**DagsHub**

# Machine Learning in Production
## What does "Production" even mean?

# Let's talk about #Buzzwords

Hyperparameter Optimization

CI/CD/CT

Feature Stores

MLOps

AGI

Data-Centric AI

Active Learning
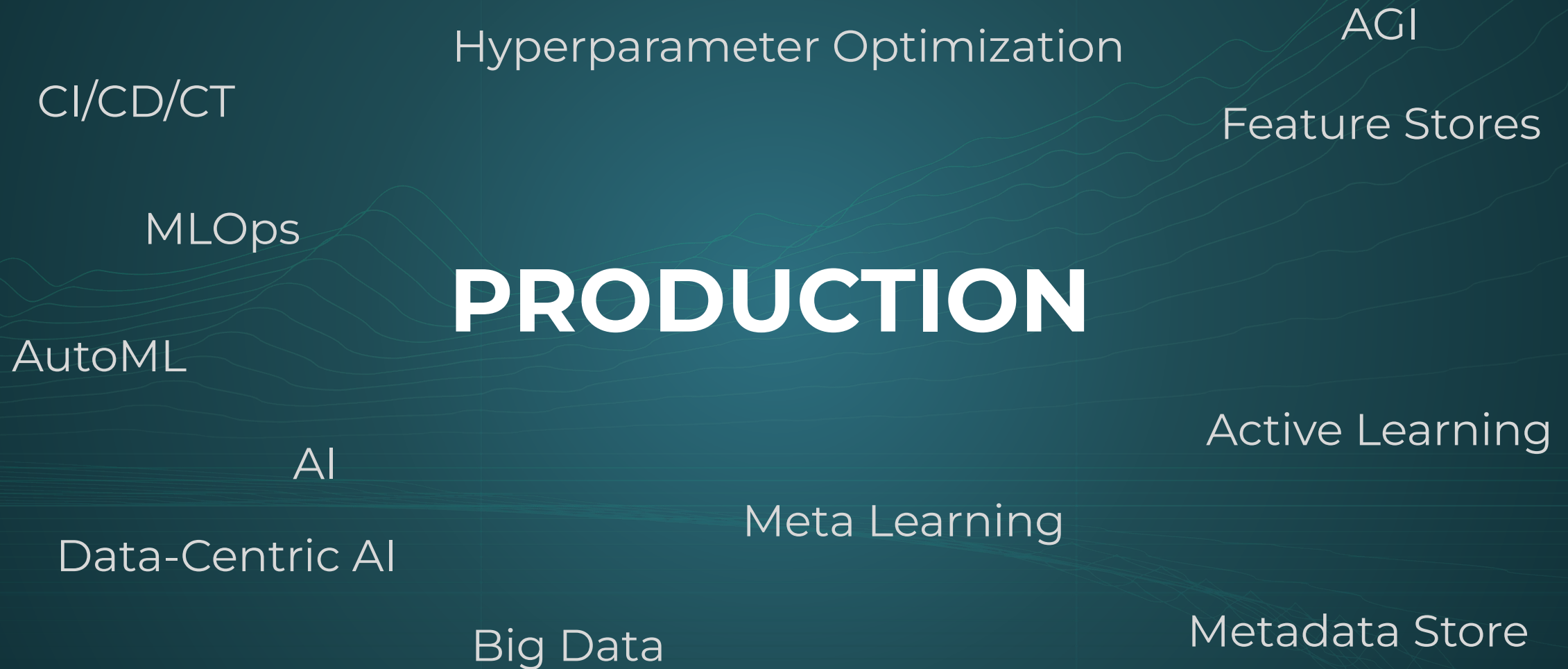
AutoML

Large Language Models
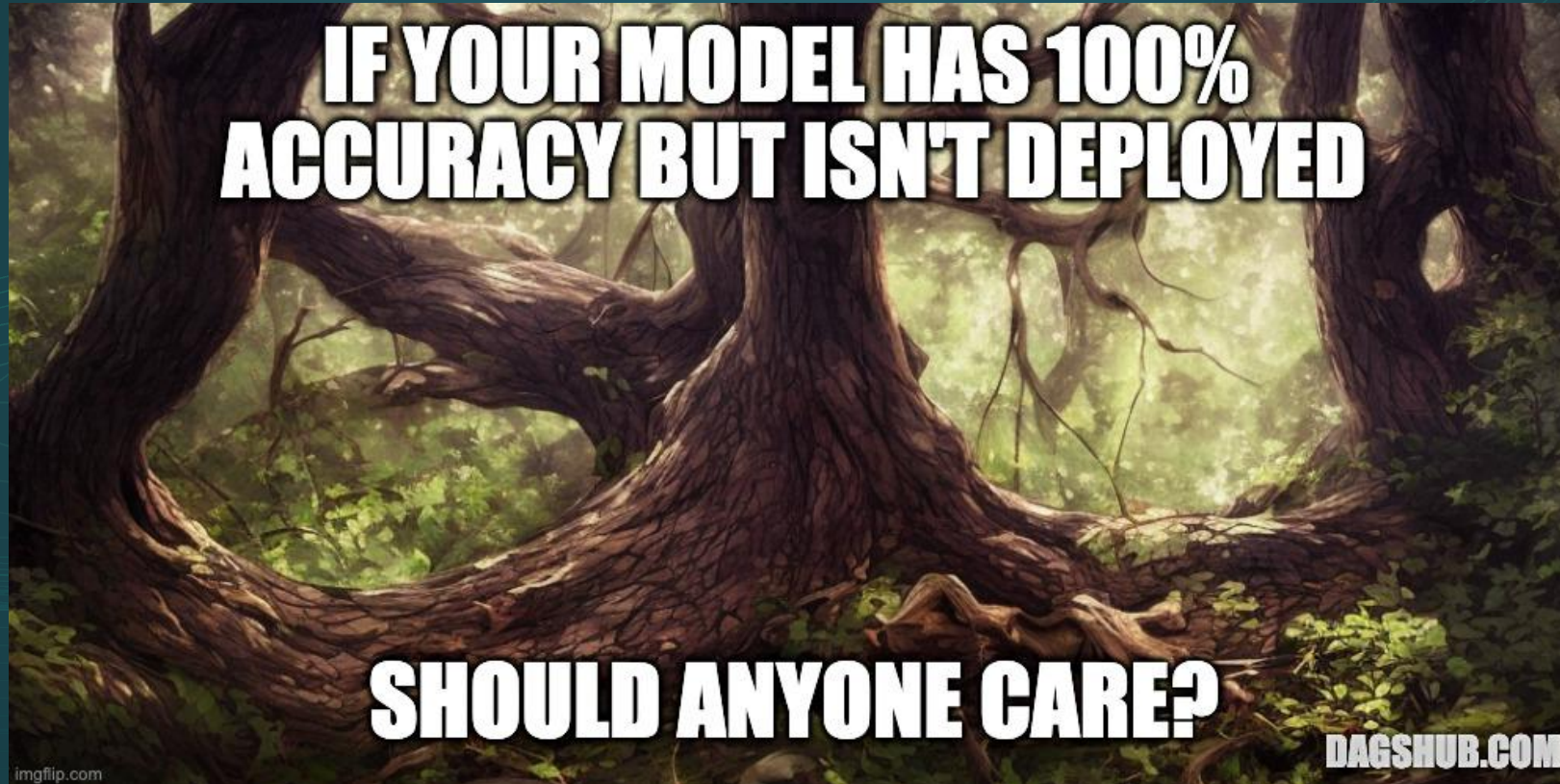
AI

Meta Learning

Metadata Store

Big Data

DagsHub

# Let's talk about #Buzzwords

Hyperparameter Optimization

AGI

CI/CD/CT

Feature Stores

MLOps

# PRODUCTION

AutoML

Active Learning

AI

Meta Learning

Data-Centric AI

Big Data

Metadata Store

DagsHub

# Defining ML in production

# About Me

## Dean Pleban

Follow me:

@DeanPlbn

/DeanPleban

DagsHub

# About Me

## Dean Pleban

Building tools for ML teamwork

Follow me:

@DeanPlbn

/DeanPleban

DagsHub

# About Me

## Dean Pleban

🛠 Building tools for ML teamwork

Strongly believe in open source

🐶 **DagsHub**

Follow me:

🐦 **@DeanPlbn**

in **/DeanPleban**

DagsHub

# About Me

## Dean Pleban

🔧 Building tools for ML teamwork

🔓 Strongly believe in open source

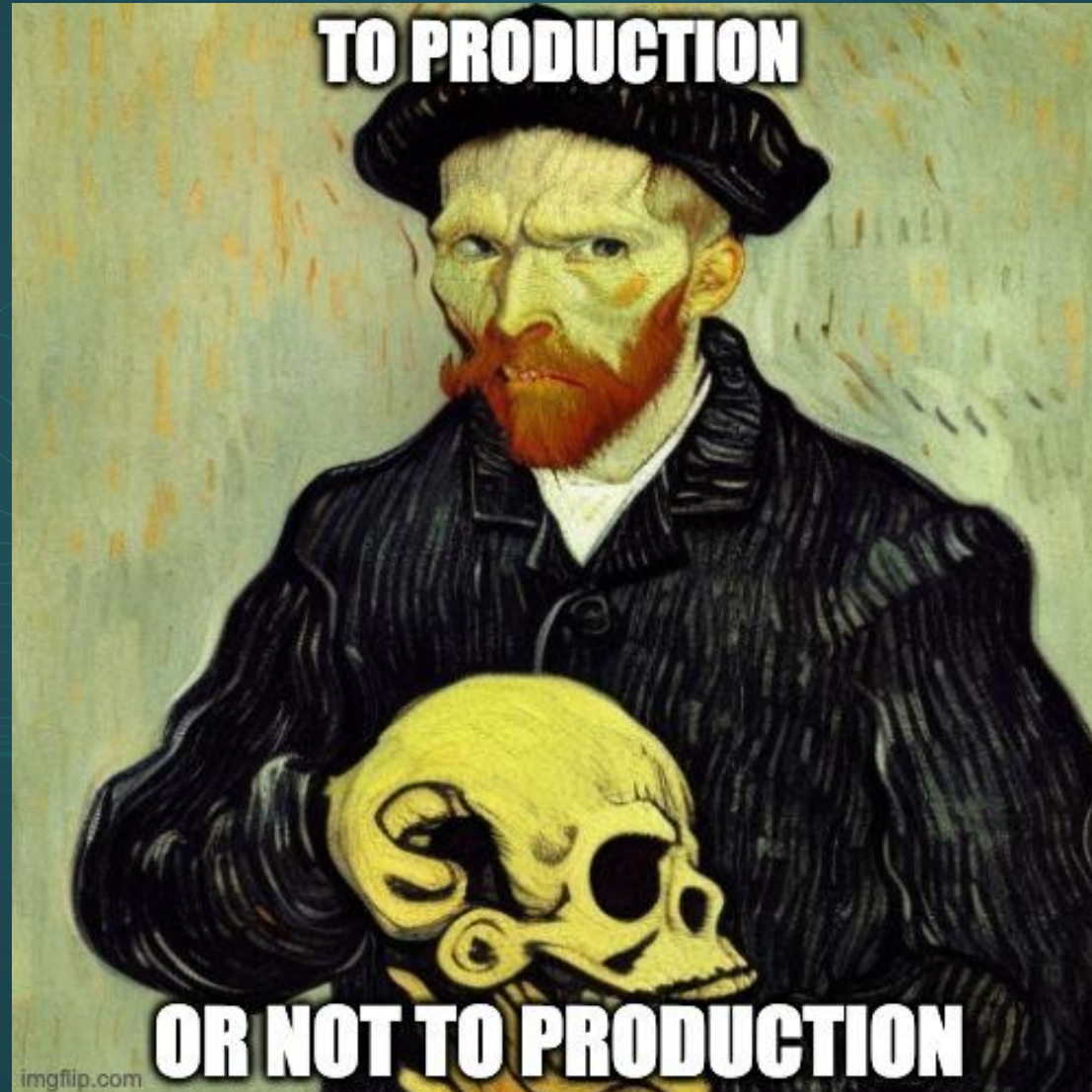🐶 DagsHub  The MLOps Podcast

Follow me:

🐦 @DeanPlbn

in /DeanPleban

DagsHub

# Defining ML in production

Remember "that" statistic...

Sponsored

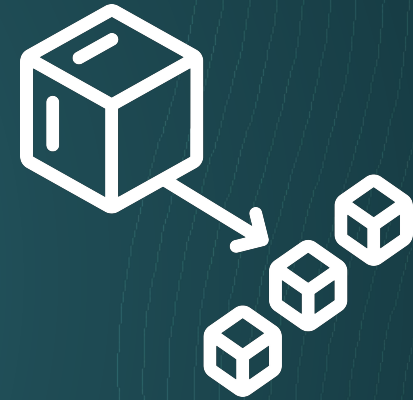Why do 87% of data science projects never make it into production?

DagsHub

# Why should you care?

Understand what production means

Framework to think about buzzwords

First-principles guide to deployment

DagsHub

# What types of production are there?



Endpoint

TheOneRavenous · 3 mo. ago

Push code to GitHub it gets reviewed and approved. Rolls server with new model weights. Then send form data to URI endpoint. Model runs inference from user queries. Then outputs the result back to the user. Microservice architecture. Just runs as a separate app so that the regular server processes I/O. Instance is always live.

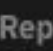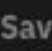DagsHub

# What types of production are there?

Endpoint

Edge Device

> **Atom_101** · 3 mo. ago
>
> Microservice is the closest one I guess. I deploy on edge. We build specialized robots. We have multiple python "servers" (just normal python scripts that receive their inputs from other processes) running, each of which control one ai model. The "main" process controlling the robot's parts is written in Java. If there is a decision that needs to be taken, the Java program sends the required information to the relevant ai process via zmq and receives a response from the model.
>
> ⬆ 31 ⬇  💬 Reply  Give Award  Share  Report  Save  Follow

DagsHub

# What types of production are there?

Endpoint

Edge Device

Dashboard

slowpush · 3 mo. ago

Create daily predictions that run overnight and push out insights and reports to the business in the morning.

⬆ 1 ⬇ 💬 Reply  Give Award  Share  Report  Save  Follow

DagsHub

# What types of production are there?



**Endpoint**          **Edge Device**          **Dashboard**          **ETL / Pipeline**

**volta_seca** · 3 mo. ago

We have a data pipeline that generate the predictions and store them in a database (we use Redis), all written in PySpark and orchestrated with Airflow. Then the API, written in Golang, accesses this database and returns predictions to the user.

⬆ 6 ⬇   💬 Reply   Give Award   Share   Report   Save   Follow

DagsHub

Models vs. Pipelines

Data

MODEL

Prediction

"Dog"

DagsHub

Models vs. Pipelines

# Models vs. Pipelines – Who cares?

DagsHub

# Models vs. Pipelines – Who cares?

End-to-end
thinking

DagsHub

# Models vs. Pipelines – Who cares?

End-to-end
thinking

Better understand
your requirements

DagsHub

# First-principles thinking



A woman thinking from **first-principles**,
lost in thought, a painting by Van Gogh
(By Stable Diffusion)

DagsHub

# What we assume



Deploying a
single model

DagsHub

# What we assume



Deploying a
single model

A simple flow
(sort of)

DagsHub

# What we assume



Deploying a single model

MODEL

"Dog"

A simple flow (sort of)

The model is trained

DagsHub

# Breaking deployment down – **The hard part**

# Breaking deployment down – **The hard part**

# Breaking deployment down – **The hard part**



"A command line, someone is typing a command"

DagsHub

# Breaking deployment down – **The hard part**

# Breaking deployment down – **The hard part**



"A website UI on a computer"



DagsHub

# Breaking deploym

# Breaking deployment down – **The hard part**

# Breaking deployment down – **The hard part**



"Two computers running the
same program"



CONTAINER

DagsHub

# Breaking deployment down – **The hard part**



CONTAINER

# Breaking deployment down – **The hard part**



DagsHub

# Breaking deployment down – **The hard part**



"Many computers connected to the cloud"



DagsHub

Breaking deployment down

1. Wrap the model in a prediction function

2. Wrap the function in an API

3. Put everything in a suitable environment

4. Provision infrastructure to host the environment

DagsHub

# 1. The prediction function

Recommended
Tools

DagsHub

# 1. The prediction function



Recommended
Tools

ONNX

Model
formats

DagsHub

# 1. The prediction function



Recommended
Tools

Model
formats

Define a class or
interface

DagsHub

# 2. The API wrapper

FastAPI

Flask
web development,
one drop at a time

Requests
http for humans

Recommended
Tools

DagsHub

# 2. The API wrapper

FastAPI

Flask
web development,
one drop at a time

Requests
http for humans

Recommended
Tools

Define the right
endpoints

DagsHub

# 2. The API wrapper

FastAPI

Flask
web development,
one drop at a time

Requests
http for humans

## Recommended Tools

## Define the right endpoints

## Authentication

DagsHub

# 3. The environment container



Recommended
Tools

DagsHub

# 3. The environment container



Recommended
Tools



Steps 2+3 in one

DagsHub

# 4. The infrastructure

Recommended
Tools

DagsHub

# 4. The infrastructure



Necessary
Recommended
Tools

DagsHub

# 4. The infrastructure



Necessary
Recommended
Tools



GPUs

DagsHub

# Further Reading

1. Building an API for ML models:
   https://towardsdatascience.com/step-by-step-approach-to-build-your-machine-learning-api-using-fast-api-21bd32f2bbdb

2. Authentication with FastAPI:
   https://fastapi.tiangolo.com/tutorial/security/

3. Docker for data science:
   https://dagshub.com/blog/setting-up-data-science-workspace-with-docker/

4. Deploy GPU Accelerated Applications with ECS and Docker:
   https://www.docker.com/blog/deploy-gpu-accelerated-applications-on-amazon-ecs-with-docker-compose/

DagsHub

# Thank You!

Follow me:

🐦 **@DeanPlbn**

in **/DeanPleban**

Follow DagsHub:

🐦 **@TheRealDAGsHub**

in **/company/dagshub/**