

A I M 3

# Build, Train, and Deploy Machine Learning Models at Any Scale | Messaging Block

Shyam Srinivasan



# The Amazon ML Stack: Broadest & Deepest Set of Capabilities

AI SERVICES

Vision

Speech

Language

Chatbots

Forecasting

Recommendations

REKOGNITION IMAGE

REKOGNITION VIDEO

TEXT TRACTOR

POLLY

TRANSCRIBE

TRANSLATE

COMPREHEND MEDICAL

LEX

FORECAST

PERSONALIZE

ML SERVICES

AMAZON SAGEMAKER

BUILD

TRAIN

DEPLOY

Pre-built algorithms & notebooks

One-click model training & tuning

One-click deployment & hosting

Data labeling (GROUND TRUTH)

Optimization (NEO)

Models without training data (REINFORCEMENT LEARNING)

Algorithms & models (AWS MARKETPLACE)

ML FRAMEWORKS & INFRASTRUCTURE

Frameworks

Interfaces

Infrastructure

TensorFlow

mxnet

PYTORCH

intel RL Coach

GLUON

K Keras

EC2 P3 & P3dn

EC2 C5

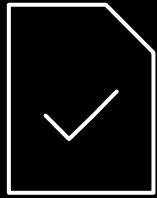
FPGA s

GREENGRASS

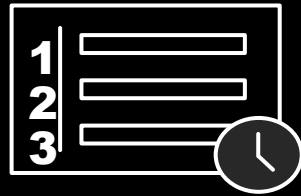
OpenVINO

ELASTIC INFERENCE

# Amazon SageMaker: Build, Train, and Deploy ML Models at Scale



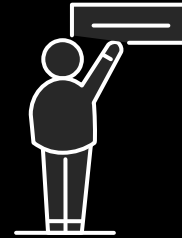
Collect and prepare  
training data



Choose and  
optimize your  
ML algorithm



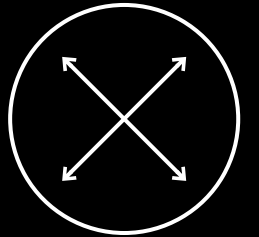
Set up and  
manage  
environments  
for training



Train and  
Tune ML Models



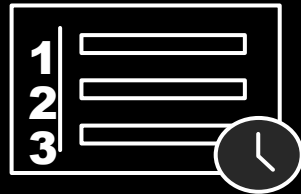
Deploy models  
in production



Scale and manage  
the production  
environment

# Amazon SageMaker: Build, Train, and Deploy ML Models at Scale

Pre-built  
notebooks  
for common  
problems



Collect and prepare  
training data

Choose and  
optimize your  
ML algorithm



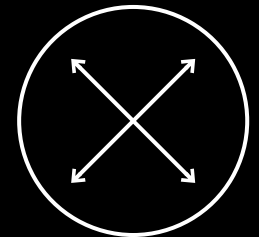
Set up and  
manage  
environments  
for training



Train and  
Tune ML Models



Deploy models  
in production



Scale and manage  
the production  
environment



# Amazon SageMaker: Build, Train, and Deploy ML Models at Scale

Pre-built  
notebooks  
for common  
problems

Collect and prepare  
training data

Built-in, high  
performance  
algorithms

Choose and  
optimize your  
ML algorithm



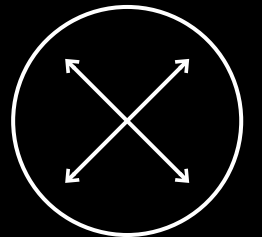
Set up and  
manage  
environments  
for training



Train and  
Tune ML Models



Deploy models  
in production



Scale and manage  
the production  
environment

# Amazon SageMaker: Build, Train, and Deploy ML Models at Scale

Pre-built  
notebooks  
for common  
problems

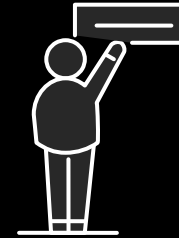
Collect and prepare  
training data

Built-in, high  
performance  
algorithms

Choose and  
optimize your  
ML algorithm

One-click  
training on the  
highest  
performing  
infrastructure

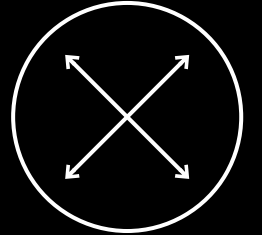
Set up and  
manage  
environments  
for training



Train and  
Tune ML Models



Deploy models  
in production



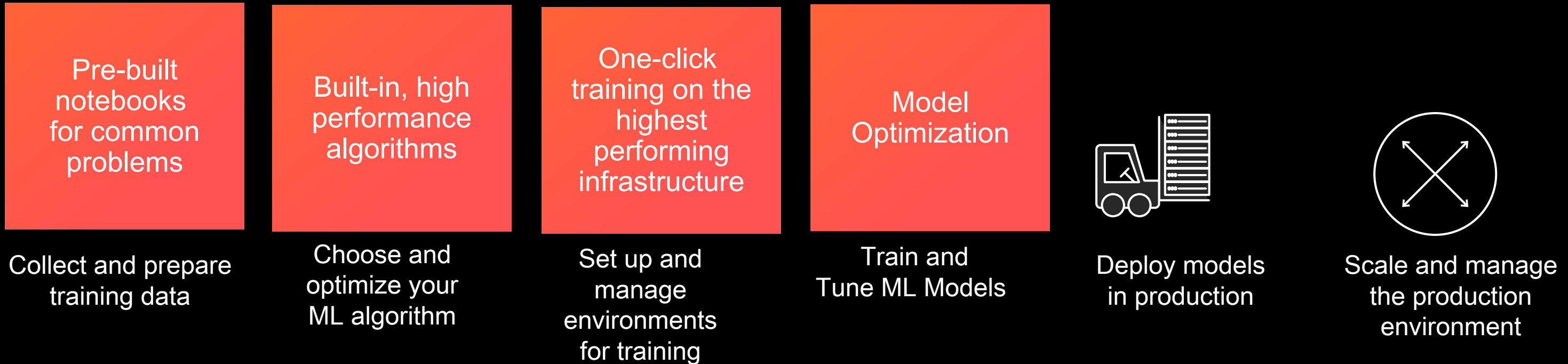
Scale and manage  
the production  
environment



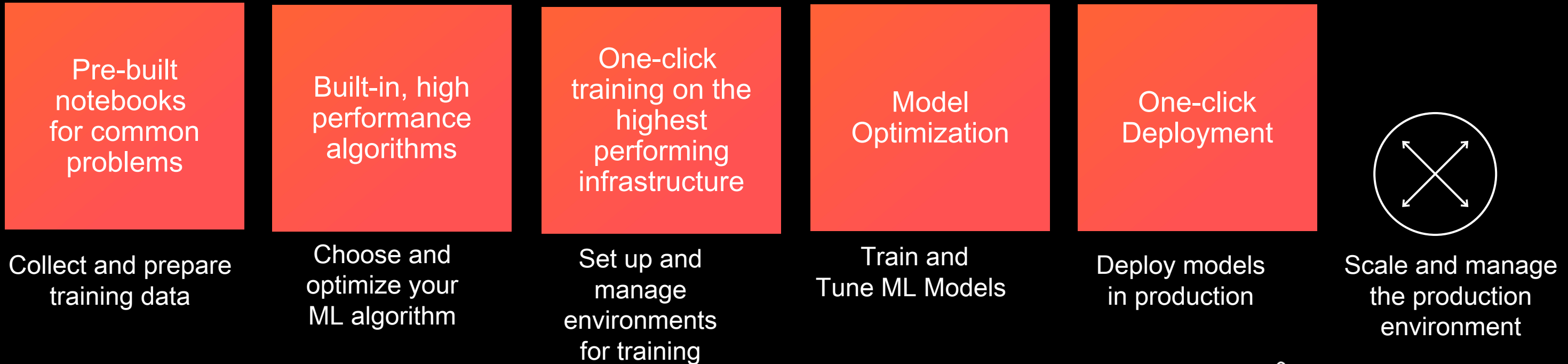
RL  
Coach



# Amazon SageMaker: Build, Train, and Deploy ML Models at Scale



# Amazon SageMaker: Build, Train, and Deploy ML Models at Scale



EC2 C5



GREENGRASS



OpenVINO



# Amazon SageMaker: Build, Train, and Deploy ML Models at Scale

Pre-built  
notebooks  
for common  
problems

Collect and prepare  
training data

**intuit.**

Built-in, high  
performance  
algorithms

Choose and  
optimize your  
ML algorithm



One-click  
training on the  
highest  
performing  
infrastructure

Set up and  
manage  
environments  
for training



**tinder**

Model  
Optimization

Train and  
Tune ML Models



One-click  
Deployment

Deploy models  
in production



Fully  
managed with  
auto-scaling  
for 75% less

Scale and manage  
the production  
environment

**CONVOY**

**SIEMENS**



DOW JONES



**SONY**



# Reducing the cost of Sentiment Analysis

Signal Labs performs nuanced sentiment analysis on billions of stories per month using AWS. Signal Labs offers solutions that analyze the entire digital media landscape to deliver instant insights for the company's Fortune 1000 customers. The company built a sentiment-analysis pipeline that uses Amazon SageMaker for machine-learning capabilities and Amazon EC2 C5 instances with Intel Xeon Scalable (Skylake) processors for faster model training and evaluation.

# 90%

Reduction of operations cost

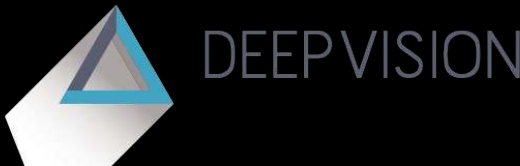


*"Using Amazon SageMaker and Amazon EC2 C5 instances, it's easier than ever for us to deliver high-quality sentiment analysis."*

Jonathan Dodson, Vice President of Engineering, Signal Labs

# Over 150 algorithms and models

## SELECTED VENDORS



## AVAILABLE ALGORITHMS & MODELS

Natural Language  
Processing

Grammar & Parsing

Text OCR

Computer Vision

Named Entity  
Recognition

Video Classification

Speech Recognition

Text-to-Speech

Speaker Identification

Text Classification

3D Images

Anomaly Detection

Text Generation

Object Detection

Regression

Text Clustering

Handwriting  
Recognition

Ranking

# Optimization is extremely complex

mxnet

TensorFlow

PYTORCH

intel

nvidia

Qualcomm

cadence

arm

xilinx

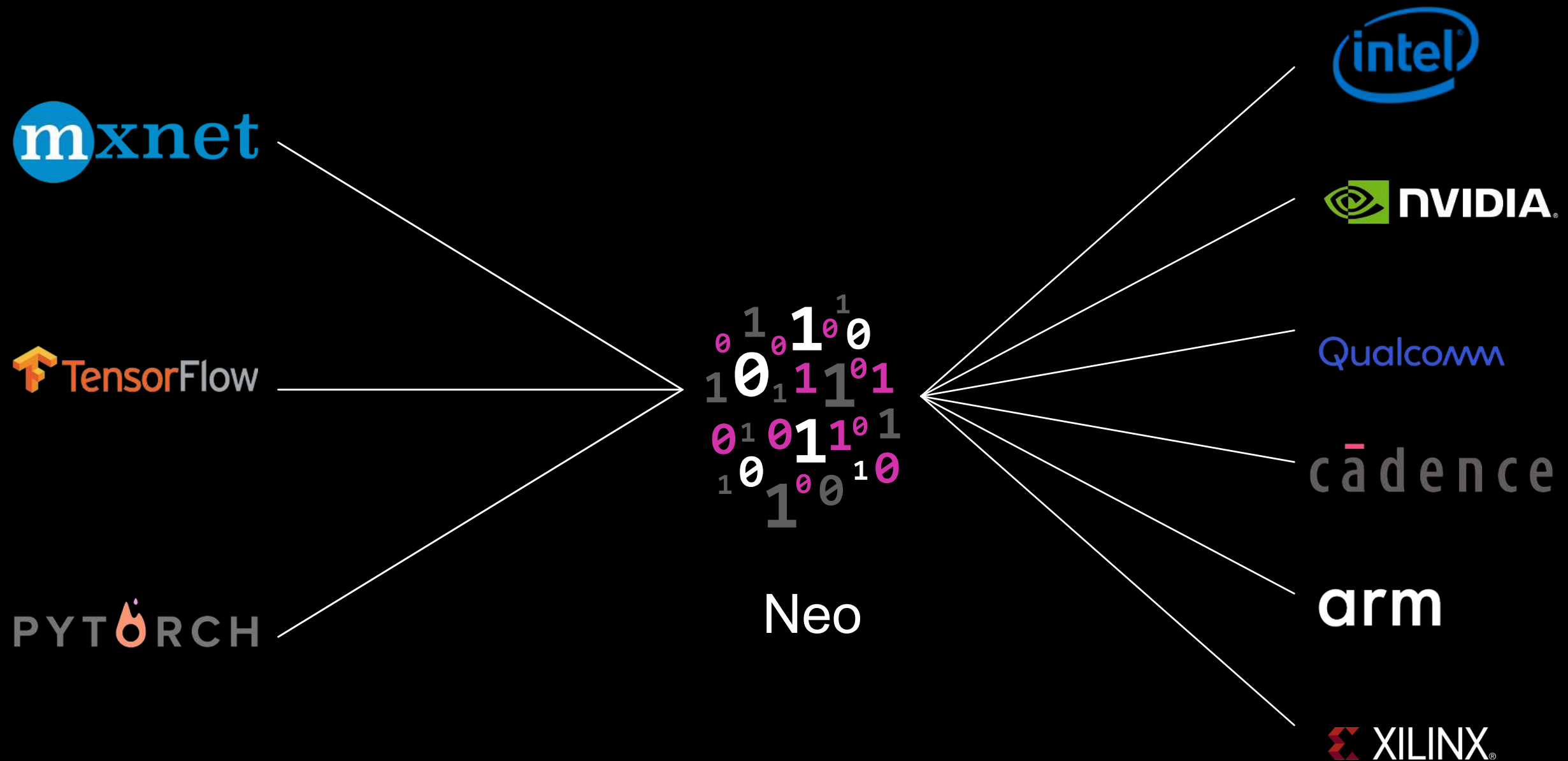
NEW

# Amazon SageMaker Neo

Train once, run anywhere with 2x the performance



# Amazon SageMaker Neo: Train once, run anywhere



# Amazon SageMaker Neo

Train once, run anywhere with 2x the performance



Get accuracy  
and performance



Automatic  
optimization



Broad framework  
support



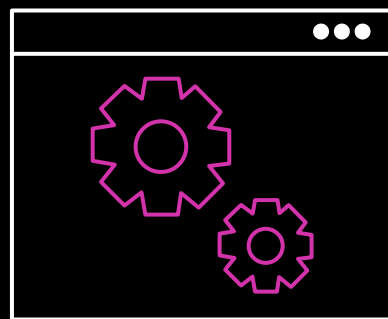
Broad hardware  
support

## KEY FEATURES

Open-source Neo-AI device runtime and  
compiler under the Apache software license;  
1/10<sup>th</sup> the size of original frameworks

# Amazon SageMaker RL

Reinforcement learning for every developer and data scientist



Fully  
managed



Broad support  
for frameworks



Broad support for simulation  
environments

## KEY FEATURES

2D & 3D physics  
environments and  
OpenGym support

Support Amazon Sumerian, AWS  
RoboMaker and the open source  
Robotics Operating System  
(ROS) project

Example notebooks  
and tutorials



SyntheticGestalt



# AWS is the best platform for Apache MXNet



Start with high quality,  
pre-trained models

- Gluon CV and Gluon NLP



Refine with fast,  
scalable training

- Keras-MXNet up to 2x faster than Keras-TensorFlow
- Near-linear scalability up to 256 GPUs
- Optimized for training on C5 Instances Intel Xeon Scalable (Skylake) Processors
- Dynamic training



Deploy using  
familiar tools

- Java/Scala APIs
- MXNet Model Server