# Lecture 4:
# LLM Application Design (Cont'd)

**SPRING 2025**

**MOHAMED FARAG**

FARAG@CMU.EDU

# Agenda

- Create Your Vertex AI Cloud Environment

- Generate Vertex AI Credentials

- FlowiseAI Components

- FlowiseAI Exercises

# Create Your Vertex AI Cloud Environment
## Hint: Enable Compute Engine and Vertex AI APIs

View: | INSTANCES | USER-MANAGED NOTEBOOKS | MANAGED NOTEBOOKS |
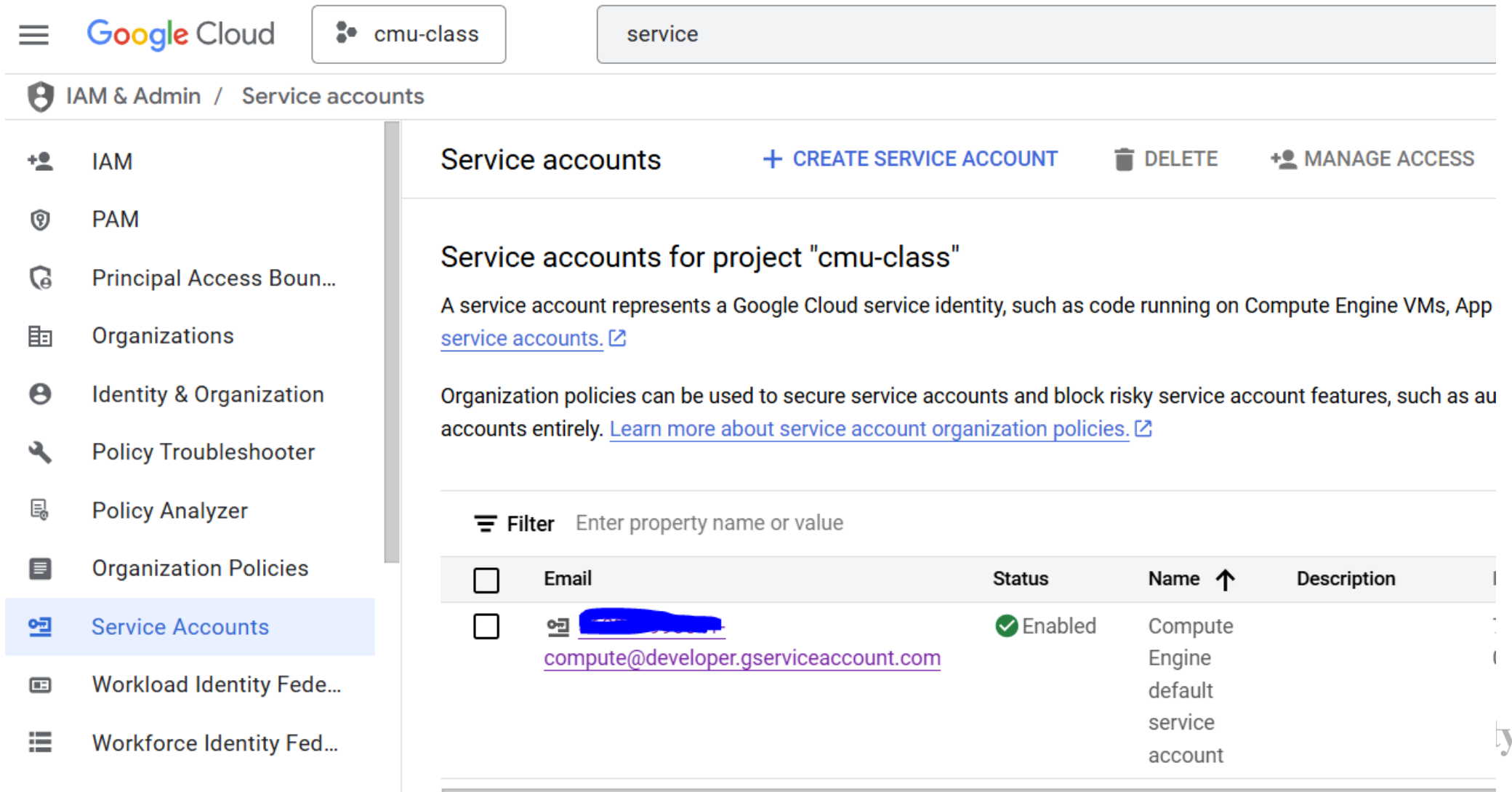
ℹ️ JupyterLab 4 is now available in Vertex AI Workbench.

Workbench Instances have JupyterLab 3 pre-installed and are configured with GPU-enabled machine learning frameworks. Learn more ↗

≡ Filter

| ☐ | ● | Instance name ↑ | | Zone | Auto upgrade | Version | Machine Type | GP |
|---|---|---|---|---|---|---|---|---|
| ☐ | ✅ | instance-20250122-223420 | OPEN JUPYTERLAB | us-central1-a | — | M127 | Efficient Instance: 4 vCPUs, 16 GB RAM | No |

# Generate Vertex API Credentials for FlowiseAI

# Create new JSON key and paste JSON file contents into your FlowiseAI credentials

# LLM Applications – General Idea



What if we would like to:

- Leverage external documents to enhance the capabilities of the LLM.
- Incorporate previous conversations into the LLM's responses.
- Integrate multiple prompts or functionalities for a cohesive experience
- Use a custom functionality that is not offered out of the box.

# FlowiseAI Components

# FlowiseAI Components (Cont'd)

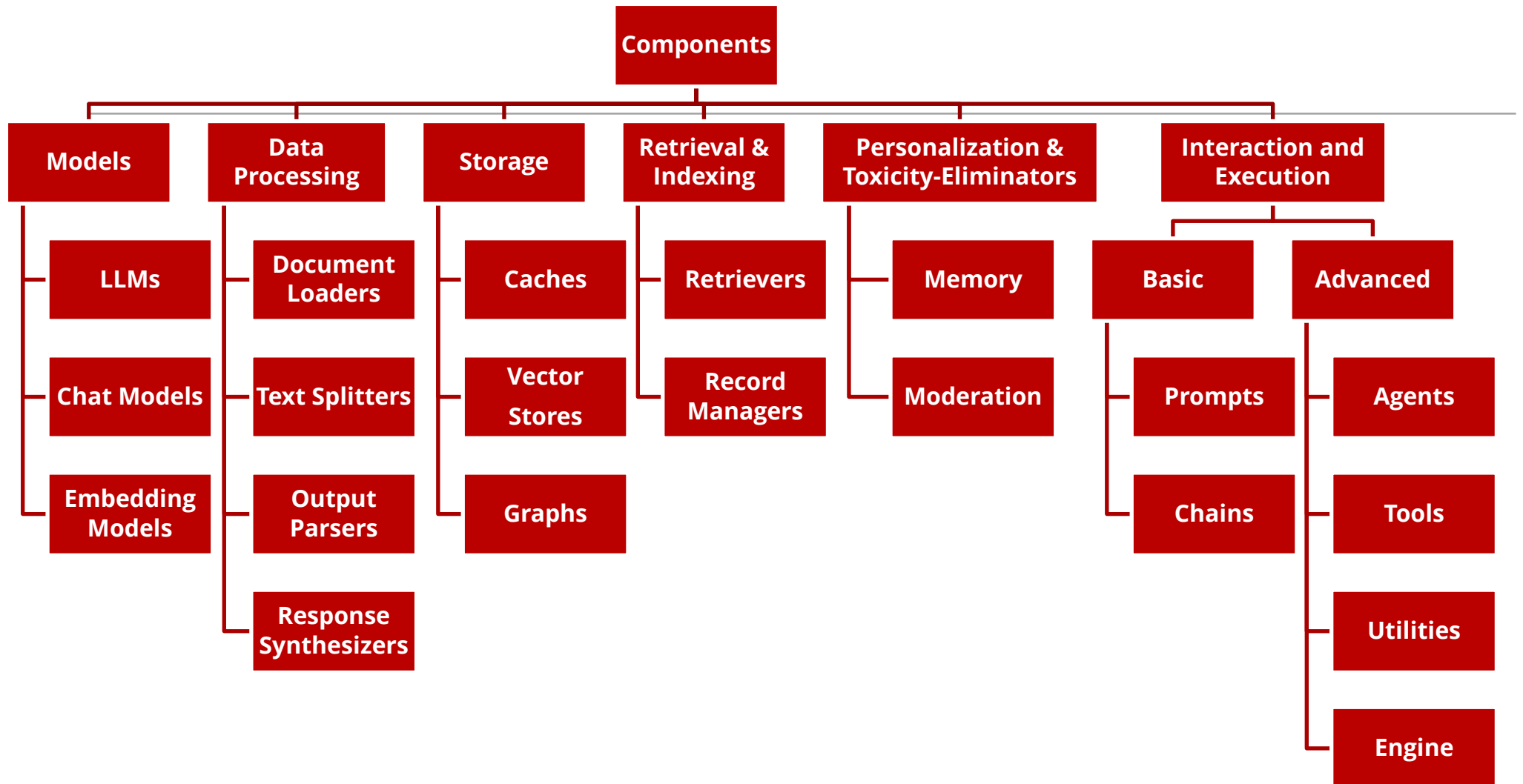| Category | Nodes | Purpose |
|---|---|---|
| Models | LLMs | Large Language Models for natural language understanding and generation |
| | Chat Models | Specialized models for conversations |
| | Embedding Models | Transforming data into vector space for machine learning |
| Data Processing | Document Loaders | Nodes for ingesting and parsing documents |
| | Text Splitters | Divide text into smaller chunks for processing |
| | Output Parsers | Convert model outputs into structured formats |
| | Response Synthesizer | Combines outputs into cohesive and structured responses (from LlamaIndex) |

# FlowiseAI Components (Cont'd)

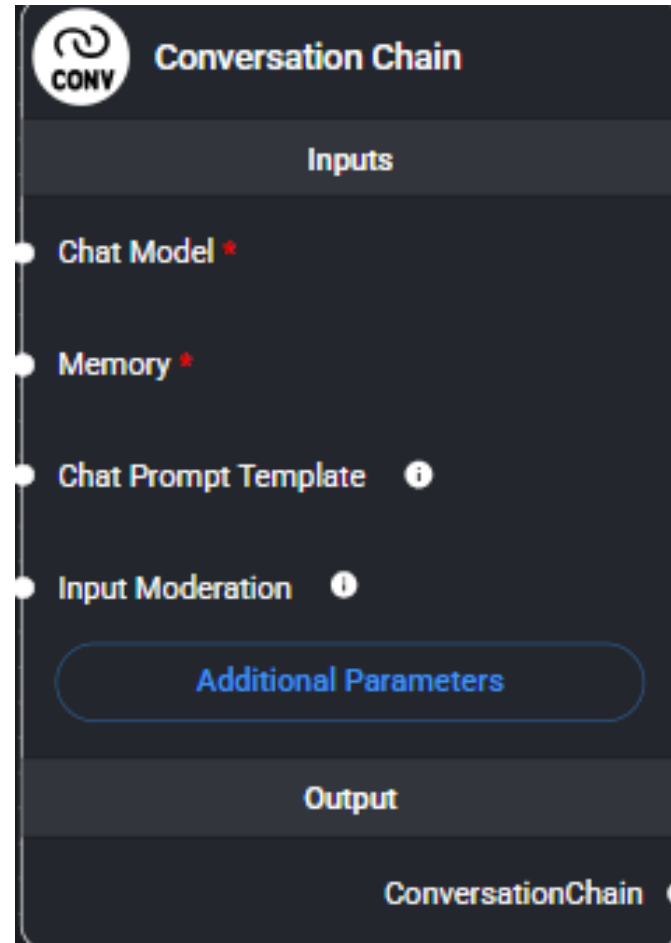| Category | Nodes | Purpose |
|---|---|---|
| Storage | Caches | Temporary storage to speed up recurring processes |
| | Vector Stores | Store and retrieve vector representations of data efficiently |
| | Graphs | Work with graph-based structures for complex relationships |
| Retrieval & Indexing | Retrievers | Fetch relevant data from storage for a query |
| | Record Managers | Organize and manage saved records in a database |
| Personalization & Toxicity-Eliminators | Memory | Retains context and past interactions |
| | Moderation | Monitor and restricts inappropriate content |

# FlowiseAI Components - Chains

Chains are leveraged to assemble modular components into versatile and reusable pipelines. Chains have different types including:

- Conversation Chain: Uses memory to keep track of previous interactions

- Retrieval QA Chain: Retrieves relevant information from a single knowledge base

- Conversational Retrieval QA Chain: Retrieves relevant information from a single knowledge base and optionally uses memory to keep track of previous interactions as well.

- Multi Retrieval QA Chain: Combines multiple retrievers.

# Chain Example – Conversation Chain

In your opinion,
why would you use a Chain?

# FlowiseAI Components - Agents

Agents are autonomous software entities designed to perform actions and execute tasks.

- What are the main differences between agents and chains?

# Agents – Cont'd

Agents are autonomous software entities designed to perform actions and execute tasks.

- What are the main differences between agents and chains?

  - Agents orchestrate chains while chains orchestrate lower-level modules.

  - Agents can take actions and perform tasks.

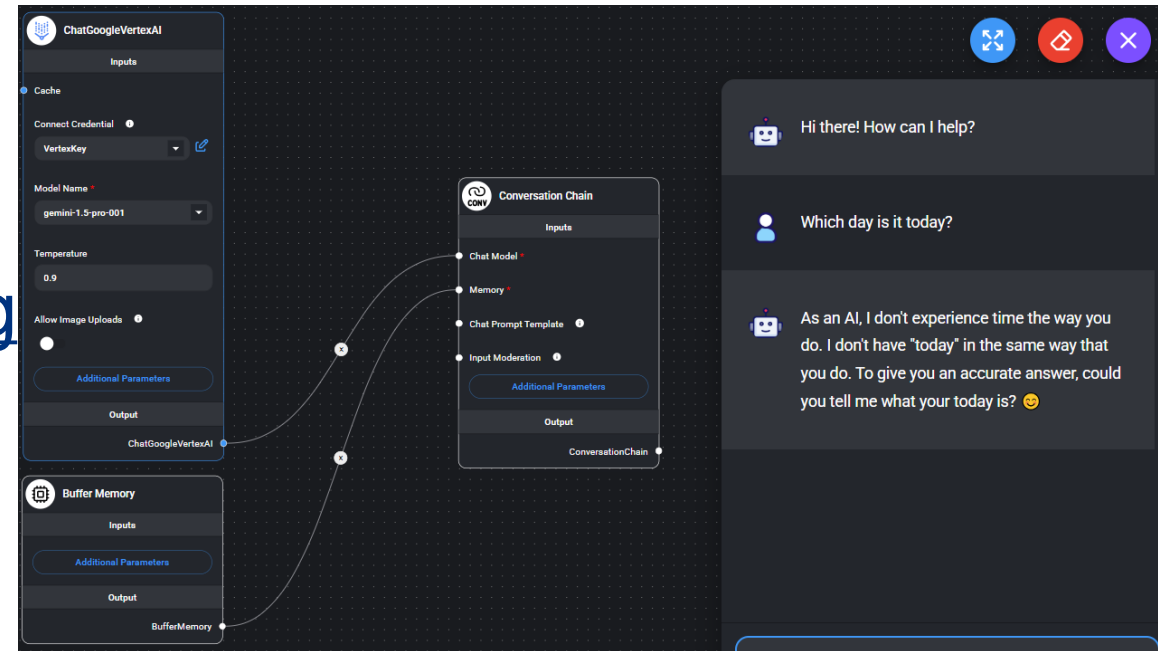Carnegie Mellon University

# Agent Example

# FlowiseAI Components - Tools

- Tools offer modular interfaces that enable <u>agents</u> to connect with external services such as databases and APIs.

- Toolkits organize tools that utilize shared resources.

- Example tools include Web Browser, Calculator, etc.

- You can create your custom tools as well.

Carnegie Mellon University

# Example

- We need to a tool to find today's date or weather since LLMs can't find this information on their own.

- Suggestions: Use
https://open-meteo.com/
OR
https://openweathermap.org
to create your custom tool
or agent.

# General LLM Application Design Guidelines

- Document Loaders → Text Splitters → Embeddings/Vector Stores

- Chat Models/LLMs → Output Parsers

- Prompts → LLM/Chat Model Chains

- Tools → Agents

- Moderation → Models, Chains, Agents, etc.

# Exercises

- Perform Q&A on PDF file.

- Build Custom Tool for retrieving weather by Lat/Long

# Readings

- Create your custom Tool in Flowise AI: https://docs.flowiseai.com/integrations/langchain/tools/custom-tool

    - Important reading for the project and HW-1.