# Self-Hosted Open Source LLMs: Best Practices and Strategies

JOSHUA ARVIN LAT & SOPHIE SOLIVEN
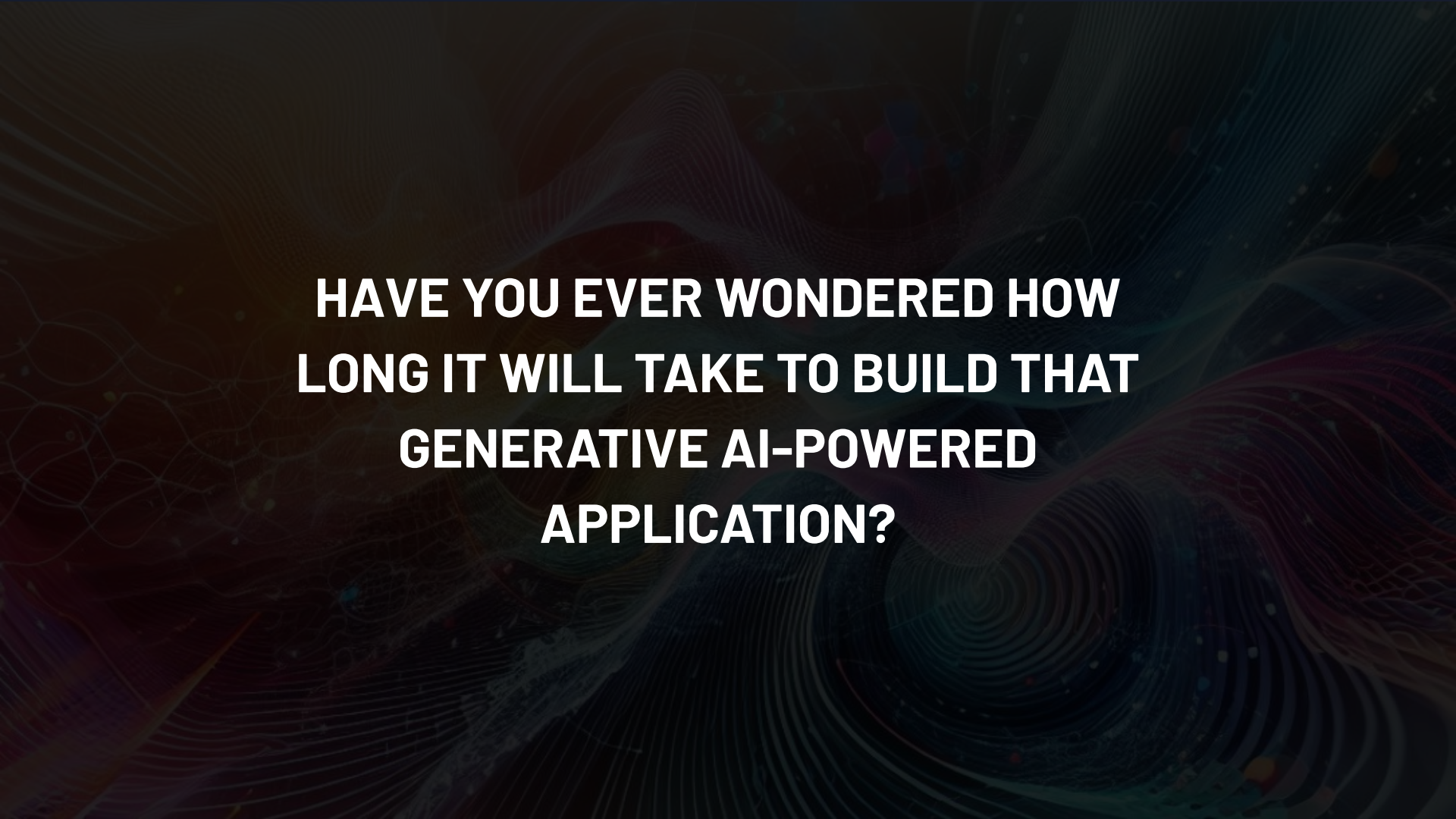
**Operations Director**
Edamama

*PREVIOUSLY*

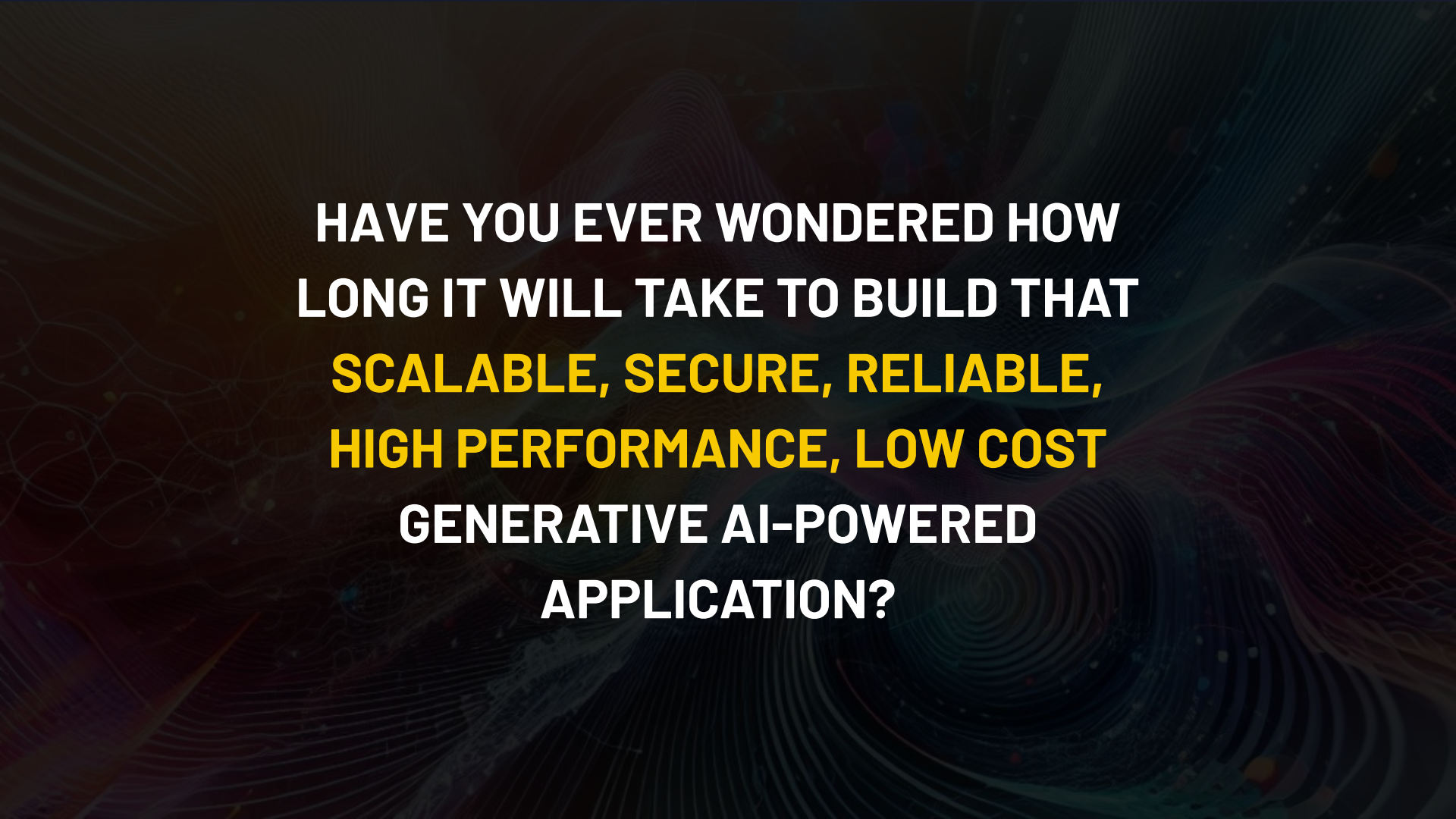**General Manager of E-Commerce Services and Dropship**

**Certifications** in:
*Cloud computing and data analytics*

**Technical Reviewer** of the book:
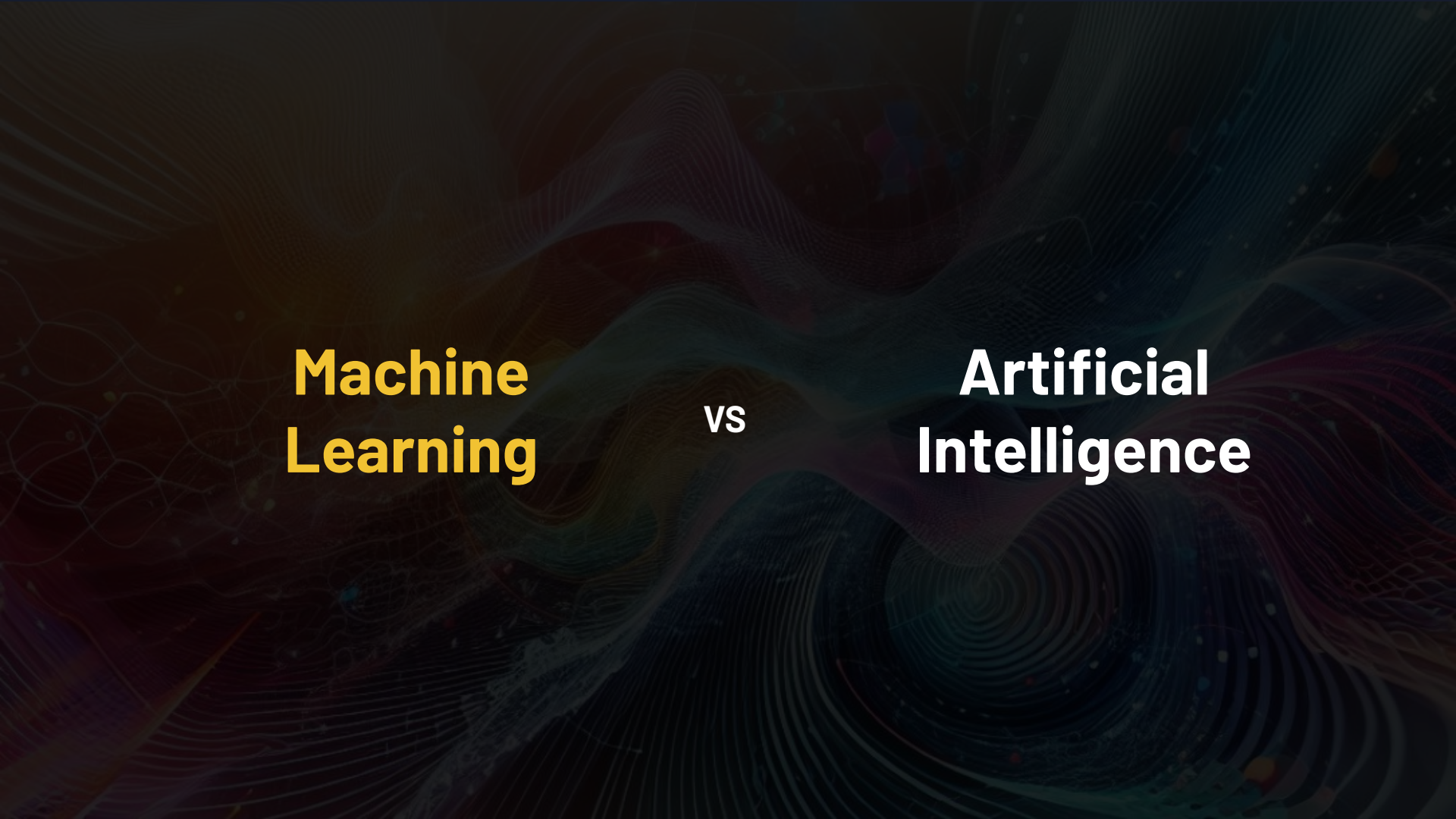*Machine Learning Engineering on AWS*

SOPHIE SOLIVEN

HAVE YOU EVER WONDERED HOW LONG IT WILL TAKE TO BUILD THAT GENERATIVE AI-POWERED APPLICATION?

HAVE YOU EVER WONDERED HOW LONG IT WILL TAKE TO BUILD THAT SCALABLE, SECURE, RELIABLE, HIGH PERFORMANCE, LOW COST GENERATIVE AI-POWERED APPLICATION?

# CONCEPTS

Machine Learning vs Artificial Intelligence

Machine Learning $\subset$ Artificial Intelligence

Artificial Intelligence

Machine Learning

## ARTIFICIAL INTELLIGENCE

| ANI<br>NARROW | AGI<br>GENERAL | ASI<br>SUPER |
|---|---|---|

## MACHINE LEARNING

| Supervised Learning | Unsupervised Learning | Reinforcement Learning |
|---|---|---|

## DEEP LEARNING

## GENERATIVE AI

| LLM<br>TEXT | IMAGES | AUDIO |
|---|---|---|

SELF-HOSTED
OPEN SOURCE LLMS

# BARRIER TO ENTRY

# INFRASTRUCTURE COST

FLEXIBILITY AND LEVEL OF CONTROL

SELF-HOSTED LLM SETUP

LLM + SAGEMAKER PYTHON SDK →

SageMaker Studio or SageMaker Notebook Instance
Data Science Environment

Large Language Model (LLM)
deployed in a SageMaker
Inference Endpoint

# SELF-HOSTED LLM SETUP



S3 Static Website Hosting (Frontend)

API Gateway

AWS Lambda

DATABASE

Large Language Model (LLM) deployed in a SageMaker Inference Endpoint

# THE END