



Using AWS Cloud for ML

Neil Mackin, AWS
nemackin@amazon.com

Our mission at AWS

Put machine learning in the hands
of every developer



Our unique approach



Customer-focused

90%+ of our ML roadmap is defined by customers



Pace of innovation

200+ new ML launches in the last year



Breadth & depth

More AI and ML services in production than any other provider



Multi-framework

Support for the most popular frameworks



Security & analytics

Deepest set of security and encryption features, with robust analytics capabilities



Embedded R&D

Customer-centric approach to advancing the state of the art

© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark



More machine learning happens on AWS than anywhere else

10,000+

customers have used machine learning on AWS

81%

of **deep learning** in the cloud runs on AWS

85%

of **TensorFlow** projects in the cloud run on AWS



AWS holds the top spots on **Stanford's** deep learning benchmark, DAWN, for fastest training time, lowest cost, lowest inference latency

© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark



Technology

Culture

© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark



Technology

Bringing AI into your digital transformation requires a new "stack" that makes it easier to put ML to work

© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark



The Amazon ML stack: Broadest & deepest set of capabilities

AI SERVICES

Easily add intelligence to applications without machine learning skills

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

ML SERVICES

Build, train, and deploy machine learning models fast

Data labeling | Pre-built algorithms & notebooks | One-click training and deployment

ML FRAMEWORKS & INFRASTRUCTURE

Flexibility & choice, highest-performing infrastructure

Support for ML frameworks | Compute options purpose-built for ML

© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark



The Amazon ML stack: Broadest & deepest set of capabilities

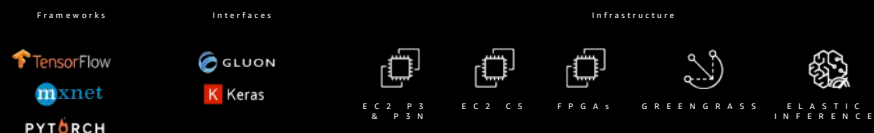
AI SERVICES



ML SERVICES



ML FRAMEWORKS & INFRASTRUCTURE



© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark



Amazon Rekognition

Deep learning-based image and video analysis



Object, Scene & Activity Recognition



Facial Recognition



Facial Analysis



Person Tracking



Unsafe Content Detection



Celebrity Recognition



Text in Images



Amazon Rekognition Images



Given this photo of Andy Jassy

...



Tell me if you find him in this image

Rekognition: Confidence level that source face is in target photo: 98 %

Amazon Rekognition Images

```

bucket = 'my-photos-bucket'
key_s = 'andy_jassy_headshot.png'
key_t = 'interview.png'

rek_client=boto3.client('rekognition', 'ap-southeast-2')

response = rek_client.compare_faces( SimilarityThreshold = 75,
    sourceImage={ 'S3Object': { 'Bucket': bucket, 'Name': key_s, } }, TargetImage={
'S3Object': { 'Bucket': bucket, 'Name': key_t, } },)

confidence = response['FaceMatches'][0]['Similarity']

print('Confidence level that source face is in target photo: ',confidence,'%')

```

© 2018, Amazon Web Services, Inc. or Its Affiliates. All rights reserved.

Amazon Rekognition Images

```

bucket = 'my-photos-bucket'
key_s = 'andy_jassy_headshot.png'
key_t = 'interview.png'

rek_client=boto3.client('rekognition', 'ap-southeast-2')

response = rek_client.compare_faces( SimilarityThreshold = 75,
    sourceImage={ 'S3Object': { 'Bucket': bucket, 'Name': key_s, } }, TargetImage={
'S3Object': { 'Bucket': bucket, 'Name': key_t, } },)

confidence = response['FaceMatches'][0]['Similarity']

print('Confidence level that source face is in target photo: ',confidence,'%')

```

© 2018, Amazon Web Services, Inc. or Its Affiliates. All rights reserved.

Build applications that understand text

APPLICATION SERVICES



AMAZON
REKOGNITION



AMAZON
REKOGNITION
VIDEO



AMAZON
POLLY



AMAZON
TRANSCRIBE



AMAZON
TRANSLATE



AMAZON
COMPREHEND



AMAZON
LEX

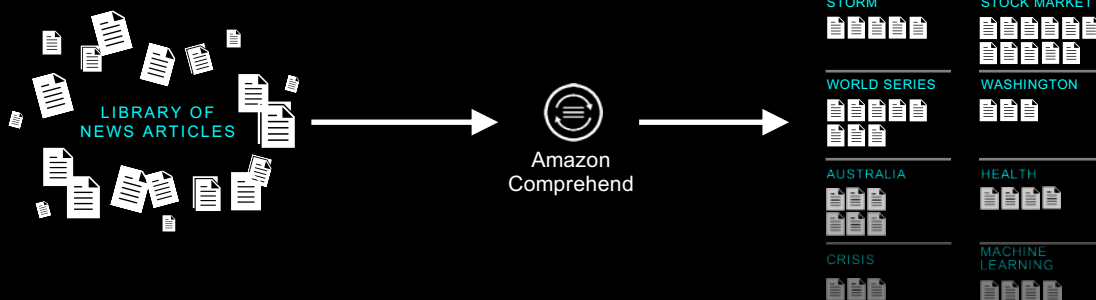
Amazon Comprehend

- Natural language processing (NLP) service that finds insights and relationships in text
- Detect sentiment & language, model documents by topic
- Identify entities, key phrases, and syntax
- Financial services, healthcare, insurance, government, document classification & search, customer analytics ++



Amazon Comprehend

Discover insights and relationships in text



Amazon Comprehend – extract insights from text

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

Key phrases

- Our customers
- books
- blenders
- great prices

Sentiment

Positive

Language

English

© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark

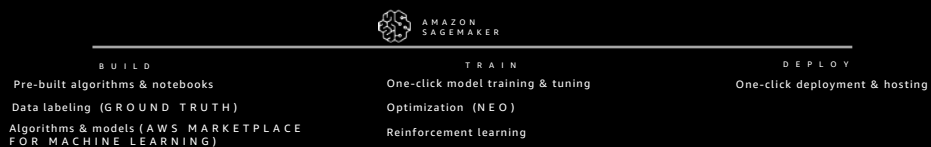


The Amazon ML stack: Broadest & deepest set of capabilities

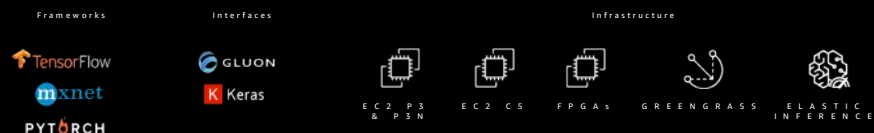
AI SERVICES



ML SERVICES



ML FRAMEWORKS & INFRASTRUCTURE



© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark



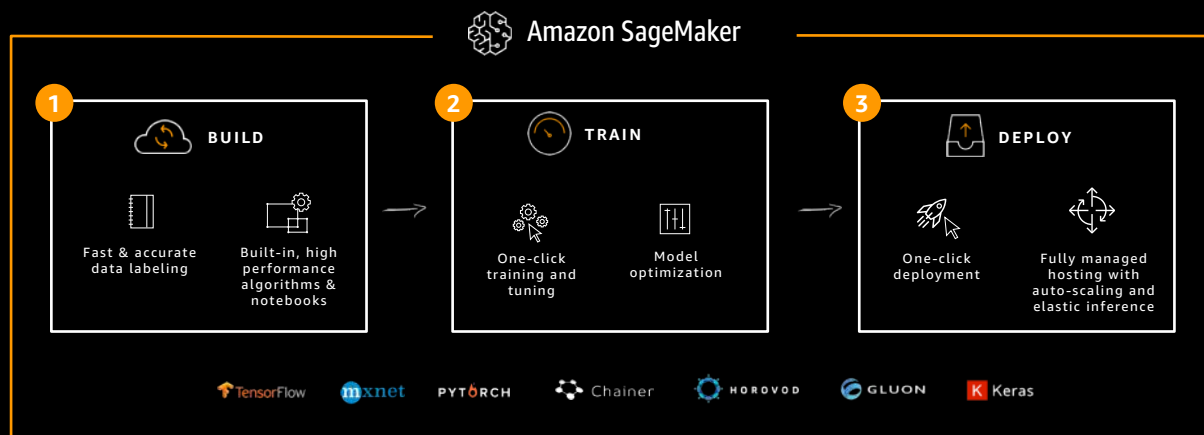
Free Data Scientists to do Data Science



© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark



Custom machine learning for your business



© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark



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

© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark



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

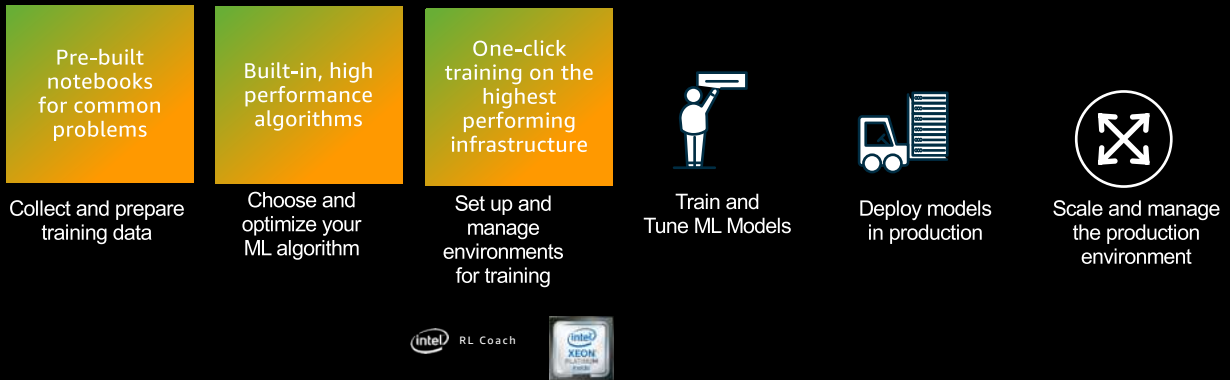


Scale and manage
the production
environment

© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark



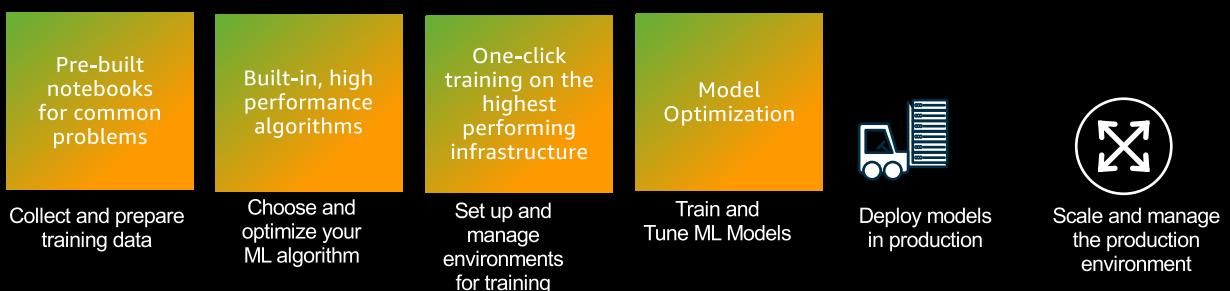
Amazon SageMaker: Build, Train, and Deploy ML Models at Scale



© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark



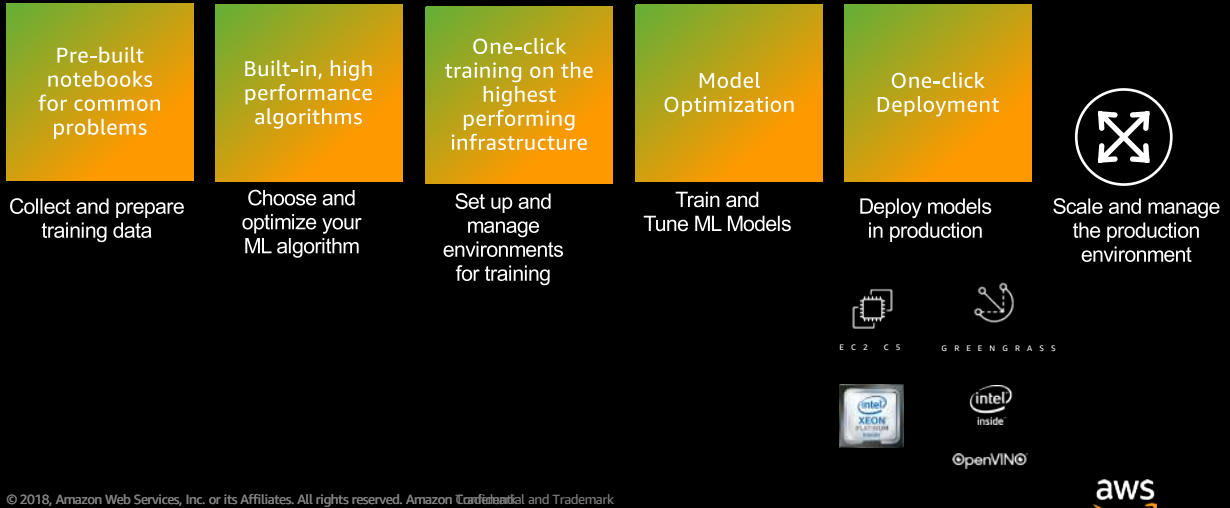
Amazon SageMaker: Build, Train, and Deploy ML Models at Scale



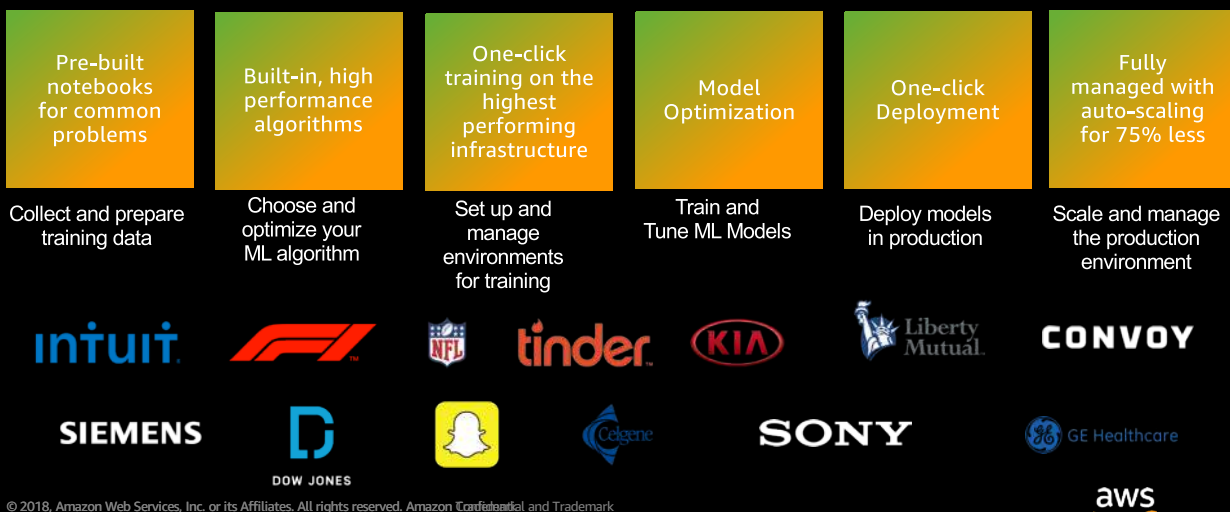
© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark



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



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



Amazon SageMaker Built-in Algorithms

- Image Classification
- Object Detection
- K-Nearest Neighbors (k-NN)
- Linear Learner
- Factorization Machines
- XGBoost
- Sequence2Sequence
- Principal Component Analysis (PCA)
- Latent Dirichlet Allocation (LDA)
- Neural Topic Model (NTM)
- DeepAR Forecasting

- Semantic Segmentation
- BlazingText
- Random Cut Forest

Three new ones:

- IP Insights: built-in algorithms for detecting suspicious IP addresses
- Object2Vec : Low dimensional embeddings for high dimensional objects
- K-means: clustering unsupervised grouping

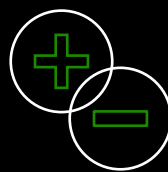
© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark



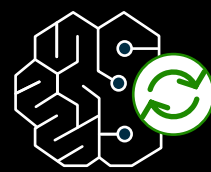
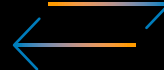
How does RL work?



Simulation environment



Scoring function



RL algorithm

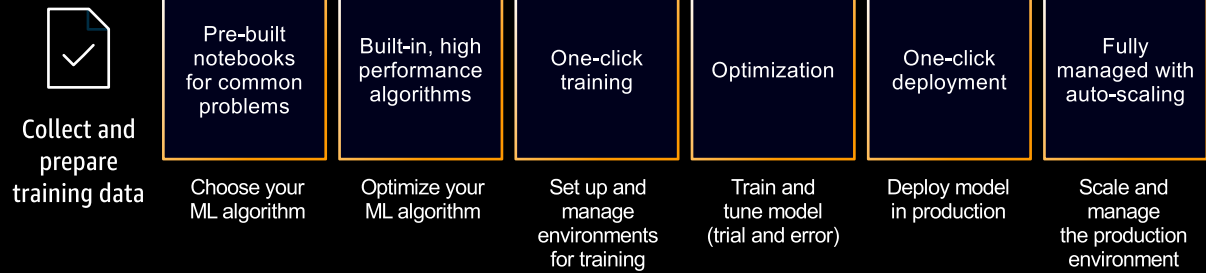
USE CASES

Supply chain simulation, manufacturing process, robot manipulation, autonomous car, drone navigation...

© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark



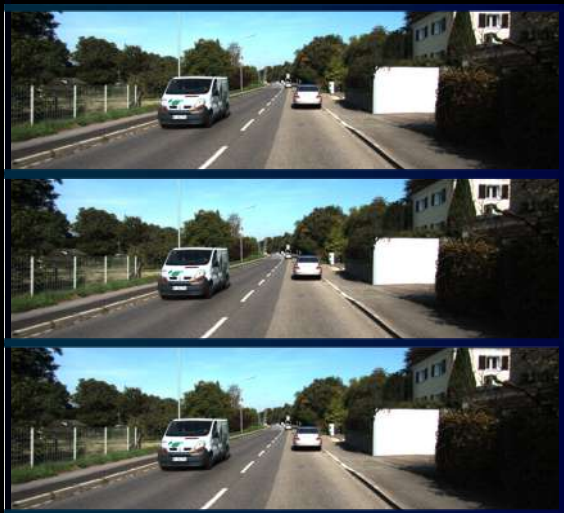
Amazon SageMaker: Build, train, and deploy ML



© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark



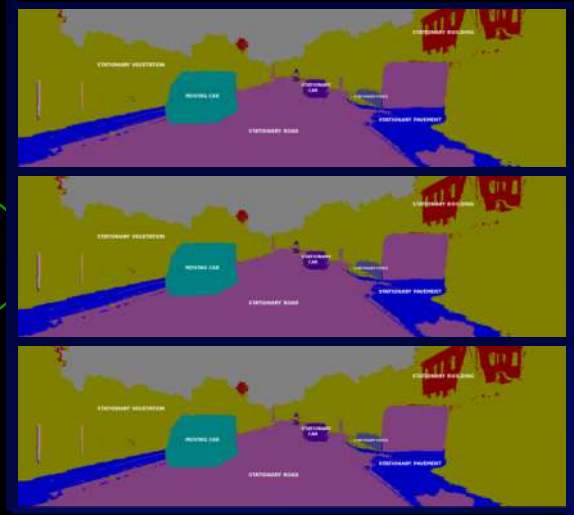
Successful models require high-quality data



© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark



Successful models require high-quality data



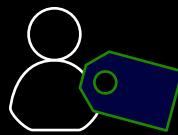
© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark



Amazon SageMaker ground truth Label machine learning training data easily and accurately



Quickly label
training data



Easily integrate
human labelers



Get accurate
results

KEY FEATURES

Automatic labeling via
machine learning

Ready-made and
custom workflows for
image bounding box,
segmentation, and text

Private and public
human workforce

Label
management

© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark



How it works



Raw Data

© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark



How it works



Raw Data

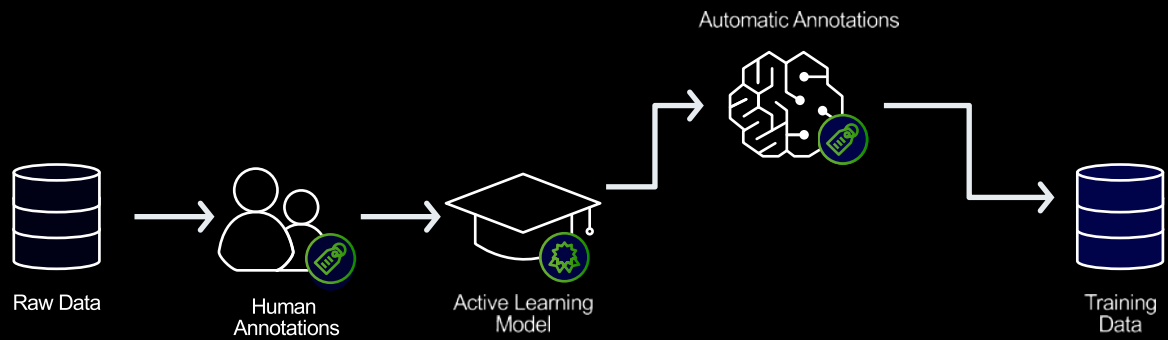


Human
Annotations

© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark



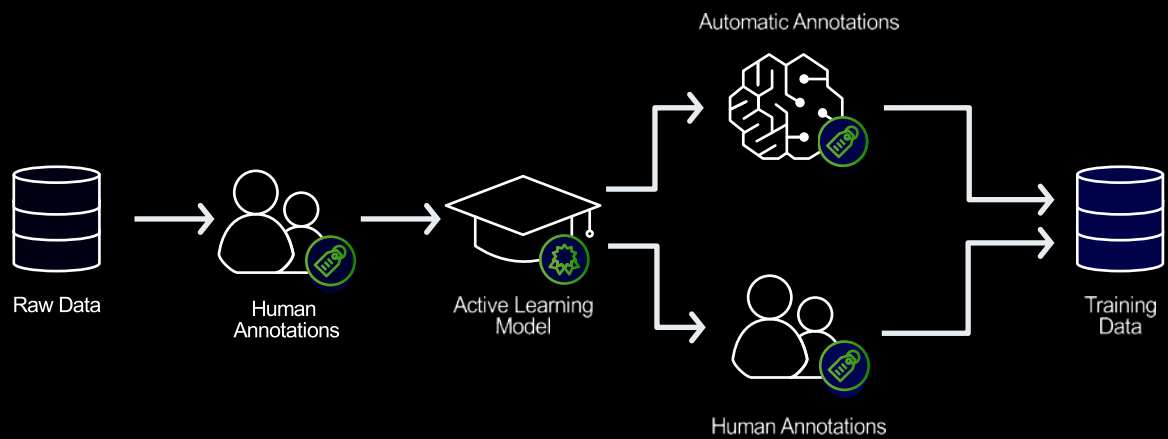
How it works



© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark



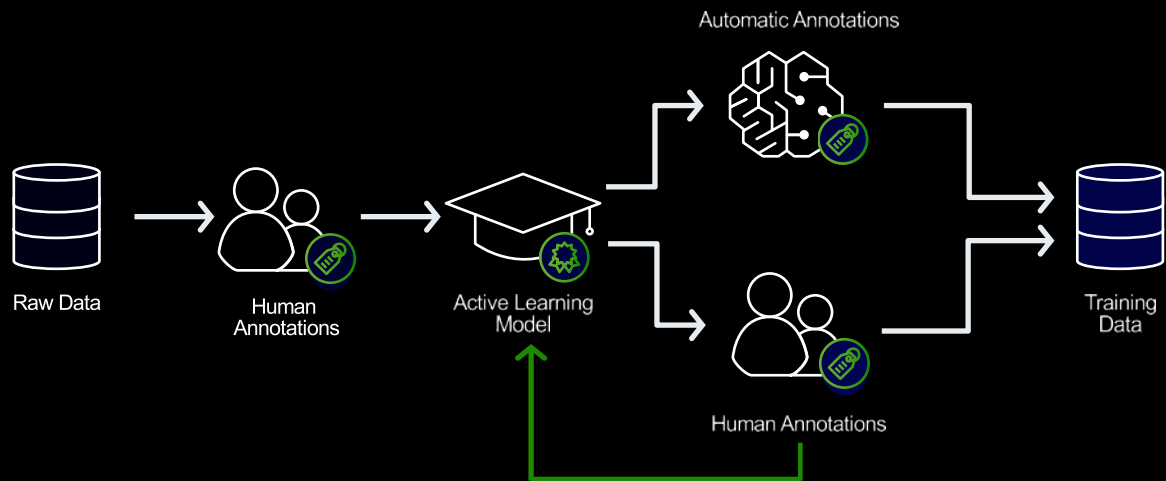
How it works



© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark



How it works



© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark



Creating training data



Mechanical
turk workers



Private labeling
workforce



Third-party
vendors



human-empowered computing



Engineer the future, today.®



© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark



Real-time Fraud Detection with Amazon SageMaker

AI and ML at Intuit have three areas of focus.

- ✓ Smart Products
- ✓ Fraud detection and prevention
- ✓ Customer Care and Expert advice

In order to keep fraudsters out of their systems and data, Intuit always stays several moves ahead by leveraging AI/ML-generated insights from data that can determine **real-time fraud detection** in TurboTax: Specifically

- Account take-over detection at login
- Identity theft detection at filing



Saving lives with Amazon SageMaker

Harnessing data and analytics across hardware, software and biotech, GE Healthcare is transforming healthcare by delivering better outcomes for providers and patients.

- Amazon SageMaker allows GE Healthcare to access powerful Artificial Intelligence tools and services to advance improved patient care.

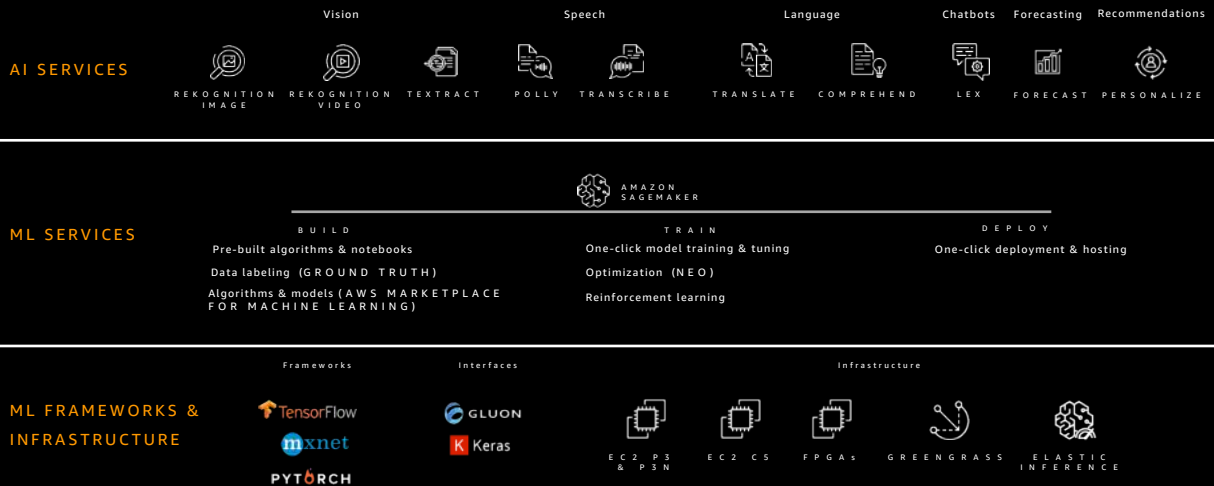


“The scalability of Amazon SageMaker, and its ability to integrate with native AWS services, adds enormous value for us. We are excited about how our continued collaboration between the GE Health Cloud and Amazon SageMaker will drive better outcomes for our healthcare provider partners and deliver improved patient care.”

- Sharath Pasupunuti, AI Engineering Leader



The Amazon ML stack: Broadest & deepest set of capabilities



© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark



Highest-performing infrastructure for your business



Build custom algorithms using the ML frameworks



Fastest and lowest-cost compute options for ML workloads



Elastic compute to provision just-right compute for your ML workloads

© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark



The best place to run TensorFlow

Amazon SageMaker is the best place to run TensorFlow in the cloud

- Fully-managed training and hosting
- Near-linear scaling across 100s of GPU
- 75% lower inference costs with Amazon Elastic Inference
- 3x faster network throughput with EC2 P3

STOCK TENSORFLOW

65%

AWS-OPTIMIZED TENSORFLOW

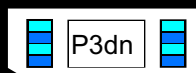
90%

Scaling efficiency with 256 GPUs

© 2018, Amazon Web Services, Inc. or its Affiliates. All rights reserved. Amazon Confidential and Trademark



Bottom Layer: Frameworks & interfaces



Model	p3dn.24xlarge
NVIDIA V100 Tensor Core GPUs	8
GPU Memory	256 GB
NVIDIA NVLink	300 GB/s
vCPUs	96
Main Memory	768 GiB
Local Storage	2 x 900 GB NVMe SSD
Network Bandwidth	100 Gbps
EBS-Optimized Bandwidth	14 Gbps



AWS Deep Learning AMI



Thank you

