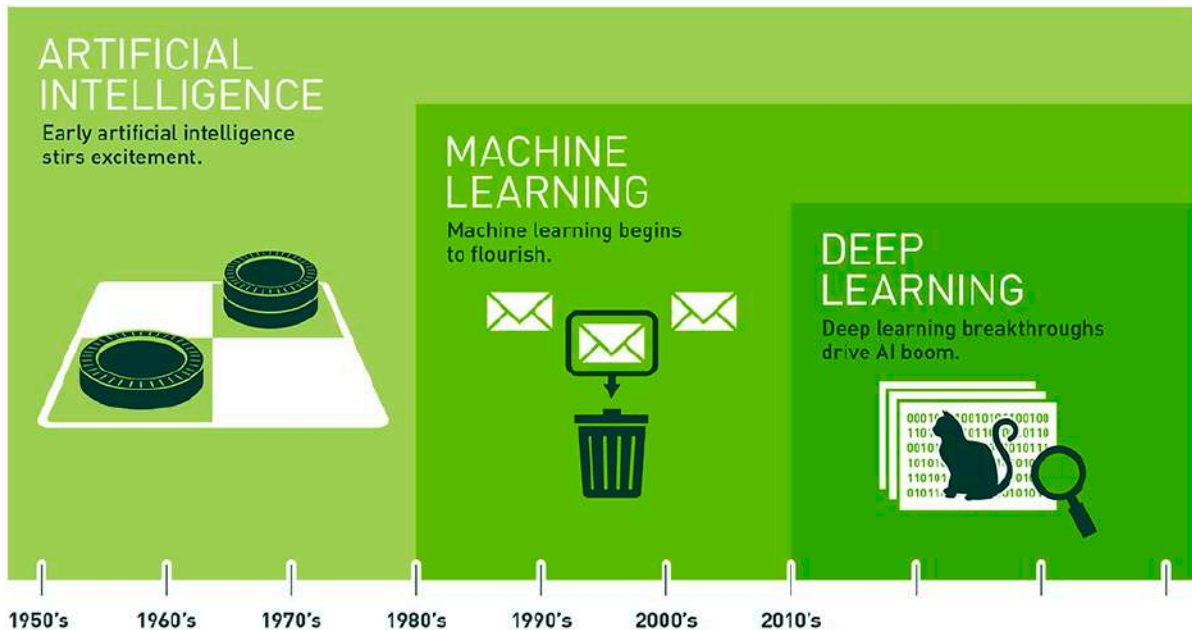# AWS | Artificial Intelligence
## Centerpiece for digital transformation

**Dylan Tong**, AI/Machine Learning Partner Solutions Architect

aws

# What is AI?

# Reinventing the Retail Experience

# AI Driven Stylist



INTRODUCING STYLE CHECK

Submit two photos to Style Check for a second opinion on which outfit looks better on you and why—based on fit, color, styling, and current trends. Through your feedback and input from our team of experienced fashion specialists, this advice gets smarter over time.

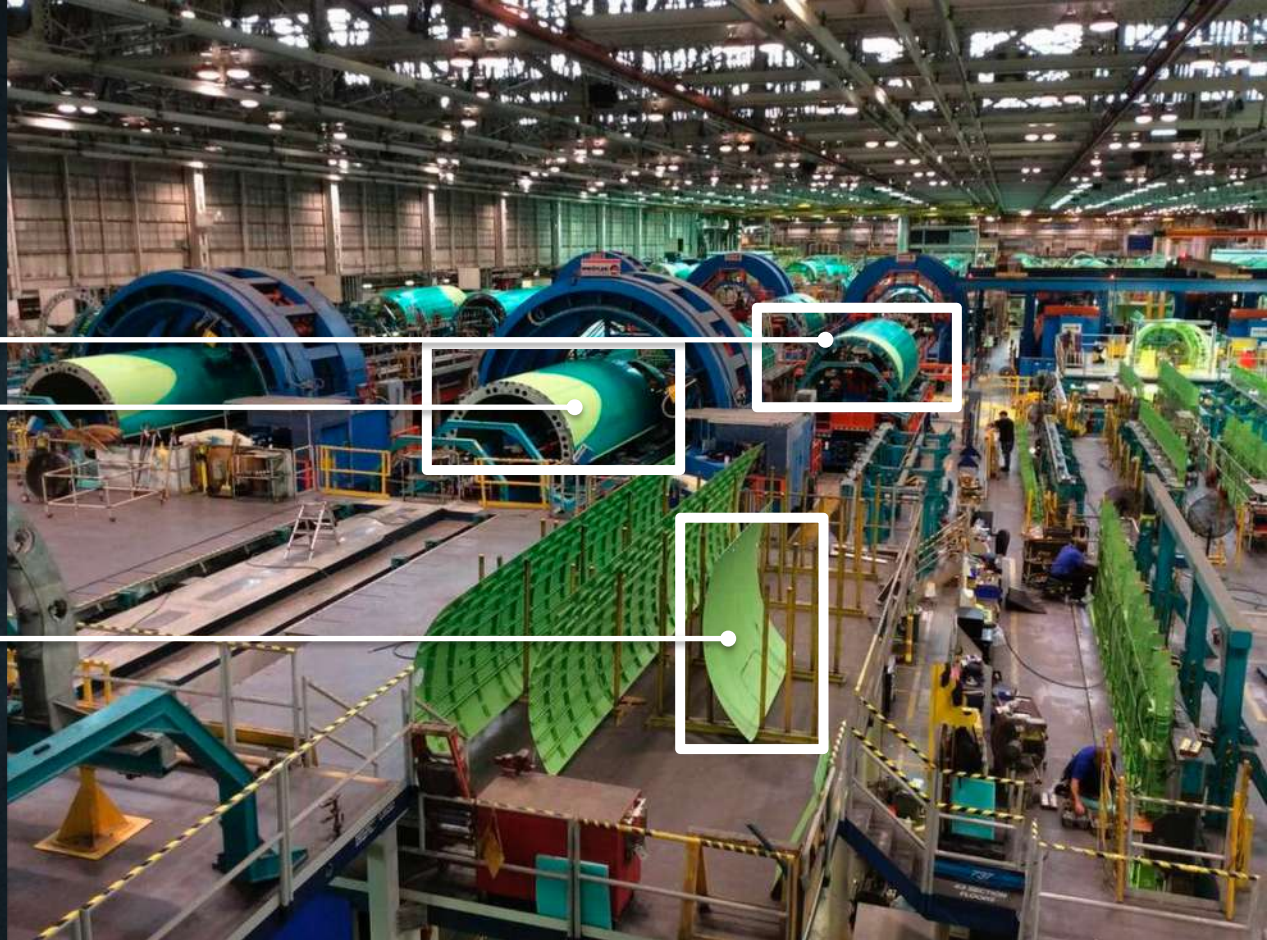72%    28%

The styling of the pieces looks better.

aws

# Smart Factory

**Inspection:**
Granular Quality
and Progress
Tracking

Asset Tracing
(SKU)

aws

# The Connected Worker



WORKER
SAFETY

● COMPLIANT

WORKER
TRACKING

JOHN ●

● DAVE

aws

# The Connected Worker

Improve worker productivity through new digital experiences:

- Voice-enabled Interfaces
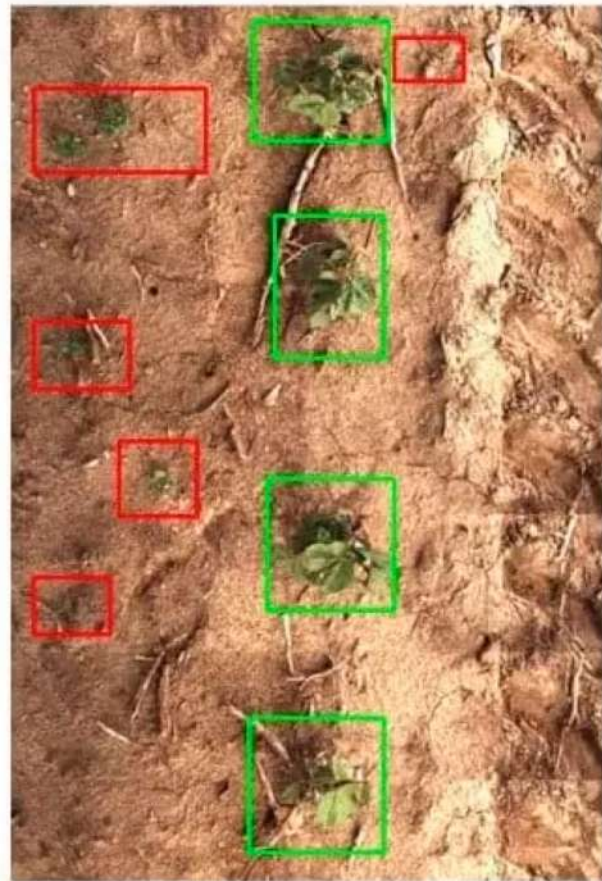
- Augmented Reality

Reduce key-strokes and lower the learning curve of complex machinery.
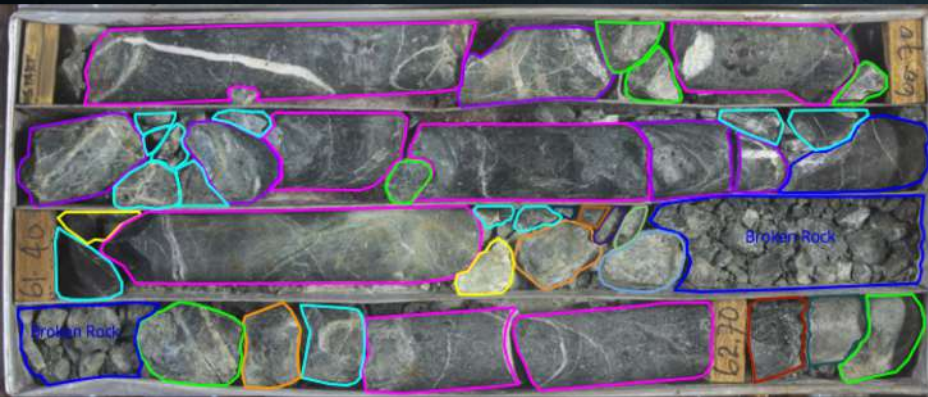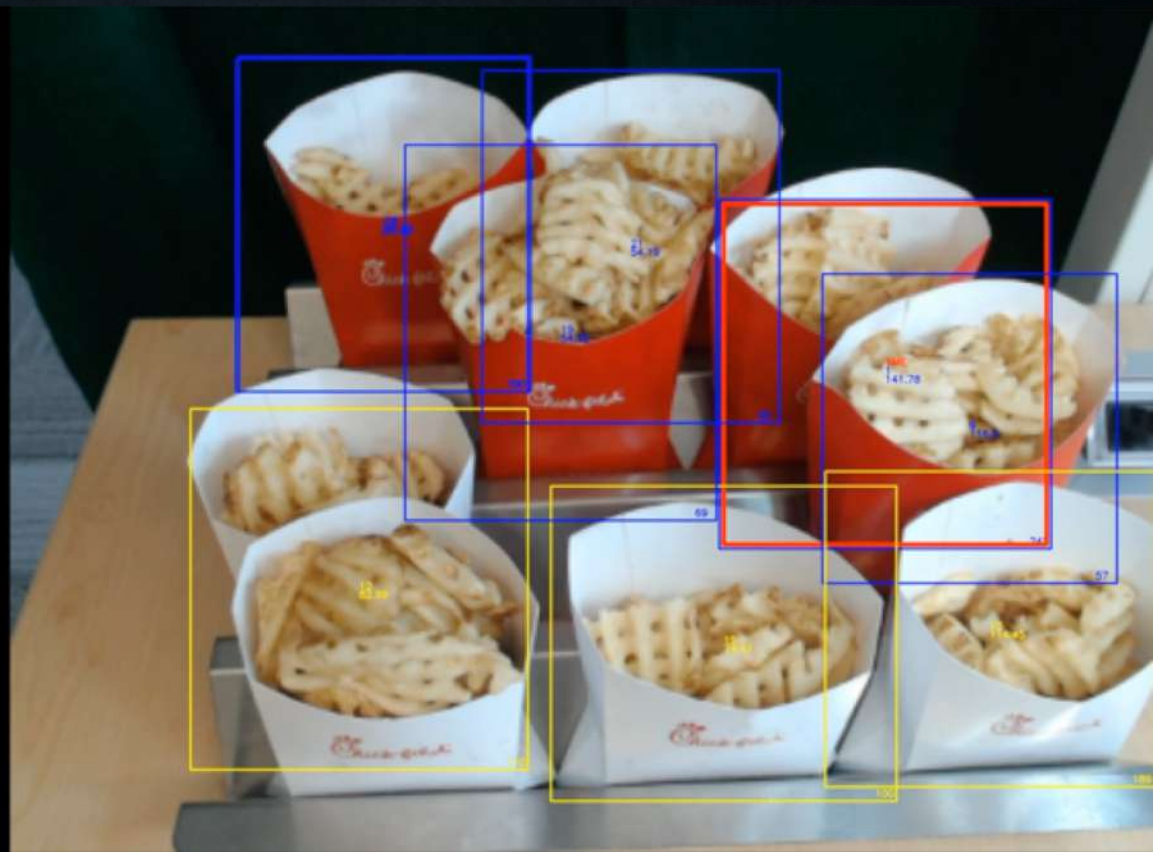
aws

## Smart Machine Era:

Reduces 90% use of herbicide through Blue River's smart sprayer. Computer Vision technology built on AWS and NVIDA GPU.
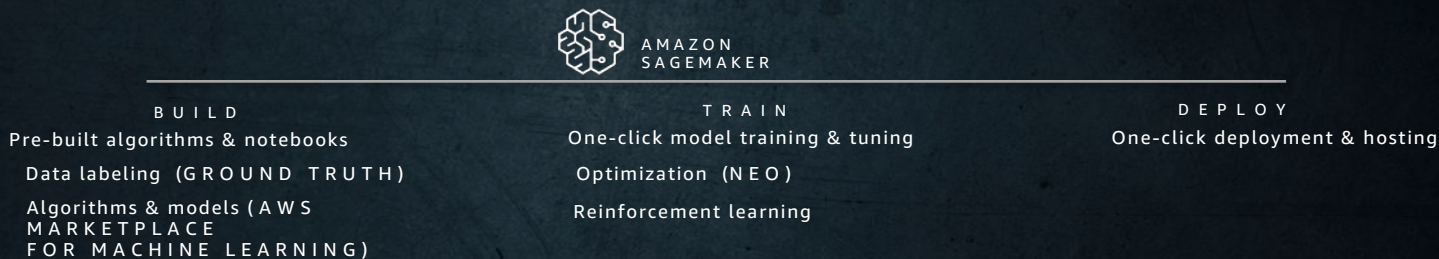
# *Technology*

**Mission:** to put machine learning in the hands of every developer.

aws

# The Amazon ML stack: Broadest & deepest set of capabilities

**AI SERVICES**

Vision | Speech | Language | Chatbots | Forecasting | Recommendations

REKOGNITION IMAGE — REKOGNITION VIDEO — TEXTRACT | POLLY — TRANSCRIBE | TRANSLATE — COMPREHEND | LEX | FORECAST | PERSONALIZE

---

**ML SERVICES**

AMAZON SAGEMAKER

BUILD
Pre-built algorithms & notebooks
Data labeling (GROUND TRUTH)
Algorithms & models (AWS MARKETPLACE FOR MACHINE LEARNING)

TRAIN
One-click model training & tuning
Optimization (NEO)
Reinforcement learning

DEPLOY
One-click deployment & hosting

---

**ML FRAMEWORKS & INFRASTRUCTURE**

Frameworks | Interfaces | Infrastructure

TensorFlow — mxnet — PYTORCH | GLUON — Keras | EC2 P3 & P3N — EC2 C5 — FPGAs — GREENGRASS — ELASTIC INFERENCE

aws

# The Amazon ML stack: Frameworks and Infrastructure



ML FRAMEWORKS & INFRASTRUCTURE

Frameworks

TensorFlow
mxnet
PYTORCH

Interfaces

GLUON
Keras

Infrastructure

EC2 P3 & P3N  EC2 C5  FPGAs  GREENGRASS  ELASTIC INFERENCE

aws

# ML Optimized Environments and Infrastructure

*More data science, less setup*



Deep Learning AMI

aws

# P3dn Instances: Optimized for Training at Scale
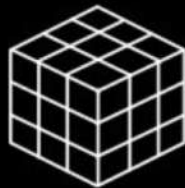
Reduce machine
learning training time

Better GPU
utilization

Support larger, more
complex models

### KEY FEATURES

**100Gbps** of networking
bandwidth
(4x more P3)

8 NVIDIA Tesla
V100 GPUs

32GB of
memory per GPU
(total 256GB,
2x more P3)

96 Intel
Skylake vCPUs
(50% more than P3)
with AVX-512

aws

# Challenges with Inference in Production



■ Prediction  ■ Training

Training 10%

Low utilization and high costs

One size does not fit all

aws

# Amazon Elastic Inference
*Reduce Deep Learning Inference costs up to 75%*

Lower inference costs

Match capacity
to demand

Available between 1 to 32 TFLOPS
per accelerator

## K E Y   F E A T U R E S

Integrated with
Amazon EC2 and
Amazon SageMaker

Support for TensorFlow, Apache
MXNet, and ONNX
with PyTorch coming soon

Single and
mixed-precision
operations

aws

# Up to 75% Reduction in Inference Costs

Inferences per second



Inception-v3

ResNet-152

SSD

c5.large + eia1.medium - $0.22/hr

c5.large + eia1.large      - $0.35/hr

c5.large + eia1.xlarge    - $0.61/hr

p2.xlarge    - $0.90/hr

p3.2xlarge - $3.06/hr

aws

# Requirements for Inference at the Edge



**BANDWIDTH**

1 billion cameras WW (2020)
10's of petabytes per day

**LATENCY**

30 images per second
200ms latency

**PRIVACY**

Confidentiality
Private cloud or on-premises storage

**AVAILABILITY**

50% of populated world < 8mbps
Bulk of uninhabited world no 3G+

aws

# AWS Greengrass

## Extend intelligence to the edge

| Local actions | Local triggers | Data and state sync | Security | Local resource access | Local ML inference |

# The Amazon ML stack: ML Services

**ML SERVICES**

AMAZON SAGEMAKER

BUILD
Pre-built algorithms & notebooks
Data labeling (GROUND TRUTH)
Algorithms & models (AWS MARKETPLACE FOR MACHINE LEARNING)

TRAIN
One-click model training & tuning
Optimization (NEO)
Reinforcement learning

DEPLOY
One-click deployment & hosting

**ML FRAMEWORKS & INFRASTRUCTURE**

Frameworks
TensorFlow
mxnet
PYTORCH

Interfaces
GLUON
Keras

Infrastructure
EC2 P3 & P3N
EC2 C5
FPGAs
GREENGRASS
ELASTIC INFERENCE

aws

# Amazon SageMaker

**BRINGING MACHINE LEARNING TO ALL DEVELOPERS**

**Business Need as a ML Problem**

**Desired Business Outcome**

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

**SIMPLIFY THE END-TO-END MACHINE LEARNING PROCESS**

aws

# Training Data for Supervised Learning

Raw Data

?

Training
Data

## Challenges

- Timely annotation of large data sets

- Managing workforces.
  - Ensuring labeling quality
  - Handling workflows

- Integration with ML development environment

aws

# Amazon SageMaker Ground Truth

*Label machine learning training data easily and accurately*

# Managed Notebooks: Exploration, Experimentation

# Your Choice. The Right Tools.

## NATIVE SUPPORT FOR MOST POPULAR FRAMEWORKS

## SAGEMAKER OPTIMIZED ALGORITHMS

- BlazingText Algorithm
- DeepAR Forecasting Algorithm
- Factorization Machines Algorithm
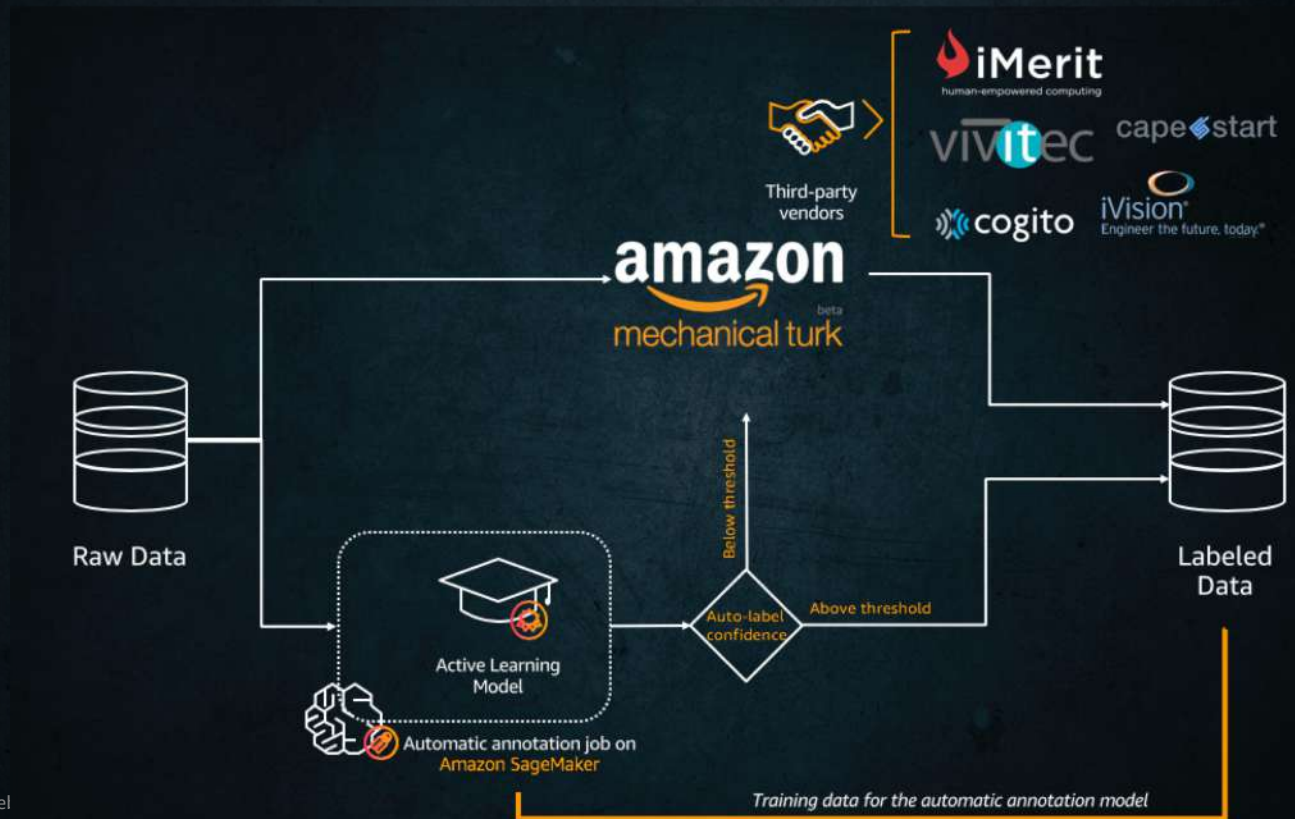- Image Classification Algorithm
- IP Insights Algorithm
- K-Means Algorithm
- K-Nearest Neighbors (k-NN) Algorithm
- Latent Dirichlet Allocation (LDA) Algorithm
- Linear Learner Algorithm
- Neural Topic Model (NTM) Algorithm
- Object2Vec Algorithm
- Object Detection Algorithm
- Principal Component Analysis (PCA) Algorithm
- Random Cut Forest (RCF) Algorithm
- Semantic Segmentation Algorithm
- Sequence-to-Sequence Algorithm
- XGBoost Algorithm

## aws marketplace INTEGRATION

# Zero Setup Training

# Automatic Model Tuning

# Amazon SageMaker Neo

*Train once, run everywhere with 2x performance and no accuracy lost*



Automate
Optimization

Neo

# Simplified Model Deployment and Hosting



**SageMaker Hosted Endpoints**

- Auto-scaling

- Performance monitoring

- A/B Testing

- Elastic Inference support

- Suited for real-time and batch workloads

aws

# Custom machine learning for your business

AMAZON SAGEMAKER

| REDUCE COSTS | INCREASE PERFORMANCE | EASE-OF-USE |
|:---:|:---:|:---:|
| **70%** | **10x** | **One-click** |
| cost reduction for data labeling using Ground Truth | better algorithm performance | model training & deployment |
| **75%** | **2x** | **Train once** |
| cost reduction for inference with Elastic Inference | performance increases from model optimization with Neo | run anywhere |

aws

# The Amazon ML stack: Broadest & deepest set of capabilities

## AI SERVICES

| Vision | | | Speech | | Language | | Chatbots | Forecasting | Recommendations |

REKOGNITION IMAGE  REKOGNITION VIDEO  TEXTRACT  POLLY  TRANSCRIBE  TRANSLATE  COMPREHEND  LEX  FORECAST  PERSONALIZE

AMAZON SAGEMAKER

## ML SERVICES

**BUILD**
Pre-built algorithms & notebooks
Data labeling (GROUND TRUTH)
Algorithms & models (AWS MARKETPLACE FOR MACHINE LEARNING)

**TRAIN**
One-click model training & tuning
Optimization (NEO)
Reinforcement learning

**DEPLOY**
One-click deployment & hosting

## ML FRAMEWORKS & INFRASTRUCTURE

Frameworks
TensorFlow
mxnet
PYTORCH

Interfaces
GLUON
K Keras

Infrastructure
EC2 P3 & P3N    EC2 C5    FPGAs    GREENGRASS    ELASTIC INFERENCE

aws

# Chatbots with Amazon Lex



Dialog Management
Session context maintenance

Deployment
One click deployment

Speech to Intent
ASR+NLU integrated into one API

Scale
Completely managed service

End to End

Text to Speech
Amazon Polly integrated into API

Business Logic
Native integration with AWS Lambda

Analytics
Monitor and improve

Security
Encrypted data in transit & at rest

aws

# Vision: Amazon Rekognition

**IMAGES**

Object and Scene Detection

Facial Analysis

Face Recognition

Unsafe Image Detection

Celebrity Recognition

Text in Image

**VIDEO**

Person Tracking

Real-time Live Stream

aws

# Vision: Amazon Rekognition

**IMAGES**

Object and Scene Detection

Facial Analysis

Face Recognition

Unsafe Image Detection

Celebrity Recognition

Text in Image

**VIDEO**

Person Tracking

Real-time Live Stream

aws

# The Connected Worker



WORKER SAFETY

🟢 COMPLIANT

WORKER TRACKING

JOHN 🔴

DAVE 🟠

aws

# Vision: Amazon Textract

*Form Extraction simplified*

| Full Name | | | Date of Birth | | | Gender |
|---|---|---|---|---|---|---|
| John | X | Doe | 01 | 01 | 1971 | Male ● |
| First | Middle | Last | MM | DD | YYYY | Female ○ |

**Output**

Full Name:
    First: John
    Middle: X
    Last: Doe

Date of Birth:
    MM: 01
    DD: 01
    YYYY: 1971

Gender:
    Male: True
    Female: False

✓ Logical groupings captured

✓ Relationships captured

✓ Glyphs captured

aws

# NLP: Amazon Comprehend

Amazon.com, Inc. is located in Seattle, WA and was founded July 5th, 1994 by Jeff Bezos. Our customers *love buying everything* from books to blenders at great prices

Named Entities
- Amazon.com: Organization
- Seattle, WA : Location
- July 5th, 1994: Date
- Jeff Bezos   : Person

Keyphrases
- Our customers
- books
- blenders
- great prices

Sentiment
- *Positive*

Language
- English

aws

# More Natural Language Processing…

**TEXT-TO-SPEECH**

Amazon
Polly

**SPEECH-TO-TEXT**

Amazon
Transcribe

**LANGUAGE
TRANSLATION**

Amazon
Translate

aws

# Solution Acceleration: Intelligent Building Blocks

# AutoML: Tailored Models Automated

**FORECASTING**          **RECOMMENDATIONS**

Amazon
Forecast

Amazon
Personalize

aws

# Computer Vision on AWS Jumpstart

aws

# Jumpstart Workshop Menu

- Object Detection
- TBD…

aws

# Object Detection Workshop

aws

# Reference Architectures

aws

# Cloud Inference on Streaming Video

aws

# AWS industrial IoT reference architecture



Industrial equipment

Protocol conversion

AWS Greengrass

ML inference

AWS IoT/AWS Greengrass/
AWS IoT Device Management/
AWS Device Defender

Protocol conversion

Industrial equipment

Amazon SNS

Email

AWS SMS

IoT rule (alerts)

IoT rule (all data)

Jupyter Notebook

AWS IoT Analytics

ML models

Amazon SageMaker

Amazon Kinesis Streams

Kinesis Data Firehose

Kinesis Data Analytics

Kinesis Data Firehose

IoT anomaly data repository

Amazon Athena

Amazon Athena

Amazon S3 Data Lake

Amazon Glacier

Amazon QuickSight

Real time and historical visualization

IAM

Amazon Cognito

IoT Cert

CloudTrail

AWS Config

CloudWatch

reInvent

aws

Workshop Objectives and Agenda

aws

# Select the Right Strategy

**Scope of Use Cases**

- **Object Detection,**
- **Segmentation,**
- **Classification…**

aws marketplace

**Amazon SageMaker Algorithms**

**Amazon Rekognition**

**Amazon SageMaker "Bring Your Own Script"**

**Amazon SageMaker "Bring Your Own Algorithm"**

**Time to Market**

Fully managed APIs. **No ML development required**

Full control over data sets, and hyperparameter tuning. **No coding.**

Full control over running TensorFlow, PyTorch, MXNet, Sklearn… scripts. **No need for managing SageMaker compliant containers**.

Integrate any algorithm **and benefit from automation.**

aws

# Right Tool for the Use Case



SEGMENTATION

CUSTOM OBJECT DETECTION

Amazon SageMaker

Amazon Rekognition

COMPLIANT

FACE SEARCH

JOHN

PEOPLE TRACKING

DAVE

aws

# Amazon SageMaker

**BRINGING MACHINE LEARNING TO ALL DEVELOPERS**

**Business Need as a ML Problem**

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

**Desired Business Outcome**

**SIMPLIFY THE END-TO-END MACHINE LEARNING PROCESS**

aws

**Lab 1:** Managing a high **quality** training set at **scale** using **SageMaker GroundTruth**

Business Need as a ML Problem

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

Desired Business Outcome

**SIMPLIFY THE END-TO-END MACHINE LEARNING PROCESS**

aws

**Lab 2:** Train, tune and deploy a custom object detector (**SSD**) with **zero coding**.

Business Need as a ML Problem

Desired Business Outcome

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
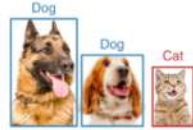
SIMPLIFY THE END-TO-END MACHINE LEARNING PROCESS

aws

**Lab 3: "Bring Your Own Script:"** train, tune and deploy a custom object detector (**YOLOv3**) on GluonCV (MXNet).

**Business Need as a ML Problem**

**Desired Business Outcome**

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

**SIMPLIFY THE END-TO-END MACHINE LEARNING PROCESS**

aws

**GLUON**

**Supported Applications**

| Application | Illustration | Available Models |
|---|---|---|
| **Image Classification:** recognize an object in an image. | Dog | 50+ models, including ResNet, MobileNet, DenseNet, VGG, ... |
| **Object Detection:** detect multiple objects with their bounding boxes in an image. | Dog Dog Cat | Faster RCNN, SSD, Yolo-v3 |
| **Semantic Segmentation:** associate each pixel of an image with a categorical label. | Background Dog Dog Cat | FCN, PSP, DeepLab v3 |
| **Instance Segmentation:** associate each pixel of an image with an instance label. | Background Dog 1 Dog 2 Cat 1 | Mask RCNN |
| **Pose Estimation:** detect human pose from images. | | Simple Pose |

**Algorithm variants :** For instance, in **Object Detection**, different algorithms offer trade-offs between **accuracy (mAP)** and **latency (fps)**
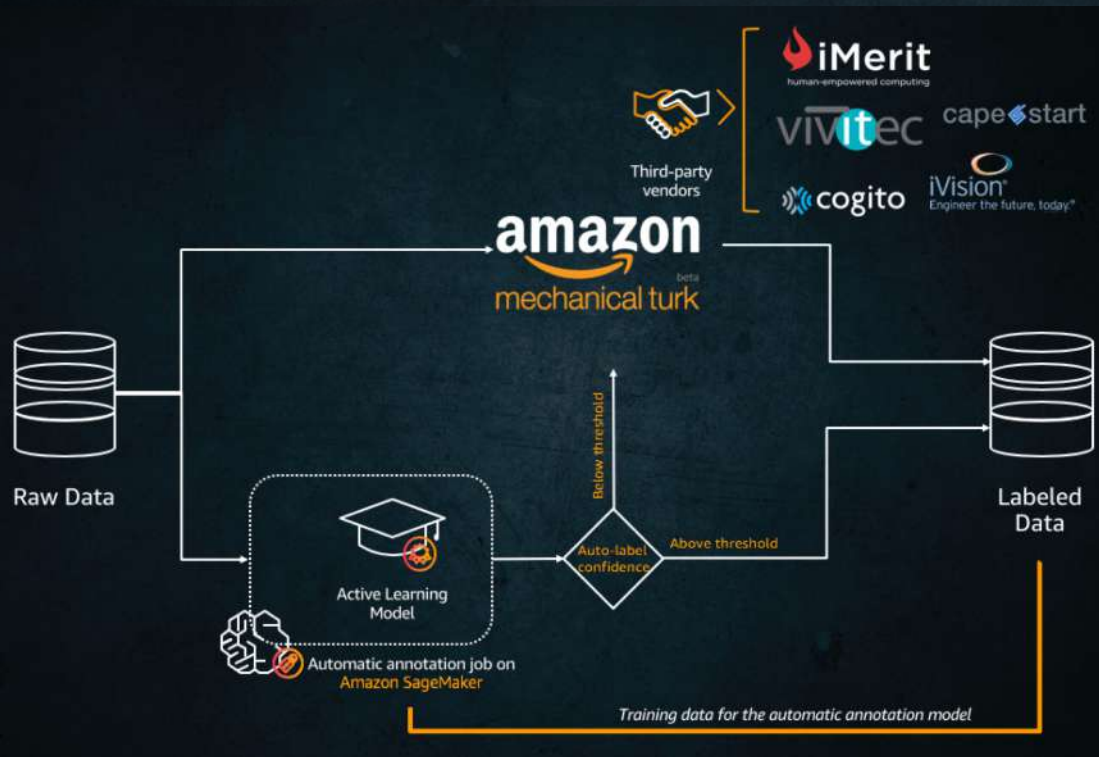
aws

# Future labs

- By feedback and demand! dylatong@amazon.com
- Recognition:
  - Face Search
  - People Tracking
- Textract solutions
- More Amazon SageMaker use cases:
  - Segmentation, Pose Estimation, Similarity Search

aws

Lab 1
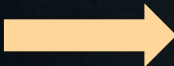
aws

# Manage a high-quality data set at scale



1. Launch a Notebook Instance.

2. Manage a private workforce.

3. Create an annotation job for Object Detection.

4. Generate a dataset and metadata compatible with Amazon SageMaker algorithms without further data wrangling!

aws

# Lab 2

aws

# Create a custom object detector with zero coding

## SAGEMAKER OPTIMIZED ALGORITHMS

- BlazingText Algorithm
- DeepAR Forecasting Algorithm
- Factorization Machines Algorithm
- Image Classification Algorithm
- IP Insights Algorithm
- K-Means Algorithm
- K-Nearest Neighbors (k-NN) Algorithm
- Latent Dirichlet Allocation (LDA) Algorithm
- Linear Learner Algorithm
- Neural Topic Model (NTM) Algorithm
- Object2Vec Algorithm
- **Object Detection Algorithm**
- Principal Component Analysis (PCA) Algorithm
- Random Cut Forest (RCF) Algorithm
- Semantic Segmentation Algorithm
- Sequence-to-Sequence Algorithm
- XGBoost Algorithm

1. Configure a hyperparameter tuning job for an Object Detection Algorithm.

2. Train on GPU

3. Deploy a managed endpoint.

4. Test and visualize!

aws

Lab 3

aws

# Bring Your Own Script and automate the ML process

## Examples

**Training:** Only modifications required is to set script certain parameters values from SageMaker container environment variables.

GluonCV YOLOv3 training script
PyTorch Siamese Network training script

**Inference:** requires overriding programmatic interface implementation.

GluonCV YOLOv3 model serving script
PyTorch Siamese Network model serving script

- **input_fn**: request format pre-processing
- **model_fn:** how to load the model
- **predict_fn:** inference logic
- **output_fn:** response format processing

1. Bring your own YoloV3 script on GluonCV

2. Prepare your data set and environment

3. Explore and prototype locally

4. Automate model tuning, and train

5. Deploy, test and visualize!

aws

https://github.com/dylan-tong-aws/aws-cv-jumpstarter

aws

aws.ai

# Appendix

aws

# How we can help…

### ML Solutions Lab

Brainstorming
Custom modeling
Training
Work side-by-side with Amazon experts

### Machine Learning Training & Certification

Practical education on ML for new & experienced practitioners
Based on the same material used to train Amazon developers

aws