



# AI and Machine Learning on AWS

**Julien Simon, Principal Evangelist, AI/ML**  
**@julsimon**



## Welcome to Amazon.com Books!

*One million titles,  
consistently low prices.*

(If you explore just one thing, make it our personal notification service. We think it's very cool!)

### SPOTLIGHT! -- AUGUST 16TH

These are the books we love, offered at Amazon.com low prices. The spotlight moves **EVERY** day so please come often.

### ONE MILLION TITLES

Search Amazon.com's [million title catalog](#) by author, subject, title, keyword, and more... Or take a look at the [books we recommend](#) in over 20 categories... Check out our [customer reviews](#) and the [award winners](#) from the Hugo and Nebula to the Pulitzer and Nobel... and [bestsellers](#) are 30% off the publishers list...

### EYES & EDITORS, A PERSONAL NOTIFICATION SERVICE

Like to know when that book you want comes out in paperback or when your favorite author releases a new title? Eyes, our tireless, automated search agent, will send you mail. Meanwhile, our human editors are busy previewing galleys and reading advance reviews. They can let you know when especially wonderful works are published in particular genres or subject areas. Come in, [meet Eyes](#), and have it all explained.

### YOUR ACCOUNT

Check the status of your orders or change the email address and password you have on file with us. Please note that you **do not** need an account to use the store. The first time you place an order, you will be given the opportunity to create an account.

# Amazon.com, 1995

# « Two Decades of Recommender Systems at Amazon.com » (2017)



G.D. Linden, J.A. Jacobi, and E.A. Benson,  
Collaborative Recommendations Using  
Item-to-Item Similarity Mappings,

US Patent 6,266,649, Amazon.com,  
Patent and Trademark Office, 2001  
(filed 1998).

<https://www.computer.org/csdl/mags/ic/2017/03/mic2017030012.html>





amazon  
fulfillment

# amazon echo



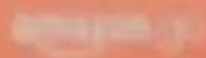




amazon



JUST  
WALK  
OUT  
SHOPPING



# Jeff Bezos' letter to Amazon shareholders

*“We are **solving problems** with **machine learning** and **artificial intelligence** that were in the realm of science fiction for the last several decades. Natural **language** understanding, machine **vision** problems, it really is an amazing renaissance.”*

<https://www.geekwire.com/2017/jeff-bezos-explains-amazons-artificial-intelligence-machine-learning-strategy/>



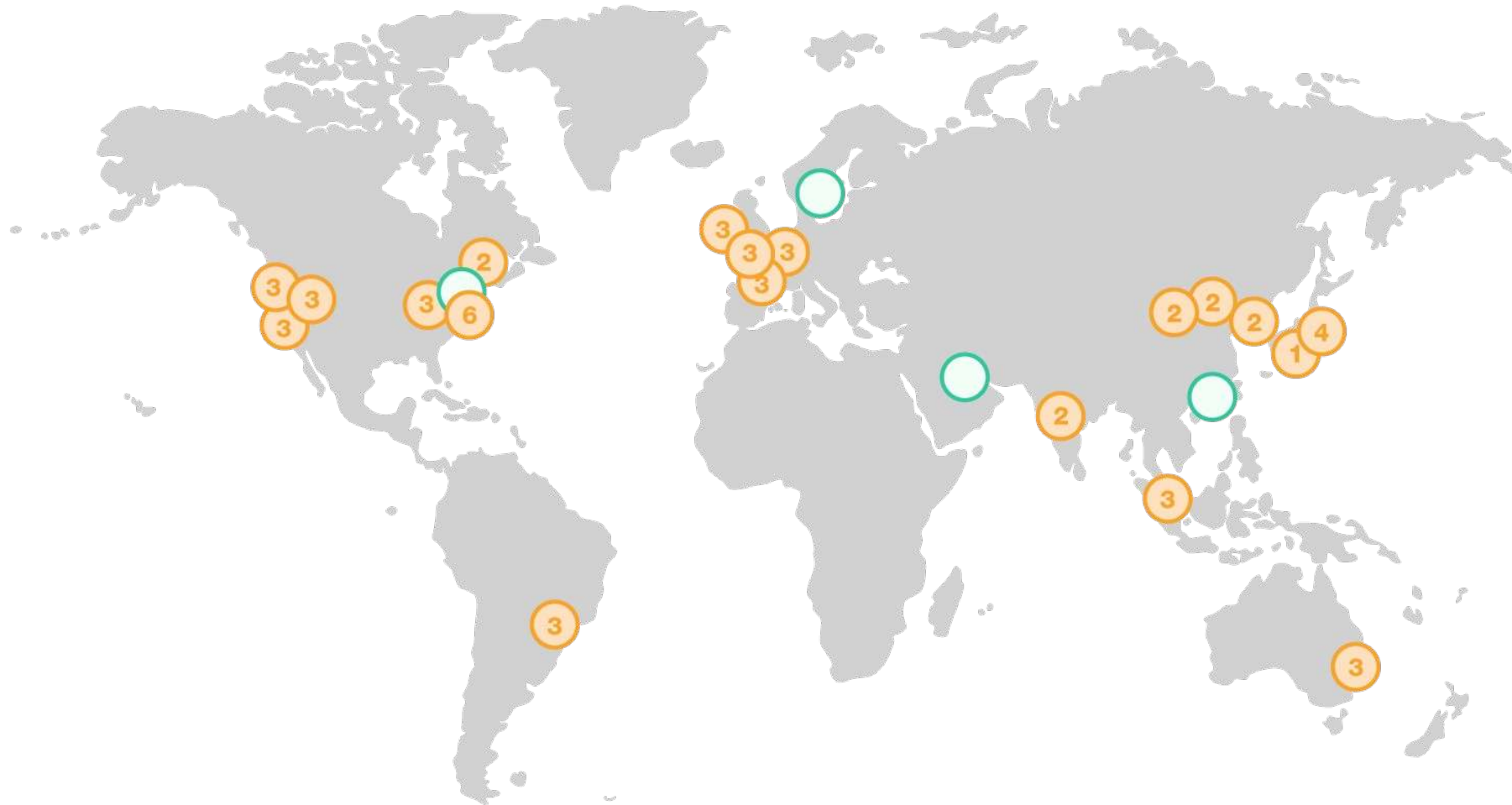
# AWS: our mission

Put AI and Machine Learning  
in the hands of every developer  
and data scientist

Machine Learning requires  
a solid infrastructure foundation

# AWS Global Infrastructure

**18** Regions – **54** Availability Zones – **103** Edge Locations





# Gartner: Cloud Infrastructure as a Service, June 2017

Figure 1. Magic Quadrant for Cloud Infrastructure as a Service, Worldwide



Source: Gartner (June 2017)

Machine Learning requires  
secure infrastructure

# AWS compliance standards

## Certifications & Attestations

Cloud Computing Compliance  
Controls Catalogue (C5)  
Cyber Essentials Plus

DE 

UK 

DoD SRG  
FedRAMP

US 

US 

FIPS

US 

IRAP

AU 

ISO 9001



ISO 27001



ISO 27017



ISO 27018  
MLPS Level 3



CN 

MTCS

SG 

PCI DSS Level 1



SEC Rule 17c-4(f)

US 

## Laws, Regulations and Privacy

CISPE

EU Model Clauses

FERPA

GLBA

HIPAA

HITECH

IRS 1075

ITAR

My Number Act

Data Protection Act – 1988

VPAT / Section 508

Data Protection Directive

Privacy Act [Australia]

Privacy Act [New Zealand]

EU 

EU 

US 

US 

US 



US 

US 

JP 

UK 

US 

EU 

AU 

NZ 

## Alignments & Frameworks

CIS (Center for Internet Security)

CJIS (US FBI)

CSA (Cloud Security Alliance)

Esquema Nacional de Seguridad

EU-US Privacy Shield

FISC

FISMA

G-Cloud

GxP (US FDA CFR 21 Part 11)

ICREA

IT Grundschutz

MITA 3.0 (US Medicaid)

MPAA

NIST



US



ES



EU



JP 

US



UK



US



DE



US



US



US



# AWS is GDPR ready

[AWS Security Blog](#)

## All AWS Services GDPR ready

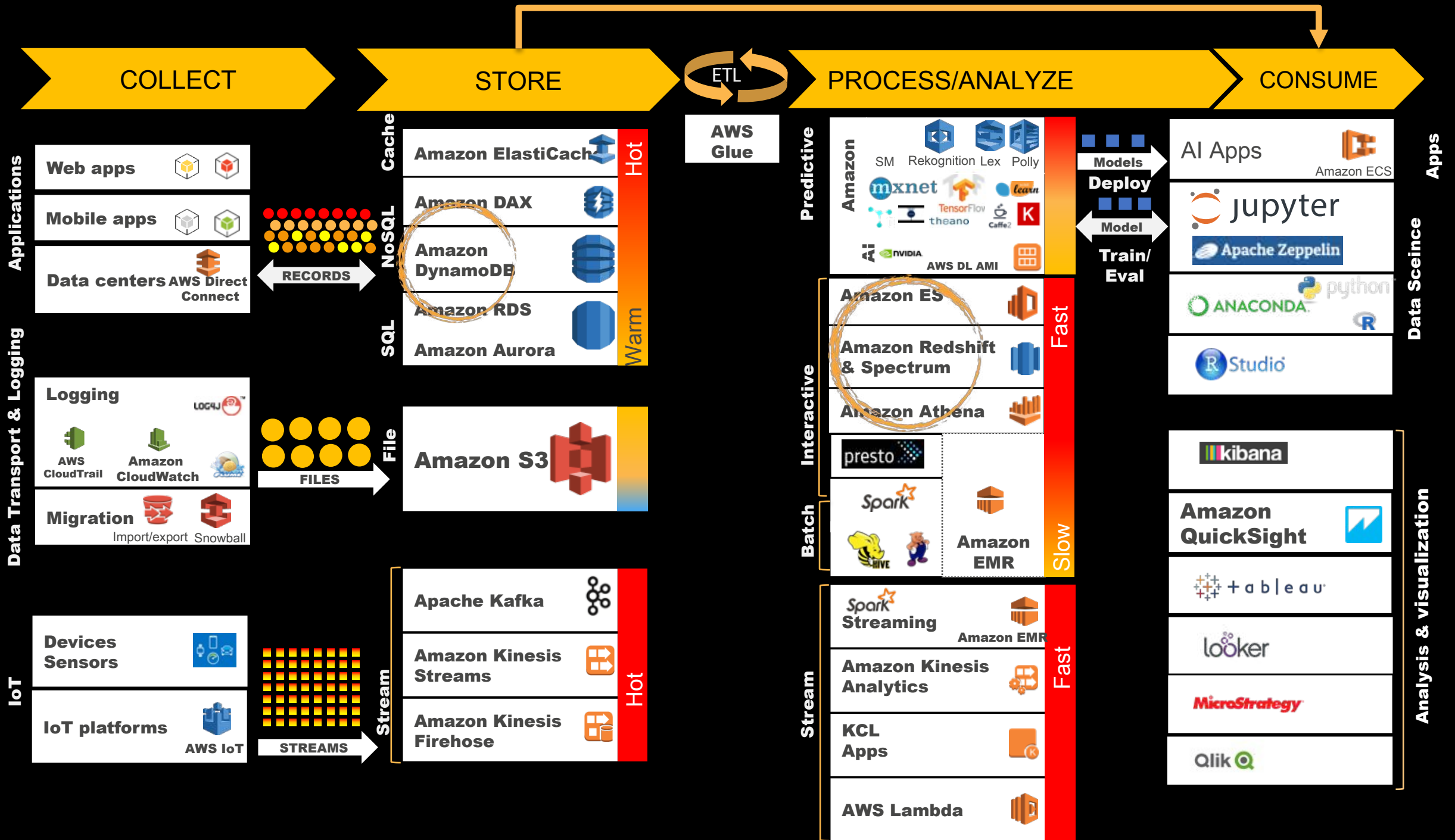
by Chad Woolf | on 26 MAR 2018 | in [Amazon GuardDuty](#), [Amazon Inspector\\*](#), [Amazon Macie\\*](#), [AWS Config\\*](#), [Compliance\\*](#), [Security, Identity, & Compliance\\*](#), [Webinars\\*](#) |

[Permalink](#) | [Comments](#) | [Share](#)

Today, I'm very pleased to announce that AWS services comply with the General Data Protection Regulation (GDPR). This means that, in addition to benefiting from all of the measures that AWS already takes to maintain services security, customers can deploy AWS services as a key part of their GDPR compliance plans.

This announcement confirms we have completed the entirety of our GDPR service readiness audit, validating that all generally available services and features adhere to the high privacy bar and data protection standards required of data processors by the GDPR. We completed this work two months ahead of the May 25, 2018 enforcement deadline in order to give customers and APN partners an environment in which they can confidently build their own GDPR-compliant products, services, and solutions.

Machine Learning requires  
extensive (big) data services





# Dynamo (2007)

## **Dynamo: Amazon's Highly Available Key-value Store**

Giuseppe DeCandia, Deniz Hastorun, Madan Jampani, Gunavardhan Kakulapati,  
Avinash Lakshman, Alex Pilchin, Swaminathan Sivasubramanian, Peter Vosshall  
and Werner Vogels

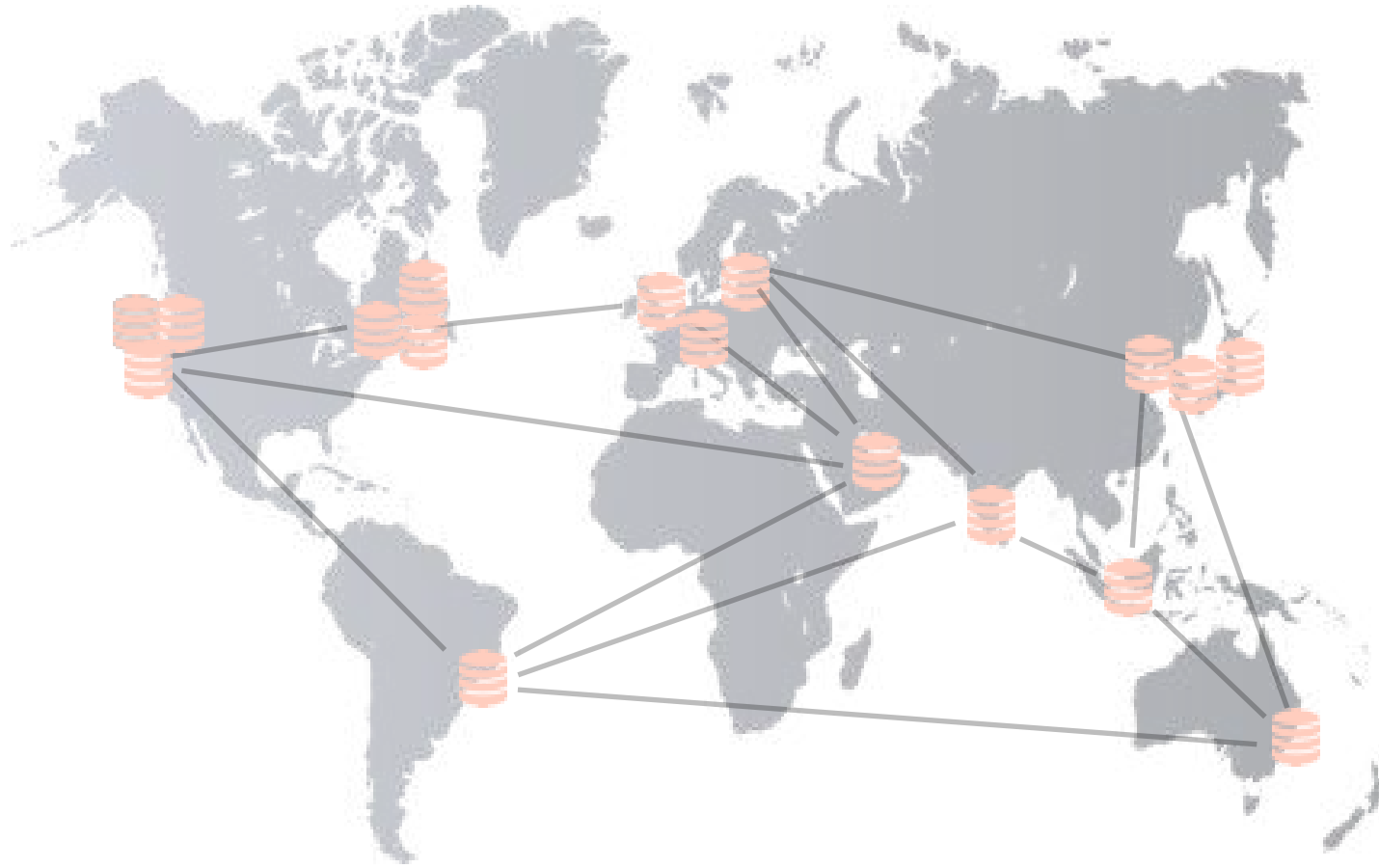
Amazon.com

Launched in 2012, Amazon DynamoDB now powers Alexa, the Amazon.com sites and the Amazon fulfillment centers.

During Prime Day 2017, it handled 3.34 trillion requests, peaking at **12.9 million requests per second**.

# DynamoDB Global Tables

First fully managed, multi-master, multi-region database



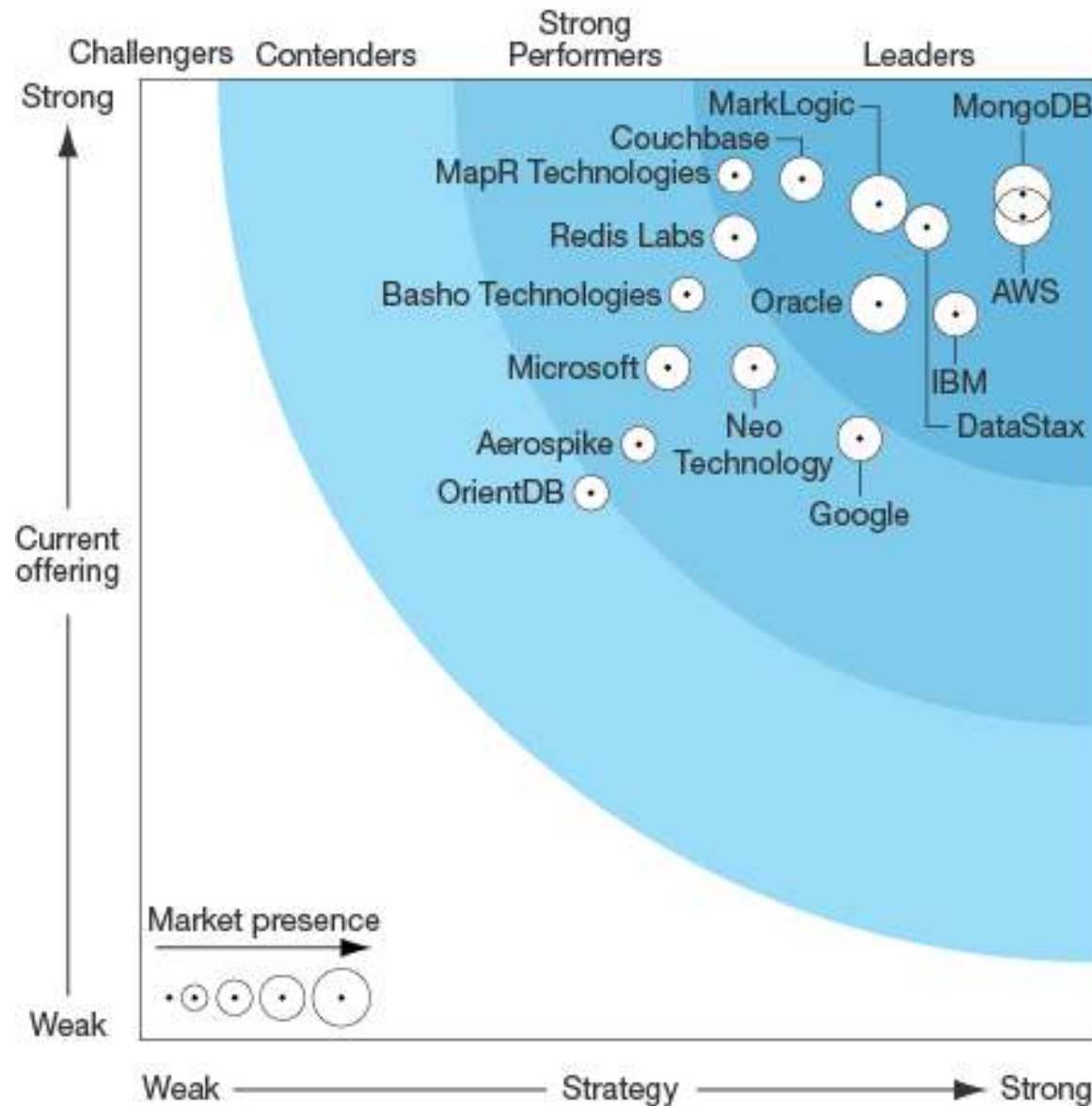
Build high performance,  
globally distributed  
applications

Low latency reads and writes  
to locally available tables

Disaster proof with multi-  
region redundancy

Easy to setup and no  
application re-writes required

# The Forrester Wave: Big Data NoSQL, Q3 2016



# The Forrester Wave: Big Data Warehouse, Q2 2017

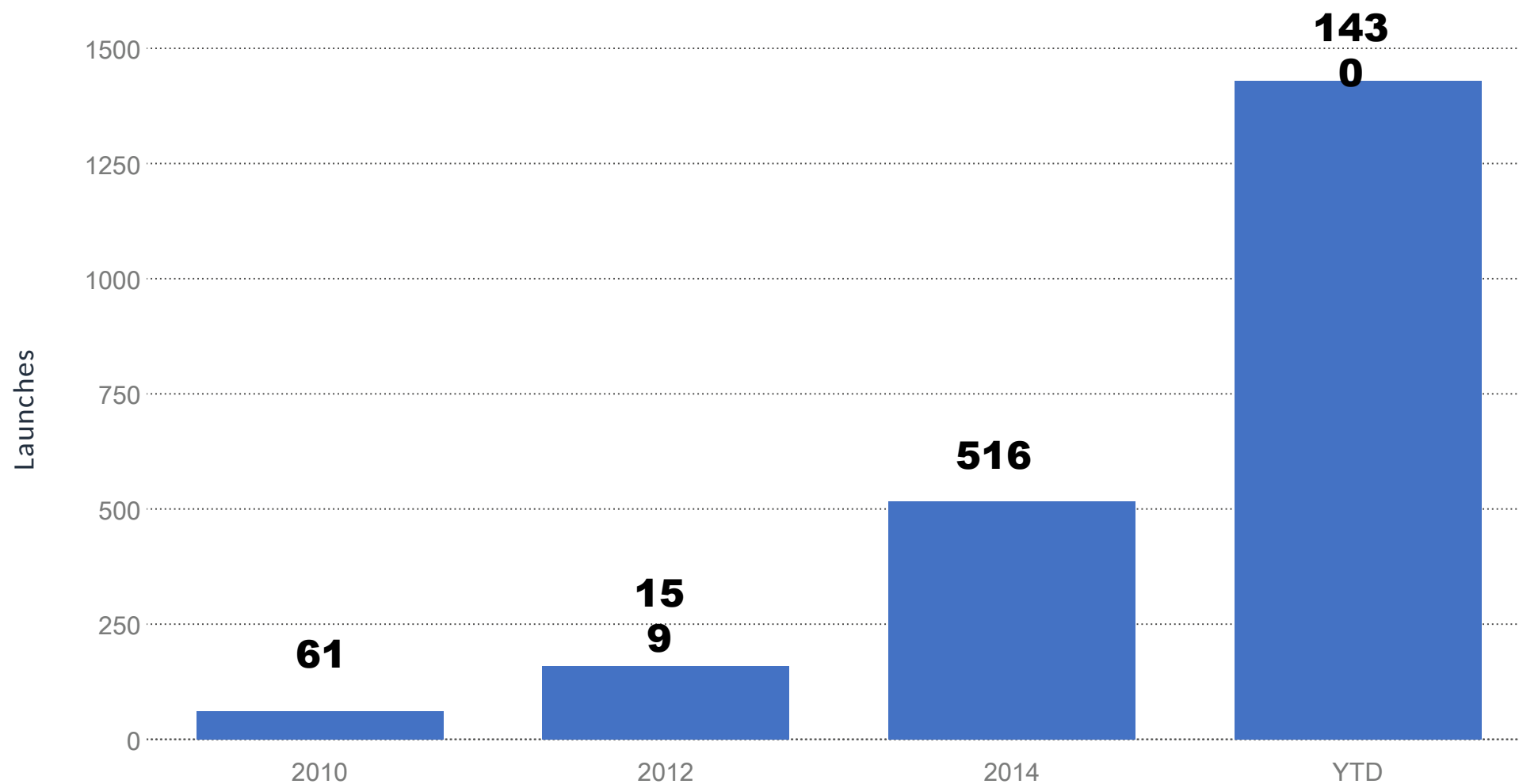
**Figure 3: Forrester Wave™: Big Data Warehouse, Q2 '17**





Machine Learning requires  
constant innovation

# AWS released 1430 features in 2017



# More AI/ML is built on AWS than anywhere else



# More AI/ML is built on AWS than anywhere else



3 times more customer references  
than anyone else

6 times more Enterprise references  
than anyone else



# AWS AI/ML Stack

## **Application Services**

API-driven services: Vision & Language Services, Conversational Chatbots

## **Platform Services**

Deploy machine learning models with high-performance machine learning algorithms, broad framework support, and one-click training, tuning, and inference.

## **Frameworks & Infrastructure**

Develop sophisticated models with any framework, create managed, auto-scaling clusters of GPUs for large scale training, or run inference on trained models.

# Application Services

## Vision Services

### Amazon Rekognition Image

*Deep learning-based image analysis*

[Learn more »](#)

### Amazon Rekognition Video

*Deep learning-based video analysis*

[Learn more »](#)



## Conversational chatbots

### Amazon Lex

*Build chatbots to engage customers*

[Learn more »](#)

## Language Services

### Amazon Comprehend

*Discover insights and relationships in text*

[Learn more »](#)



### Amazon Translate

*Fluent translation of text*

[Learn more »](#)



### Amazon Transcribe

*Automatic speech recognition*

[Learn more »](#)

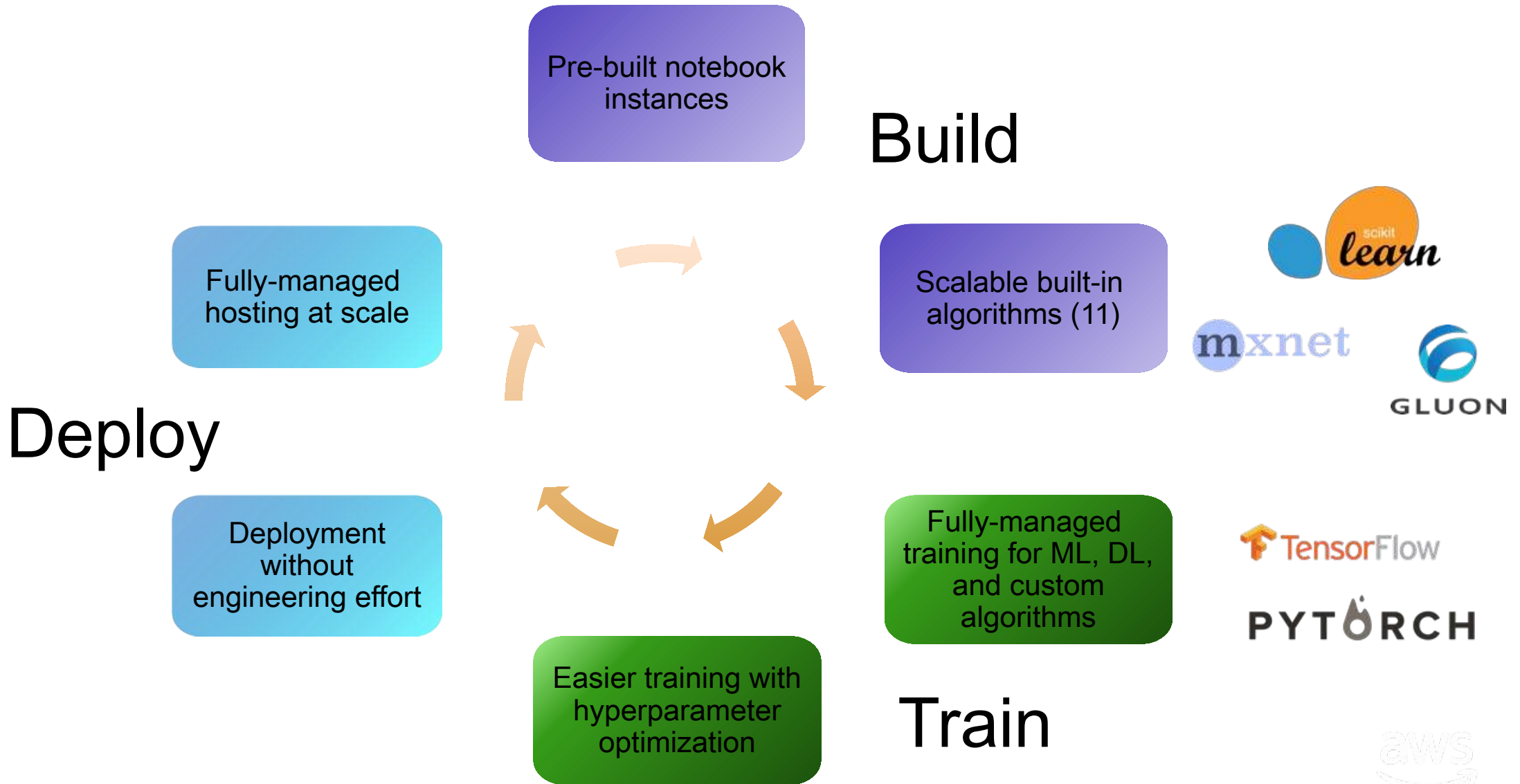


### Amazon Polly

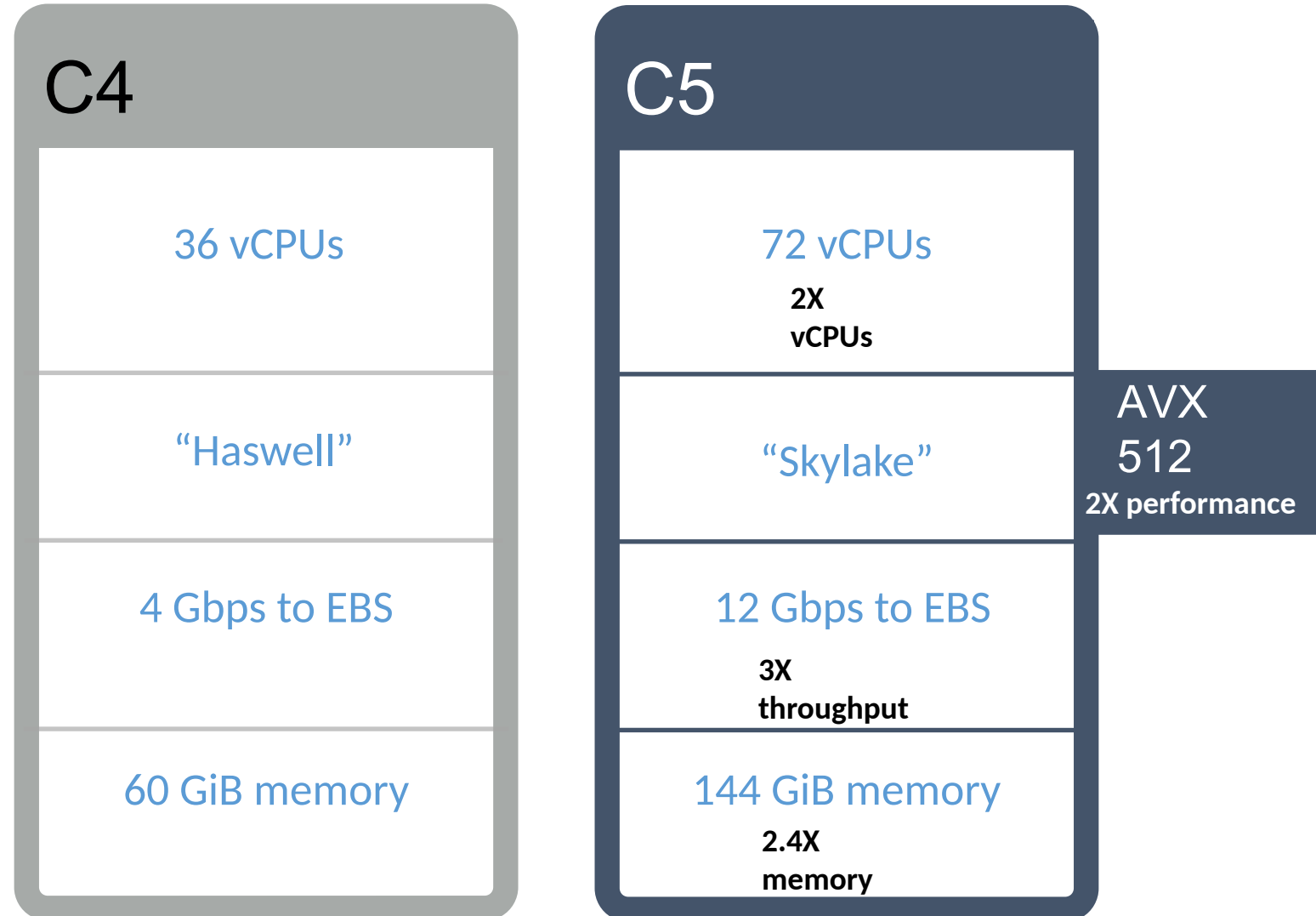
*Natural sounding text to speech*

[Learn more »](#)

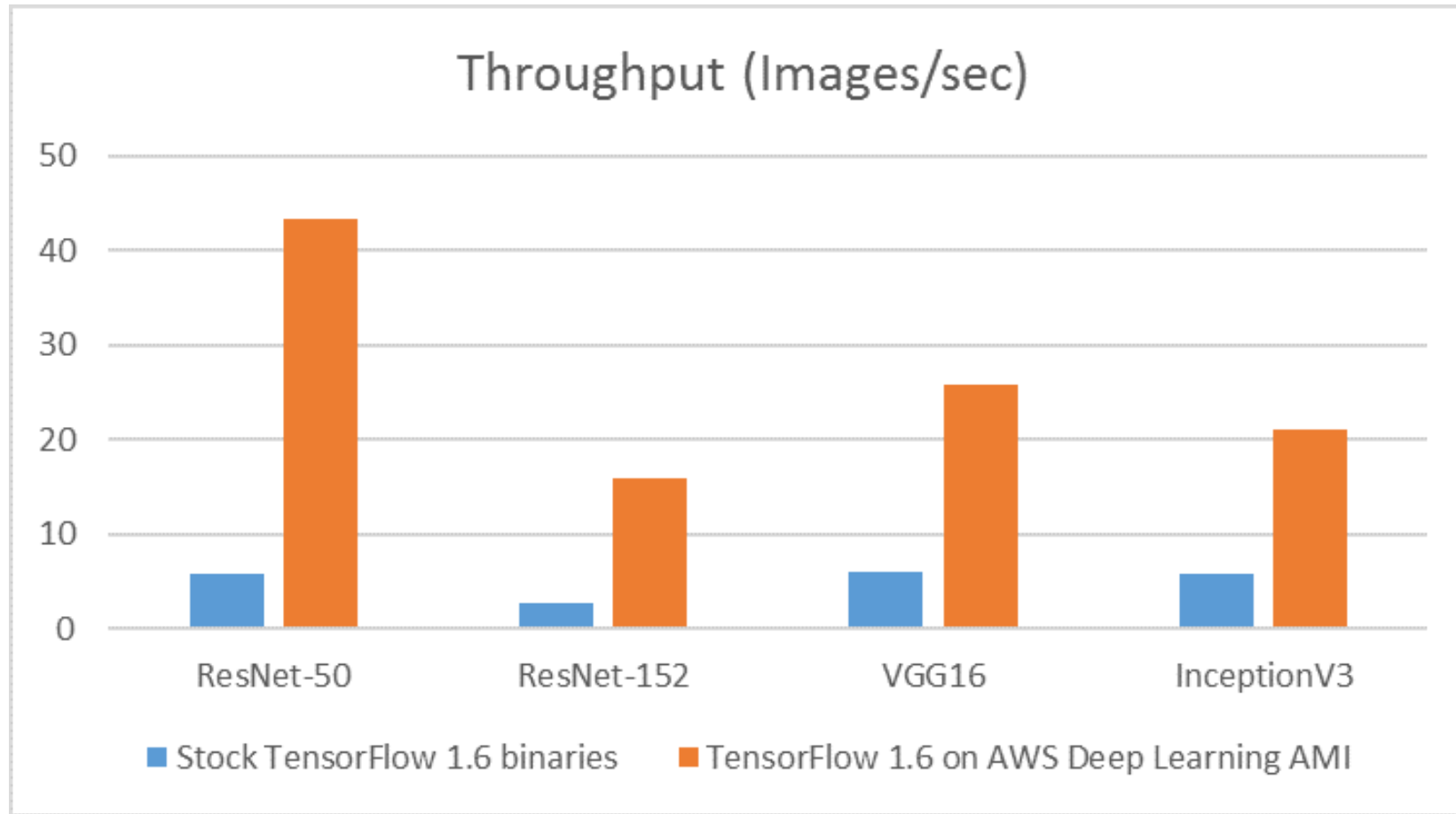
# Amazon SageMaker



# Infrastructure – Amazon EC2 C5



# Train TensorFlow 7x faster on AWS

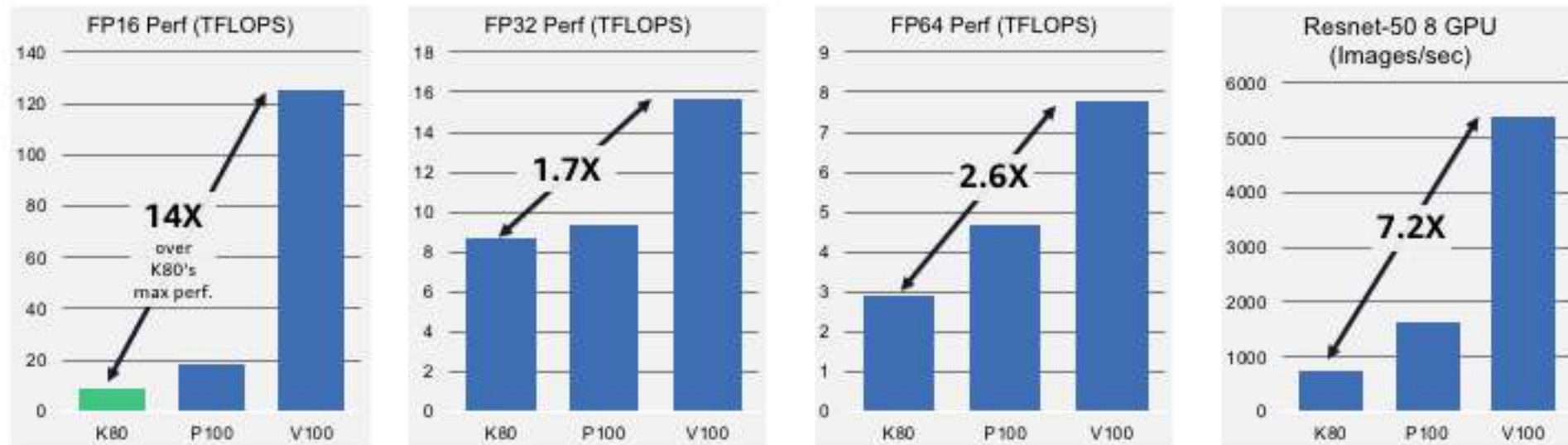


<https://aws.amazon.com/blogs/machine-learning/faster-training-with-optimized-tensorflow-1-6-on-amazon-ec2-c5-and-p3-instances/>



# Infrastructure - GPU instances

- NVIDIA GPU Architectures: Kepler > Maxwell > Pascal > Volta
- P2 Instances use K80 Accelerator (Kepler Architecture)
- P3 Instances use V100 Accelerator (Volta Architecture)



P3 family launched 25/10/2017: up to 8 GPUs in the same server (46,000+ cores)  
Now available in 7 regions (US, EU, APAC)

# NVIDIA V100 availability

		K80	P100	V100
GCE		Available	Available	-
Azure		Available	Available	Available 6/3, 1 US region, 4 GPUs
Oracle		N/A	Available	Available 27/3, 1 US region
IBM		Available	Available	Available 31/1, bare metal, 2 GPUs
NVIDIA		N/A	Available	Available

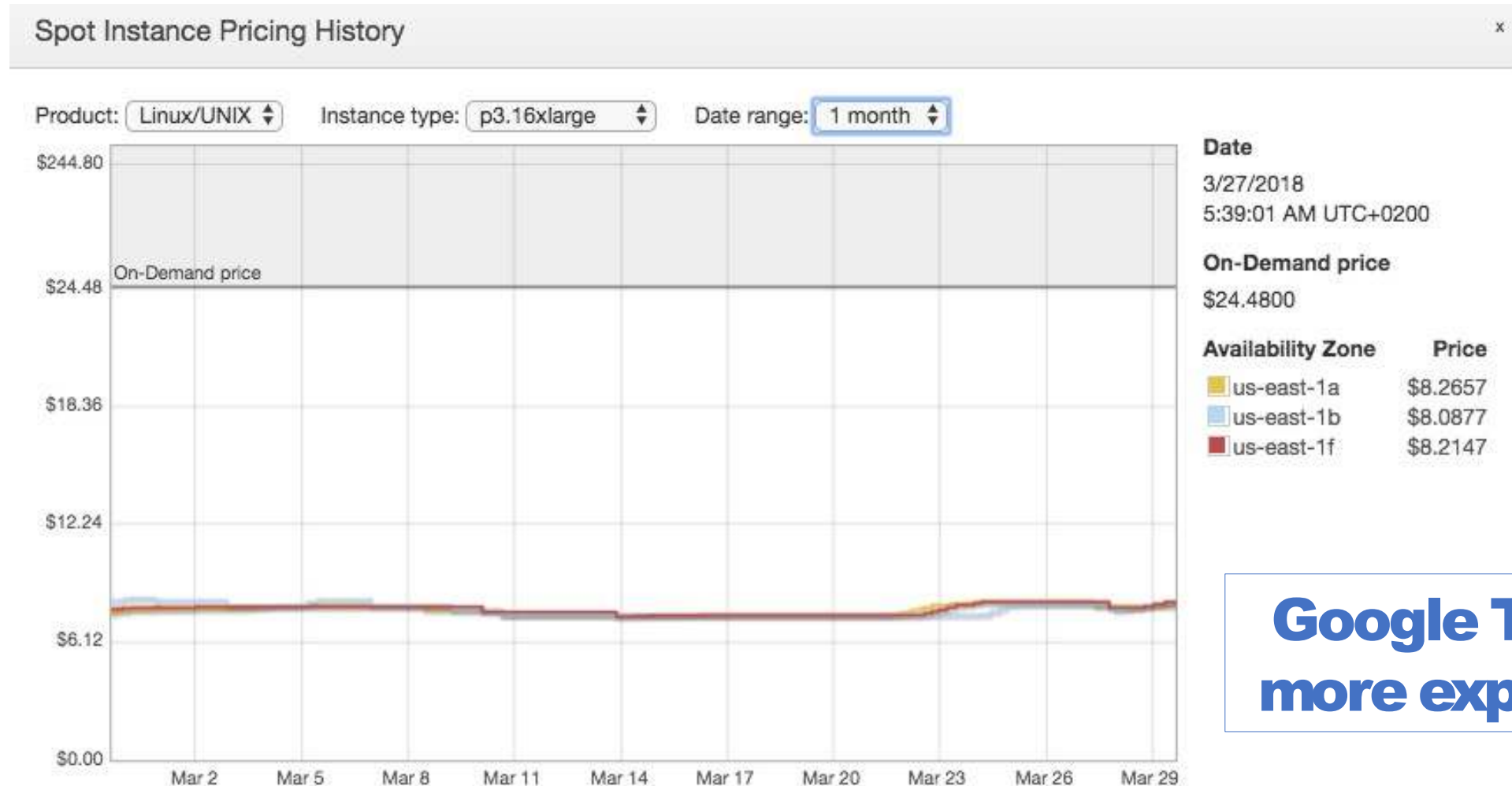
## « Google Announces Expensive Cloud TPU Availability »

	Cloud TPU (4x TPU2)	AWS P3.2xlarge (1xV100)	AWS P3.8xlarge (4xV100)	AWS P3.16xlarge (8xV100)
RN50 Image/Second	1369	819	3242	6309
RN50 Time to Train (hr)	23.4	39.1	9.9	5.1
\$/Hr	\$6.88	\$3.06	\$12.24	\$24.48
\$ to Train ImageNet, 90 epochs	\$161	\$120	\$121	\$124
Availability	Beta	Now	Now	Now

« Net it out, and the *Google part costs ~33% more to do the same work.* »

<https://www.forbes.com/sites/moorinsights/2018/02/13/google-announces-expensive-cloud-tpu-availability/>

# EC2 Spot Instances: 67% discount on p3



**Google TPU 4x  
more expensive**

# AWS DeepLens

## The world's first deep learning enabled video camera for developers

AWS DeepLens helps put deep learning in the hands of developers, literally, with a fully programmable video camera, tutorials, code, and pre-trained models designed to expand deep learning skills.

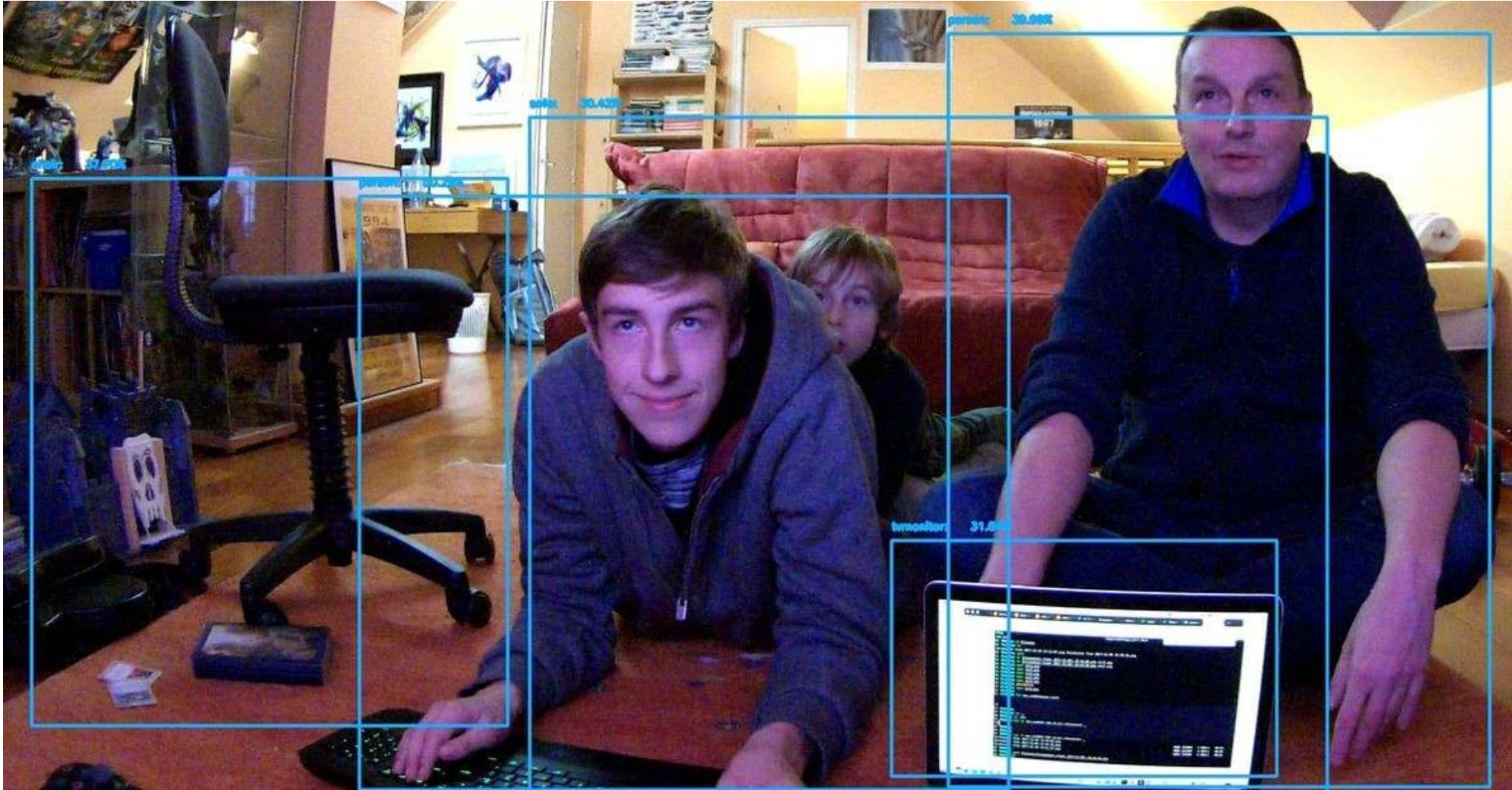
[Pre-order](#)

[Get started with your DeepLens](#)





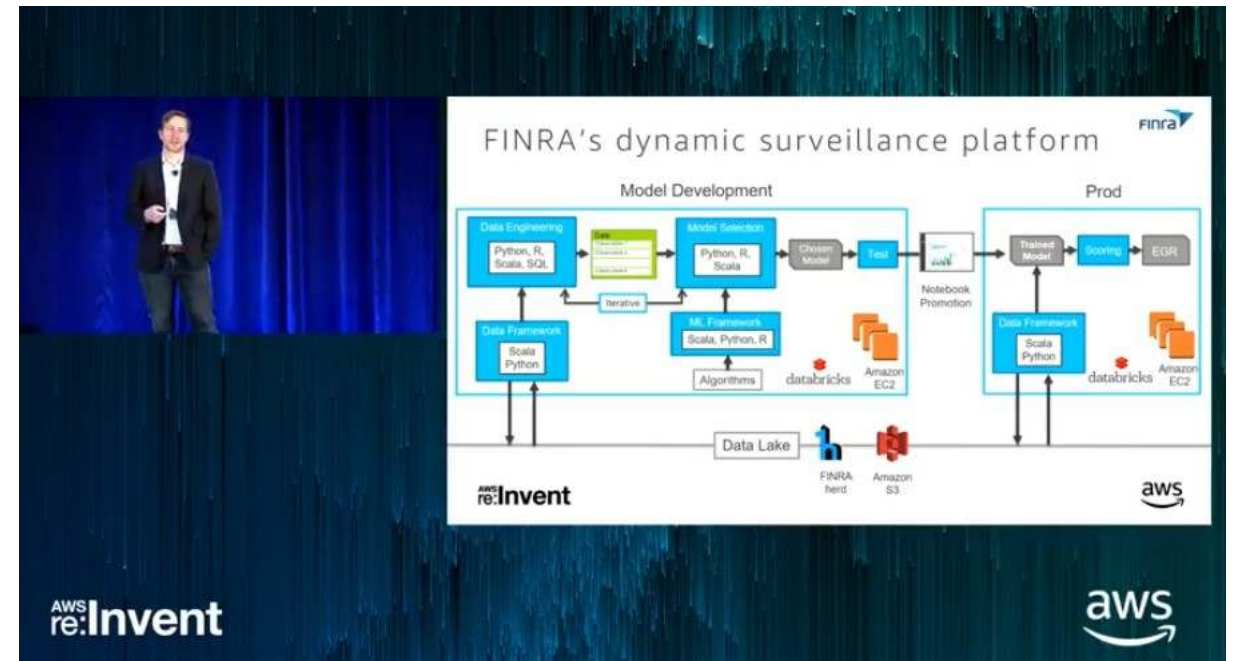
# Object detection with AWS DeepLens



Machine Learning at scale

# FINRA

- FINRA is the primary regulatory agency for stock brokers in the US.
- They handle about **75 billion market events every day** and run hundreds of surveillance algorithms.
- A typical complex query that was taking an hour now takes five to ten seconds.
- Cost savings amount up to **\$20 million** annually.



<https://www.youtube.com/watch?v=gdKmMO5PVOY>



# Autodesk



- Autodesk demonstrates the power of **generative design** to explore all the possible permutations of a solution quickly.
- Autodesk relies on high performance computing on AWS to scale its solutions.



<https://www.youtube.com/watch?v=A31A8KDC9S4>

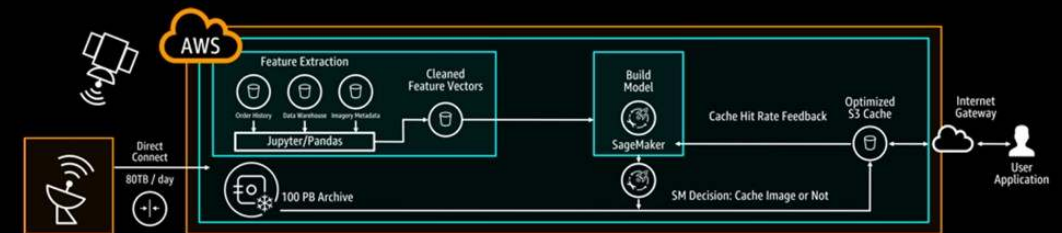
# Digital Globe



- In the last 18 years DigitalGlobe has been operating Earth imaging satellites, they have collected over 100 PB of imagery.
- They make extensive use of Machine Learning on [Amazon SageMaker](#) to extract information.
- Working with the [AWS ML Lab](#), Digital Globe also built a predictive model that will [reduce cloud storage costs by 50%](#).



## USING AMAZON SAGEMAKER TO CUT CLOUD STORAGE COSTS IN HALF



<https://aws.amazon.com/solutions/case-studies/digitalglobe-machine-learning/>

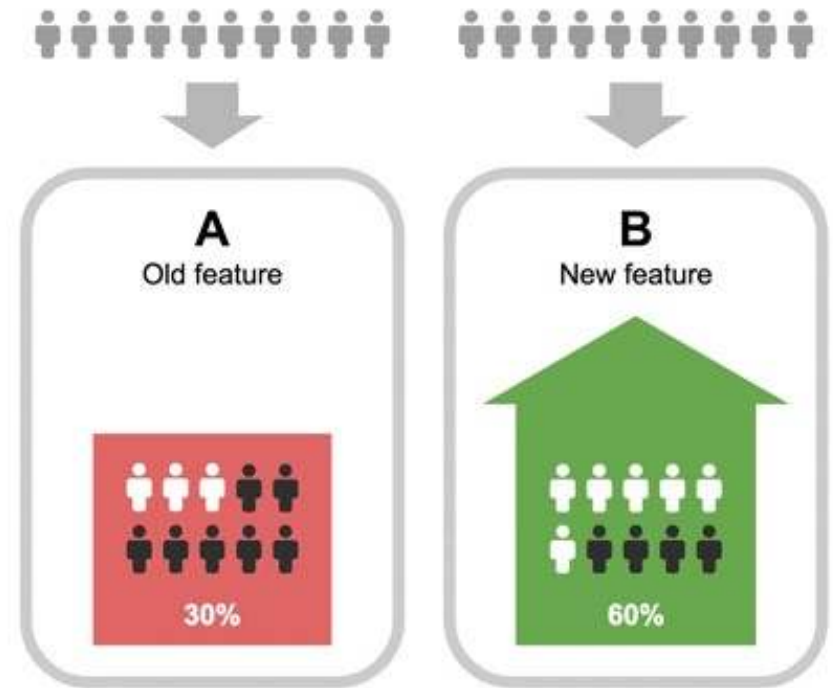
Speech



# Duolingo



- Duolingo is the most popular language-learning platform and the most downloaded education app in the world, with more than 170 million users.
- They have run six A/B tests, testing an [Amazon Polly](#) voice against a voice from other TTS providers.
- For all of these experiments, the winning condition was the Amazon Polly voice



<https://aws.amazon.com/blogs/machine-learning/powering-language-learning-on-duolingo-with-amazon-polly/>

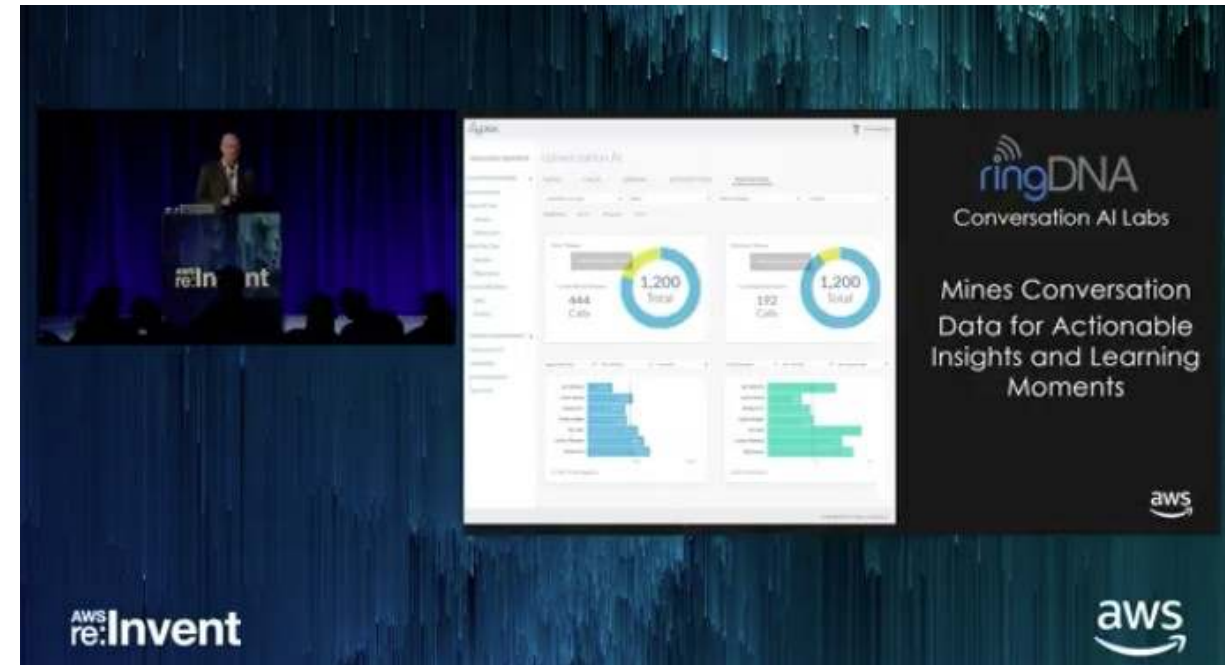
# ringDNA

- RingDNA is an end-to-end communications platform for sales teams.
- Hundreds of enterprise organizations use RingDNA to dramatically increase productivity, engage in smarter sales conversations, gain predictive sales insights and improve their win rate.

## Speech to Text

"A critical component of RingDNA's Conversation AI requires best of breed speech-to-text to deliver transcriptions of every phone call. RingDNA is excited about [Amazon Transcribe](#) since it provides high-quality speech recognition at scale, helping us to better transcribe every call to text "

Howard Brown, CEO & Founder, RingDNA



[https://www.youtube.com/watch?v=1ZJ\\_f1bDdog](https://www.youtube.com/watch?v=1ZJ_f1bDdog)

# Natural Language Processing

# ClearView Social

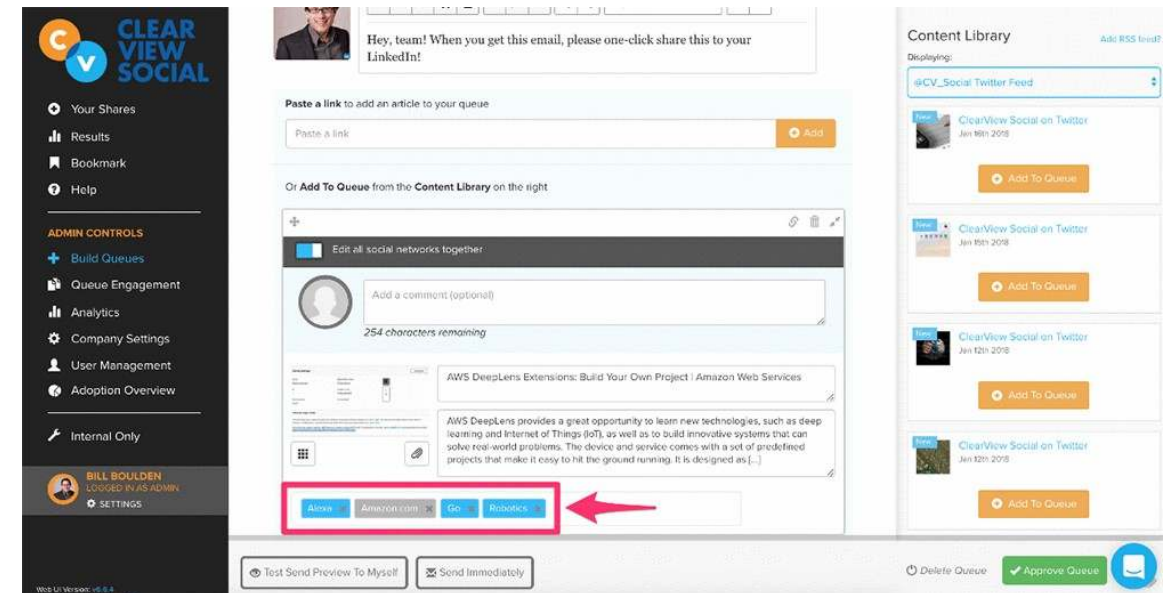


- ClearView Social enables a company's employees to share approved content on LinkedIn, Twitter, and other social networks.
- To eliminate the reliance on manual tagging, ClearView Social turned to **Amazon Comprehend** to find **insights** and **relationships** in text.

## Natural Language Processing

"We use **Amazon Comprehend** to read an article and extract topics, which are automatically tagged using machine learning. This automatic tagging helps customers easily estimate the market value of their engagement according to the current bid prices from the Google AdWords API"

Bill Boulden, CTO, ClearView Social



<https://aws.amazon.com/blogs/machine-learning/clearview-social-uses-amazon-comprehend-to-measure-the-impact-of-social-sharing/>

# Digital Genius

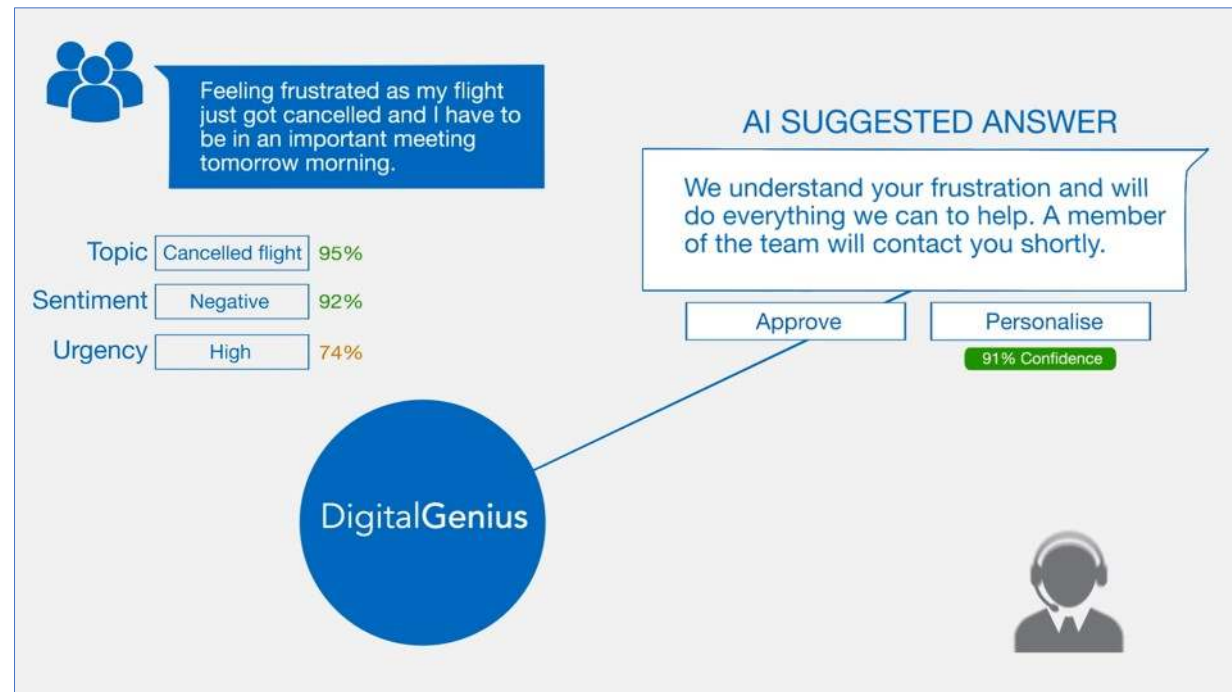
Digital Genius is an AI-powered customer service solution that monitors social media, texts and email for negative customer sentiment, identifies cases that need more attention, and suggests ideal responses to queries.

## Scaling on the cloud

« We're unlocking the intelligence value of historical data while helping customer service agents deliver a faster and more accurate experience for their consumers »

« NVIDIA GPUs make it possible for us to train very large neural nets with millions of parameters in a matter of hours rather than days »

Mikhail Naumov, President, DigitalGenius



<https://blogs.nvidia.com/blog/2017/01/27/faster-customer-service-with-ai/>

# Computer Vision



# GumGum

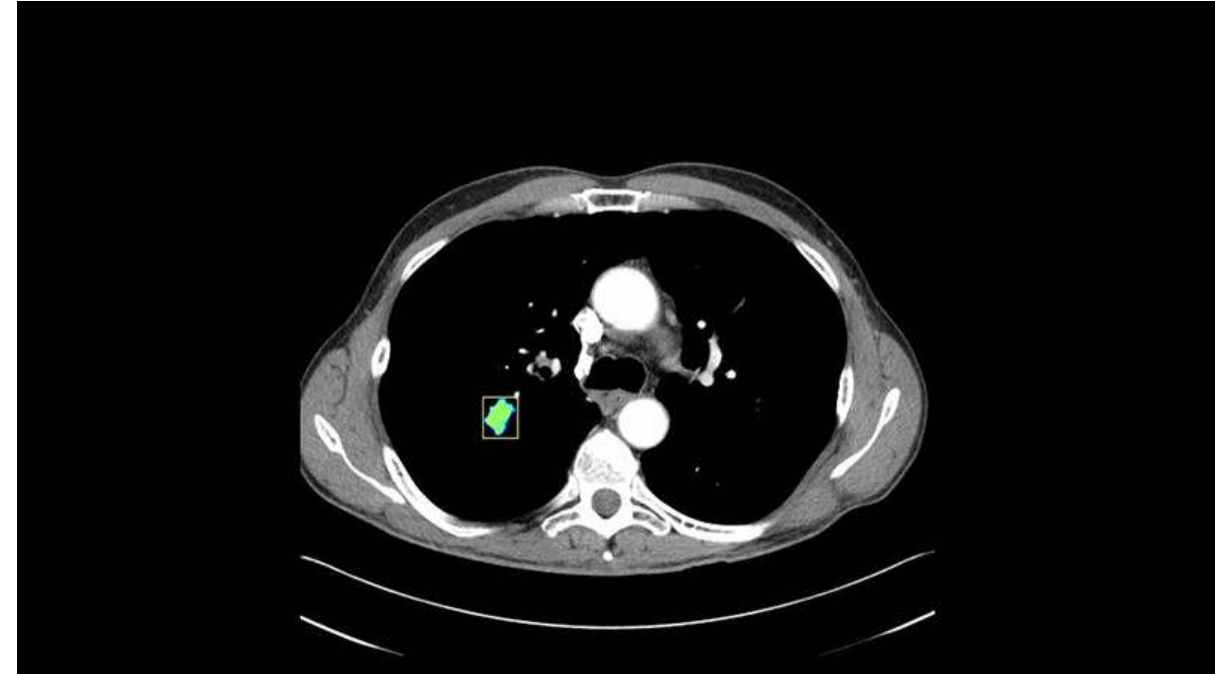
- GumGum unlocks the value of images in advertising & professional sports.
- Manti, GumGum's real-time social media listening tool can analyze the more than **1.8 billion social images** posted daily, where 80% of them lack text to help marketers find them.
- Using **Deep Learning** with **Apache MXNet** and **GPU instances**, the company reduced model training time from 21 hours to 6 hours.



<https://aws.amazon.com/solutions/case-studies/gumgum/>

# Matrix Analytics

- Matrix Analytics is helping to save lives. The Colorado-based startup uses **Deep Learning** on AWS (**Deep Learning AMI, TensorFlow, GPU instances**) to track disease progression for patients diagnosed with pulmonary nodules in their lungs.
- The Matrix Analytics tools are able to outperform previous methods in their ability to diagnose cancer from a CT image.
- The software automates follow-up care to ensure that each patient follows through with recommendations in order to monitor changes in their condition.



## Deep Learning

« Using the convenience of the [Deep Learning] AMI on AWS gives us the opportunity to offer up different business models, which allows us to become excellent technology partners as the market evolves at an ever-increasing pace.

# PowerScout

- PowerScout is changing the sales model for solar-powered homes by using **Deep Learning** to identify households likely to embrace solar panels.
- Moving away from the door-to-door sales approach, the company's deep learning models analyze satellite imagery to evaluate solar-worthy factors for each home.
- PowerScout uses **GPU instances**, the CUDA parallel processing platform and the **cuDNN** deep neural network library.



<https://blogs.nvidia.com/blog/2016/12/27/ai-solar-powered-homes/>



# TuSimple

- TuSimple, a leader in self-driving technology, uses **Deep Learning** to build sophisticated algorithms for computer vision and driving simulation.
- They rely on **Apache MXNet** to teach computers how to recognize and track objects and to make decisions to avoid collisions and prioritize safety.
- They simulated a billion miles of road driving with a wide range of variables and driving conditions—the largest simulation of its kind in history.



<https://www.oreilly.com/ideas/self-driving-trucks-enter-the-fast-lane-using-deep-learning>

More AI/ML is built on AWS than anywhere else

<https://aws.amazon.com/machine-learning/>



# Thank you

**Julien Simon, AI Evangelist, EMEA**  
**@julsimon**

