

| | |
|-----------------------|---------------------------|
| COURSE TITLE | AGENTIC AI AND LLM |
| COURSE CODE | 01AI0609 |
| COURSE CREDITS | 4 |

Objective:

- 1 To equip students with the knowledge and practical skills to design, build, and deploy intelligent agentic systems and LLM-powered applications using advanced AI frameworks and tools.

Course Outcomes: After completion of this course, student will be able to:

- 1 Understand the architecture and functioning of Agentic AI, AI Agents, and LLMs. (Understand)
- 2 Build AI-driven conversational systems with memory and planning. (Apply)
- 3 Analyze retrieval-augmented generation for grounded knowledge use. (Analyze)
- 4 Use LangChain to create agent workflows integrating tools and API and assess its performance using different evaluation metrics. (Evaluate)
- 5 Apply Agentic AI and LLMs algorithms/tools (LangGraph, CrewAI, and AutoGen) in solving real-world problems (Apply)

Pre-requisite of course: Foundational knowledge of Machine Learning, Deep Learning, and Natural Language Processing, Probability and Statistics, Proficiency in Python Programming.

Teaching and Examination Scheme

| Theory Hours | Tutorial Hours | Practical Hours | ESE | IA | CSE | Viva | Term Work |
|--------------|----------------|-----------------|-----|----|-----|------|-----------|
| 3 | 0 | 2 | 50 | 30 | 20 | 25 | 25 |

| Contents : Unit | Topics | Contact Hours |
|-----------------|--|---------------|
| 1 | Introduction to LLMs Transformer Architecture – self attention mechanism, types of transformer architecture, Prompting and prompt engineering, Pre-training and fine tuning- single task, multi-task, Tokenization, Evaluation metrics - Perplexity, BLEU, ROUGE, METEOR, and BERTScore. | 4 |
| 2 | Conversational AI Dialogue Systems - task-oriented vs open domain, rule-based vs neural conversational agents, Context Management – dialogue history tracking, memory strategies, context window length and truncation techniques, Intent Recognition – classification-based models, zero-short/few shot intent detection, named entity recognition. | 8 |

| Contents : Unit | Topics | Contact Hours |
|------------------------|--|----------------------|
| 3 | Retrieval-Augmented Generation (RAG) Fundamentals of RAG – architecture, dense vs sparse retrieval, fusion-in decoder vs fusion-in-encoder, Retrieval techniques – vector embeddings, BM25, DPR, ColBERT, RAG pipeline. | 5 |
| 4 | LangChain Chains – simpleChain, sequentialChain, LLMChain, ConversationChain, and RetrievalQAChain, Pipeline and modular component – PromptTemplate, LLM, OutputParser, Memory, Input/output schema handling, Error Correction - Retry chains, Exception handling, LangSmith and logging/debugging tools, Context Handling - vectorStores, document loaders, memory management and history buffers | 8 |
| 5 | Function Calling OpenAI-style function calling, Defining functions via JSON schema, Tool execution and argument parsing, Tool integration - External APIs, code execution, database queries, Responsiveness - Multi-turn reasoning with intermediate tool use, Streaming responses, Error Handling. | 4 |
| 6 | Agentic Systems Introduction to Agents – reflection, long term memory, autonomy vs interactivity, LangGraph, Architecture – Planner-executer loop, tool chaining and task resolution, Planning and Decision Making – dynamic routing, conditional logic and retry mechanism, logging and introspection. | 4 |
| 7 | Advanced Agentic Tools CrewAI – Defining roles, tasks, communication protocols between agents, AutoGen – multi agent chat-based orchestration, system-agent-human interaction loops, human-in-the-loop design, Task Management – decomposition and delegation, prioritization and concurrent execution, memory sharing and results aggregation | 6 |
| 8 | Project Development and Deployment Final Integration - Combining RAG, LangChain, agentic tools, Handling API dependencies and configuration, Testing | 3 |
| Total Hours | | 42 |

Suggested List of Experiments:

| Contents : Unit | Topics | Contact Hours |
|------------------------|--|----------------------|
| 1 | Practical 1 Use BERT/GPT tokenizer to tokenize input text and visualize token embeddings using PCA or t-SNE. | 2 |
| 2 | Practical 2 Implement a mini-transformer model using PyTorch or TensorFlow to understand self-attention and positional encoding. | 2 |

Suggested List of Experiments:

| Contents : Unit | Topics | Contact Hours |
|--------------------|---|------------------|
| 3 | Practical 3 Summarize Dialogue with an Instruction Prompt (Zero shot, one shot, few shot inference) | 2 |
| 4 | Practical 4 Use pre-trained models (e.g., GPT-4, LLaMA2) for text generation, summarization, or question answering. | 2 |
| 5 | Practical 5 Load a pre-trained sentiment classifier and predict sentiment on movie reviews or tweets. | 2 |
| 6 | Practical 6 Use spaCy or transformers to extract named entities (like persons, organizations, dates) from text. | 2 |
| 7 | Practical 7 Build a chatbot using LangChain's ConversationChain with memory support. Track and update context over turns. | 2 |
| 8 | Practical 8 Load a sample data file using LangChain's DocumentLoader and build a simple question-answer system over it. | 2 |
| 9 | Practical 9 Develop an LLM agent that can use multiple tools based on user intent. | 2 |
| 10 | Practical 10 Set up a multi-agent chat system to iteratively refine a solution to a user query. | 2 |
| 11 | Practical 11 Create a small FAISS index using sentence embeddings and perform basic similarity search over short documents. | 2 |
| 12 | Practical 12 Containerize and deploy the LLM application using FastAPI and Docker. Host on Render or Hugging Face Spaces. | 2 |
| Total Hours | | 24 |

Textbook :

- 1 Speech and Language Processing (3rd ed. draft), Jurafsky, D., & Martin, J. H. , Stanford University, 2023
- 2 Natural Language Processing with PyTorch: Build Intelligent Language Applications Using Deep Learning, Delip Rao, B., & McMahan, B., O'Reilly Media, 2018

References:

- 1 Practical Retrieval-Augmented Generation: Building RAG Applications with Haystack and Transformers, Practical Retrieval-Augmented Generation: Building RAG Applications with Haystack and Transformers, Haystack Team, -, 2023

References:

- 2 A Concise Introduction to Multiagent Systems and Distributed Artificial Intelligence., A Concise Introduction to Multiagent Systems and Distributed Artificial Intelligence., Wooldridge, M., Wiley, 2021

Suggested Theory Distribution:

The suggested theory distribution as per Bloom’s taxonomy is as follows. This distribution serves as guidelines for teachers and students to achieve effective teaching-learning process

| Distribution of Theory for course delivery | | | | | |
|--|------------|-------|---------|----------|----------------------------------|
| Remember / Knowledge | Understand | Apply | Analyze | Evaluate | Higher order Thinking / Creative |
| 10.00 | 25.00 | 30.00 | 15.00 | 10.00 | 10.00 |

Instructional Method:

- 1 The course delivery method will depend upon the requirement of content and need of students. The teacher in addition to conventional teaching method by black board, may also use any of tools such as demonstration, role play, Quiz, brainstorming, MOOCs etc.
- 2 The internal evaluation will be done on the basis of continuous evaluation of students in the laboratory and class-room.
- 3 Practical examination will be conducted at the end of semester for evaluation of performance of students in laboratory.
- 4 Students will use supplementary resources such as online videos, NPTEL videos, e courses, Virtual Laboratory

Supplementary Resources:

- 1 <https://www.coursera.org/learn/generative-ai-with-llms>
- 2 https://onlinecourses.nptel.ac.in/noc25_cs45/