

# Generative AI & RAG Systems on NVIDIA

**Duration: 02 days**

## **Course Prerequisites:**

- Introductory deep learning knowledge, with comfort with PyTorch and transfer learning preferred.
- Intermediate Python experience, including object-oriented programming and libraries.

=====  
=====

## **Module 01: Introduction to Generative AI**

### **About this Module**

Generative AI describes technologies that are used to generate new content based on a variety of inputs. In recent time, Generative AI involves the use of neural networks to identify patterns and structures within existing data to generate new content. In this course, you will learn Generative AI concepts, applications, as well as the challenges and opportunities in this exciting field.

### **Learning Objectives**

Upon completion, you will have a basic understanding of Generative AI and be able to more effectively use the various tools built on this technology.

### **Topics Covered**

This no coding course provides an overview of Generative AI concepts and applications, as well as the challenges and opportunities in this exciting field.

### **Course Outline**

- Define Generative AI and explain how Generative AI works
- Describe various Generative AI applications
- Explain the challenges and opportunities in Generative AI

=====  
=====

## **Module 02: Building RAG Agents with LLMs**

### **About this Module**

The evolution and adoption of large language models (LLMs) have been nothing short of revolutionary, with retrieval-based systems at the forefront of this technological leap. These models are not just tools for automation; they are partners in enhancing productivity, capable of holding informed conversations by interacting with a vast array of tools and documents. This course is designed for those eager to explore the potential of these systems, focusing on practical deployment and the efficient implementation required to manage the considerable demands of both users and deep learning models. As we delve into the intricacies of LLMs, participants will gain insights into advanced orchestration techniques that include internal reasoning, dialog management, and effective tooling strategies.

### **Learning Objectives**

- Compose an LLM system that can interact predictably with a user by leveraging internal and external reasoning components.
- Design a dialog management and document reasoning system that maintains state and coerces information into structured formats.
- Leverage embedding models for efficient similarity queries for content retrieval and dialog guardrailing.
- Implement, modularize, and evaluate a RAG agent that can answer questions about the research papers in its dataset without any fine-tuning.

### **Topics Covered**

The workshop includes topics such as LLM Inference Interfaces, Pipeline Design with LangChain, Gradio, and LangServe, Dialog Management with Running States, Working with Documents, Embeddings for Semantic Similarity and Guardrailing, and Vector Stores for RAG Agents.

### **Course Outline**

- Introduction to the workshop and setting up the environment.
- Exploration of LLM inference interfaces and microservices.
- Designing LLM pipelines using LangChain, Gradio, and LangServe.
- Managing dialog states and integrating knowledge extraction.
- Strategies for working with long-form documents.
- Utilizing embeddings for semantic similarity and guardrailing.
- Implementing vector stores for efficient document retrieval.