



# Building with Large Language Models

## Leveraging LLMs in the Age of Generative AI

### Your Learning Transformation

By the end of this course, you will be able to **build** and **deploy** complex AI applications that produce **high-quality** outputs using the latest in open-source LLM technology.

### Learning Outcomes

By the end of this workshop you will:

1. Learn best-practices of prompt engineering and fine-tuning
2. Learn how to fine-tune and instruct-tune LLMs using data-centric approaches
3. Create a ChatGPT-like interface for your own data using indexing and chaining
4. Build and deploy an end-to-end Generative AI application each week!

### Project-Based Learning

In this workshop, you will deploy three applications in three weeks, including your own unique LLM capstone project in Week 3. You'll leverage the most powerful LLMs to build and deploy personalized applications in your domain and for your use cases. You'll also fine-tune LLMs to improve classic ML algorithm performance, teach LLMs new structure, and build a question-answering tool for your documents.

### Technical Requirements and Prerequisites

This workshop is intended for software engineers, data scientists, ML engineers, or others who have strong Python skills and have working experience deploying applications. Before the workshop, you should:

- Create an account on Hugging Face and be familiar with Hugging Face Datasets, Models, and Spaces. We recommend checking out our [Zero to Deployment Adventure Guide](#) from our recent Building Generative AI Applications Workshop, and checking out the [Hugging Face Course](#) for more information.
- Create a [Google Colab Pro](#) account, which is what we will be using for demos during the workshop.



## Detailed Schedule

Module	Topics	Build Activities
<b>Week 1</b> Fine-Tuning and Deploying LLM Applications	<b>Introduction and Networking</b> <ul style="list-style-type: none"> <li>• Program Overview</li> <li>• A Brief History of GPT Models</li> </ul> <b>Prompt Engineering</b> <ul style="list-style-type: none"> <li>• Prompt Engineering Best Practices</li> <li>• Zero-Shot vs. Few-Shot Learning</li> </ul> <b>LLM Fine-Tuning</b> <ul style="list-style-type: none"> <li>• Fine-Tuning Input-Output Schema</li> <li>• Instruction-Tuning</li> </ul> <b>Building with LLMs</b> <ul style="list-style-type: none"> <li>• The Role of LLM Developer APIs</li> <li>• Leading 2023 LLM Models</li> <li>• The Hugging Face Stack (Datasets, Models, Spaces)</li> </ul>	<ul style="list-style-type: none"> <li>• <b>Keeping up with LLMs!</b></li> <li>• <b>Building with gpt-3.5-turbo and OpenAI's API</b></li> <li>• Interactive End-to-End BLOOMZ Data Curation, Fine-Tuning, and Deployment Demo on Hugging Face</li> <li>• <b>Homework: Fine-Tuning BLOOMZ with Google Colab</b></li> <li>• <b>Homework: Capstone Project Commitments</b></li> </ul>
<b>Week 2</b> Indexing & Document Question Answering	<b>Technologies for “Building ChatGPT For Your Data”</b> <ul style="list-style-type: none"> <li>• <i>Indexing</i>, or Structuring Documents for LLMs to Interact with Them</li> <li>• <i>Vector Databases</i>, or Search and Retrieval using Vector Embeddings</li> <li>• <i>Chaining</i>, or Building Complex AI Applications with LLMs</li> <li>• <i>Agents</i>, or Building with Chains that depend on user inputs</li> <li>• Emerging Tools on the Market (e.g., LangChain, LlamaIndex, HayStack)</li> </ul> <b>Applications Considerations</b> <ul style="list-style-type: none"> <li>• From One Document to Many</li> <li>• From Qualitative QA to Quantitative</li> <li>• Chatbots versus Virtual Assistants</li> </ul>	<ul style="list-style-type: none"> <li>• <b>Answering Questions from a Single Document</b></li> <li>• Interactive End-to-End Document Question Answering Deployment Demo on Hugging Face</li> <li>• <b>Homework: Indexing and Querying Multiple Documents using Indexing with LangChain and Open AI's GPT 3.5</b></li> <li>• <b>Homework: Capstone Project MVPs (Dataset, Model Card, Space)</b></li> </ul>
<b>Week 3</b> Validating LLM Application	<b>How Do I Know If My Output is Good?</b> <ul style="list-style-type: none"> <li>• Challenges with Assessing Quality</li> <li>• Early LLM Metrics (Perplexity,</li> </ul>	<ul style="list-style-type: none"> <li>• Interactive Demos of how to leverage Stanford's HELM and Eluether AI's Test</li> </ul>



Outputs and Demo Day!	Burstiness, etc.) <ul style="list-style-type: none"><li>Emerging Metrics &amp; Model Card Best-Practices</li></ul> <b>Benchmark Frameworks &amp; Test Harnesses</b> <ul style="list-style-type: none"><li>A Brief History of NLP Benchmarks</li><li>LLM Trade Offs: Size and Capability vs. Inference Cost</li><li>Stanford's Holistic Evaluation of Large Language Models (HELM)</li><li>Eleuther AI's Test Harness</li></ul>	Harness <ul style="list-style-type: none"><li><b>PROJECT: Curate Data, Build a Fine-Tuned AI Application. Deploy on Hugging Face; include a framework for measuring quality of outputs</b></li></ul>
-----------------------	--	--

\* Activities in **Bold Pink** will be working sessions during the class!

[Register for Building with LLMs](#)

### About FourthBrain

FourthBrain trains engineers, developers, data scientists, and leaders to make an impact in the Artificial Intelligence field, with our flexible, accessible education programs. We are training a new generation of engineers and leaders who have more than just technical ability; they have an awareness and mindset of what is needed to succeed with AI. We are part of the AI Fund, founded by Andrew Ng.