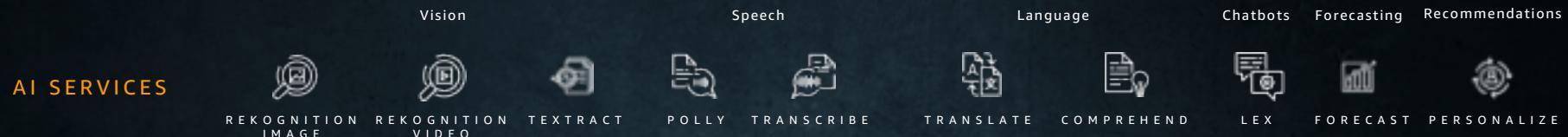
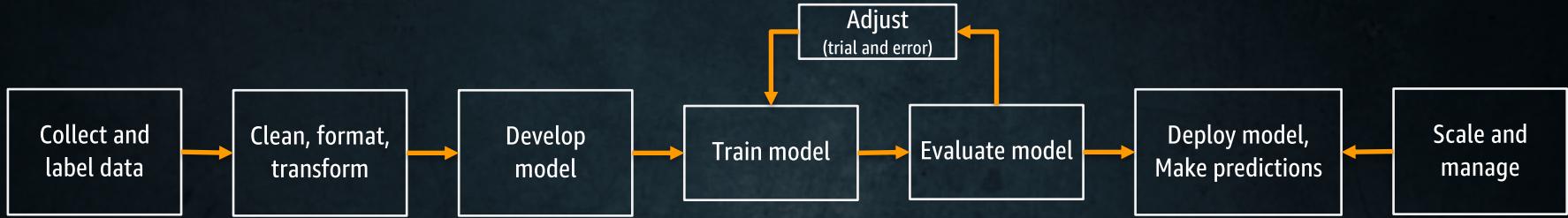


# AWS | Introduction to Amazon SageMaker

# The Amazon ML stack

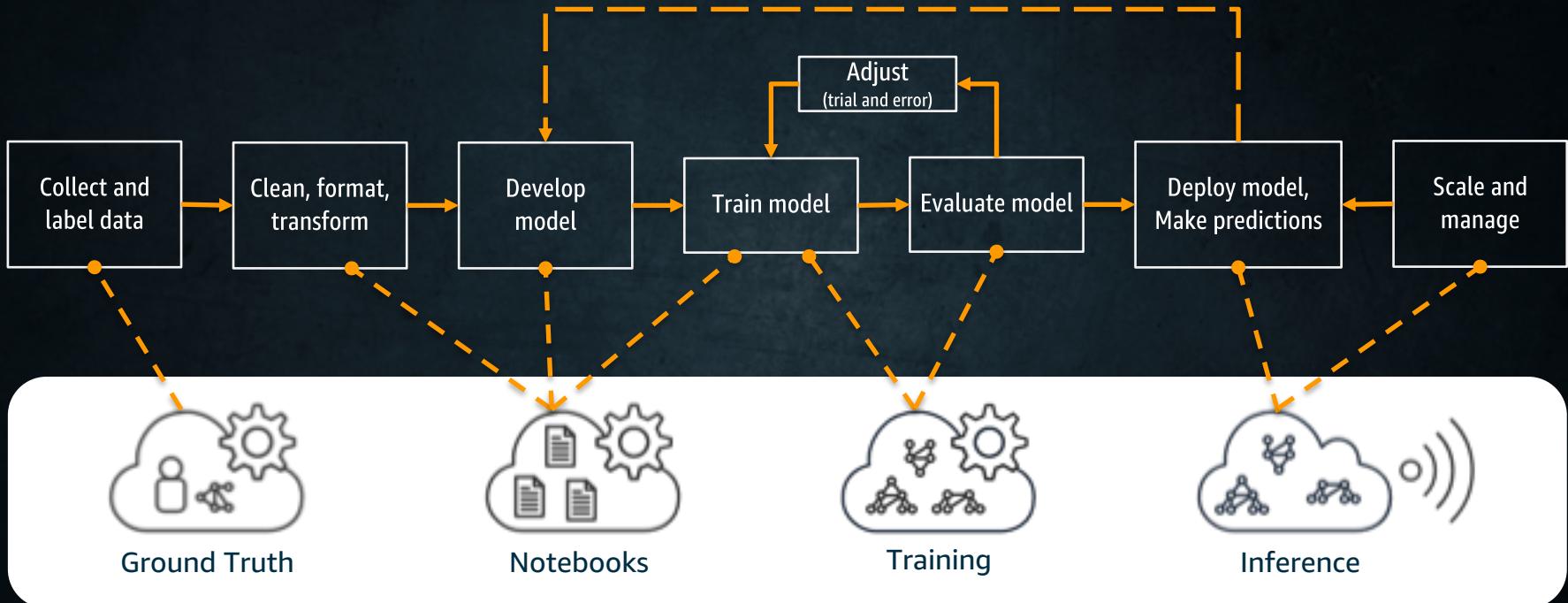


# ML was too complicated for everyday developers



# Amazon SageMaker

Build, train, and deploy machine learning models quickly



# Amazon SageMaker | Ground Truth



Ground Truth



Notebooks



Training

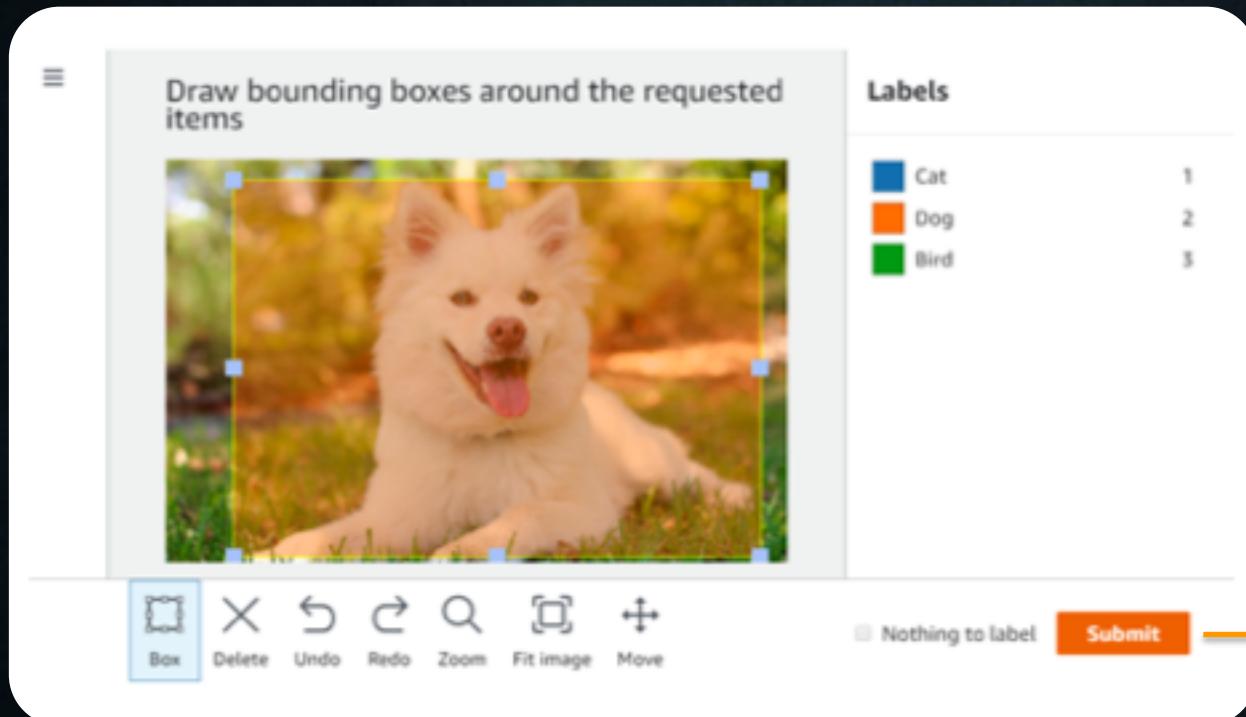


Inference

# Amazon SageMaker | Ground Truth



# Amazon SageMaker | Ground Truth



```
{  
  "boundingBox": {  
    "boundingBoxes": [  
      {  
        "height": 291,  
        "label": "dog",  
        "left": 59,  
        "top": 86,  
        "width": 417  
      }  
    ],  
    "inputImageProperties": {  
      "height": 480,  
      "width": 640  
    }  
  }  
}
```

# Amazon SageMaker | Ground Truth

## Creating training data



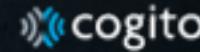
Mechanical  
Turk workers



Private labeling  
workforce



Third-party  
vendors



# Amazon SageMaker | Notebooks



Ground Truth



Notebooks



Training



Inference

# Amazon SageMaker | Notebooks

- Jupyter notebooks
- Support Jupyter Lab
- Multiple built-in kernels
- Install external libraries and kernels
- Integrate with Git
- Sample notebooks

The screenshot shows a Jupyter notebook titled "jupyter linear\_learner\_mnist (unsaved changes)". The notebook has a "Not Trusted" status and a "conda\_python3" kernel. The first cell contains a section titled "Permissions and environment variables" with instructions for specifying S3 bucket and IAM role. The second cell contains Python code for defining an S3 bucket and IAM role:

```
bucket = '<your_s3_bucket_name_here>'  
prefix = 'sagemaker/0000-Linear-Mnist'  
  
# Define IAM role  
import boto3  
import re  
from sagemaker import get_execution_role  
  
role = get_execution_role()
```

The third cell starts with "Data ingestion" and describes reading a dataset from an online URL. The fourth cell begins with "In [ ]: `vitime`". A tooltip for "vitime" lists several options: R, Sparkmagic (PySpark), Sparkmagic (PySpark3), Sparkmagic (Spark), Sparkmagic (SparkR), conda\_chainer\_p27, conda\_chainer\_p36, conda\_mxnet\_p27, conda\_mxnet\_p36, conda\_python2, conda\_python3, conda\_pytorch\_p27, conda\_pytorch\_p36, conda\_r\_jupyter\_systemenv, conda\_tensorflow\_p27, and conda\_tensorflow\_p36.

# Amazon SageMaker | Training



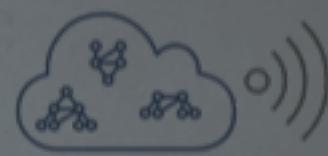
Ground Truth



Notebooks

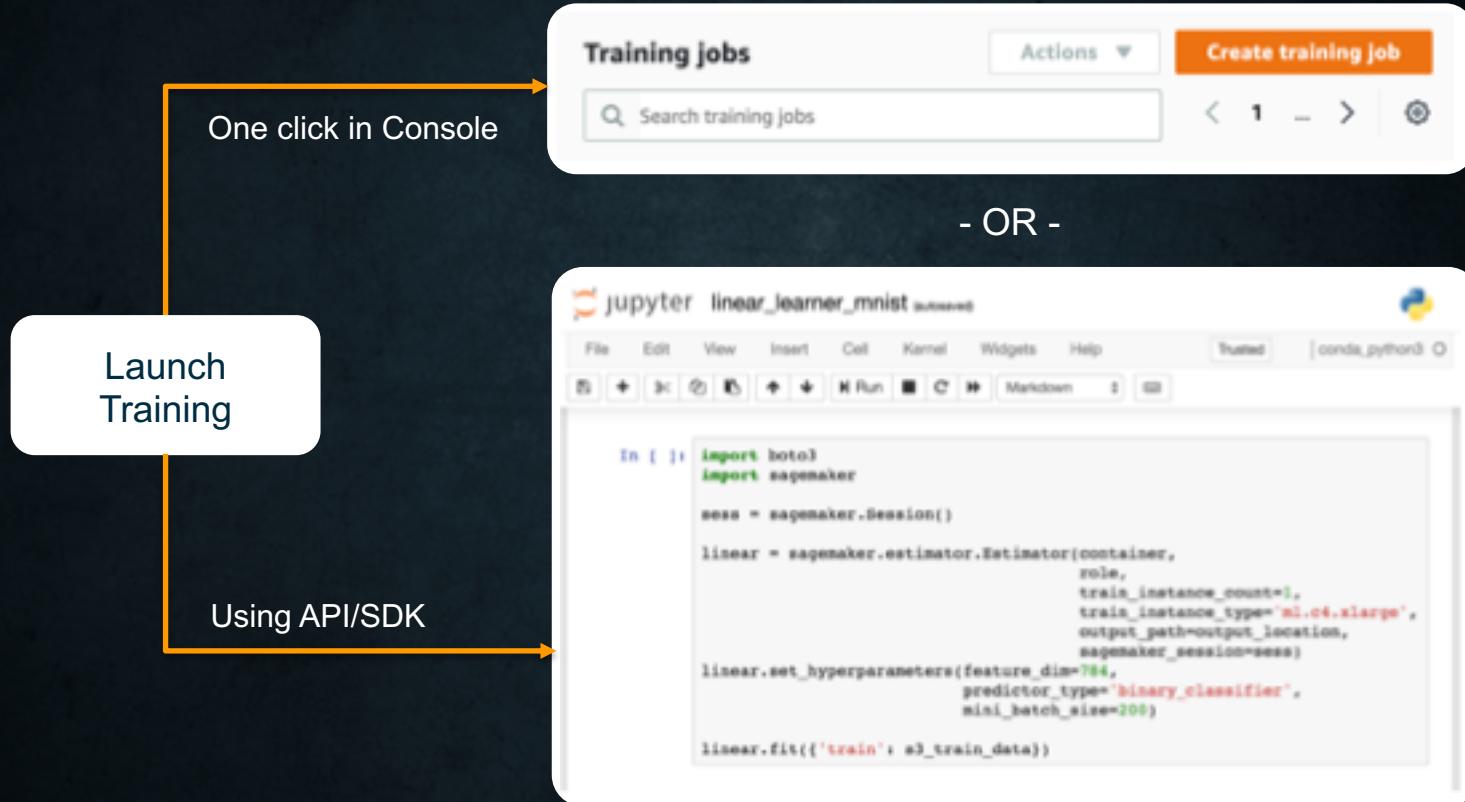


Training

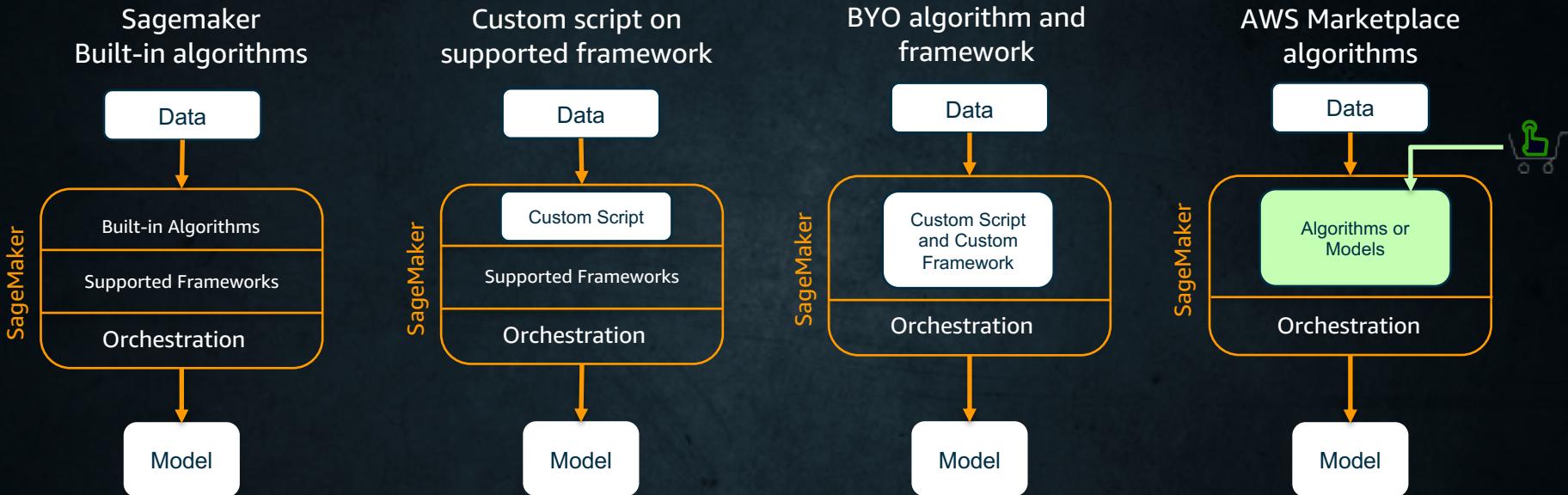


Inference

# Amazon SageMaker | Training



# Amazon SageMaker | Training



17 Built in high performance  
algorithms

Supported Frameworks: Apache  
MXNet, TensorFlow , Scikit-learn,  
PyTorch, Chainer

Docker containers with your own  
algorithms and frameworks

3<sup>rd</sup> party algorithms and models

# Amazon SageMaker | Training

## Built-in algorithms



XGBoost, FM, Linear,  
and Forecasting for  
supervised learning



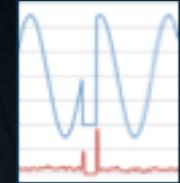
Kmeans, PCA, and  
BlazingText (Word2Vec)  
for clustering and pre-  
processing



Image classification and  
object detection with  
convolutional neural  
networks



LDA and NTM for topic  
modeling, seq2seq for  
translation



Random Cut Forest  
for anomaly  
detection

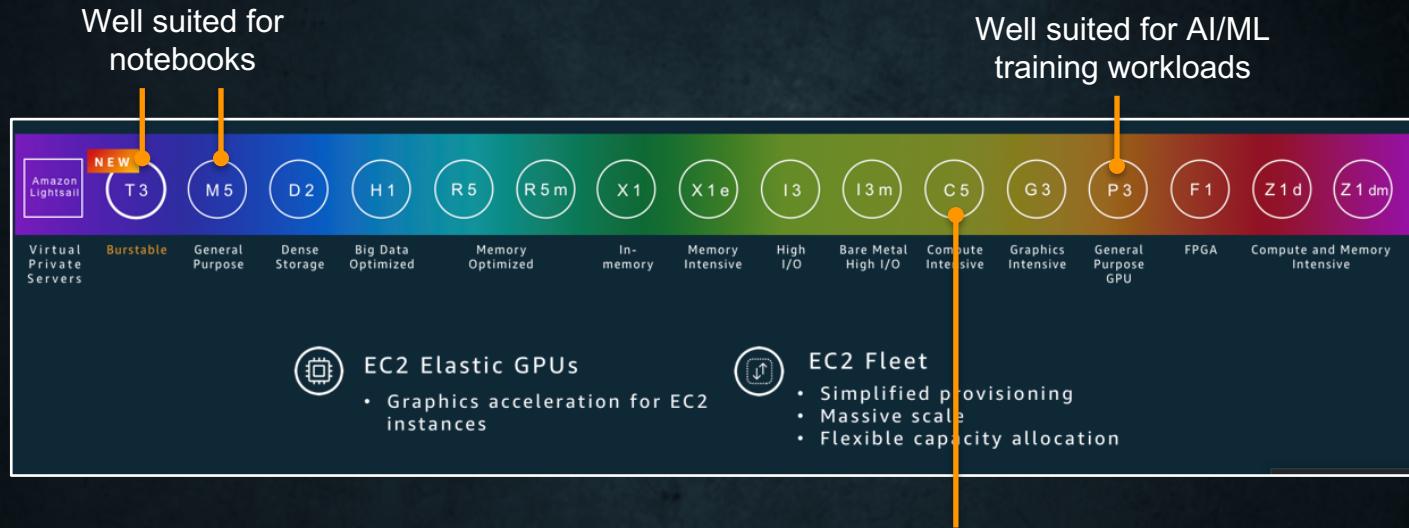
# Amazon SageMaker | Training

## Hyperparameter Optimization



- Define Metrics
- Hyperparameter ranges/scaling
- Stop tuning job early
- Use warm start

# Amazon SageMaker | Training



# Amazon SageMaker | Inference



Ground Truth



Notebooks

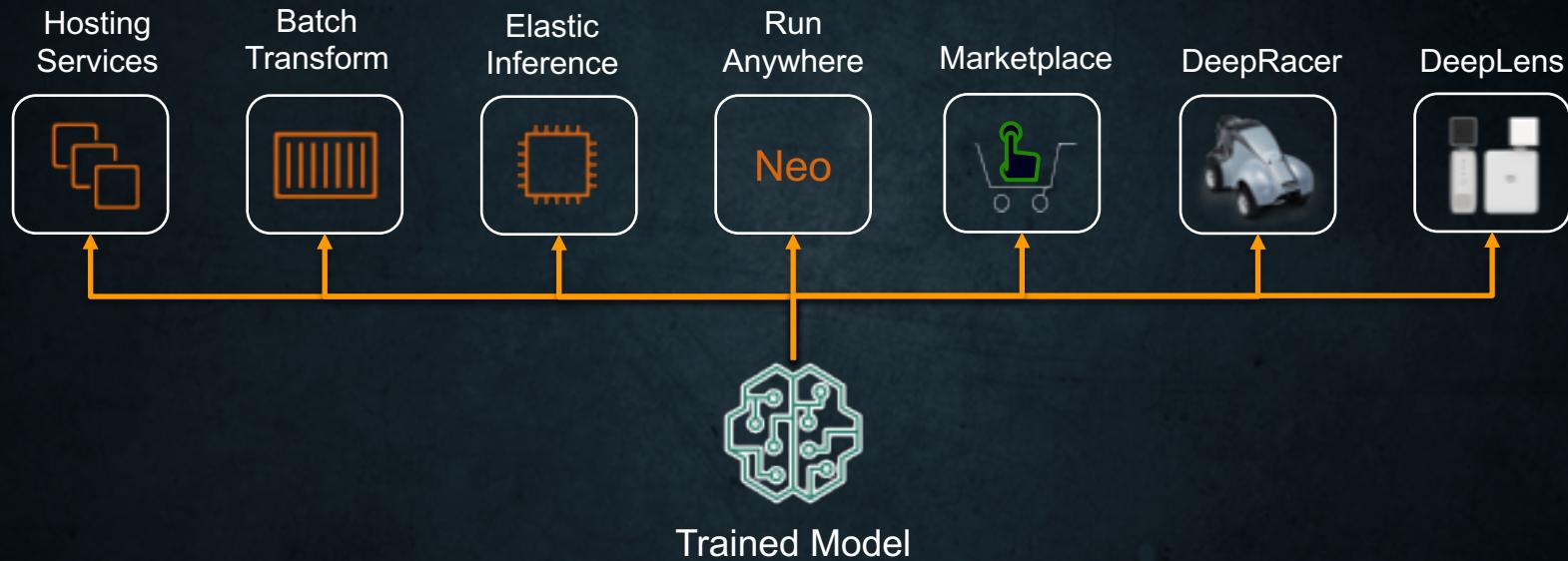


Training



Inference

# Amazon SageMaker | Inference

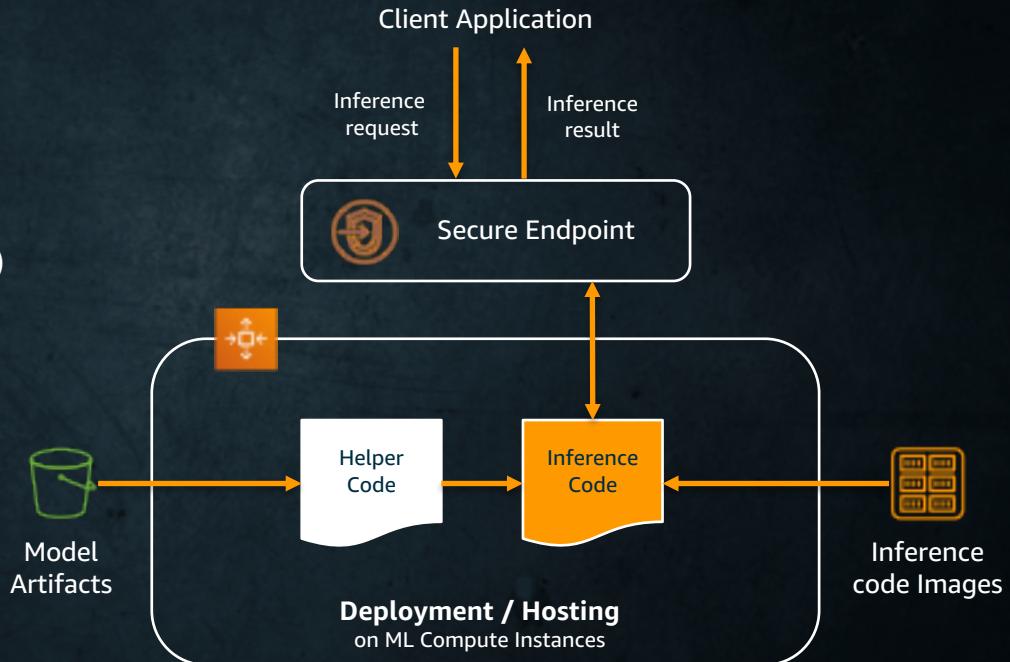


# Amazon SageMaker | Inference



## Hosting Services

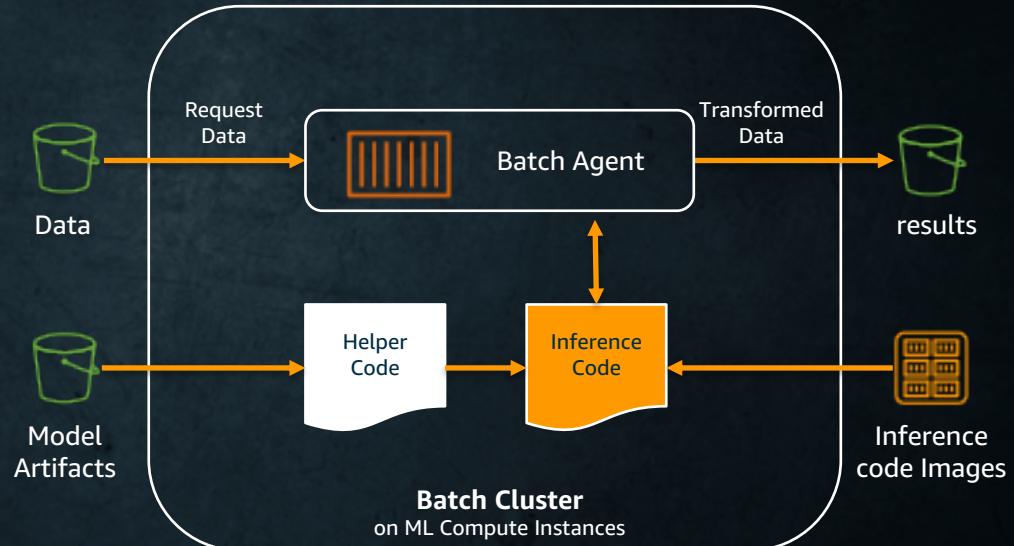
- Secure HTTPS endpoint
- Elastic configuration (scales with traffic)
- Multiple versions of the a model (A/B testing)
- Multi-AZ support
- Sub-second latencies
- Persistent deployments



# Amazon SageMaker | Inference

## Batch Transform

- Predictions for an entire dataset
- Transient resources (instances provisioned and terminated once job is done)
- No infrastructure to manage

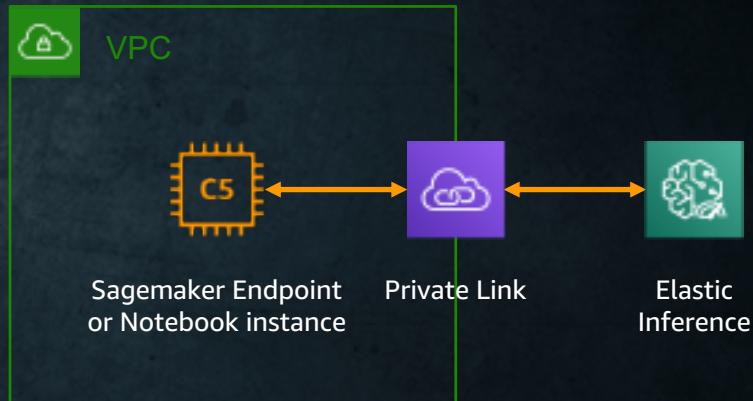


# Amazon SageMaker | Inference



## Elastic Inference

- Inference acceleration at a fraction of full GPU instance cost
- Add accelerators to CPU instances
- Works with inference as well as notebook instances
- Supports TensorFlow and MXNet frameworks, any other can be used via ONNX

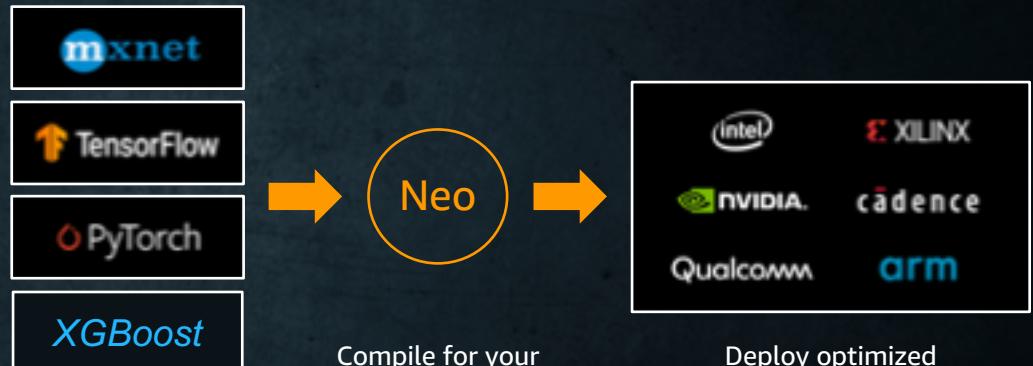


# Amazon SageMaker | Inference

## Neo



- Train your model once and deploy anywhere
- Optimization for target platform
- Deploy on EC2, Edge, IoT
- Open Source
- Up to 2X performance increase
- 1/10<sup>th</sup> of the size of the original framework

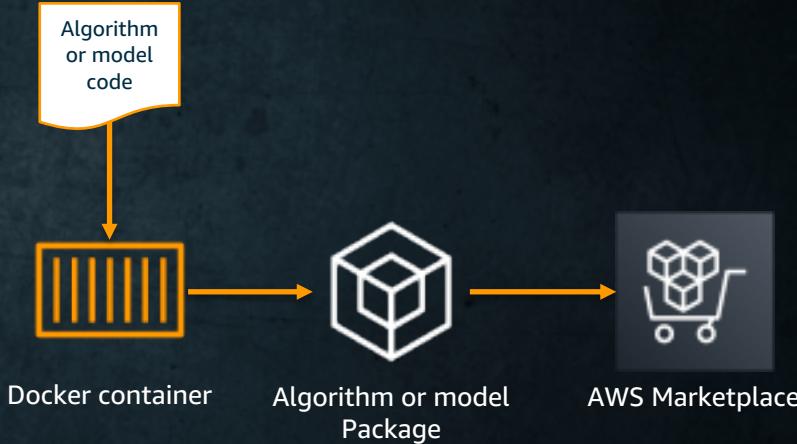


# Amazon SageMaker | Inference



## Marketplace

- Develop your algorithm or model, and package it in a Docker container.
- Create an algorithm or model package resource in Amazon SageMaker.
- Register as a seller on AWS Marketplace and list your algorithm or model package on AWS Marketplace.



# Amazon SageMaker | Inference

## DeepLens



Learn computer vision through projects, tutorials, and real-world hands-on exploration with the world's first deep learning enabled video camera for developers.

- Inference at the edge
- SageMaker trained models
- Simplified deployment
- Sample notebooks

## DeepRacer



A fully autonomous, 1/18th scale race car, packed full of everything you need to learn about reinforcement learning through autonomous driving.

# Amazon SageMaker | >10k users



# To learn more:



AWS Machine Learning University

---

Practical education on ML for new & experienced practitioners

Based on the same material used to train Amazon developers

# ml.aws