

# LangGraph RAG Python Code Assistant

## AI-Powered Coding Companion with Memory

Doaa Said

December 7, 2025

# Key Libraries

- **Streamlit:** Interactive Web UI
- **LangGraph:** State-machine workflow for AI tasks
- **LangChain + ChatOllama:** LLM integration
- **Chroma:** Vector database for RAG
- **HuggingFaceEmbeddings:** Sentence embeddings
- **datasets:** HumanEval dataset for code examples

# LLM Setup

## ChatOllama Setup

```
1 llm = ChatOllama(  
2     model="llama3.2",  
3     temperature=0.6,  
4     num_ctx=2048,  
5     base_url="http://localhost:11434",  
6 )
```

- **temperature=0.6** → controls creativity
- **num\_ctx=2048** → context window
- **base\_url** → local LLaMA server

# Dataset Vector Store

## Code Example

```
1 dataset = load_dataset("openai/openai_humaneval"  
    , split="test")  
2 documents = [  
3     Document(page_content=f"Problem:\n{item['  
        prompt']}\n\nSolution:\n{item['  
        canonical_solution']}")  
4     for item in dataset  
5 ]  
6  
7 embeddings = HuggingFaceEmbeddings(model_name="  
    sentence-transformers/all-MiniLM-L6-v2")  
8 vector_store = Chroma(collection_name="  
    py_examples", embedding_function=embeddings)  
9 vector_store.add_documents(documents)
```

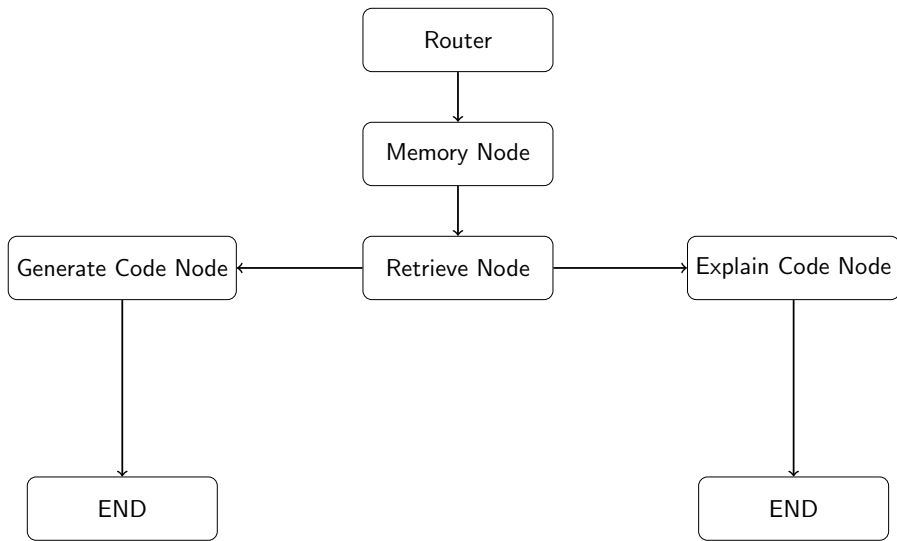
# Memory Module

## Conversation Memory

```
1 memory_buffer = deque(maxlen=5)
2
3 def add_to_memory(role, content):
4     memory_buffer.append({"role": role, "content": content})
5
6 def generate_memory_prompt():
7     return "\n".join([f"{m['role']}: {m['content']}" for m in memory_buffer])
```

- Maintains last 5 conversation turns
- Provides context to LLM

# LangGraph Workflow



# Key Features

- Integrates RAG + LLM for code generation
- Conversational memory (last 5 turns)
- Intent detection: Generate or Explain
- Interactive web interface via Streamlit
- Modular extensible workflow

# Conclusion

- Combines RAG + LLM + Memory for interactive coding assistant
- Modular workflow using LangGraph
- Real-time Streamlit UI
- Scalable for future datasets LLMs