⚙    👤 **tliu** ▾

**Data Science on AWS**    ✕

# Introduction

This workshop consists of three parts. Each part can be run independantly.

## Part 1: Distributed Data Preprocessing

Learn how to utilize distributed data processing in your data science projects that supports both interactive serverless development and ephemeral jobs for large scale processing.

More specifically, you will learn how to do the following:

- Register parquet data with AWS Glue and run interactive Spark queries and visualizations against this data in SageMaker Studio, all in an easy to use serverless Spark cluster.
- Use SageMaker processing jobs to deploy a large scale data processing job.
- Simplify your Spark workloads with SageMaker including EMR, Glue, SageMaker processing, and more.

## Part 2: Fine-tuning LLMs with Amazon SageMaker

Learn how to build generative AI and large language models (LLMs) for your data science projects that support both interactive development and ephemeral training

More specifically, you will learn how to do the following:

- Use SageMaker Studio and SageMaker training jobs to build generative AI models interactively and at-scale.
- Build and fine-tune a state-of-the-art generative model with SageMaker and HuggingFace.
- Use SageMaker endpoints to deploy a model as a real-time application endpoint.

## Part 3: Automating Fine-Tuning workflows with SageMaker Pipelines

Learn how to build repeatable MLOps pipelines for your data science projects including prompt engineering, generative AI, and large-language models (LLMs).

More specifically, you will learn how to do the following:

- Use SageMaker Studio and SageMaker Pipelines to create a repeatable pipeline for data preparation, model training, and model deployment.
- Engineer a set of features, fine-tune a state-of-the-art language model, and deploy the model as an HTTPS endpoint.
- Use ephemeral SageMaker processing jobs, training jobs, and pipelines to automate the entire model lifecycle and minimize cost.