



EMEA HEALTHCARE & LIFE SCIENCES WORKSHOPS

# Machine Learning on AWS

Yegor Tokmakov  
Sr. Solutions Architect  
AWS Health

Dmitri Laptev  
Sr. Solutions Architect  
AWS Startups

Aamna Najmi  
Data Scientist  
AWS Professional Services

# Agenda

- 13:00 - 13:15 Welcome and introductions
- 13:15 - 13:45 Amazon SageMaker Introduction
- 13:45 - 14:30 Lab - Train, Tune and Deploy model using SageMaker Built-in Algorithm
- 14:30 - 15:00 Create model, prediction and inference
- 15:00 - 15:45 Lab - Bring your custom model
- 15:45 - 16:00 Coffee break
- 16:00 - 16:30 Advanced ML: Transformers, LLMs, HuggingFace
- 16:30 - 17:15 Lab - Summarize Scientific Documents with Amazon SageMaker and HuggingFace
- 17:15 - 17:30 Conclusion, QA

# Workshop materials and instructions

- <https://catalog.us-east-1.prod.workshops.aws/workshops/63069e26-921c-4ce1-9cc7-dd882ff62575/en-US>
- <https://catalog.workshops.aws/hcls-aiml/en-US>

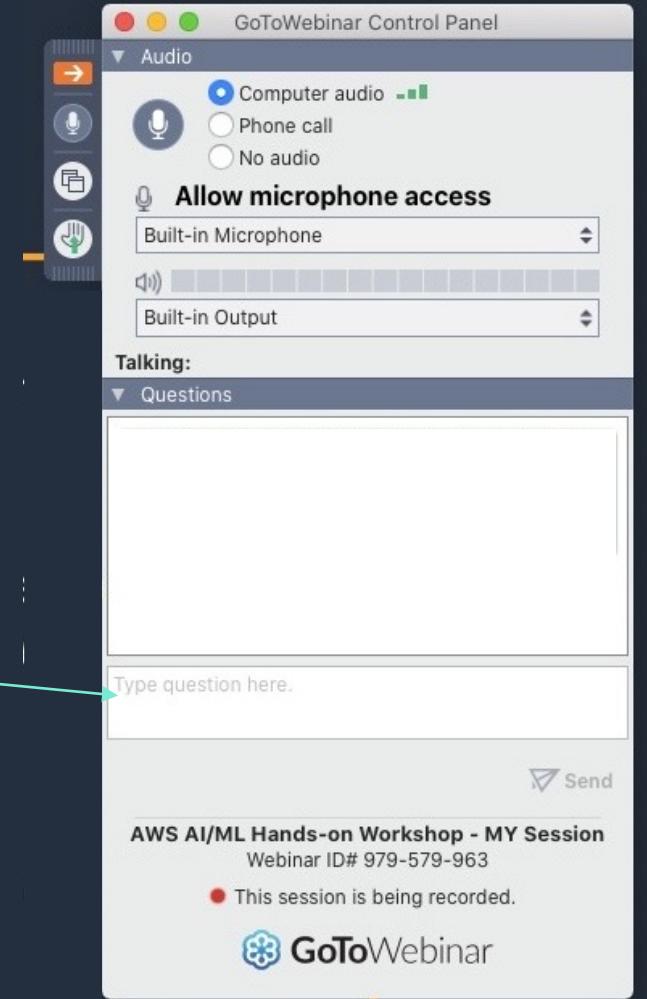
# Questions & answers

If you have any questions or encounter issues during the workshop, our support team is online.

You can submit your query in the GoToWebinar Questions function. To submit questions, select "Send"



Type your question here



# Amazon SageMaker Introduction

# The reach of ML is growing



## INCREASED SPENDING

By 2025, global spending on artificial intelligence will reach \$204 billion

—IDC

IDC, "Worldwide Spending on Artificial Intelligence,"  
<https://bit.ly/3y7hDoP>

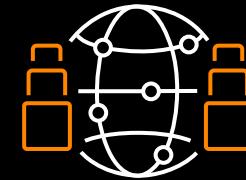


## FROM PILOTING TO OPERATIONALIZING

By the end of 2024, 75% of enterprises will shift from piloting to operationalizing AI

—Gartner

Gartner, "Gartner Identifies Top 10,"  
<https://gtnr.it/3Bln3uU>



## AI TRANSFORMATION

57% said that AI would transform their organization in the next three years

—Deloitte

Deloitte, "Thriving in the Era of Pervasive AI"  
<https://bit.ly/3JSeCi6>

# AI/ML with AWS

## Innovation, choice, and flexibility

**100,000+**

customers have used  
machine learning (ML)  
on AWS

**250+**

new capabilities for  
ML and artificial  
intelligence (AI) in just  
the last 12 months

**92%**

of deep learning (DL) in the cloud runs on  
AWS

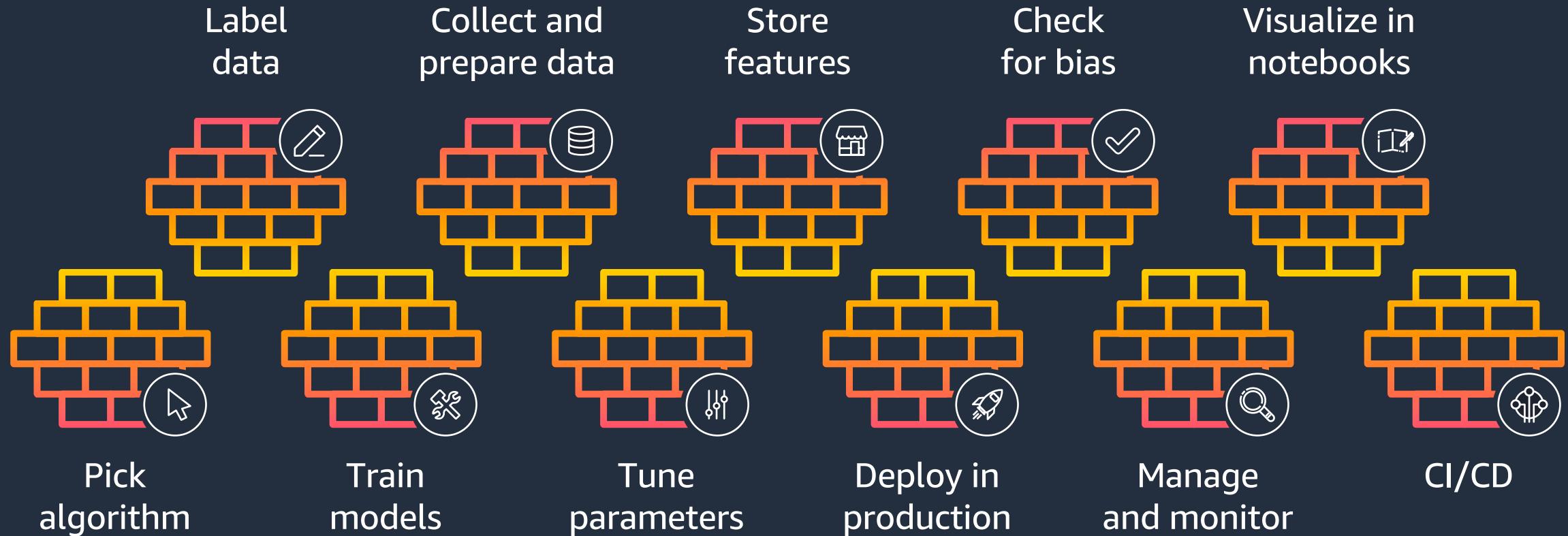
**91%**

of cloud-based PyTorch runs on  
AWS

### **AWS ML SOLUTIONS**

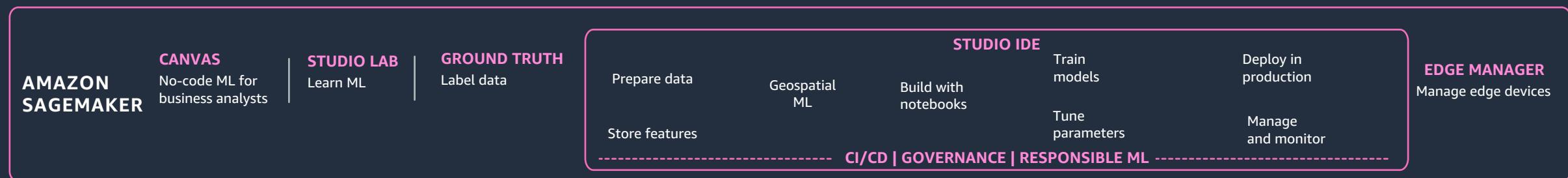
Reduce training time by 50%  
Provide 90% scaling efficiency  
Deliver 3x faster network throughput  
Improve price and performance by 25%

# Machine learning development is complex and costly



# The AWS AI/ML stack

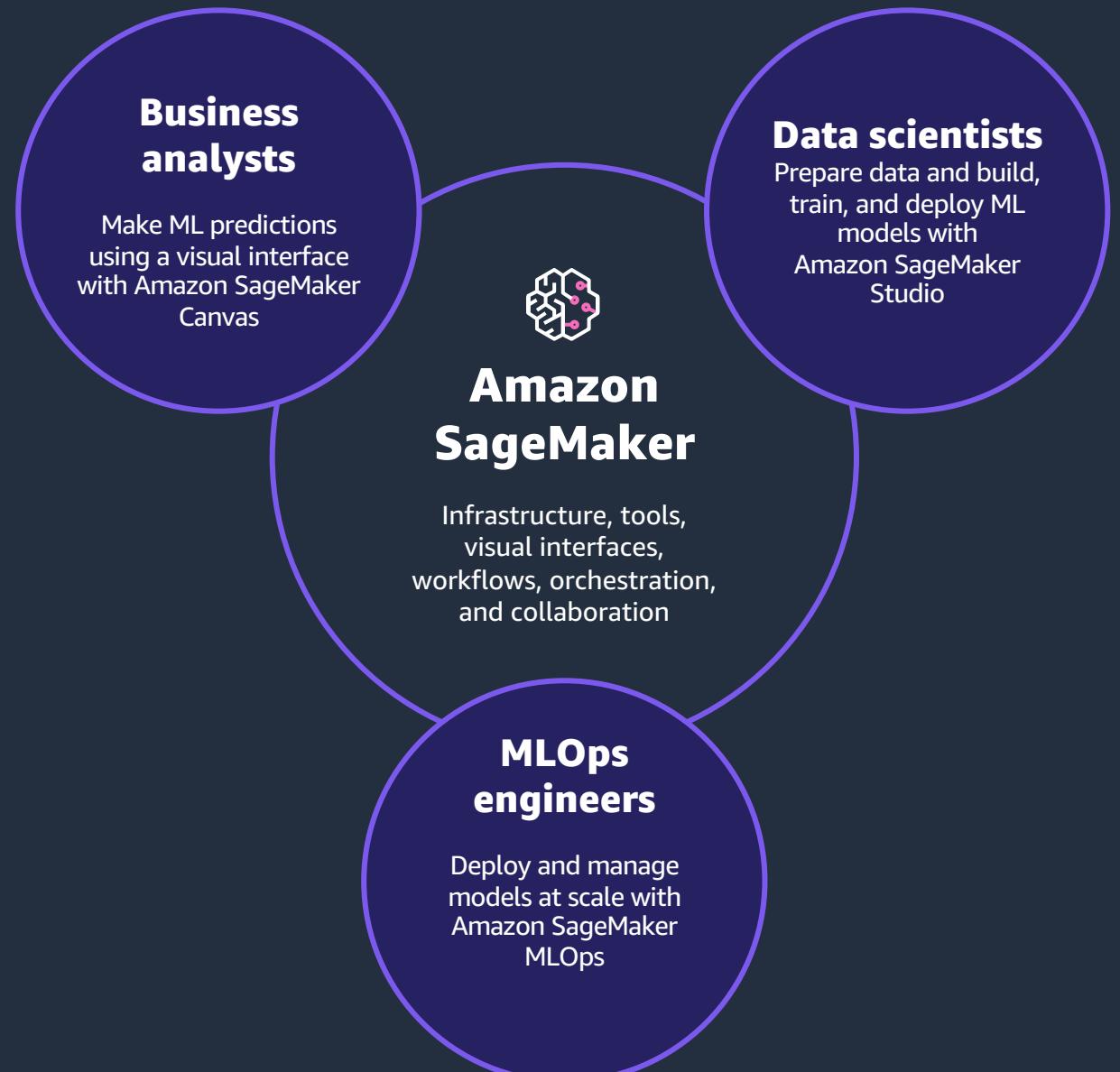
BROADEST AND MOST COMPLETE SET OF MACHINE LEARNING CAPABILITIES





# What is SageMaker?

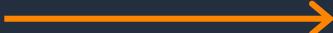
# Amazon SageMaker helps organizations harness ML



# Overcoming the barriers to ML



**Not enough ML builders**



## No-code ML tools

Make ML predictions regardless of ML experience



**Access, process, and label massive volumes of data for ML**



## Purpose-built data preparation tools

Access, process, and label data for ML



**Disparate data science tools**



## Integrated ML tools in a single interface

Build, train, and deploy models using IDEs



**Tedious, manual ML operations**



## Built-in MLOps

Automate and standardize MLOps practices



**Challenging to govern ML projects efficiently**



## Out-of-box ML governance tools

Simplify access control and enhance transparency across ML lifecycle

# Amazon SageMaker feature tour

PREPARE DATA AND BUILD, TRAIN, AND DEPLOY ML MODEL FOR ANY USE CASE

## PREPARE →

### Geospatial

Visualize geospatial data

### Ground Truth

Create high quality datasets for ML

### Data Wrangler

Aggregate and prepare data for ML

### Processing

Built-in Python, BYO R/Spark

### Feature Store

Store, catalog, search, and reuse features

### Clarify

Detect bias and understand model predictions

## BUILD →

### Studio Notebooks & Notebook Instances

Fully managed Jupyter notebooks with elastic compute

### Studio Lab

Free ML development environment

### Built-in Algorithms

Integrated tabular, NLP, and vision algorithms

### JumpStart

UI based discovery, training, and deployment of models, solutions, and examples

### Autopilot

Automatically create ML models with full visibility

### Bring Your Own

Bring your own container and algorithms

### Local Mode

Test and prototype on your local machine

## TRAIN & TUNE →

### Fully Managed Training

Broad hardware options, easy to setup and scale

### Distributed Training Libraries

High performance training for large datasets and models

### Training Compiler

Faster deep learning model training

### Automatic Model Tuning

Hyperparameter optimization

### Managed Spot Training

Reduce training cost by up to 90%

### Debugger and Profiler

Debug and profile training runs

### Experiments

Track, visualize, and share model artifacts across teams

### Customization Support

Integrate with popular open source frameworks and libraries

## DEPLOY & MANAGE →

### Fully Managed Deployment

Ultra low latency, high throughput inference

### Real-Time Inference

For steady traffic patterns

### Serverless Inference

For intermittent traffic patterns

### Asynchronous Inference

For large payloads or long processing times

### Batch Transform

For offline inference on batches of large datasets

### Multi-Model Endpoints

Reduce cost by hosting multiple models per instance

### Multi-Container Endpoints

Reduce cost by hosting multiple containers per instance

### Shadow Testing

Validate model performance in production

### Inference Recommender

Automatically select compute instance and configuration

### Model Monitor

Maintain accuracy of deployed models

### Kubernetes Operators & Components

Manage and monitor models on edge devices

### Edge Manager

Manage and monitor models on edge devices

### MLOps: Pipelines | Projects | Model Registry

Workflow automation, CI/CD for ML, central model catalog

### Canvas

Generate accurate machine learning predictions—no code required

### Studio | RStudio

Integrated development environment (IDE) for ML

### Governance

Model Cards | Dashboard | Permissions



# Amazon SageMaker Data Wrangler

EXPLORE, PREPARE, AND PROCESS  
DATA WITH LITTLE TO NO CODE



**Import data from multiple sources**

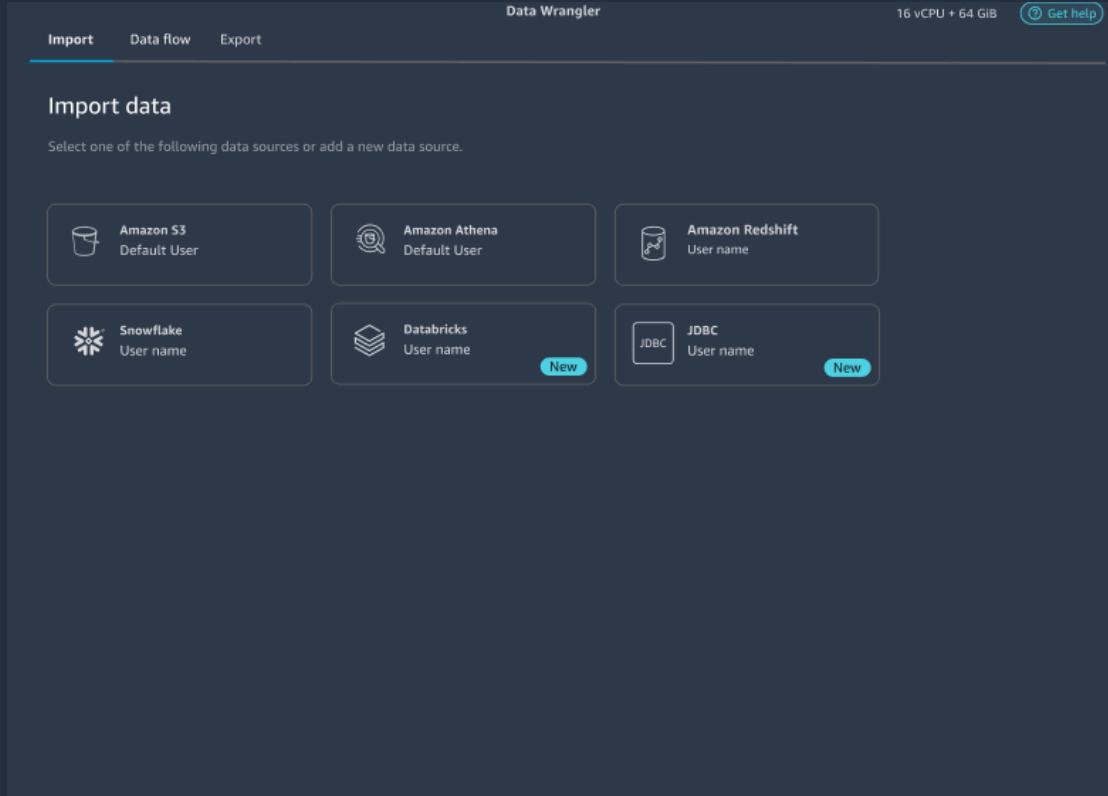
**Get insights on data and data quality**

**Visually explore, analyze, and prepare data**

**Quickly perform feature engineering**

**Automate ML data preparation workflows**

# Quickly connect and query data (1/3)



## Connect to multiple data sources:

Amazon S3

Amazon Athena

Amazon Redshift

Snowflake

Databricks DeltaLake (NEW)

SageMaker Feature Store & more coming

## Support for common file formats:

CSV files

Parquet files

JSON & JSONL (NEW)

ORC file (NEW)

Database tables

# Quickly connect and query data (2/3)

The screenshot shows the AWS Data Catalog interface. On the left, there's a search bar with 'Enter an S3 URL' and a 'Go' button. Below it, a list shows an object named 'airports.csv' located at 'S3 / sagemaker-us-east-2-562522975874 / airports'. The object has a size of 209.09KB and was last modified on 2020-11-05 17:17:30+00:00. On the right, a 'DETAILS' panel shows the file name as 'airports.csv' and the file type as 'csv'. It also includes options to preview the data on S3 or a database, schema, tables, and objects in Snowflake. At the bottom, there's a preview section titled 'PREVIEW • airports.csv (first 100 rows shown)' displaying the first few rows of the CSV file.

| iata | airport              | city             | state | country | lat         | long         |
|------|----------------------|------------------|-------|---------|-------------|--------------|
| 00M  | Thigpen              | Bay Springs      | MS    | USA     | 31.95376472 | -89.23450472 |
| 00R  | Livingston Municipal | Livingston       | TX    | USA     | 30.68586111 | -95.01792778 |
| 00V  | Meadow Lake          | Colorado Springs | CO    | USA     | 38.94574889 | -104.5698933 |
| 01G  | Perry-Warsaw         | Perry            | NY    | USA     | 42.74134667 | -78.05208056 |
| 01J  | Hilliard Airpark     | Hilliard         | FL    | USA     | 30.6880125  | -81.90594389 |
| 01M  | Tishomingo County    | Belmont          | MS    | USA     | 34.49166667 | -88.20111111 |
| 02A  | Gragg-Wade           | Clanton          | AL    | USA     | 32.85048667 | -86.61145333 |
| 02C  | Capitol              | Brookfield       | WI    | USA     | 43.08751    | -88.17786917 |
| 02G  | Columbiana County    | East Liverpool   | OH    | USA     | 40.67331278 | -80.64140639 |

Visually browse data sources like objects on S3, or database, schema, tables and objects in Snowflake

Preview & sample top rows

Join data from multiple sources

Support for VPC, KMS, CMK, and AWS Secrets Manager

# Quickly connect and query data (3/3)

Query Amazon Athena

Write a SQL query to extract data

Data catalog  
AwsDataCatalog

Database  
airlines

Advanced configuration

Workgroup

Location of query results  
s3://sagemaker-us-east-2-562522975874/athena

This value cannot be modified

Enable sampling

Run Cancel query

```
1 select * from airlines_2008 where origin = 'SFO'
```

Results

| year | month | dayofmonth | dayofweek | deptime | crsdeptime |
|------|-------|------------|-----------|---------|------------|
| 2008 | 3     | 17         | 1         | 1117    | 1121       |
| 2008 | 3     | 17         | 1         | 1841    | 1640       |
| 2008 | 3     | 17         | 1         | 1519    | 1455       |
| 2008 | 3     | 17         | 1         | 1959    | 2000       |

Query data using Amazon Athena

Browse data catalog, databases, & tables

Write queries & preview before importing data

Includes support for Athena Workgroups & custom output bucket

# Easily transform data for ML with 300+ built-in transforms

300+ built-in data transformations (no code) for common data prep needs and ML specific needs

Built by data scientists for data scientists

**ML specific transforms such as:**

One hot encoding  
Balance data  
Time series transforms

The image displays two identical-looking user interface panels titled "ADD TRANSFORM" with a dark header bar. Below the header, there are two columns of transformation options, each with a title, a brief description, and a "Learn more" link.

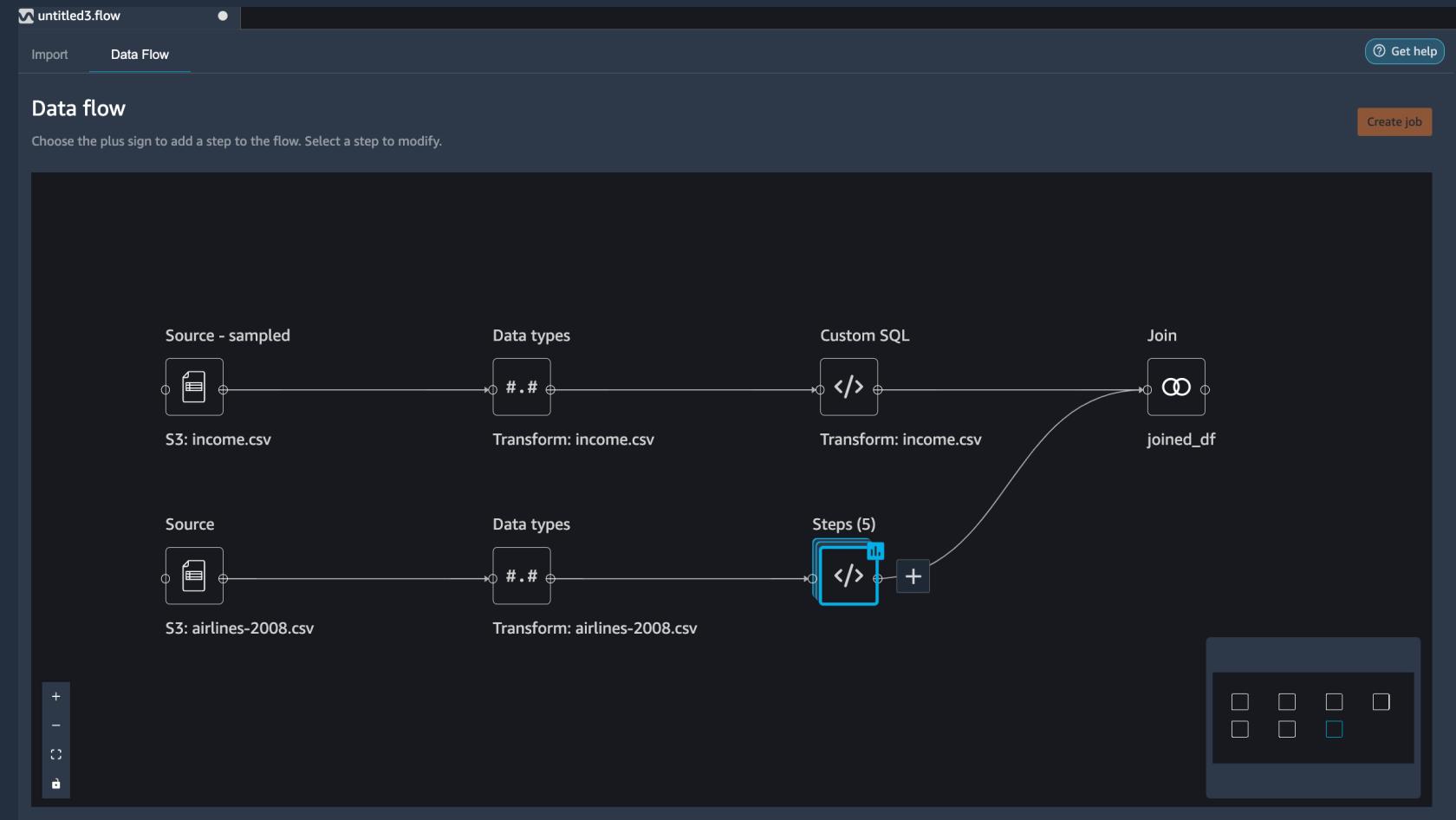
| Category                  | Transformation Type      | Description   | Link                       |
|---------------------------|--------------------------|---|----------------------------|
| Common Data Prep          | Custom transform         | Use Pyspark, Pandas, or Pyspark (SQL) to define custom transformations...                         | <a href="#">Learn more</a> |
|                           | Balance data             | Balance the data for binary classification problems using random oversampling...                  | <a href="#">Learn more</a> |
|                           | Custom formula           | Define a new column using a Spark SQL expression to query data in the dataset.                    | <a href="#">Learn more</a> |
|                           | Encode categorical       | Convert categorical variables to numeric or vector representations. <a href="#">Learn more</a>    | <a href="#">Learn more</a> |
|                           | Featurize date/time      | Encode date/time values to numeric and vector representations. <a href="#">Learn more</a>         | <a href="#">Learn more</a> |
|                           | Featurize text           | Generate vector representations from natural language text. <a href="#">Learn more</a>            | <a href="#">Learn more</a> |
|                           | Format string            | Clean and prepare strings using standard string formatting operations. <a href="#">Learn more</a> | <a href="#">Learn more</a> |
|                           | Group by                 | Add an aggregated column after group by as a new column.  |                            |
|                           | Handle missing           | Replace, drop, or add indicators for missing values. <a href="#">Learn more</a>                   | <a href="#">Learn more</a> |
|                           | Handle outliers          | Remove or replace outlier numeric and categorical values. <a href="#">Learn more</a>              | <a href="#">Learn more</a> |
| Machine Learning Specific | Handle outliers          | Remove or replace outlier numeric and categorical values. <a href="#">Learn more</a>              | <a href="#">Learn more</a> |
|                           | Handle structured column | Flatten JSON and perform other operations on structured data                                      |                            |
|                           | Manage columns           | Move, drop, duplicate or rename columns in the dataset. <a href="#">Learn more</a>                | <a href="#">Learn more</a> |
|                           | Manage rows              | Sort, shuffle or drop duplicate rows.   |                            |
|                           | Manage vectors           | Expand or create vector columns. <a href="#">Learn more</a>                                       | <a href="#">Learn more</a> |
|                           | Parse column as type     | Cast a column to a new data type. <a href="#">Learn more</a>                                      | <a href="#">Learn more</a> |
|                           | Process numeric          | Transform numeric values to improve machine learning model performance.                           |                            |
|                           | Search and edit          | Find, replace, split, and otherwise transform input string values using regular expressions.      |                            |
| Time Series               | Time Series              | Transformers to preprocess and manipulate time series. <a href="#">Learn more</a>                 | <a href="#">Learn more</a> |
|                           | Validate string          | Validate the format of string values using standard string functions. <a href="#">Learn more</a>  | <a href="#">Learn more</a> |

# Easily visualize the steps of your data processing pipeline

Data Wrangler records all the steps of data prep workflow in a data flow graph

Visualize the order of transformations, join and concatenate operators

Easily navigate data transformation flow, and modify and delete steps iteratively



# Amazon SageMaker Feature Store

SECURELY STORE, DISCOVER,  
AND SHARE FEATURES FOR ML



**Online and off-line**



**Millisecond latency**



**Consistent features**



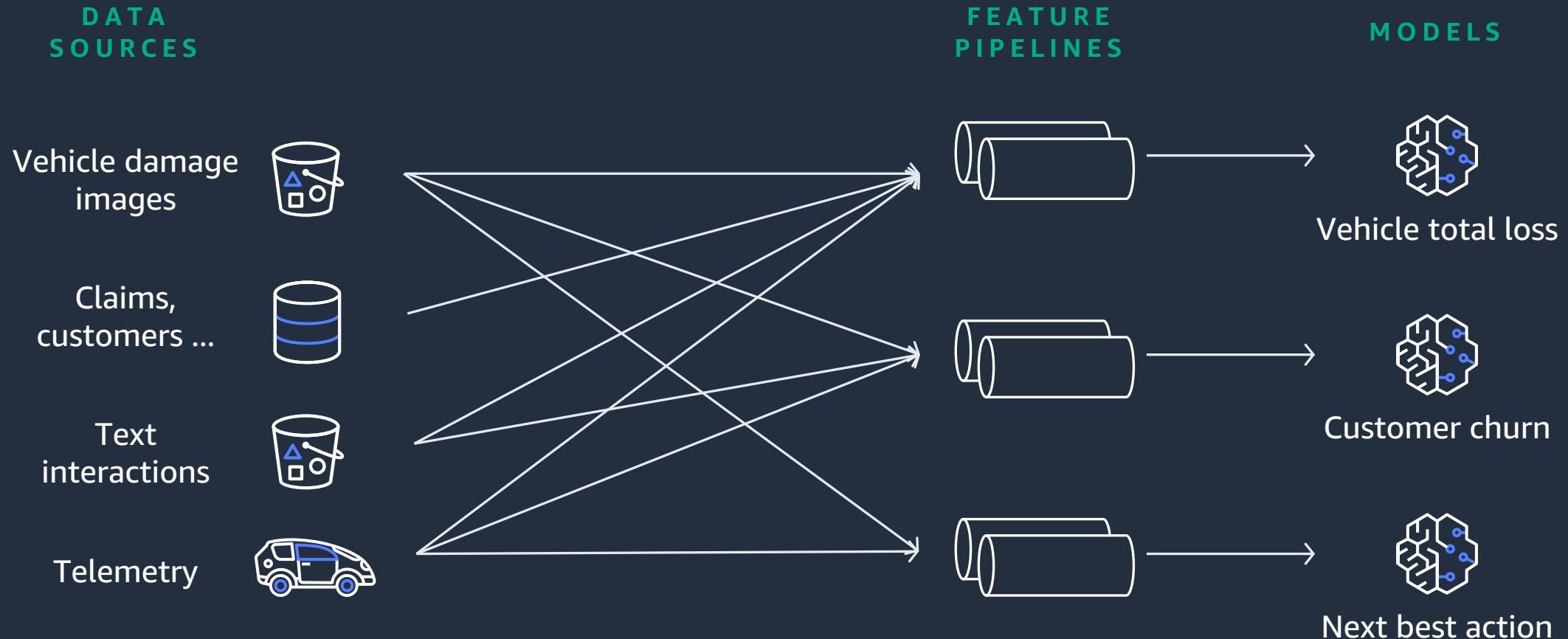
**Visual search**



**Sharing and collaboration**

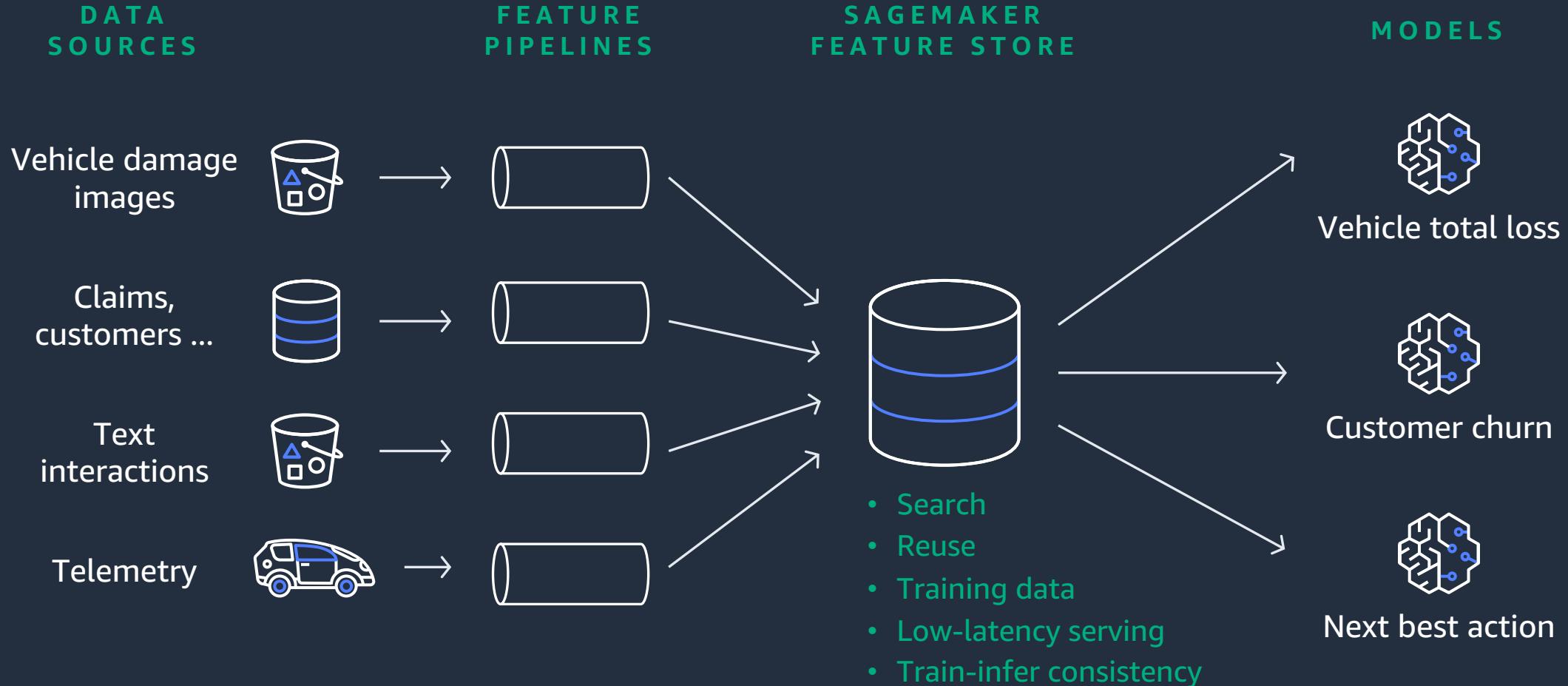
# Without a feature store ...

Standalone feature engineering for each new model



# With SageMaker Feature Store ...

Build features once, and reuse them across teams and models



# Support for separate feature stores



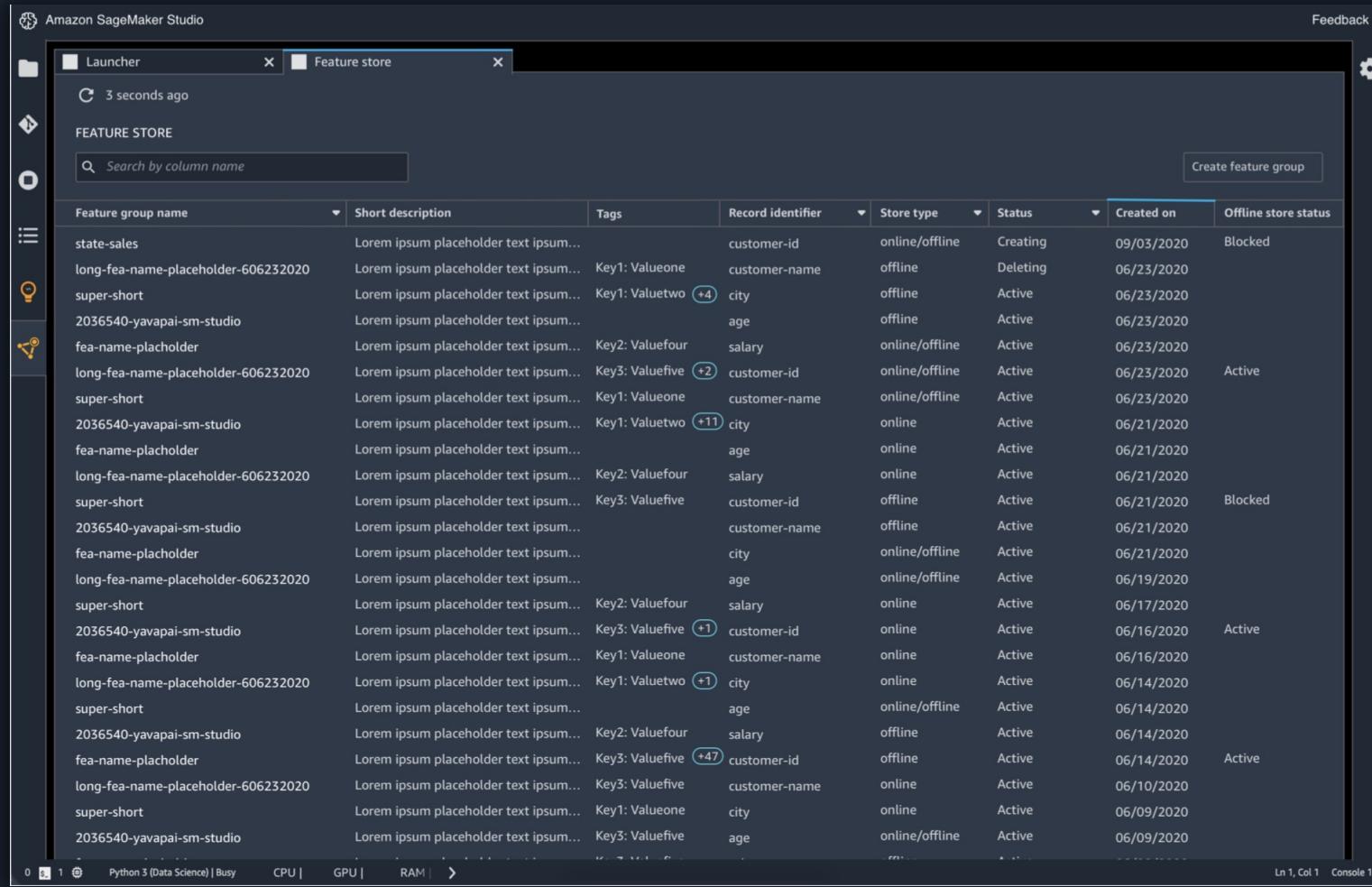
## Online feature store

- Primarily used for real time predictions
- Use cases such as real-time fraud detection
- Latest copy of feature data
- High throughput writes
- Low millisecond latency reads

## Offline feature store

- Primarily used for batch predictions and model training
- Historical record of feature data
- High throughput writes
- <15 minutes read after write consistency

# Search and discover features using Feature Store

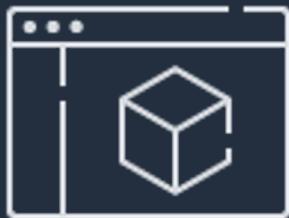


The screenshot shows the Amazon SageMaker Studio interface with the 'Feature store' tab selected. On the left, there's a sidebar with various icons and a search bar labeled 'Search by column name'. The main area is titled 'FEATURE STORE' and contains a table with the following columns: Feature group name, Short description, Tags, Record identifier, Store type, Status, Created on, and Offline store status. The table lists numerous feature groups, many of which have placeholder text in their descriptions and tags. Some entries show 'customer-id' as the record identifier, while others show 'age' or 'city'. The 'Status' column indicates various states like 'Creating', 'Deleting', 'Active', and 'Blocked'. The 'Created on' and 'Offline store status' columns provide dates and additional status information. A 'Create feature group' button is located in the top right corner of the table area.

- Search features individually or by groups visually with SageMaker Studio
- Discover features by name, description, tags, and other metadata
- Understand how features are grouped relevant to ML applications

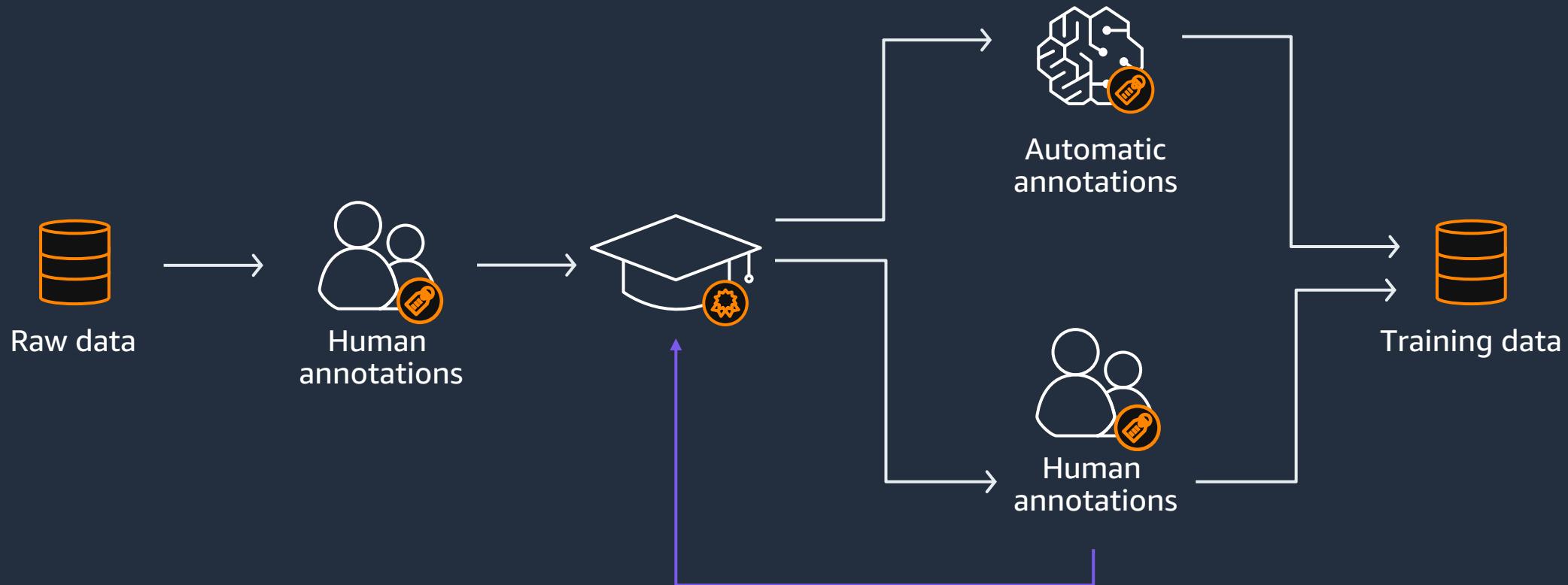
# Amazon SageMaker Ground Truth

Build highly accurate training datasets using machine learning



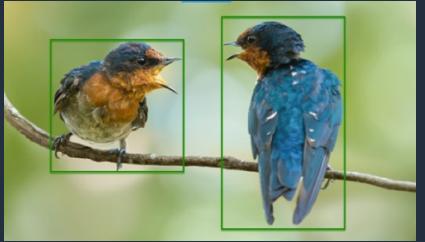
- Reduce data labeling costs by up to 70%
- Access labelers through Amazon Mechanical Turk, Amazon approved vendors, or use private human labelers
- Achieve accurate results quickly

# How Amazon SageMaker Ground Truth Works



# Turnkey data labeling service

70+ workflows



Bounding boxes



Image classification



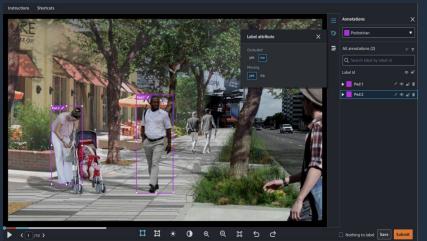
Segmentation



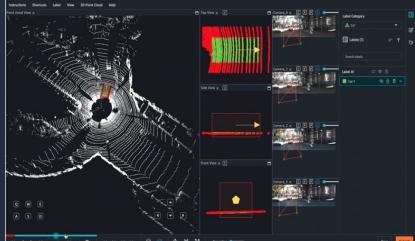
Label verification



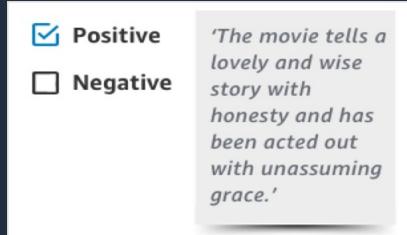
Custom template



Video



LIDAR 3D Point Cloud



Text classification

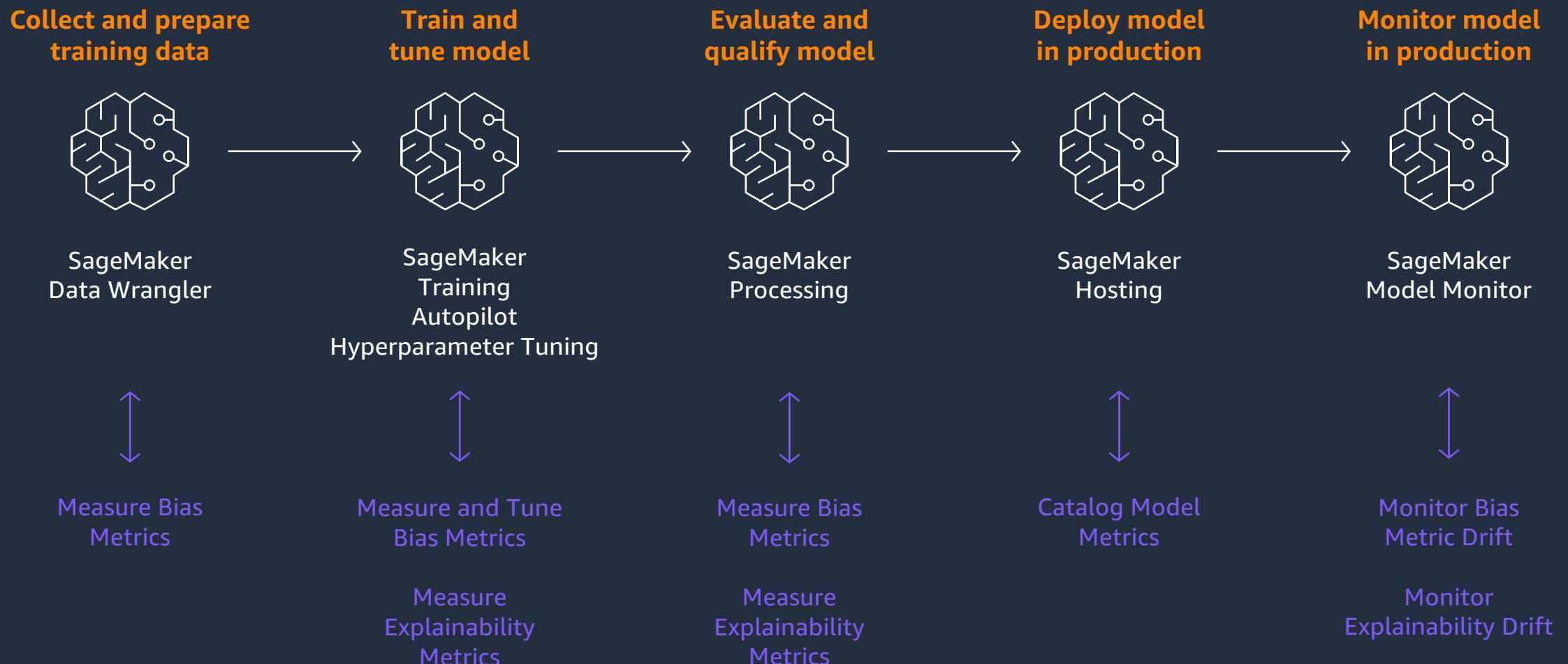


Named entity recognition

| COMPANY     | SAN JOAQUIN FACILITY MANAGEMENT   | STATE                          | CALIFORNIA |
|-------------|---|--------------------------------|------------|
| WELL        | RED RIBBON RANCH 90   |                                |            |
| FIELD/BLOCK | FRUITVALE   |                                |            |
| COUNTY      | KERN  |                                |            |
| API No.     | 04-030-55890  |                                |            |
| Location    | 631W AND 237N FROM THE SE CORNER OF SEC 27 - T29S - R27E MGR&M<br>NAD 83<br>LAT: 36.37503<br>LONG: -119.05665 |                                |            |
| STATE       |   |                                |            |
| COMPANY     |   | Other Services:<br>RWH<br>SFTI |            |
| WELL        |   |                                |            |
| FIELD/BLOCK |   |                                |            |
| COUNTY      |   |                                |            |
| Sect.       | 27  | Twp.                           | 29S        |
|             |   | Rge.                           | 27E        |

OCR, form, table

# SageMaker Clarify works across the ML lifecycle





**MLOps: Pipelines | Projects | Model Registry**  
Workflow automation, CI/CD for ML, central model catalog

**Canvas**  
Generate accurate machine learning predictions—no code required

**Studio | RStudio**  
Integrated development environment (IDE) for ML

**Governance**  
Model Cards | Dashboard | Permissions

# Model options



Training code



AWS Marketplace for  
Machine Learning



Amazon SageMaker  
AutoPilot

- XGBoost - Gradient Boosted Trees
- Matrix Factorization
- Regression
- Principal Component Analysis
- K-Means Clustering
- And More!

Built-in Algorithms  
No ML coding required



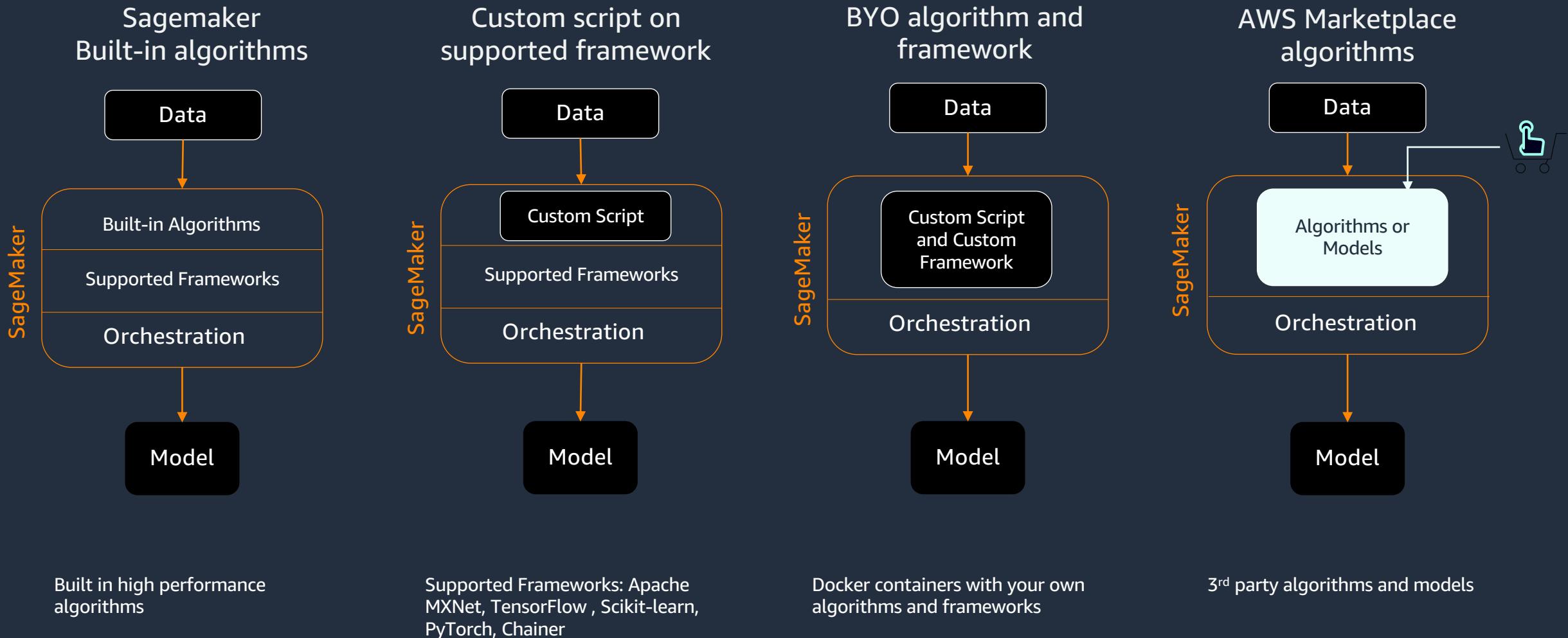
Bring Your Own Script  
Amazon SageMaker builds the container  
Open source containers



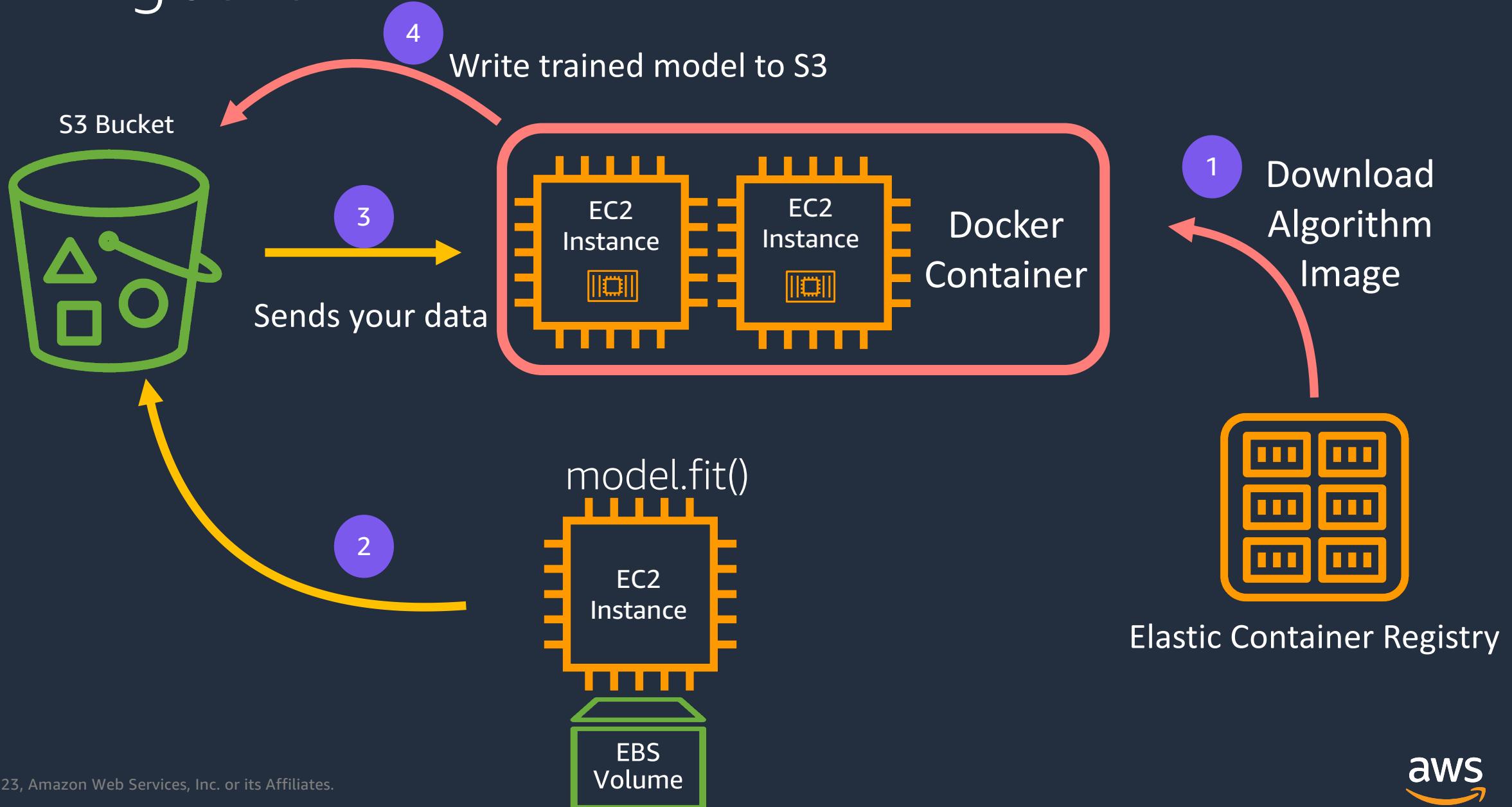
Bring Your Own Container  
Full control, you build the container  
R, C++, etc

Fully Managed, Distributed, Auto-Scaled, Secured

# Amazon SageMaker | Training



# Training Jobs



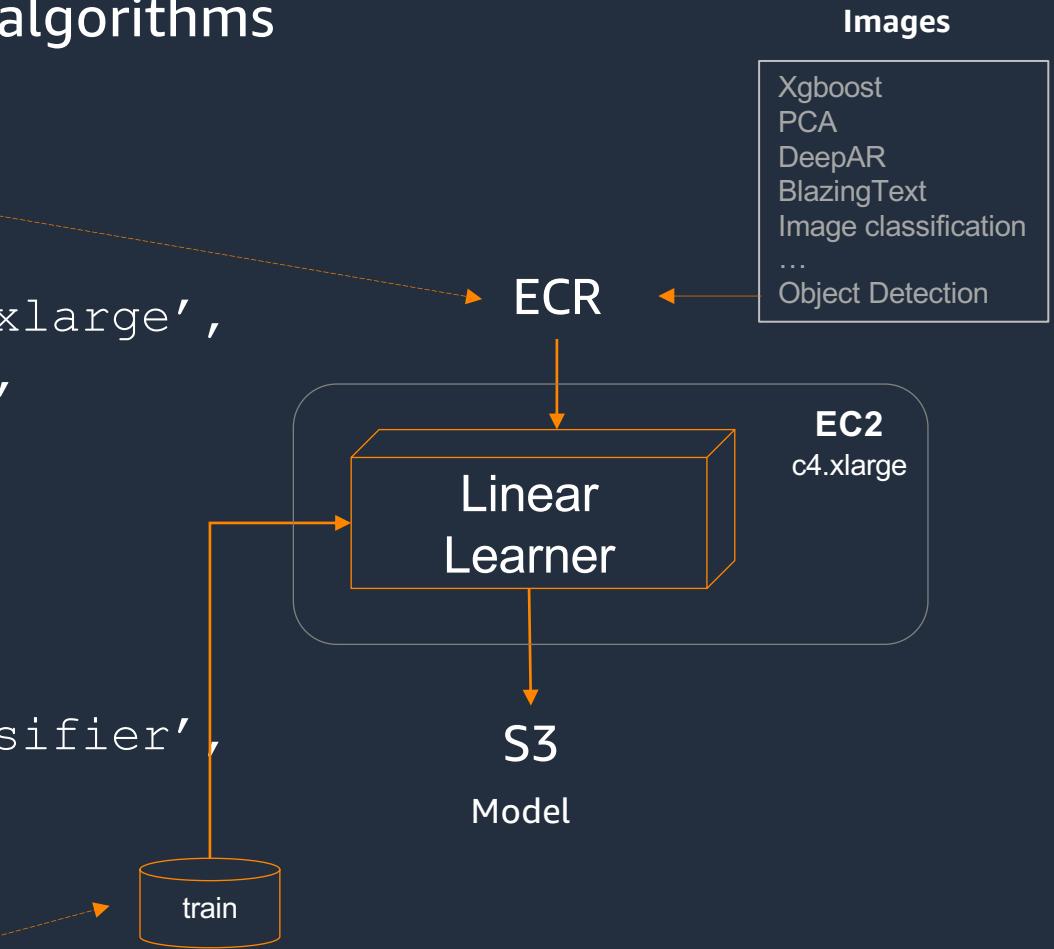
# Amazon SageMaker | Training

Use built-in algorithms

```
linear = Estimator('linear-learner',
                    train_instance_count=1,
                    train_instance_type='ml.c4.xlarge',
                    output_path=output_location,
                    sagemaker_session=sess)

linear.set_hyperparameters(
    feature_dim=784,
    predictor_type='binary_classifier',
    mini_batch_size=200)

linear.fit({'train': s3_train_data})
```



# Training – Sample Code

```
s3_input_train = sagemaker.inputs.TrainingInput(s3_data='s3://{}:{} /train'.format(bucket, prefix), content_type='csv')
s3_input_validation = sagemaker.inputs.TrainingInput(s3_data='s3://{}:{} /validation/'.format(bucket, prefix), content_type='csv')

sess = sagemaker.Session()

xgb = sagemaker.estimator.Estimator(container,
                                      role,
                                      instance_count=1,
                                      instance_type='ml.m4.xlarge',
                                      output_path='s3://{}:{} /output'.format(bucket, prefix),
                                      sagemaker_session=sess)

xgb.set_hyperparameters(max_depth=5,
                        eta=0.2,
                        gamma=4,
                        min_child_weight=6,
                        subsample=0.8,
                        silent=0,
                        objective='binary:logistic',
                        num_round=100)

xgb.fit({'train': s3_input_train, 'validation': s3_input_validation})
```

```
In [ ]: bt_model = sagemaker.estimator.Estimator(container,
                                                role,
                                                train_instance_count=1,
                                                train_instance_type='ml.c4.4xlarge',
                                                train_volume_size = 30,
                                                train_max_run = 360000,
                                                input_mode= 'File',
                                                output_path=s3_output_location,
                                                sagemaker_session=sess)
```

Number of EC2 instances

Type of EC2 instances

Disk space

## Algorithm Container

```
In [ ]: bt_model = sagemaker.estimator.Estimator(container,
                                                role,
                                                train_instance_count=1,
                                                train_instance_type='ml.c4.4xlarge',
                                                train_volume_size = 30,
                                                train_max_run = 360000,
                                                input_mode= 'File',
                                                output_path=s3_output_location,
                                                sagemaker_session=sess)
```

## SageMaker Estimator

## Execution Role

Cluster comes online

Logs to CloudWatch

Monitor via console or  
notification stream

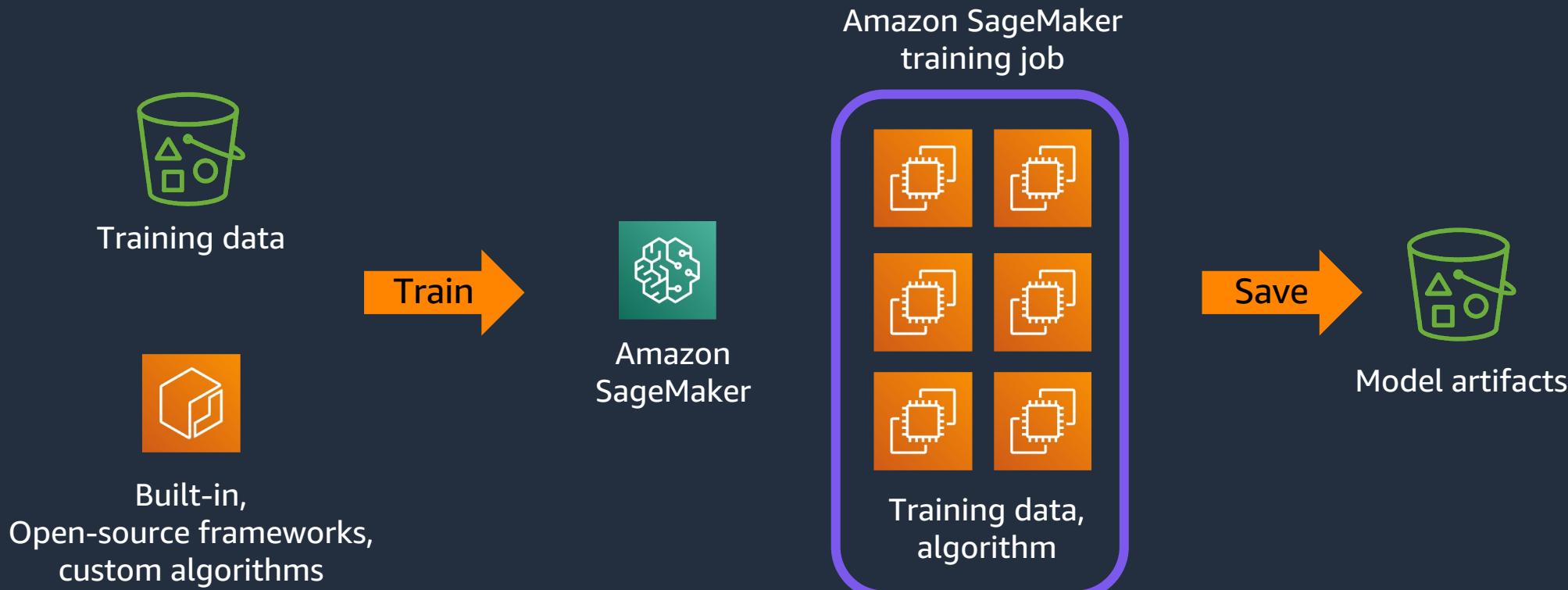


# Loading training data from Amazon S3

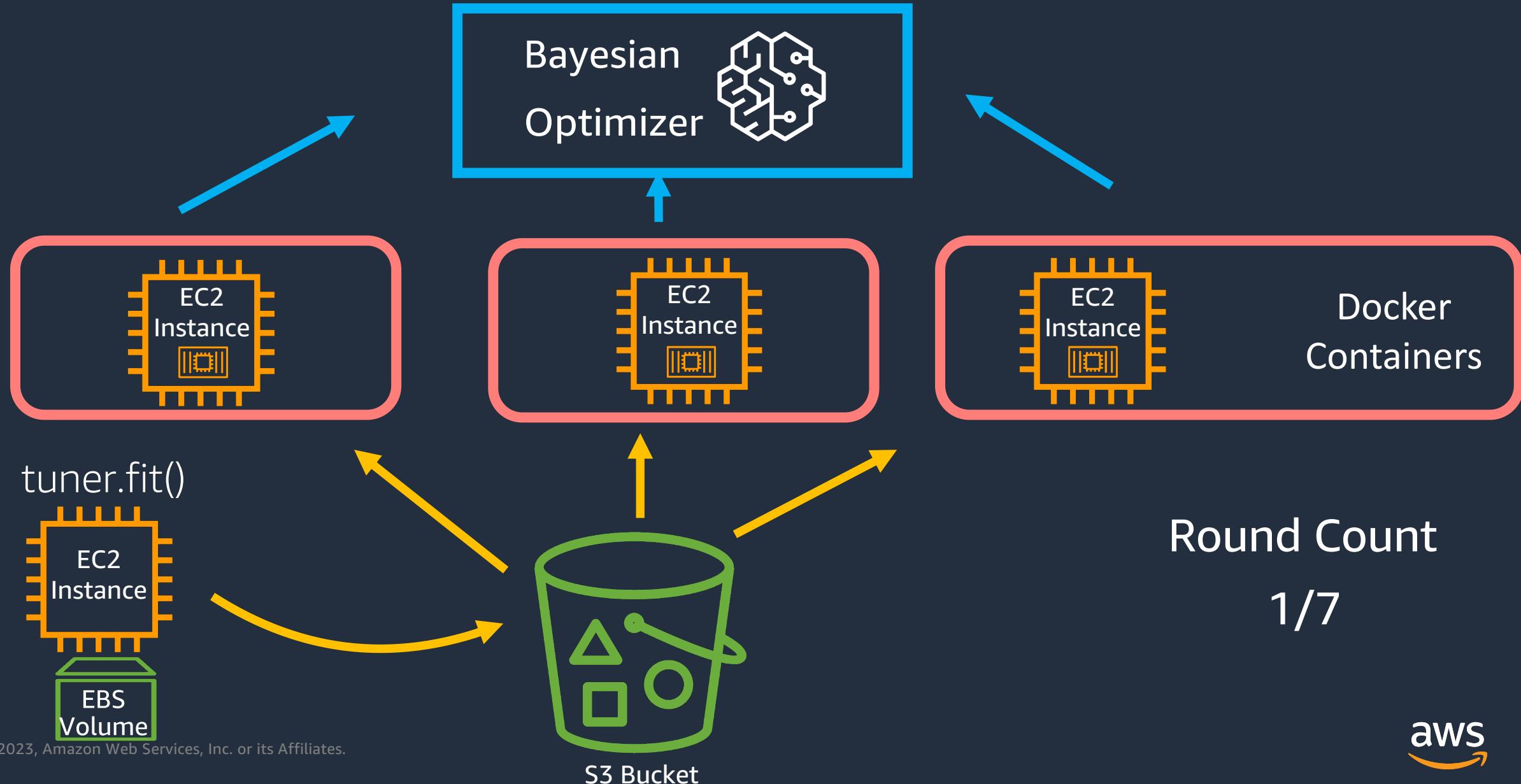
- Two modes: File Mode and Pipe Mode
  - *input\_mode* parameter in *sagemaker.estimator.estimator*
- File Mode **copies** the dataset to training instances
  - You need to provision enough storage
  - *S3DataSource* object
  - *S3DataDistributionType*: *FullyReplicated* | *ShardedByS3Key*
  - Different data formats are supported: CSV, protobuf, JSON, libsvm (check algo docs!)
- Pipe Mode **streams** the dataset to training instances
  - This allows you to process infinitely large datasets
  - Training starts faster
  - This mode is supported by some built-in algos as well as TensorFlow
  - Your dataset must be in **recordIO-encoded protobuf** format

# Fully managed training

## Billable by the second, 100% utilization



# Hyperparameter Tuning Jobs



## 1. Pick hyperparameters and ranges

```
: hyperparameter_ranges = {'eta': ContinuousParameter(0, 1),  
                           'min_child_weight': ContinuousParameter(1, 10),  
                           'alpha': ContinuousParameter(0, 2),  
                           'max_depth': IntegerParameter(1, 10)}
```

## 2. Pick objective metric

```
: objective_metric_name = 'validation:auc'
```

## 3. Pick job parameters

```
: tuner = HyperparameterTuner(xgb,  
                               objective_metric_name,  
                               hyperparameter_ranges,  
                               max_jobs=20,  
                               max_parallel_jobs=3)
```

# Deployment flexibility

## SageMaker Real Time Inference

- Creates a long-running microservice as an endpoint
- Effective for low latency for small data batches
- Can be deployed as a Multi-Model Endpoint

## SageMaker Batch Transform

- Bulk data transformation
- Periodic arrival of big chunks of data
- Configurable inference compute infrastructure

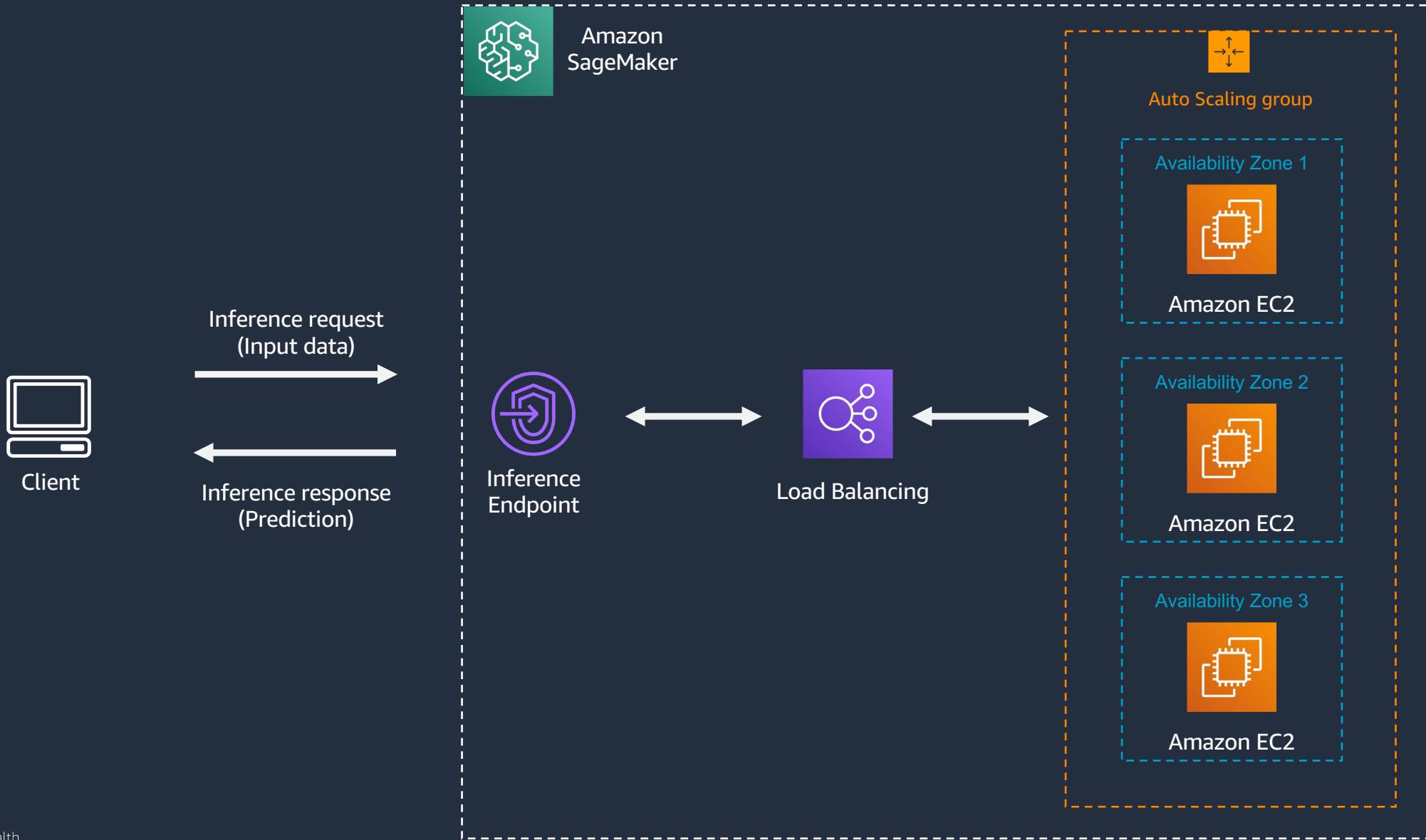
## SageMaker Serverless Inference

- Deploy without managing infrastructure
- Ideal for intermittent or unpredictable traffic

## SageMaker Asynchronous Inference

- Queues incoming requests
- Ideal for large payloads and long processing times
- Can scale down to zero instances, saving costs

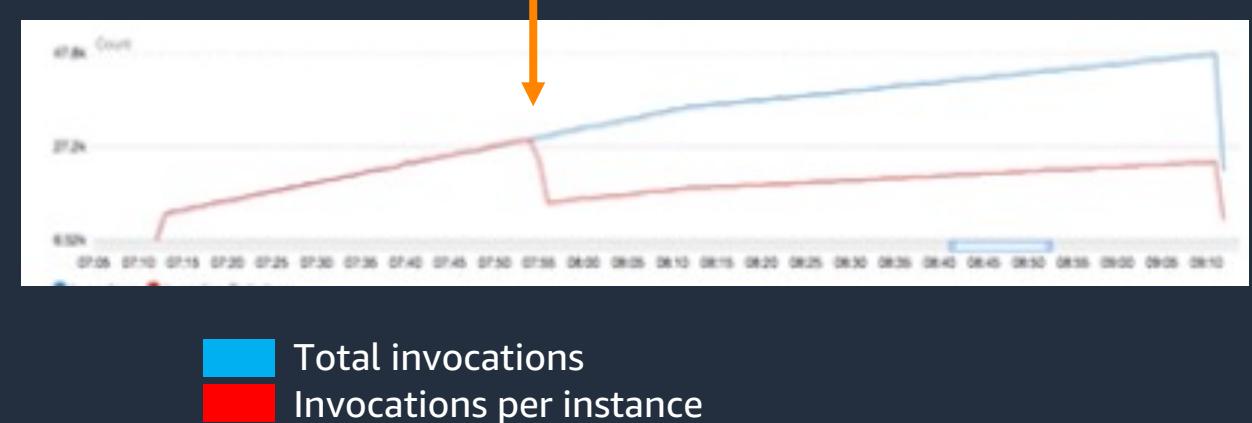
# Amazon SageMaker Real-time endpoint



# Autoscaling for real-time inference endpoints

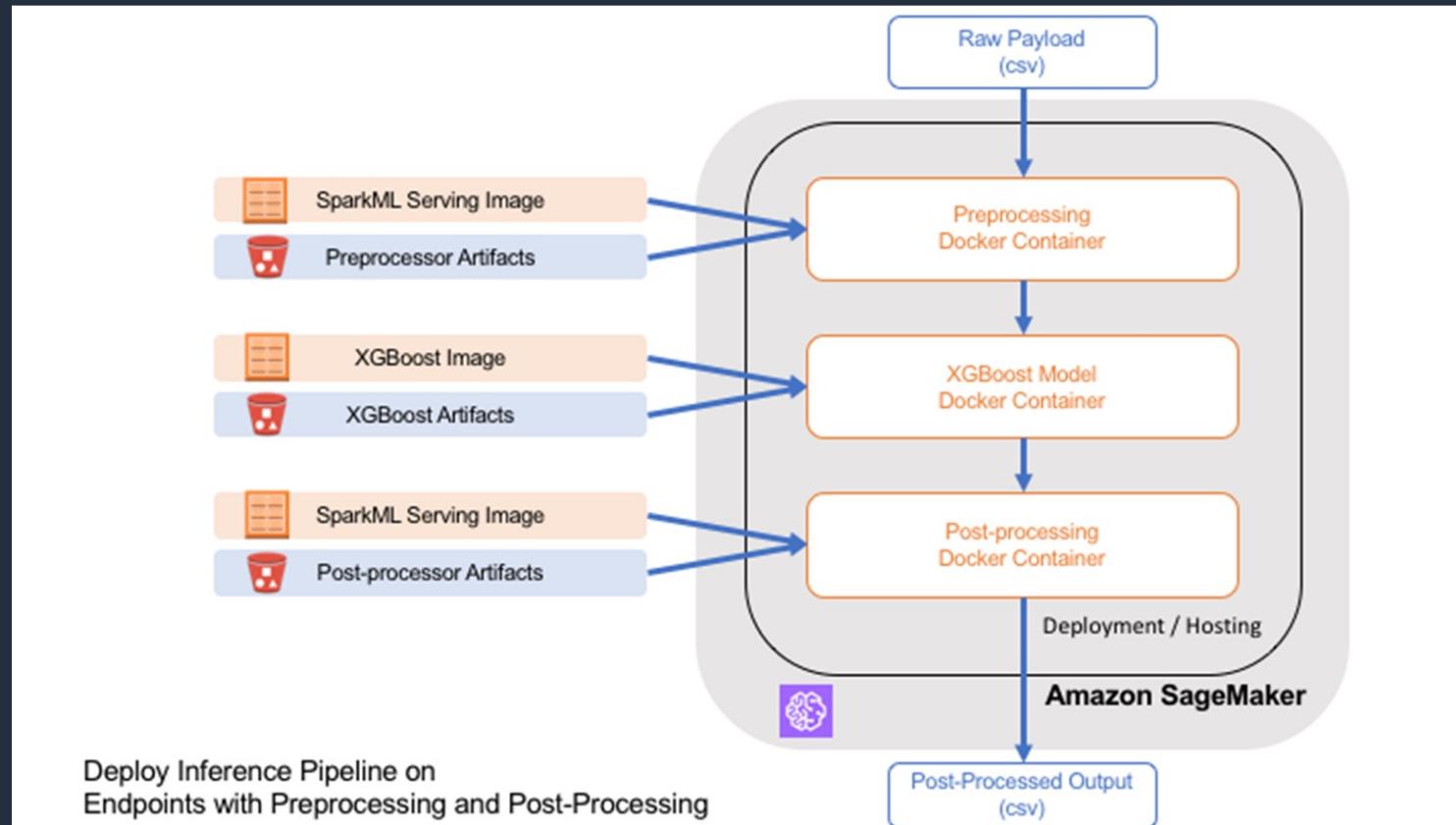
- Provision steady state capacity
- Set minimum, maximum instance counts, scaling criteria
- SageMaker auto-scales up to meet demand, and scales back down
- Significant cost savings

After scaling, invocations are split across more instances



# Inference Pipelines for sequential execution of models

Execute data processing on inference requests. Maintain single copy of data processing code for training and inference.

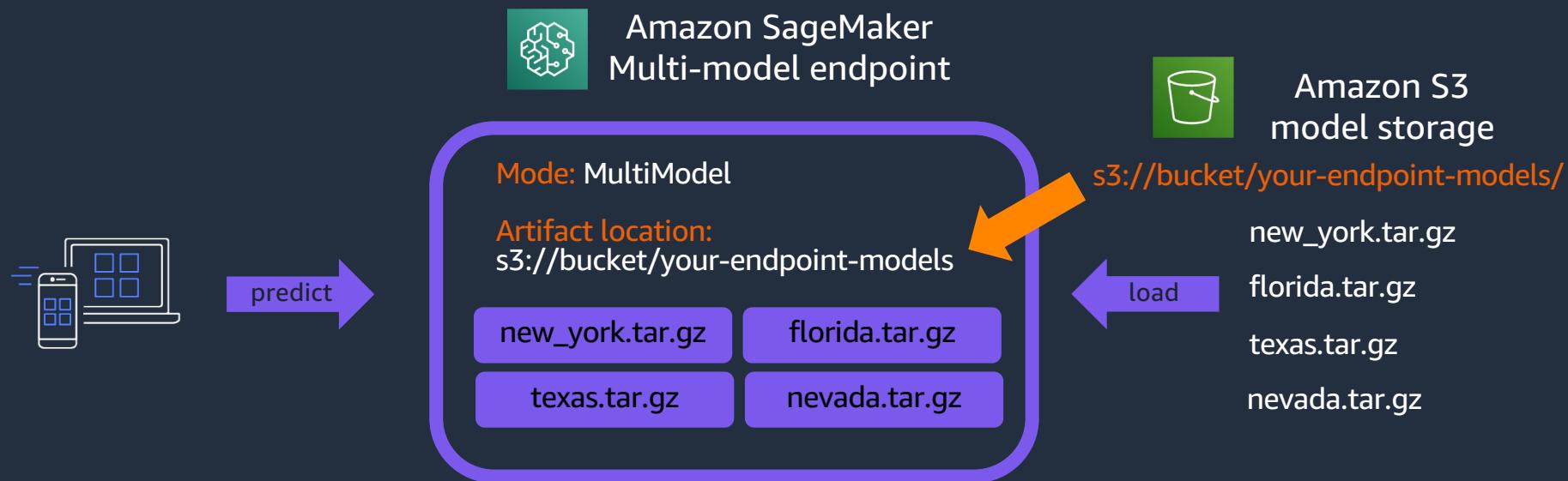


**Built-in containers –**  
Scikit-Learn and  
Apache Spark MLLib

Add Up to 5 containers;  
execute sequentially

Containers co-located on  
instances for low latency

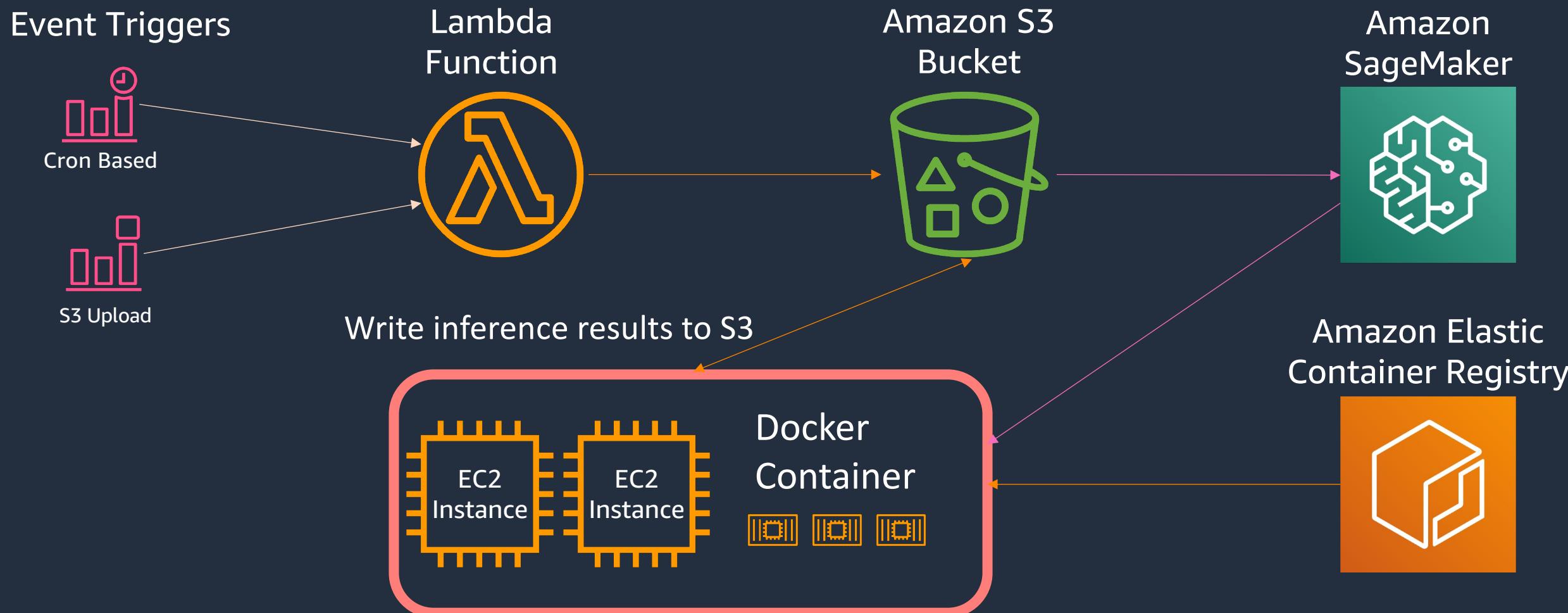
# Multi-Model Endpoints



## Key Capabilities

- Deploy tens to tens of thousands of models; target model for each inference request
- Two modes for model caching: Caching Enabled, Caching Disabled
- Support for built in algorithms, SageMaker frameworks, and custom models
- Support for inference pipelines and condition keys to restrict access to models

# Batch Transform Common Design Pattern



# Serverless Deployment - Fully managed offering



**Managed  
infrastructure**

Security  
Monitoring  
Logging  
Built in availability  
and fault tolerance



**Serverless**

No need to select instance types or provision capacity  
Choose memory options based on inference processing needs



**Automatically scale  
out, in, and down to 0**

No need to set scaling policies

# Pro Tips

- Turn your endpoints off when not in use with Lambda
- Use inference pipelines for pre and post processing
- Bring more models in your own Docker container
- Consider the size of data hitting a single endpoint
- You can train your model elsewhere and host in SageMaker

# Lab: Train, Tune and Deploy model using SageMaker Built-in Algorithm

Prerequisites + “Lab 2. Train, Tune and Deploy XGBoost”

Q&A and chat are open!

Until 14:30

# Bring Your Own Custom Models

# Model Options

## In this session: Bring Your Own Script/Container



Training code



AWS Marketplace for  
Machine Learning



Amazon SageMaker  
AutoPilot

- XGBoost - Gradient Boosted Trees
- Matrix Factorization
- Regression
- Principal Component Analysis
- K-Means Clustering
- And More!

Built-in Algorithms (17)  
No ML coding required



Bring Your Own Script  
Amazon SageMaker builds the container  
Open source containers



Bring Your Own Container  
Full control, you build the container  
R, C++, etc

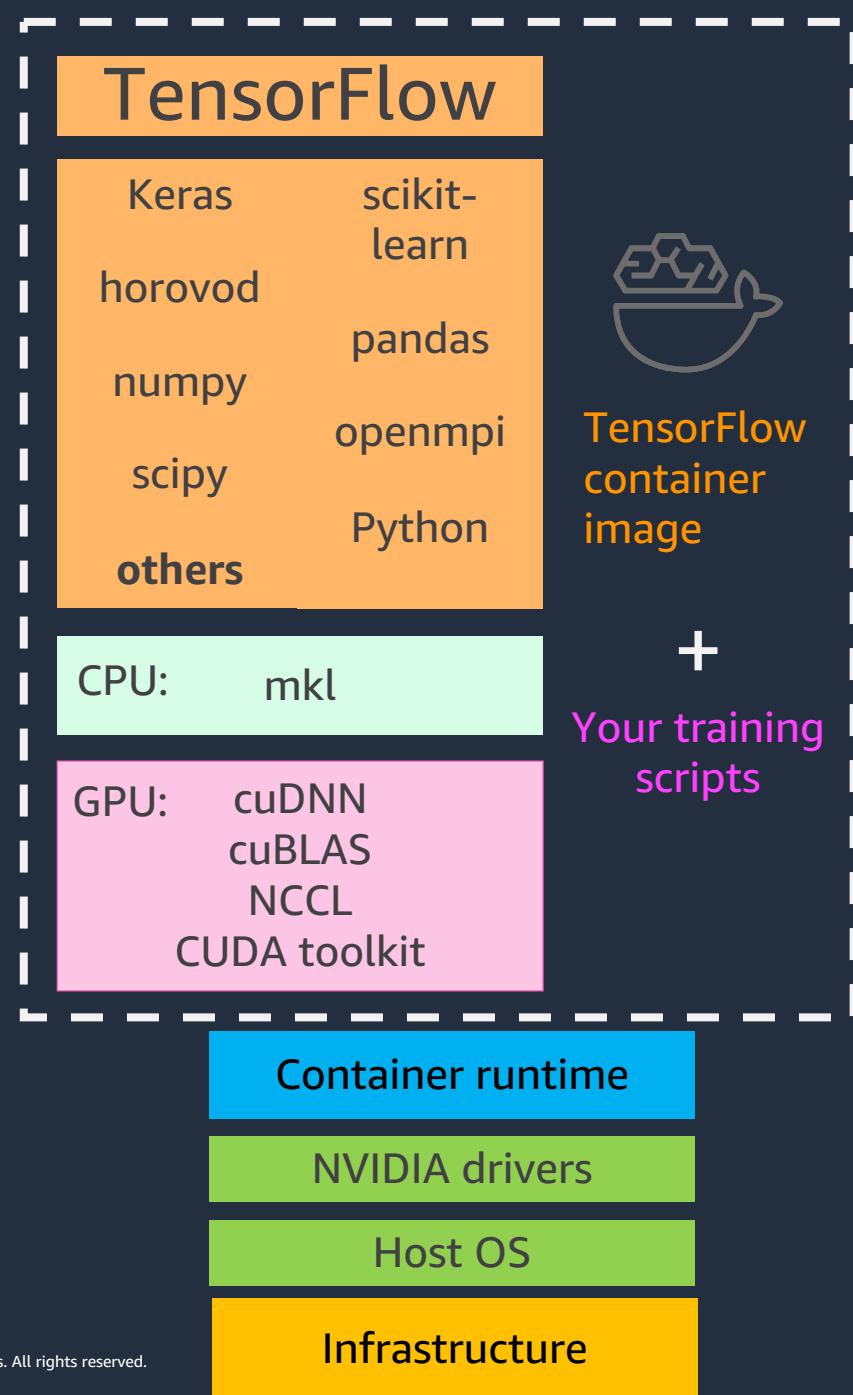
Fully Managed, Distributed, Auto-Scaled, Secured



# Containers for Machine Learning



# Containers for machine learning



ML environments that are:

- Lightweight
- Portable
- Scalable
- Consistent

Packages:

- Training code
- Dependencies
- Configurations

# TensorFlow

Keras      scikit-  
horovod    learn  
numpy      pandas  
scipy      openmpi  
**others**     Python

CPU: mkl

GPU: cuDNN  
      cuBLAS  
      NCCL  
      CUDA toolkit



TensorFlow  
container  
image

+

Your training  
scripts

Container runtime

NVIDIA drivers  
Host OS



Development system

push

pull



Amazon ECR

# TensorFlow

Keras      scikit-  
horovod    learn  
numpy      pandas  
scipy      openmpi  
**others**     Python

CPU: mkl

GPU: cuDNN  
      cuBLAS  
      NCCL  
      CUDA toolkit



TensorFlow  
container  
image

+

Your training  
scripts

Container runtime

NVIDIA drivers  
Host OS

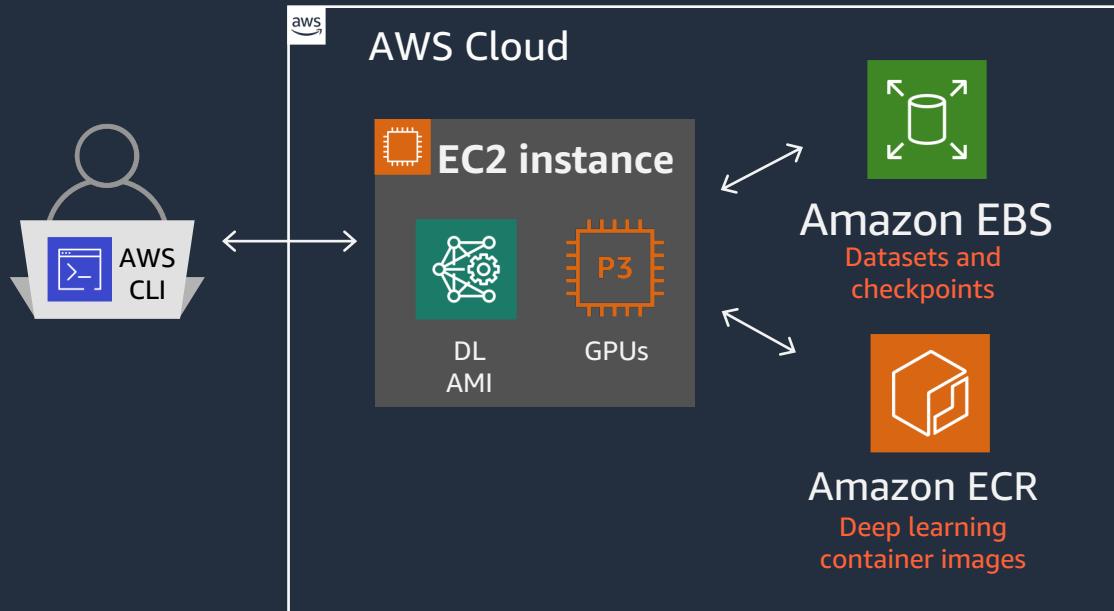
Training cluster

Container  
registry

# AWS Deep Learning Containers

- Prepackaged Docker container images fully configured and validated
- Optimized for performance with latest NVIDIA driver, CUDA libraries, and Intel libraries
- Consistent and reproducible deployment and lightweight
- Optimized for distributed machine learning
- Runs on Amazon ECS, Amazon EKS and **Amazon SageMaker**

# Machine Learning Setups on AWS Today



Code and dependencies

Collaborative development

Performance optimizations

Scaling

Infrastructure management

# ML Infrastructure and Cluster Management

## ML services

Fully managed service that covers the entire machine learning workflow

### Amazon SageMaker



## Management

Deployment, scheduling, scaling, and management of containerized applications



Amazon Elastic Container Service  
(Amazon ECS)



Amazon Elastic Kubernetes Service  
(Amazon EKS)



Kubeflow

## Image registry

Container image repository



Amazon Elastic Container Registry  
(Amazon ECR)

## Compute

Where the containers run



Amazon EC2

# 3 Ways to Train Using Amazon SageMaker

## Use built-in algorithms

(Bring Your Own Data)

- K-Means Clustering
- Principal Component Analysis
- Neural Topic Modelling
- Factorization Machines
- Linear Learner (Regression)
- BlazingText
- Reinforcement learning
- XGBoost
- Topic Modeling (LDA)
- Image Classification
- Seq2Seq
- Linear Learner (Classification)
- DeepAR Forecasting



© 2023, Amazon Web Services, Inc. or its affiliates. All rights reserved.

## Use deep learning frameworks

(Bring Your Own Data)

(Bring your own training script)



## Use custom containers

(Bring your own data)

(Bring your own container)

### Custom container

```
import tensorflow as tf
import numpy as np
import argparse
from os import path
from tensorflow import keras
from tensorflow.keras.layers import Input
from tensorflow.keras.models import Model
from tensorflow.keras.utils import multi_gpu_model
from tensorflow.keras.optimizers import Adam
HEIGHT = 32
WIDTH = 32
NUM_DEPTHS = 3
NUM_CLASSES = 10
```

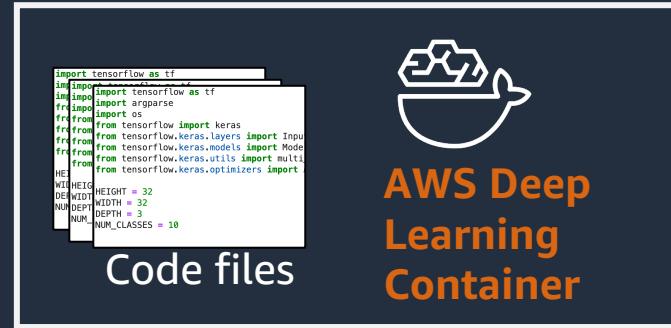
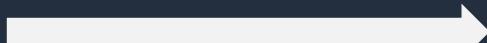
Code files



# Amazon SageMaker

## Bring your own Script

# High Level Workflow



```
import tensorflow as tf
import tensorflow.keras
from tensorflow import keras
from tensorflow.keras import layers
from tensorflow.keras import Input
from tensorflow.keras.models import Model
from tensorflow.keras.optimizers import Adam
from tensorflow.keras.utils import multi_gpu_model
HEIGHT = 32
WIDTH = 32
DEPTH = 32
NUM_DEP = 3
NUM_CLASSES = 10
```

Code files



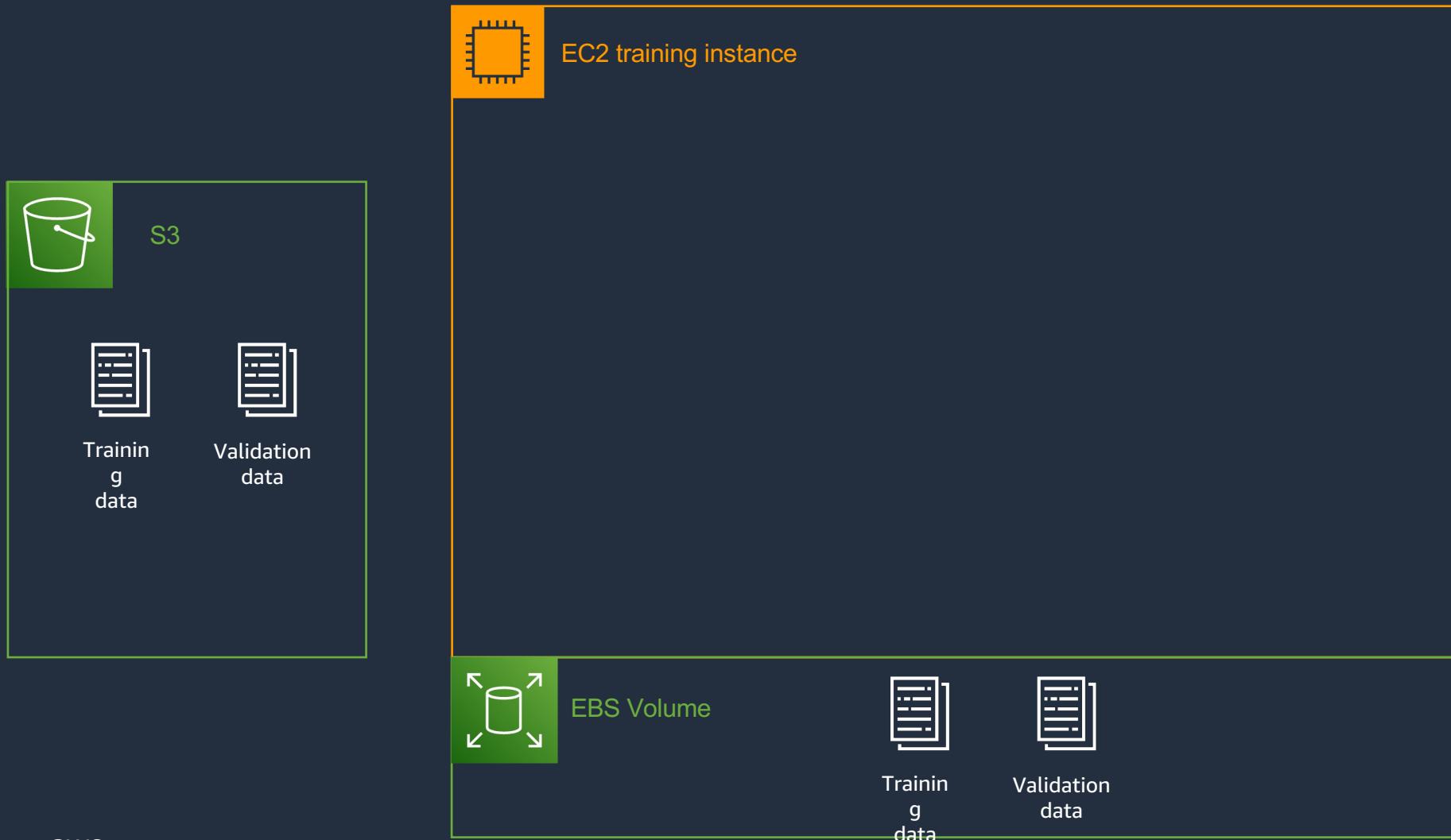
Container registry

```
from sagemaker.tensorflow import TensorFlow

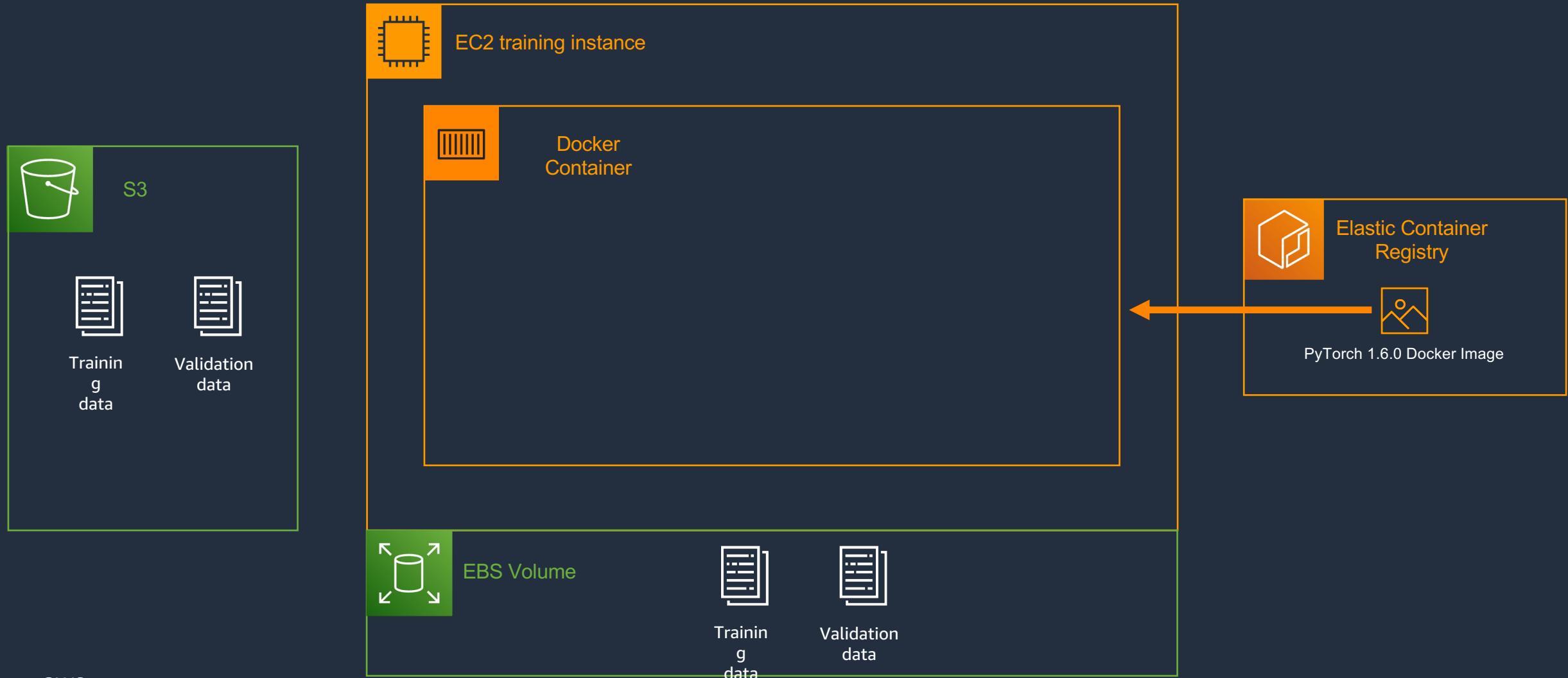
cifar10_estimator = TensorFlow(entry_point='source/cifar10.py',
                                role=role,
                                framework_version='1.14',
                                train_instance_count=1,
                                train_instance_type='ml.p3.xlarge')
```



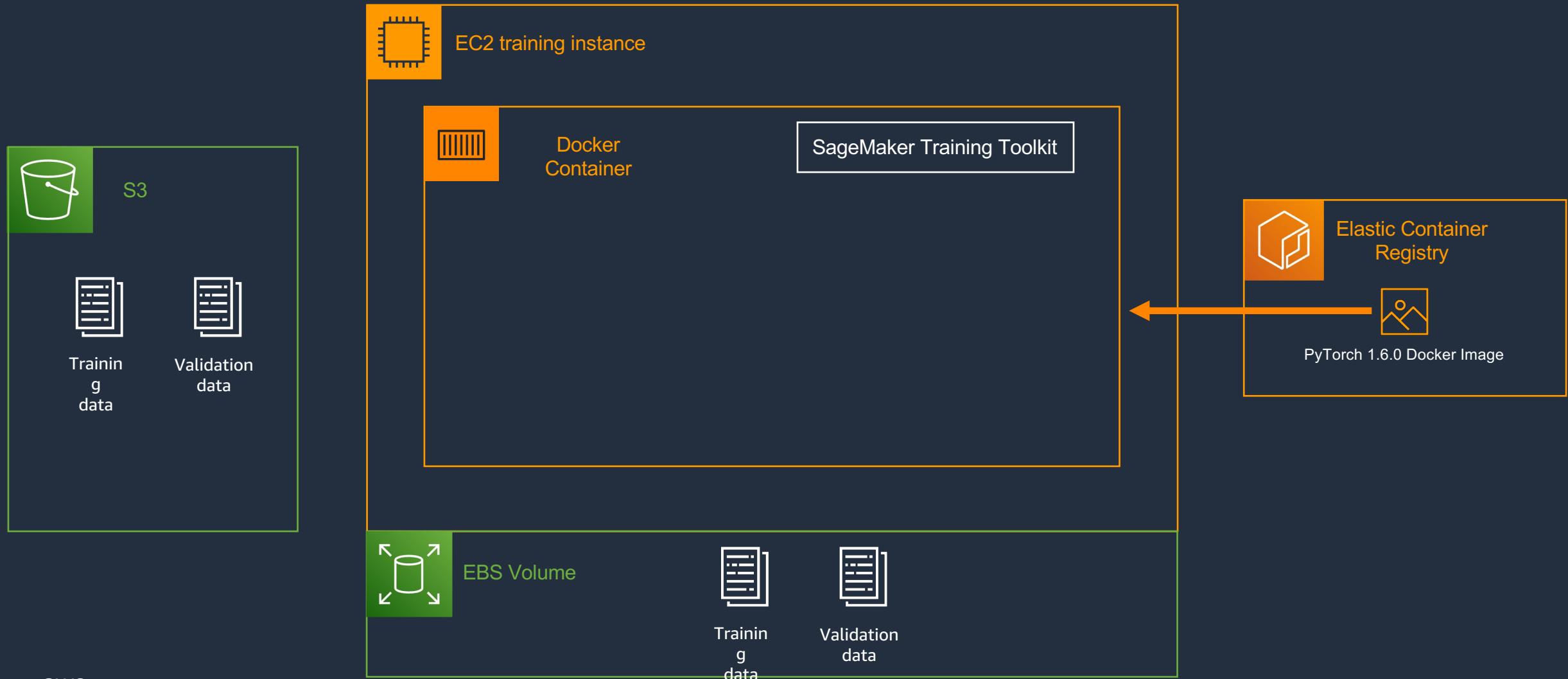
# Copy the data from S3 to EBS volume of the EC2 instance



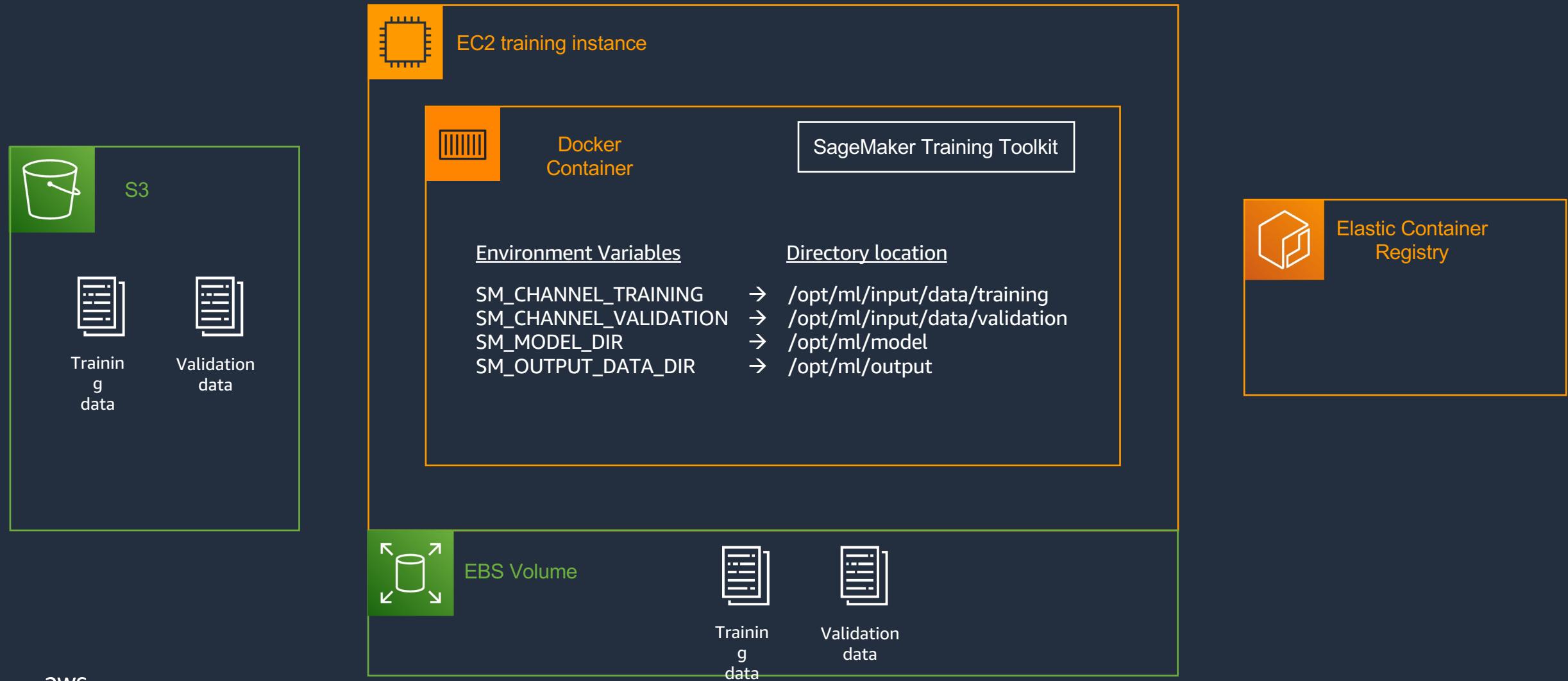
# Docker Container



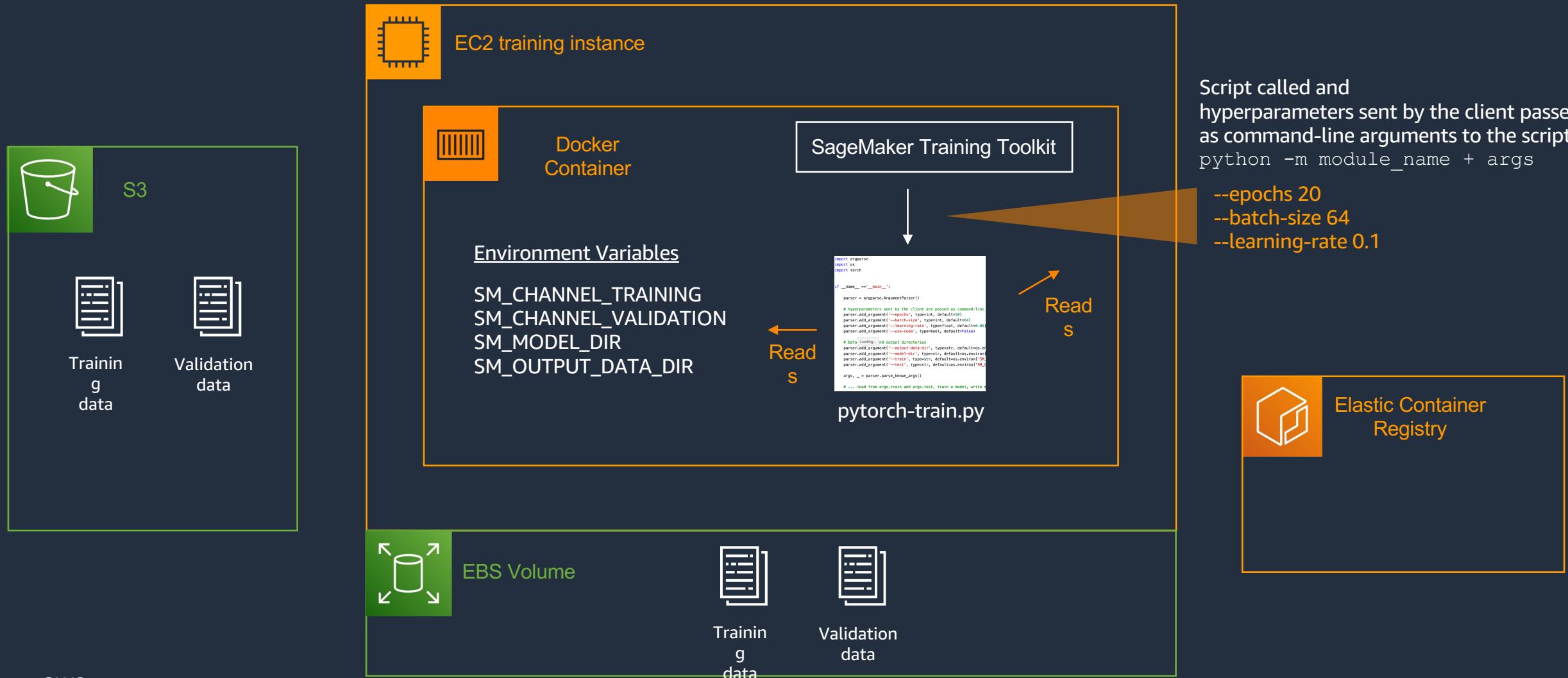
# SageMaker Training Toolkit



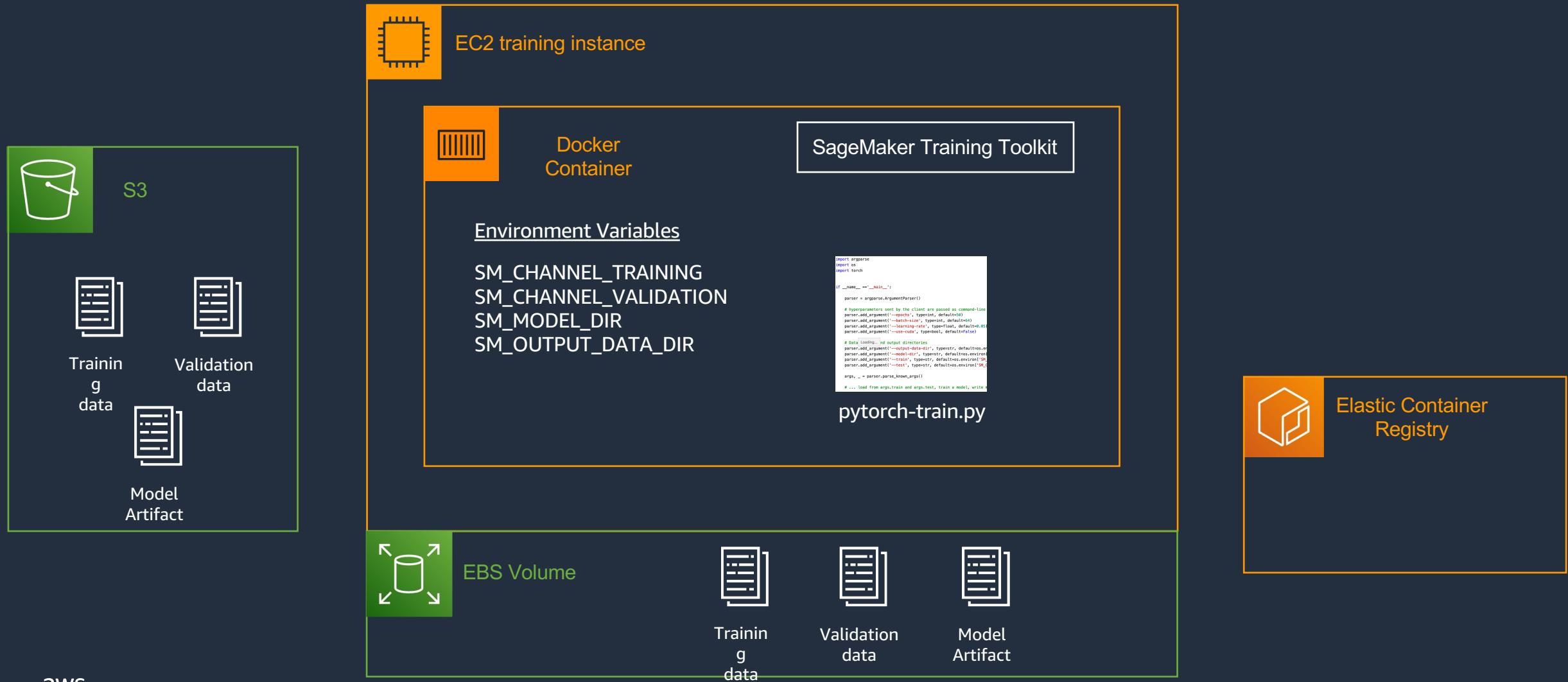
# Environment variables and directory location



# Your script is called with arguments present in the Estimator



# Copy of the Model Artifact to S3 after training





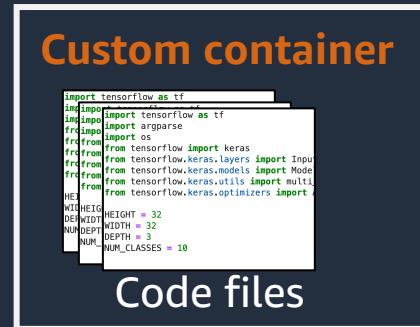
# Amazon SageMaker

## Bring your own Container

# High Level Workflow



Docker build



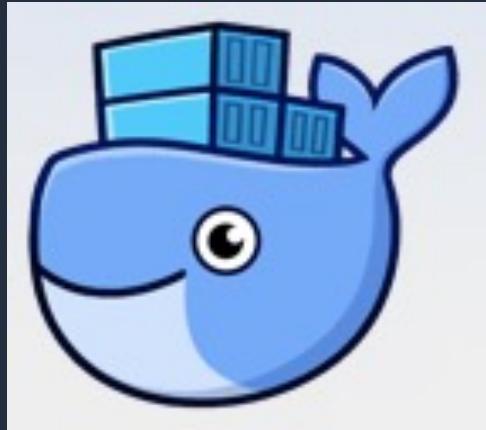
Container registry

```
ic = sagemaker.estimator.Estimator(training_image,
                                     role,
                                     train_instance_count=1,
                                     train_instance_type='ml.p3.xlarge',
                                     train_volume_size = 50,
                                     train_max_run = 360000,
                                     input_mode= 'File',
                                     output_path=s3_output_location,
                                     sagemaker_session=sess)
```



Amazon Simple Storage Service

# Bring Your Own Docker File

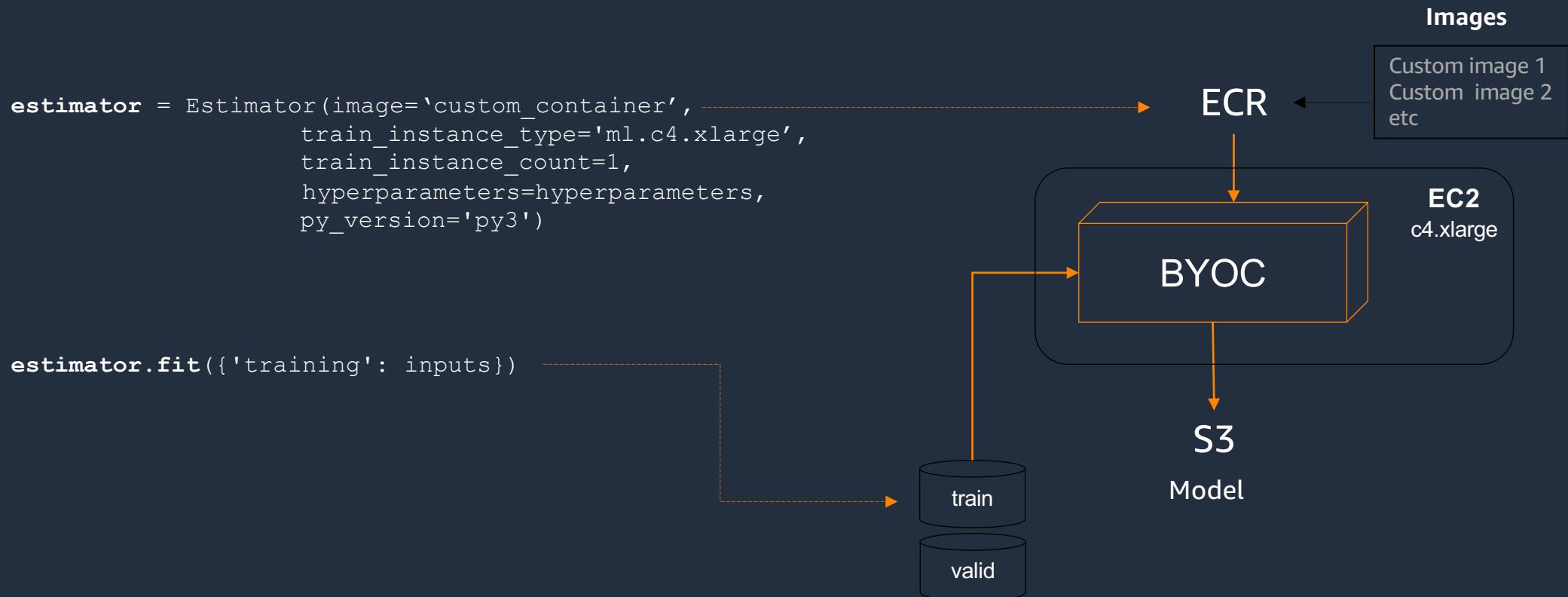


Customer Managed

1. Write your model however you please
2. Point to your model within your Docker file
3. Register your container on ECR
4. Point to your container's address in ECR
5. Don't forget to implement a serve() function!

# Amazon SageMaker | Training

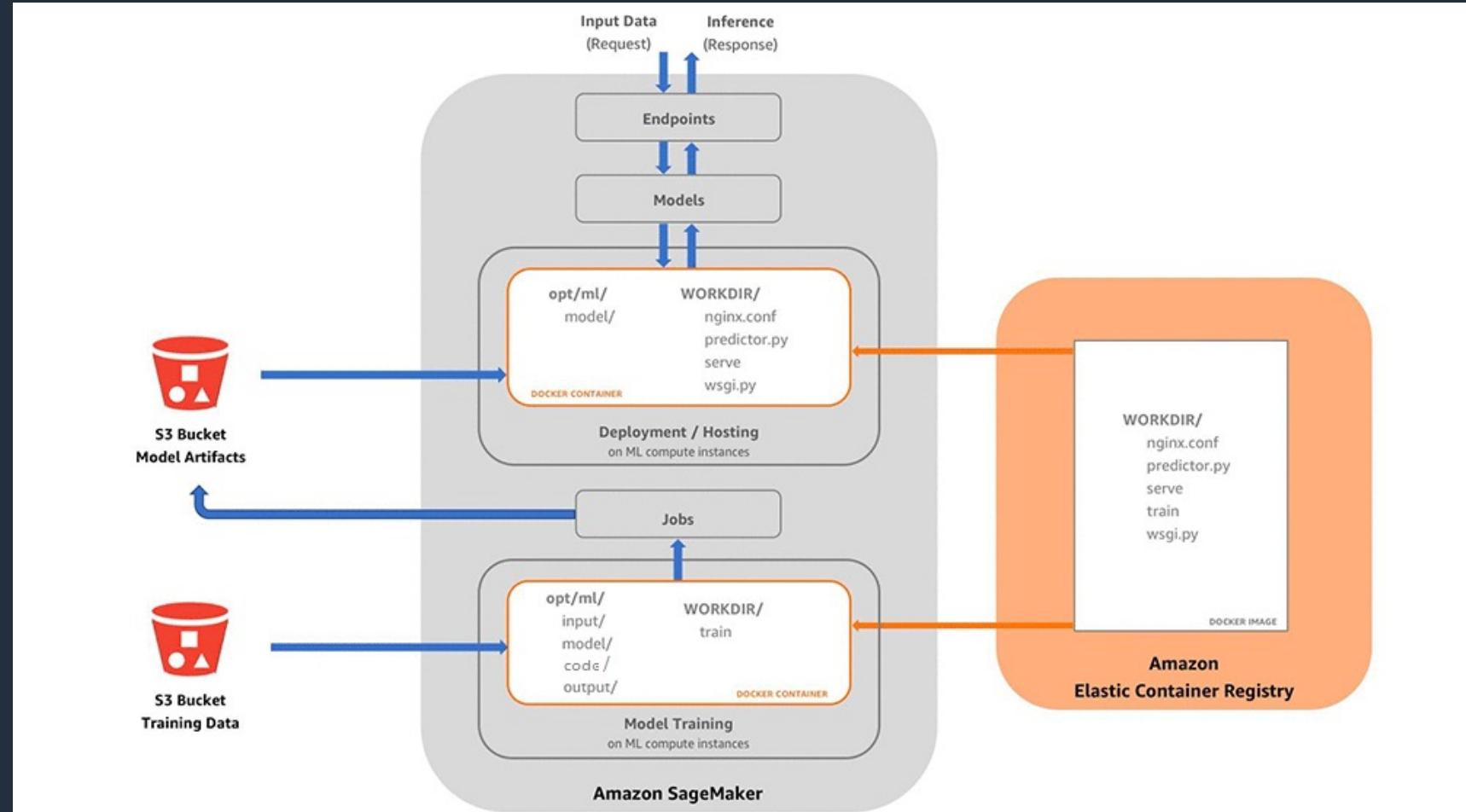
Bring your own container



# Bring Your Own Container

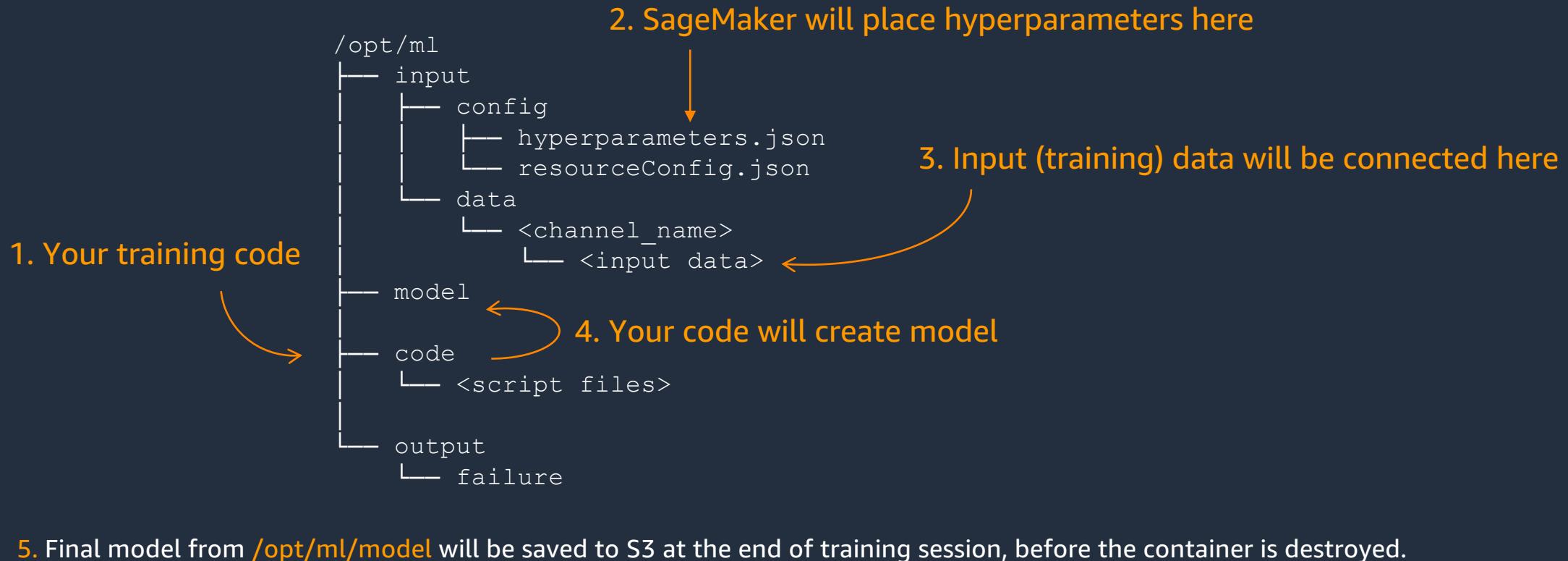
Custom Docker image:

- Training image
- Inference (serving) image



# Amazon SageMaker | Training

Bring your own container



# ML Marketplace

## Algorithms

- You can train on your data

## Models

- Use a pretrained model artifact

- Subscription model
- Free tier!

230+ solutions!

Categories for algorithms and models for Amazon SageMaker



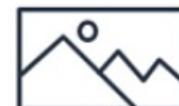
Computer Vision



Natural Language Processing



Speech Recognition



Image



Text



Structured



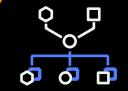
Audio



Video

## Model parallelism

Automated and efficient model partitioning

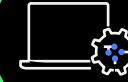


## Data parallelism

Reduced training time



Minimal code change



Optimized for AWS network



Optimized for training large models



Support for popular ML framework APIs

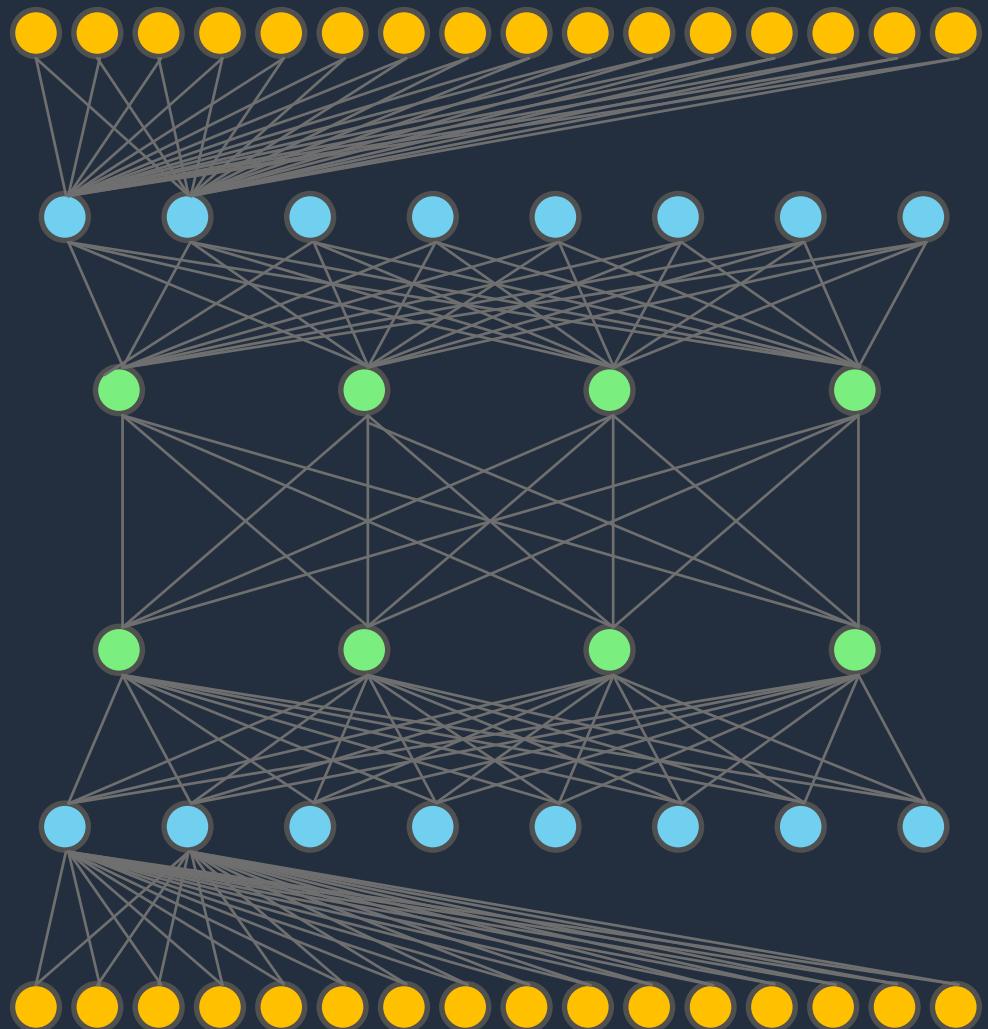


# Distributed training on Amazon SageMaker

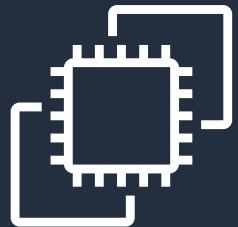
<https://aws.amazon.com/sagemaker/distributed-training/>

# DataParallel in SageMaker

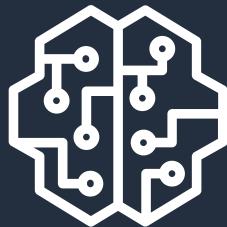
- Library for distributed training of deep learning models in TensorFlow and PyTorch
- Accelerates training for network-bound workloads
- Built and optimized for AWS network topology and hardware
- 20%–40% faster and cheaper; best performance on AWS



# Model parallelism on SageMaker



Efficient  
pipelined training



Automated  
model partitioning

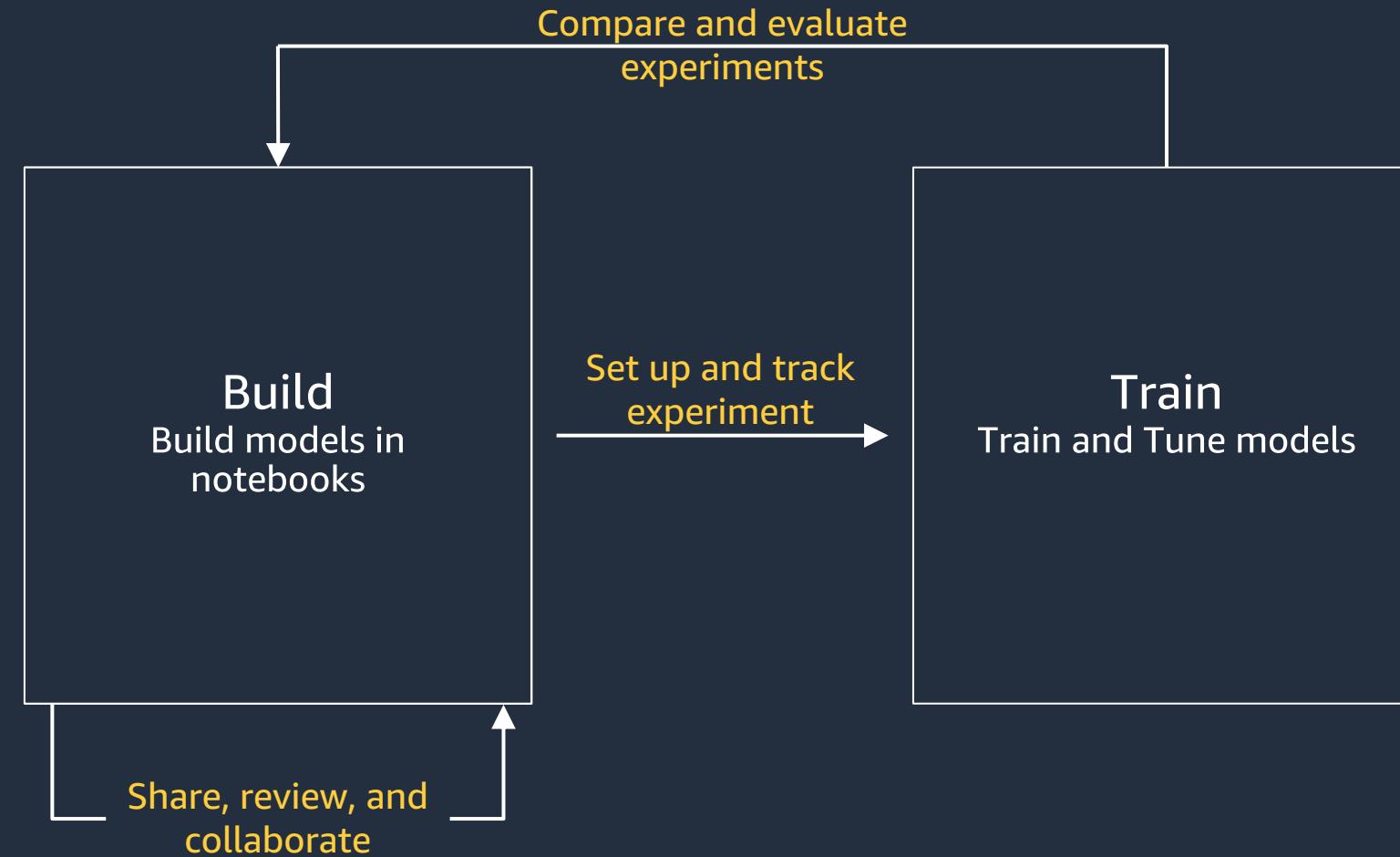


Managed  
SageMaker training



Tight framework  
integration

# Manage Experiments



# Amazon SageMaker Experiments

Organize, track, and compare training experiments



## Tracking at scale

Track parameters and metrics across experiments and users

## Custom organization

Organize experiments by teams, goals, and hypotheses

## Visualization

Easily visualize experiments and compare

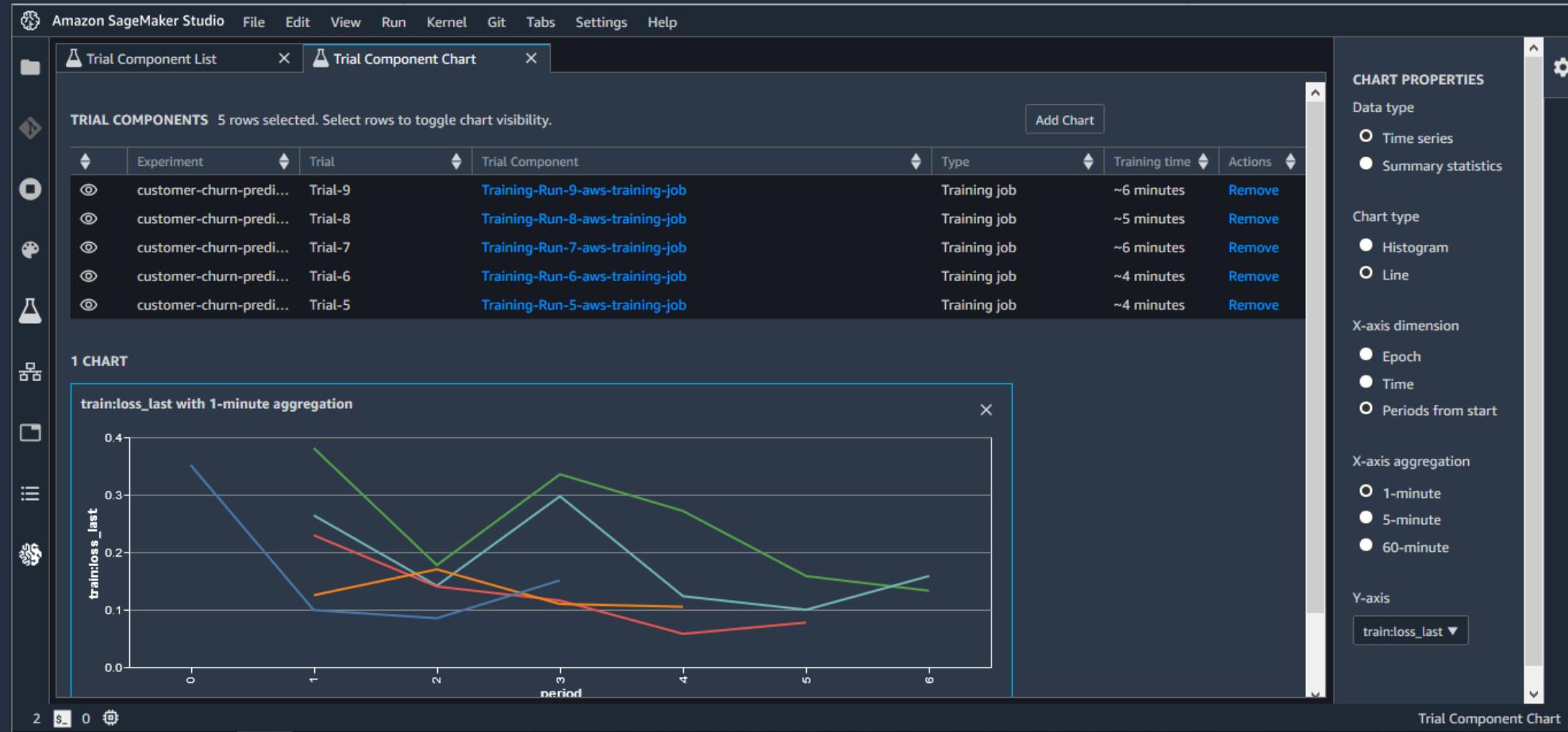
## Metrics and logging

Log custom metrics using the Python SDK and APIs

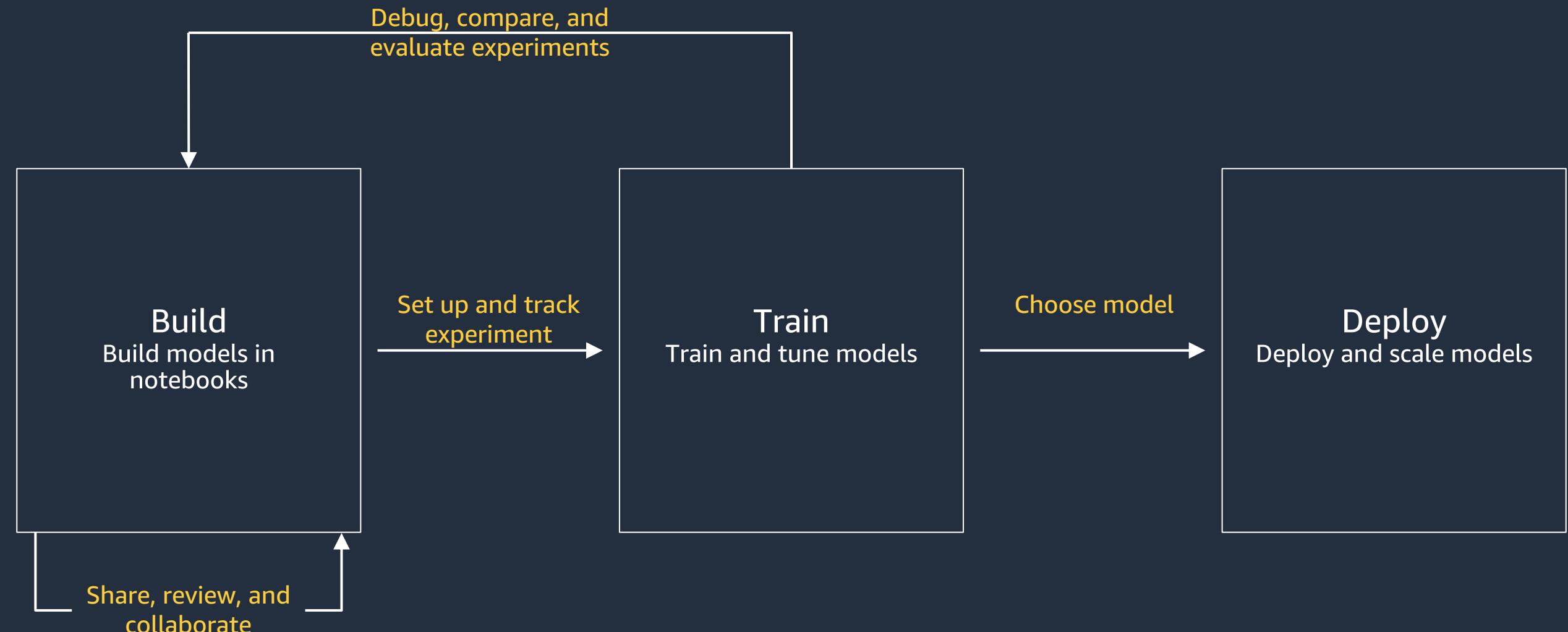
## Fast iteration

Quickly go back and forth, and maintain high-quality

# Use Amazon SageMaker Experiments to track and manage thousands of experiments



# Deploy models



# Pro Tips

|      | <b>Script Mode</b>           | <b>BYO Docker</b>   | <b>ML Marketplace</b>                   |
|------|------------------------------|---------------------|---|
| Pros | Quickly train your own model | Most flexible       | Huge variety: easy to add value quickly |
| Cons | Limited to managed options   | More time consuming | Less insight into solution              |

SageMaker reads from /opt/ml

Find an example to modify

# Takeaways

Code and dependencies

- Containers let you build lightweight, portable and consistent ML environments
- AWS DL containers include frameworks optimized by experts to deliver the best performance on CPUs and GPUs
- Leverage Amazon ECR along with git for collaborative development
- Leverage Amazon SageMaker to manage large-scale ML workloads

Scaling

Infrastructure management

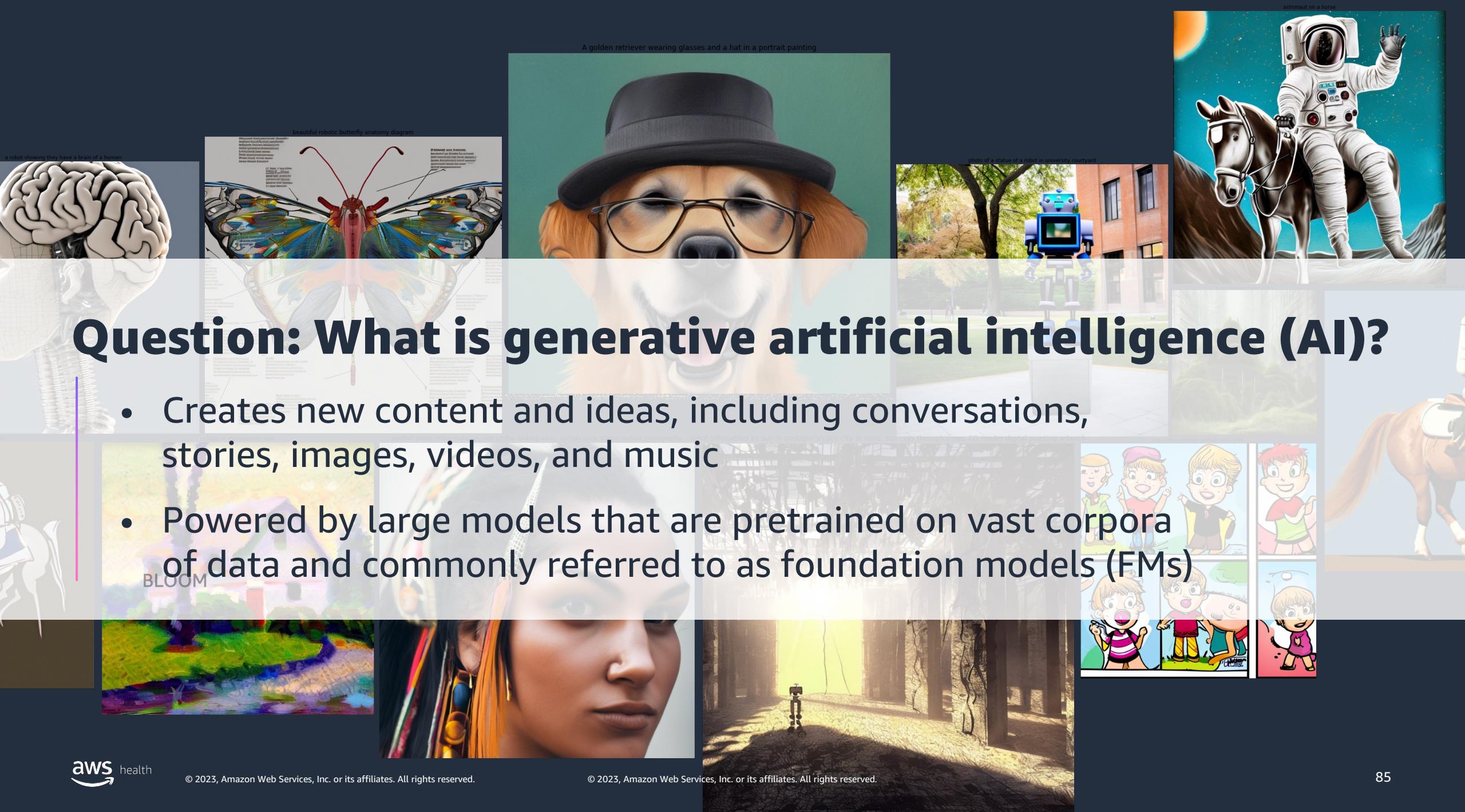
# Lab: Bring your custom model

**Lab 3. Bring your own model**

**Q&A and chat are open!**

Until 15:45

# Advanced ML & Generative AI



# Question: What is generative artificial intelligence (AI)?

- Creates new content and ideas, including conversations, stories, images, videos, and music
- Powered by large models that are pretrained on vast corpora of data and commonly referred to as foundation models (FMs)

# Where does generative AI fit?



## Artificial intelligence (AI)

Any technique that allows computers to mimic human intelligence using logic, if-then statements, and machine learning



## Machine learning (ML)

A subset of AI that uses machines to search for patterns in data to build logic models automatically



## Deep learning (DL)

A subset of ML composed of deeply multi-layered neural networks that perform tasks like speech and image recognition



## Generative AI

Powered by large models that are pretrained on vast corpora of data and commonly referred to as foundation models (FMs)

# Common generative AI use cases



Text  
generation



Q&A



Text  
summarization



Text  
extraction



Paraphrase  
rephrase



Search



Code  
generation



Image  
generation



Image  
classification

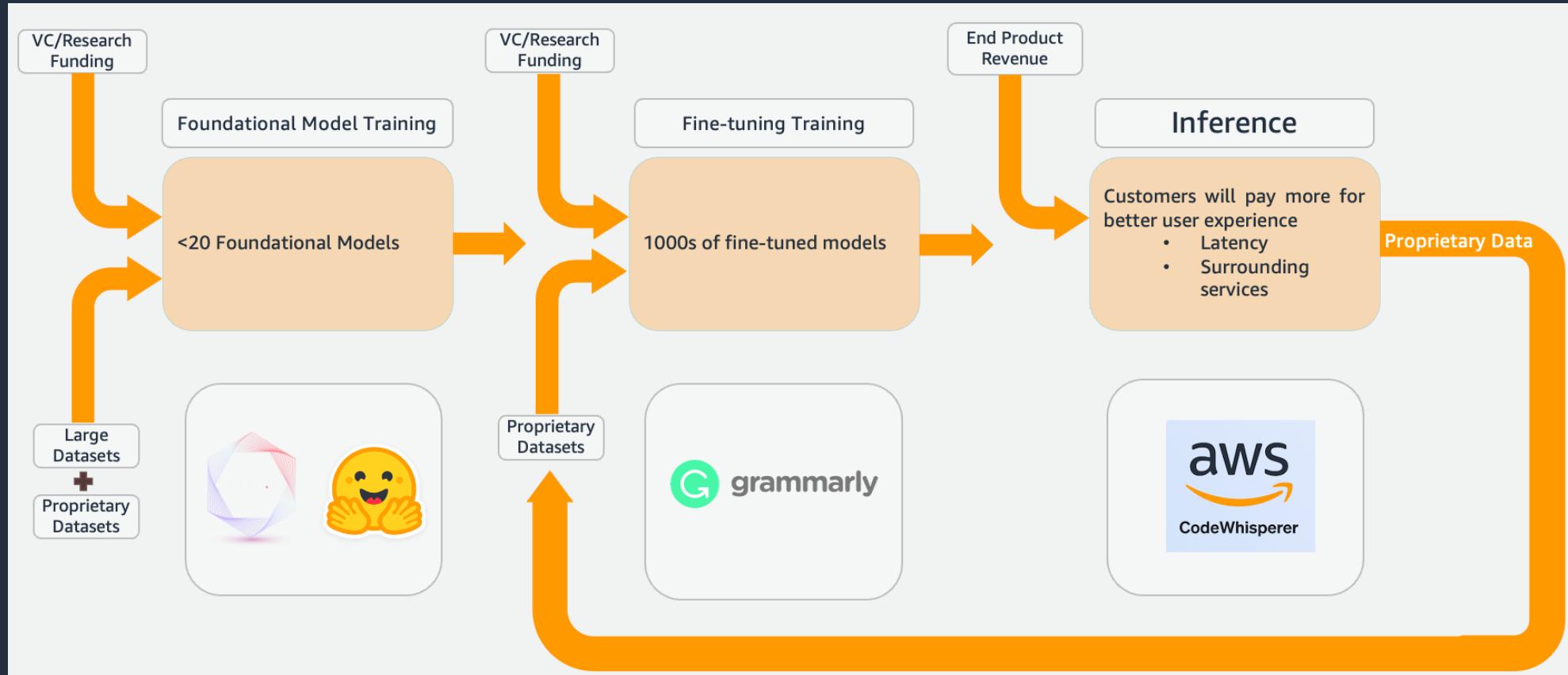


Audio  
generation

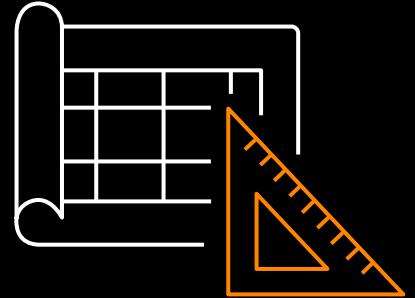


Video  
generation

# GenAI Workload Flow



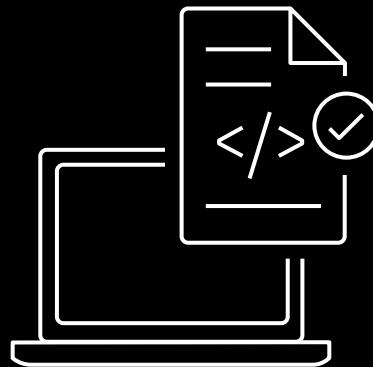
# GenAI Workloads



**Build your own foundation model from scratch**

Expensive, time consuming and requires deep expertise

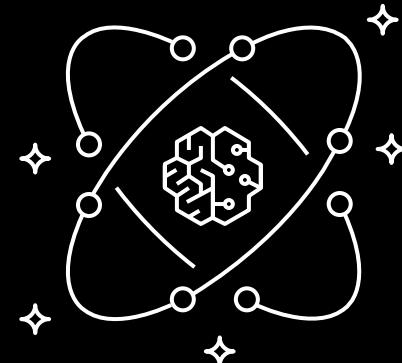
Model Provider



**Use publicly available foundation models and Fine-Tune them**

Substantial undifferentiated work needed to operationalize

Model Tuner

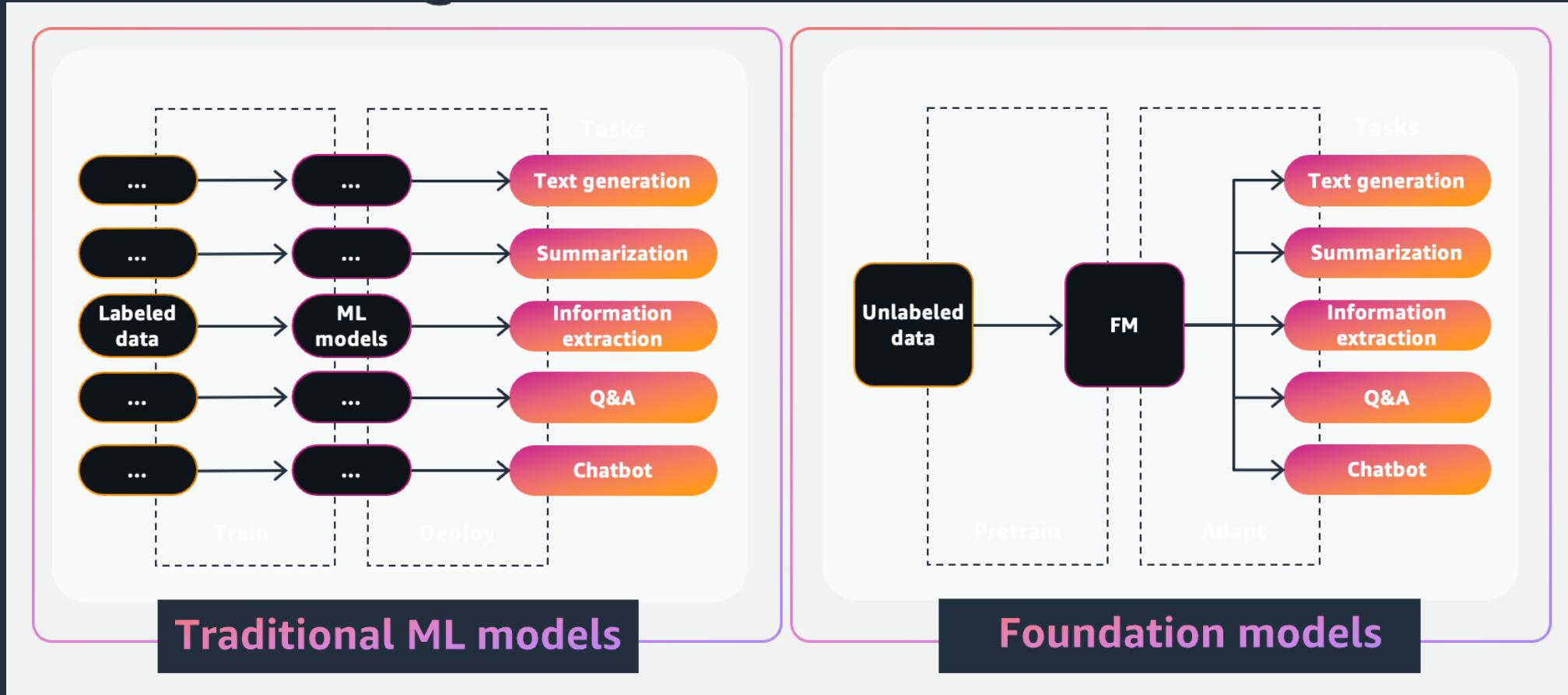


**Use generative AI services or APIs offered by foundation model vendors**

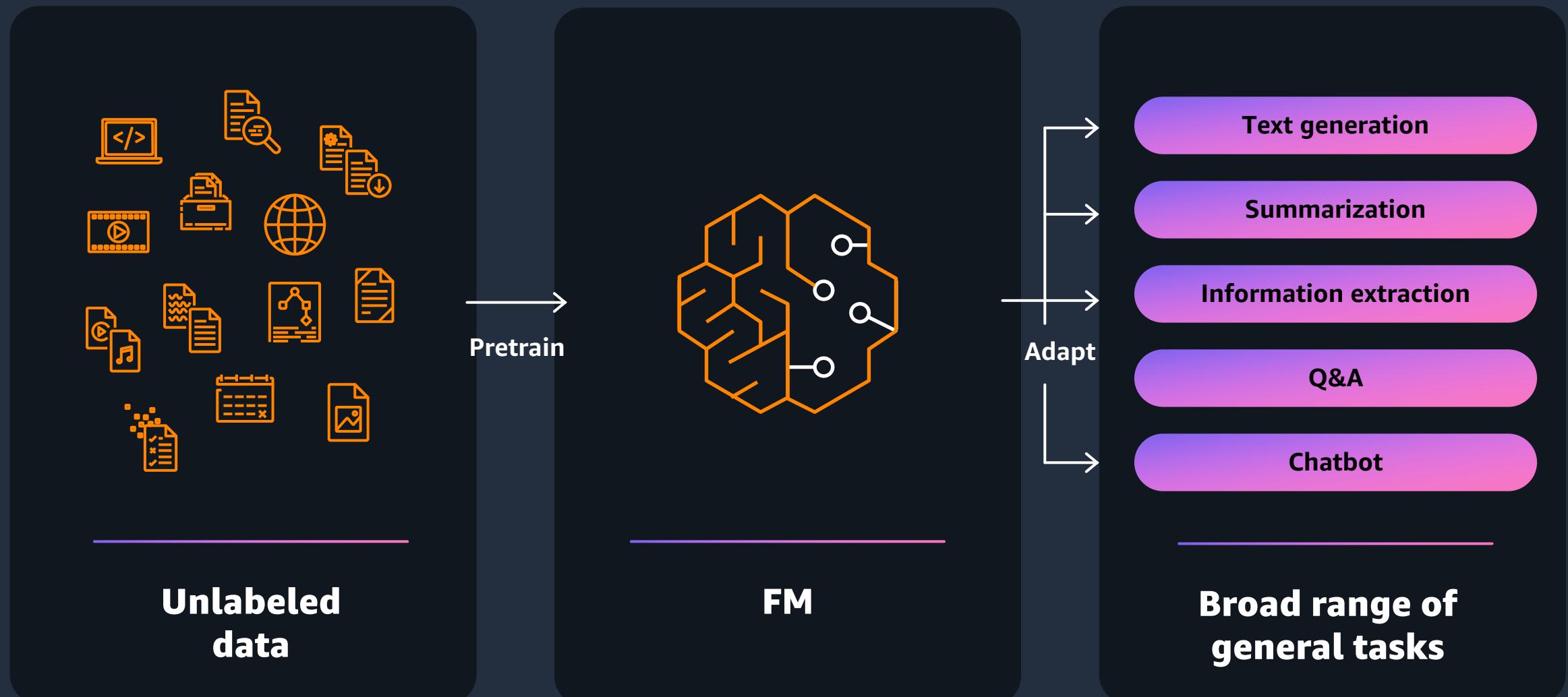
No control over data, costs, and no customization support

Model Consumer

# Generative AI @ AWS: Foundational models



# How foundation models work



# Types of foundation models

**Input**

"Summarize the articles on impact  
of walking on heart health"

"hand soap"

"a photo of an astronaut  
riding a horse on mars"

**FM**

## **Text-to-text**

Generate text from simple natural-language  
prompts for various applications

## **Text-to-embeddings**

Generate numerical representation of text  
for applications like search and finding  
similarities between documents

## **Multimodal**

Generate and edit images from  
natural-language prompts

**Output**

"Ten thousand steps per day  
is optimum for maintaining  
a healthy heart"

Numerical representation of  
"Hand soap refills"  
"Hand soap dispenser"  
"Hand soap antibacterial"



NEW

# Amazon Bedrock

**The easiest way to build and  
scale generative AI  
applications with FMs**

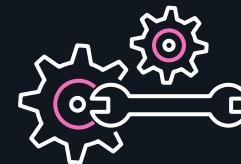
# Amazon Bedrock key benefits



Accelerate development of generative AI applications using FMs through an API, without managing infrastructure



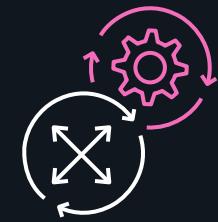
Choose FMs from AI21 Labs, Anthropic, Stability AI, and Amazon to find the right FM for your use case



Privately customize FMs using your organization's data



Enhance your data protection using comprehensive AWS security capabilities



Use AWS tools and capabilities that you are familiar with to deploy scalable, reliable, and secure generative AI applications

# Bedrock supports a wide range of foundation models

## FMs from Amazon



Titan Text



Titan  
Embeddings

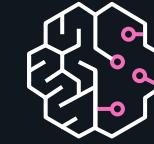
## FMs from AI21 Labs, Anthropic, and Stability AI



Jurassic-2



Claude



Stable  
Diffusion

# Amazon Titan

INNOVATE RESPONSIBLY WITH HIGH-PERFORMING FMs FROM AMAZON



Titan Text  
focused on  
NLP tasks



Titan Embeddings  
for enterprise tasks  
such as search and  
personalization

## Benefits

- Built with 20+ years of Amazon ML experience
- Automate language tasks such as summarization and text generation with Amazon Titan Text FM
- Enhance search accuracy and improve personalized recommendations with Amazon Titan Embeddings FM
- Support responsible use of AI by reducing inappropriate or harmful content

# Foundation models from top AI startups



**AI21 labs**

## Jurassic-2

Multilingual LLMs for text generation in Spanish, French, German, Portuguese, Italian, and Dutch



**ANTHROPIC**

## Claude

LLM for conversations, question answering, and workflow automation based on research into training honest and responsible AI systems



**stability.ai**

## Stable Diffusion

Generation of unique, realistic, high-quality images, art, logos, and designs

# Privately customize foundation models using your organization's data



## Fine-tune

**PURPOSE**

Maximizing accuracy for specific tasks

**DATA NEED**

Small number of labeled examples

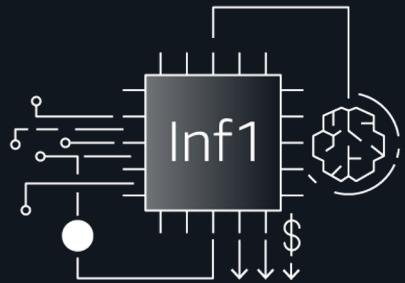
NEW

# Amazon EC2 Trn1n and Amazon EC2 Inf2

Purpose-built ML accelerators for  
the best price performance for  
training and inference in the cloud

# Purpose-built accelerators for generative AI

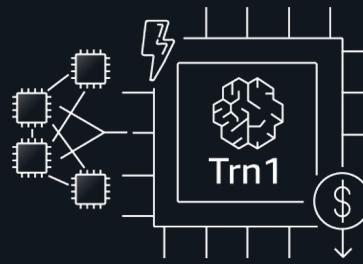
## AWS Inferentia



Lowest cost per inference  
in the cloud for running  
deep learning (DL) models

Up to 70% lower  
cost per inference  
than comparable  
Amazon EC2 instances

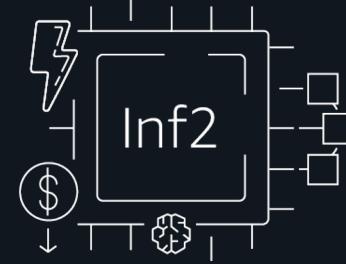
## AWS Trainium



The most cost-efficient, high-  
performance training of  
LLMs and diffusion models

Up to 50% savings  
on training costs  
over comparable  
Amazon EC2 instances

## AWS Inferentia2



High performance at the  
lowest cost per inference for  
LLMs and diffusion models

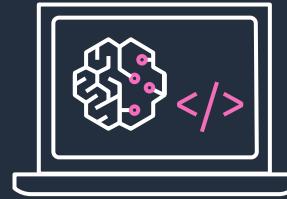
Up to 40% better  
price performance  
than comparable  
Amazon EC2 instances

NEW

# Amazon CodeWhisperer

Build applications faster and  
more securely with an AI coding  
companion

# Amazon CodeWhisperer: Now generally available and free to use for individual developers!



---

Generate code suggestions in real time



---

Scan code for hard-to-find vulnerabilities



---

Flag code that resembles open-source training data or filter by default

During preview Amazon ran a productivity challenge, and participants who used Amazon CodeWhisperer were **27% more likely to complete tasks successfully and did so an average of 57% faster** than those who did not use CodeWhisperer.

# Healthcare and Life Sciences



---

Accelerate drug discovery and research using foundation models to design vaccines, antibodies, and enzymes

---

Improve customer support with faster agent assisted claim and document processing

---

Deliver better care with personalized medicine

---

Keep patient data safe using synthetic data generation for research

# Lab: Advanced ML & Generative AI

Until 17:20



EMEA HEALTHCARE & LIFE SCIENCES WORKSHOPS

<https://aws-experience.com/>

- June 27th: AWS Cloud Fundamentals
- June 28th: Genomic Data Analysis with Amazon Omics
- July 4th: Machine Learning
- July 5th: AWS for Pharma Manufacturing
- July 6th: Security, Encryption & Data Protection Immersion Day
- July 11th: Sustainability
- July 13th: High-performance computing
- July 18th: Compliance in the Cloud



Please complete the  
workshop survey

**<http://bit.ly/3CVzhO0>**





# Thank you!

**Yegor Tokmakov**

Sr. Solutions Architect  
AWS Health

**Dmitri Laptev**

Sr. Solutions Architect  
AWS Startups

**Aamna Najmi**

Data Scientist  
AWS Professional Services