

Week 7: Local LLM + RAG + Streamlit

Full Stack RAG with Local LLM

Semester 2/2568

Agenda

1 Ollama

2 RAG Pipeline

3 Streamlit

4 Lab 6

What is Ollama?

Local LLM Runtime

- Run LLMs locally
- No API key
- Free unlimited
- Data stays local

We use: qwen2.5:7b

Ollama Commands

```
# Pull model  
ollama pull qwen2.5:7b  
  
# Run  
ollama run qwen2.5:7b  
  
# API at http://localhost:11434
```

Call from Python

```
import requests

response = requests.post(
    "http://localhost:11434/api/generate",
    json={
        "model": "qwen2.5:7b",
        "prompt": "What is RAG?",
        "stream": False
    }
)
print(response.json()["response"])
```

Question → Embed → Search → Context → LLM → Answer

- ① Embed question (bge-m3)
- ② Search OpenSearch
- ③ Build context
- ④ Send to Ollama
- ⑤ Return answer

Streamlit UI

```
import streamlit as st

st.title("RAG Q&A")

if prompt := st.chat_input("Ask..."):
    st.chat_message("user").write(prompt)
    response = call_api(prompt)
    st.chat_message("assistant").write(response)
```

Lab 6: Complete RAG (3.75%)

Tasks:

- ① Setup Ollama
- ② Test LLM
- ③ Run complete RAG
- ④ Run Streamlit UI

Questions? Complete RAG System!