# Build Intelligent Applications with Machine Learning on AWS

Shyam Srinivasan
Sr. Product Marketing Manager, AWS

aws

# Machine Learning at Amazon: A long heritage

**Personalized Recommendations**

**Fulfillment automation & Inventory Management**

**Drones**

**Voice driven Interactions**

**Inventing entirely new Customer experiences**

aws

# ML @ AWS: OUR MISSION

Put machine learning in the hands of every developer and data scientist

aws

# Tens of thousands of customers running ML on AWS

aws

# The Amazon Machine Learning Stack

## APPLICATION SERVICES

Rekognition

Rekognition Video

Polly

Transcribe

Translate

Comprehend

Lex

## PLATFORMS

Amazon SageMaker

AWS DeepLens

## FRAMEWORKS

TensorFlow · mxnet · PYT🔥RCH · Caffe2 · Chainer · HOROVOD · GLUON · Keras

aws

# Complete Control over Frameworks & Infrastructure

*For the data scientist, AI researcher, or advanced ML practitioner*

## FRAMEWORKS & INTERFACES



TensorFlow    mxnet    PYTORCH    Caffe2    Chainer    HOROVOD    GLUON    Keras

## INFRASTRUCTURE (GPU)

| P3 | | | |
|---|---|---|---|
| **NVIDIA**<br>**Tesla V100 GPUs**<br>*(14x faster than P2)* | 5,120 Tensor cores | 1 Petaflop of compute | **Machine Learning AMIs** |
| | 128GB of memory | NVLink 2.0 | |

## INFRASTRUCTURE (CPU)

| C5 | | | |
|---|---|---|---|
| **Intel Xeon**<br>**3.0 GHz Skylake CPU**<br>*(25% better perf/price than C4)* | 72 vCPUs | AVX 512 | **Machine Learning AMIs** |
| | 144 GB of memory | Nitro Hypervisor | |

aws

# AWS Deep Learning AMI

Get started quickly with easy-to-launch tutorials

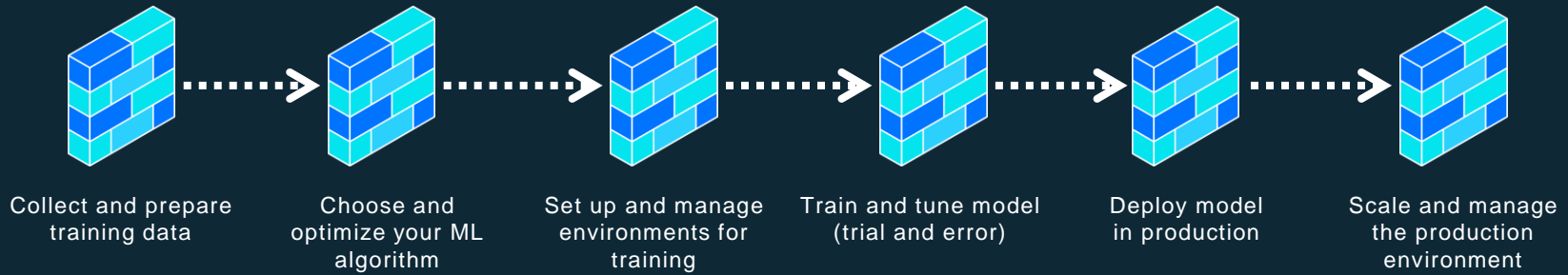Hassle-free setup and configuration

Pay only for what you use – no additional charge for the AMI

Accelerate your model training and deployment

Support for popular deep learning frameworks

# The Machine Learning Process



Collect and prepare training data → Choose and optimize your ML algorithm → Set up and manage environments for training → Train and tune model (trial and error) → Deploy model in production → Scale and manage the production environment

aws

# Amazon SageMaker



Pre-built notebooks for common problems

**BUILD**

Choose and optimize your ML algorithm

Set up and manage environments for training

Train and tune model (trial and error)

Deploy model in production

Scale and manage the production environment

aws

# Amazon SageMaker

| ALGORITHMS | | |
|---|---|---|
| | K-Means Clustering | XGBoost |
| | Principal Component Analysis | Latent Dirichlet Allocation |
| | Neural Topic Modelling | Image Classification (ResNet) |
| | Factorization Machines | Sequence2Sequence |
| | Linear Learner – Regression | Linear Learner – Classification |
| | DeepAR | BlazingText Word2Vec |
| FRAMEWORKS | Apache MXNet<br>TensorFlow | Caffe2, CNTK, PyTorch,<br>Torch, Chainer |

**Pre-built notebooks for common problems**

**Built-in, high performance algorithms**

Set up and manage environments for training

Train and tune model (trial and error)

Deploy model in production

Scale and manage the production environment

## BUILD

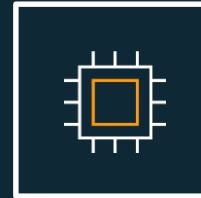aws

# Designing Better Algorithms

Algorithms for "infinite scale"

Distributed by default

Train on a data stream

01010

1

Single pass training

Not memory bound

Checkpoint for re-training

aws

# Amazon SageMaker

**BUILD**

Pre-built notebooks for common problems

Built-in, high performance algorithms

**TRAIN**

One-click training

Train and tune model (trial and error)

Deploy model in production

Scale and manage the production environment

aws

# Amazon SageMaker

Pre-built notebooks for common problems

Built-in, high performance algorithms

**BUILD**

One-click training

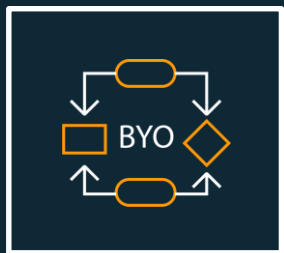Automatic Model Tuning

**TRAIN**

Deploy model in production

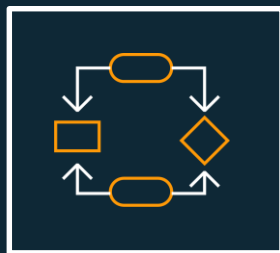Scale and manage the production environment

aws

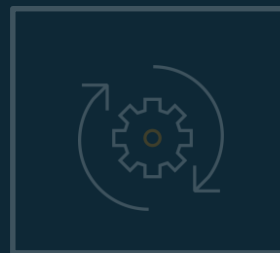# Amazon SageMaker:
# Custom ML Models Using Your Own Data



Train with your own algorithms

Train with a framework of your choice

Train with SageMaker algorithms

Optimize your models

Host models

A/B test

aws

"How can we make our algorithms
as fast as those in SageMaker?"

aws

# SageMaker Streaming For Custom Algorithms

**ACCELERATE YOUR OWN ALGORITHMS BY STREAMING LARGE VOLUMES OF TRAINING DATA FROM AMAZON S3**

NEW

Stream data to your own algorithm

Quicker time to start training

Faster training

Lower cost training

TensorFlow

Additional frameworks coming soon

aws

# Amazon SageMaker

**Pre-built notebooks for common problems**

**Built-in, high performance algorithms**

**One-click training**

**Automatic Model Tuning**

**One-click deployment**

Scale and manage the production environment

BUILD

TRAIN

DEPLOY

aws

# Amazon SageMaker



**BUILD**

Pre-built notebooks for common problems

Built-in, high performance algorithms

**TRAIN**

One-click training

Automatic Model Tuning

**DEPLOY**

One-click deployment
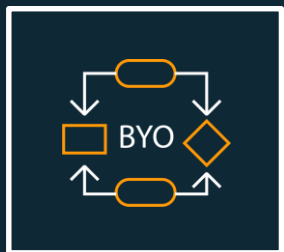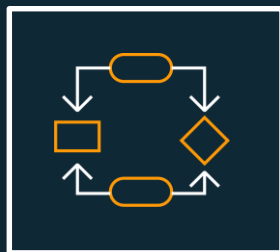
Fully managed hosting with auto-scaling

aws

# Amazon SageMaker:
# Custom ML Models Using Your Own Data



Train with
your own
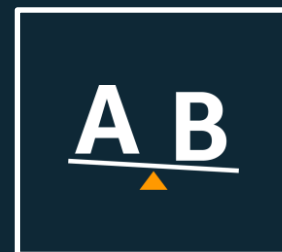algorithms

Train with
a framework
of your choice

Train with
SageMaker
algorithms

Optimize
your models

Host models

A/B test

aws

# Amazon SageMaker:
# Elastic Hosting For Custom Models

Easy to integrate API endpoint

Low latency

Auto-scaling

Fault tolerant, multi-AZ

But…

What if you need to batch process?

What if your files are big?

aws

# Amazon SageMaker Batch Transform

RUN FULLY MANAGED, HIGH-THROUGHPUT BATCH TRANSFORM JOBS WITH A SIMPLE API CALL

NEW

Process data dumps in a batch

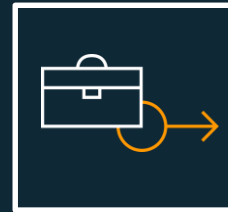Process large files more easily

Batch test models before hosting

Test models before deployment at the edge

Use same model for batch and real-time

Reuse pre-processing pipelines between training and prediction

Fully managed batch transforms

aws

# Amazon SageMaker Demo

## Build, Train and Deploy Machine Learning Models Quickly & Easily

aws

# The Amazon Machine Learning Stack

## APPLICATION SERVICES

Rekognition

Rekognition Video

Polly

Transcribe

Translate

Comprehend

Lex

## PLATFORMS

Amazon SageMaker

AWS DeepLens

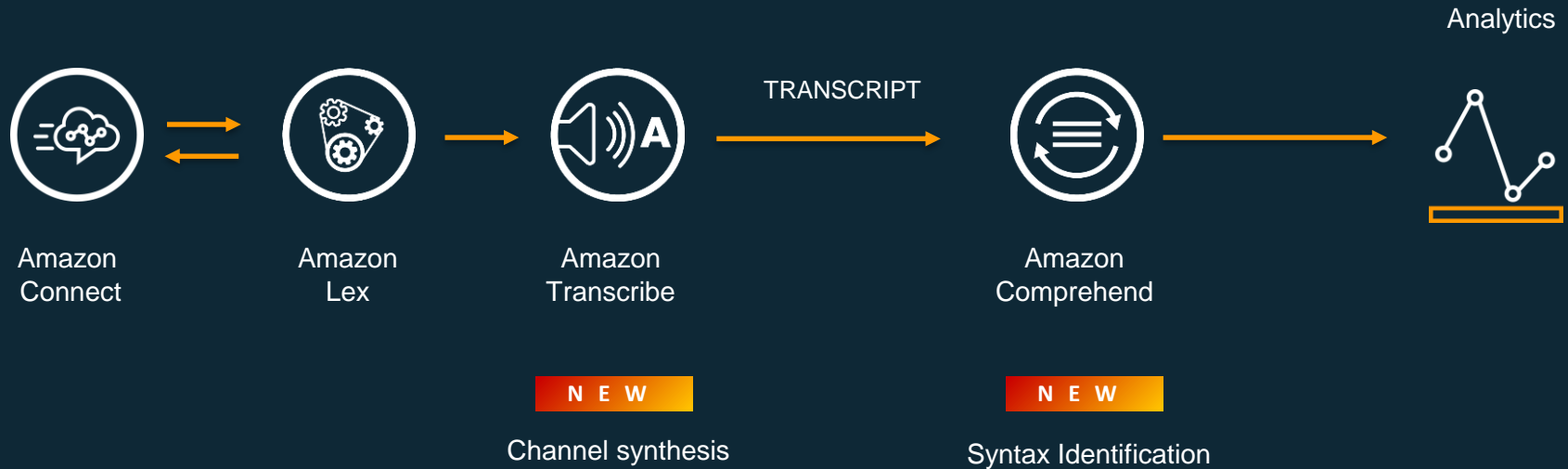## FRAMEWORKS

TensorFlow    mxnet    PYTORCH    Caffe2    Chainer    HOROVOD    GLUON    Keras
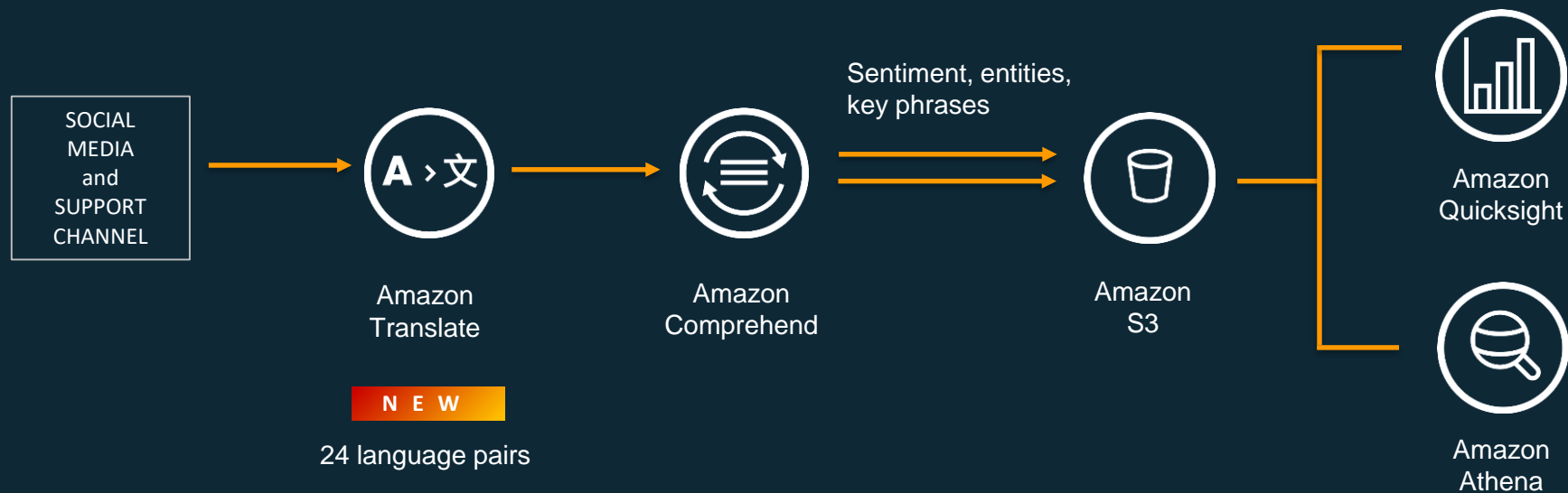
aws

# Improving Contact Centers With Artificial Intelligence

Analytics

Amazon
Connect

Amazon
Lex

Amazon
Transcribe

TRANSCRIPT

Amazon
Comprehend

**NEW**

Channel synthesis

**NEW**

Syntax Identification

aws

# Improving Voice Of Customer Analytics With Artificial Intelligence

SOCIAL MEDIA and SUPPORT CHANNEL

Amazon Translate

**NEW**

24 language pairs

Amazon Comprehend

Sentiment, entities, key phrases

Amazon S3

Amazon Quicksight

Amazon Athena

aws

# The Amazon Machine Learning Stack

## APPLICATION SERVICES

Rekognition · Rekognition Video · Polly · Transcribe · Translate · Comprehend · Lex

## PLATFORMS

Amazon SageMaker · AWS DeepLens

## FRAMEWORKS

TensorFlow · mxnet · PYTORCH · Caffe2 · Chainer · HOROVOD · GLUON · Keras

aws

# Thank You!

## ml.aws

aws