com/using-flowise/agentflows/sequential-agents