

# Week 1: Introduction to RAG + Local LLM

## Full Stack RAG with Local LLM

Semester 2/2568

# Agenda

- 1 Course Overview
- 2 What is RAG?
- 3 Why Local LLM?
- 4 Technology Stack
- 5 Assessment
- 6 Setup

# Course Information

Item	Detail
Course Code	CSI403
Course Name	Full Stack RAG with Local LLM
Credits	3 (2-3-5)
Semester	2/2568

**"Build Your Own AI, Not Just API Calls"**

# What Makes This Course Different?

## Traditional Course:

- OpenAI API Key
- Pay per Token \$\$\$
- Cloud Dependency
- Data sent to Cloud

## This Course:

- Local Ollama LLM
- Free & Unlimited
- Self-Hosted
- Data Stays Local

# The Problem with LLMs

**Large Language Models have limitations:**

- Knowledge cutoff date
- No access to your private documents
- Can hallucinate (make up facts)
- Generic responses

**Solution: RAG (Retrieval-Augmented Generation)**

## Retrieval-Augmented Generation

**RAG = Search your documents + Use LLM to generate answers**

- ① Retrieval:** Find relevant documents
- ② Augmented:** Add context to prompt
- ③ Generation:** LLM creates answer

# Cost Comparison

	Service	Cost	Type
<b>Cost per 1 Million Tokens</b>	OpenAI GPT-4	\$30.00	Cloud API
	OpenAI GPT-3.5	\$2.00	Cloud API
	Ollama (Local)	<b>\$0.00</b>	<b>Self-Hosted</b>

**Local LLM = Free Forever!**

# Benefits of Local LLM

## Cost & Freedom:

- No API fees
- Unlimited usage
- No rate limits

## Privacy & Control:

- Data stays local
- Offline capable
- Full customization

# Tech Stack (100% Self-Hosted)

Component	Technology	Type
Frontend	Streamlit	Local
Backend	FastAPI	Local
Embedding	HuggingFace (bge-m3)	Local
Vector DB	OpenSearch	Local
<b>LLM</b>	<b>Ollama (qwen2.5:7b)</b>	<b>Local</b>

No API Key Required!

# Assessment Overview

<b>Category</b>	<b>Total</b>
Attendance	10%
Quiz (4 times)	20%
Lab (8 times)	30%
Project	40%
<b>Total</b>	<b>100%</b>

# Software to Install

Please install before Week 2:

- ① **Miniconda** - [Python 3.10+](#)
- ② **Git** - [git-scm.com](#)
- ③ **VS Code** - [code.visualstudio.com](#)
- ④ **Docker Desktop** - [docker.com](#)
- ⑤ **Ollama** - [ollama.ai](#)

Questions? Welcome to Full Stack RAG with Local LLM! "Build Your

Own AI, Not Just API Calls"