

Cloud Operated Open Source AI Robots

Adrian Cockcroft @adrianco



- Applications of Deep Learning
- Apache MXNet Overview
- Apache MXNet API
- Code and DIYrobocars
- Tools and Resources

Applications of Deep Learning

The Challenge For AI: Scale



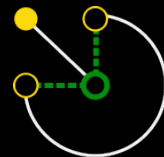
Data



Training



Inference



At the Edge

The Challenge For AI: Scale



Data

PBs of existing data
New data created
on AWS



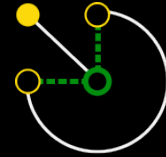
Training

Tons of GPUs
Elastic capacity
Pre-built images



Inference

Tons of GPUs and
CPUs
Serverless



At the Edge

IoT and mobile deployment
Mobile optimization
IoT device optimization

AI On AWS Today

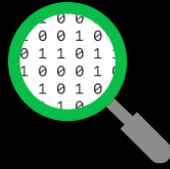


Thousands Of Amazon Engineers Focused On Machine Learning



|

Fulfilment &
logistics



|

Search &
discovery



|

Existing
products



|

New
products



|

At AWS

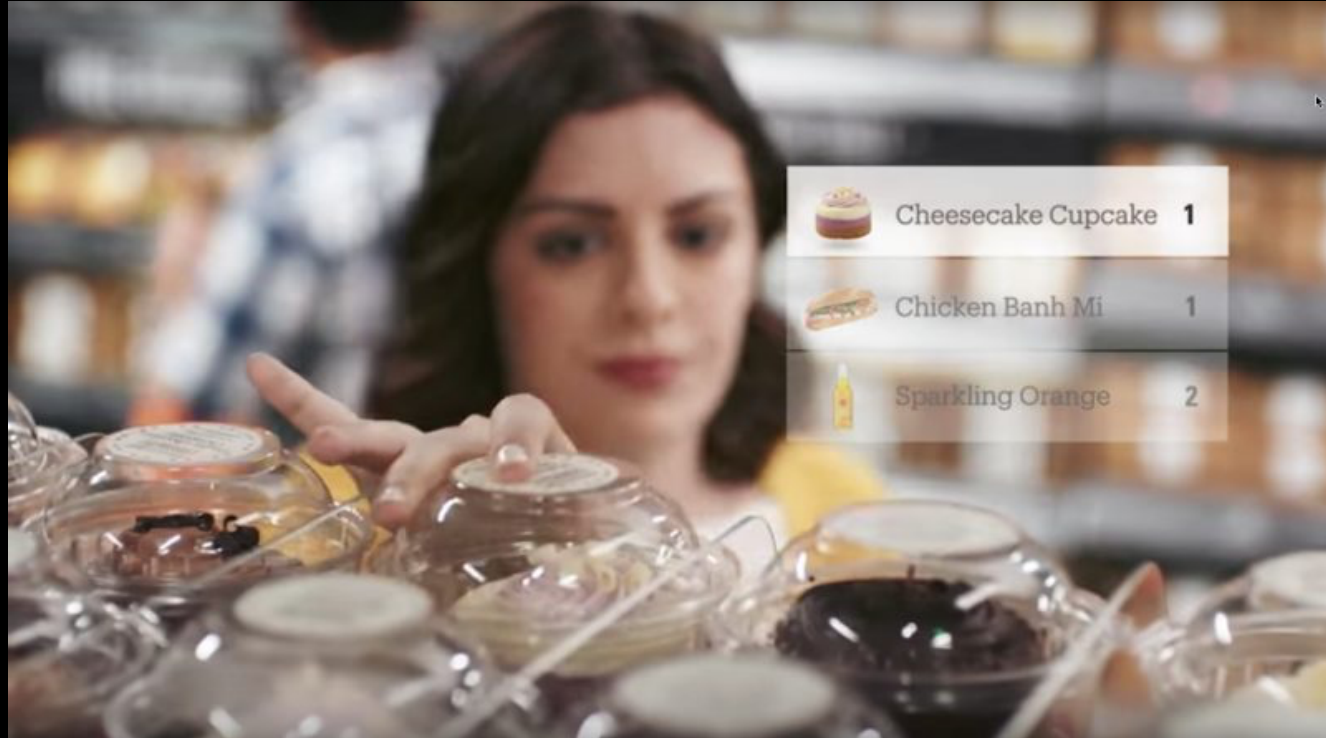
Deep Learning using MXNet @Amazon

- Applied Research
- Core Research
- Alexa
- Demand Forecasting
- Risk Analytics
- Search
- Recommendations
- AI Services | Rek, Lex, Polly
- Q&A Systems
- Supply Chain Optimization
- Advertising
- Machine Translation
- Video Content Analysis
- Robotics
- Lots of Computer Vision..
- Lots of NLP/U..

*Teams are either actively evaluating, in development, or transitioning to scale production

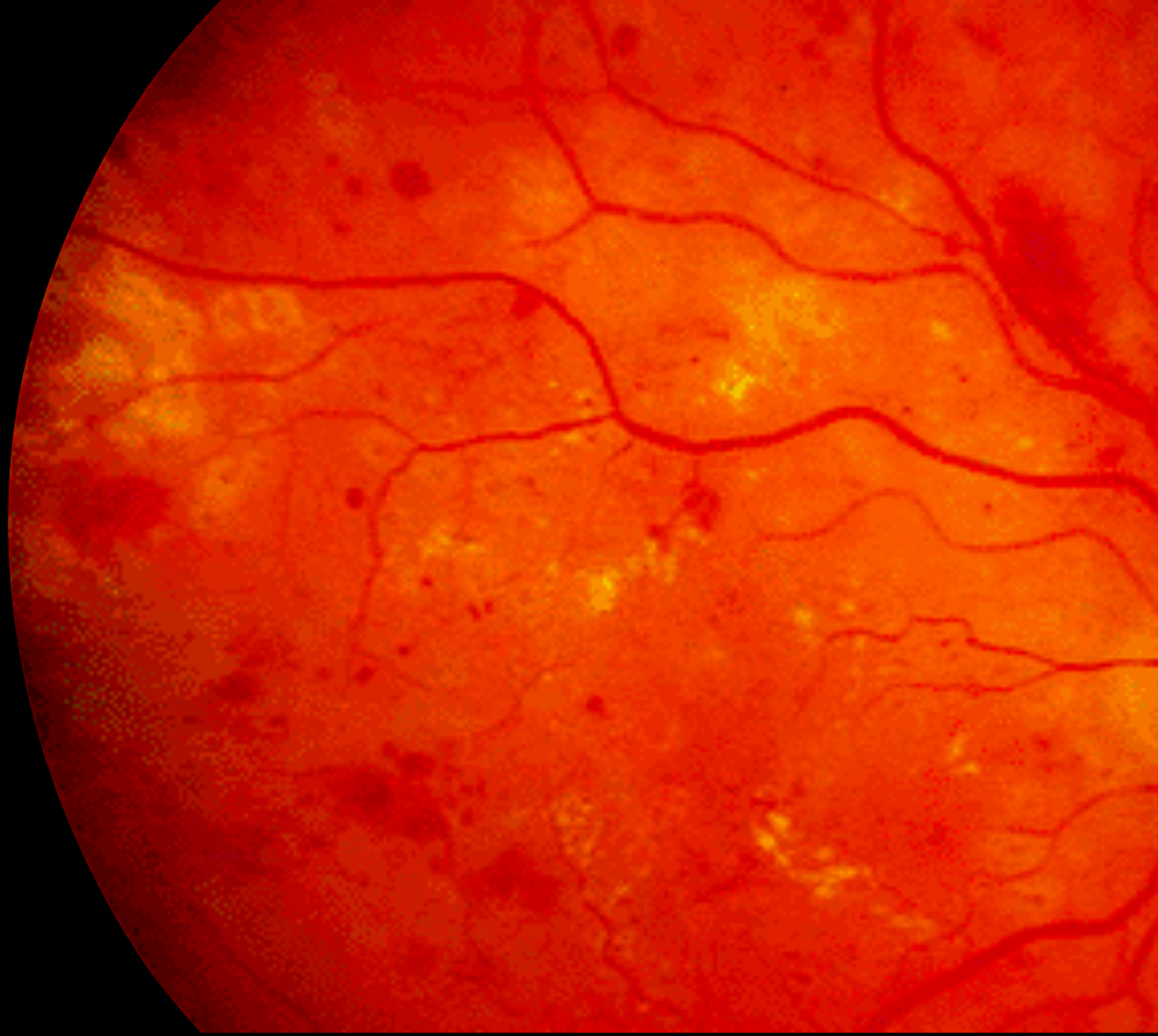


amazon go



Stanford

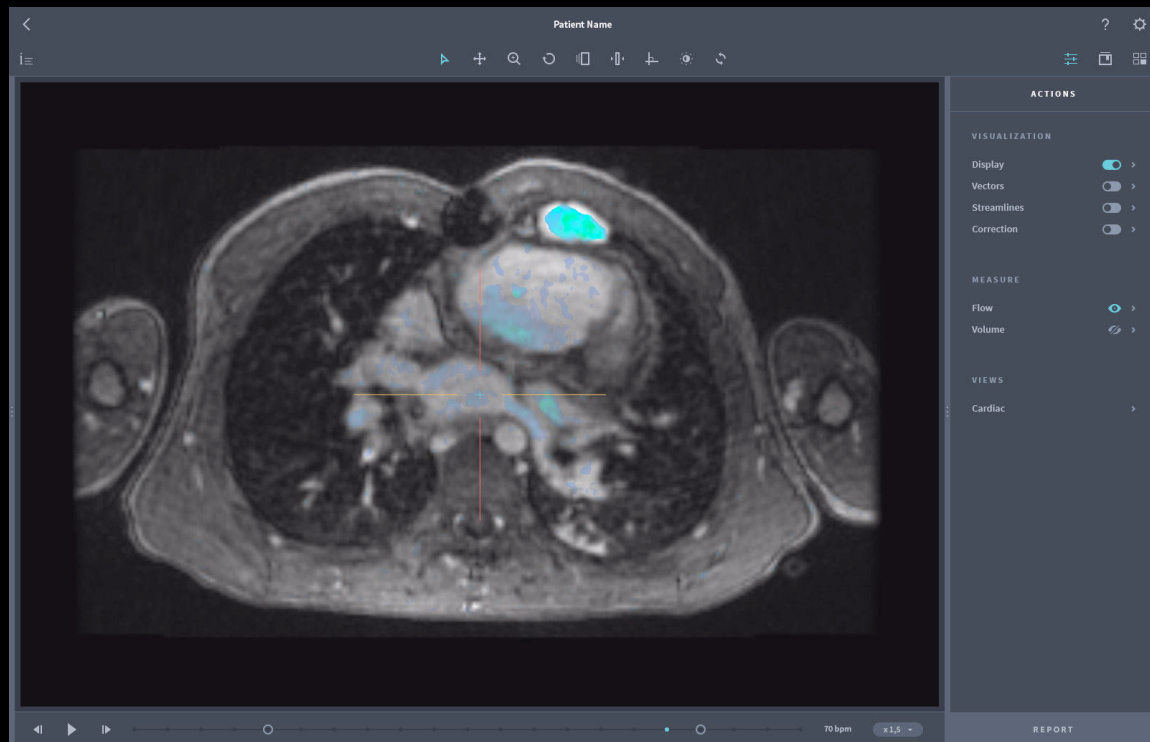
Early detection of
diabetic complications

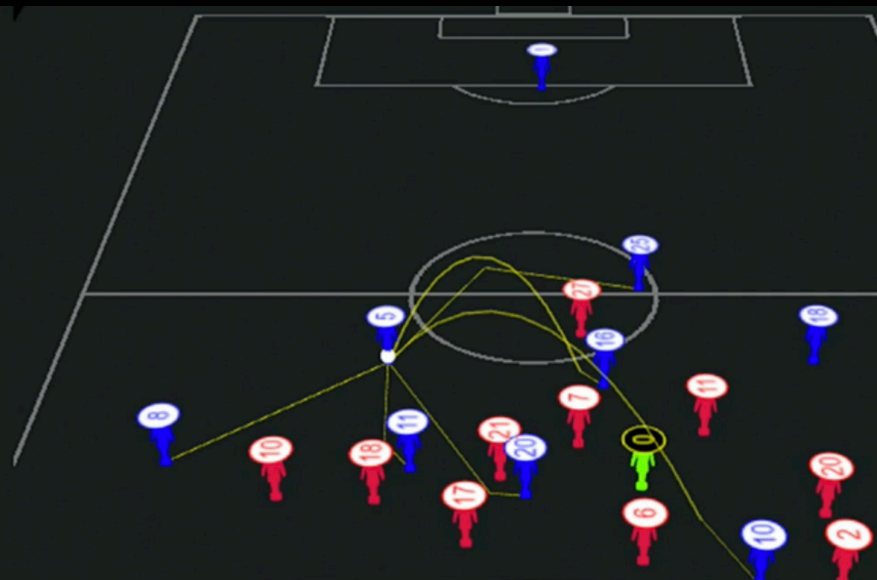
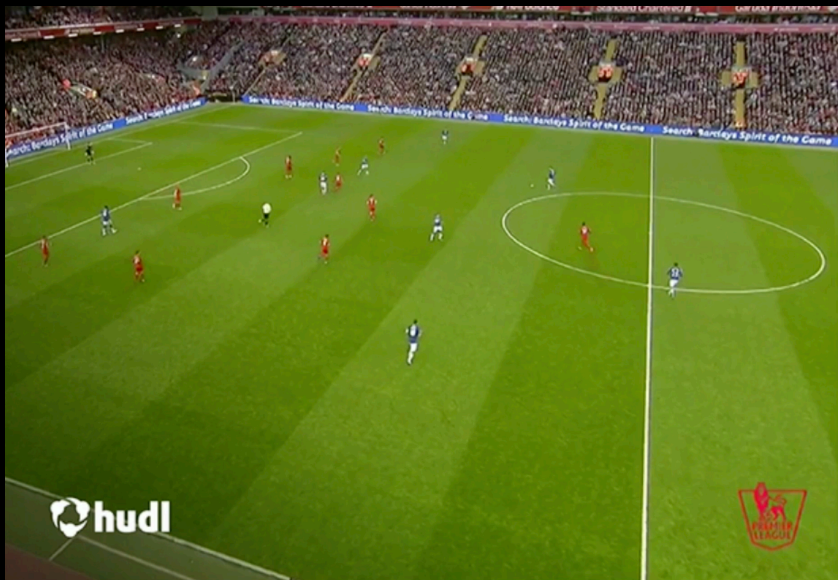





ARTERYS



FDA-approved
medical imaging














 **WolframAlpha**[®] computational knowledge engine.

who recorded pet sounds?

 Web Apps  Examples  Random

Assuming "pet sounds" is a music album | Use as [a music work](#) instead

Input interpretation:

Pet Sounds (music album) artist

Result:

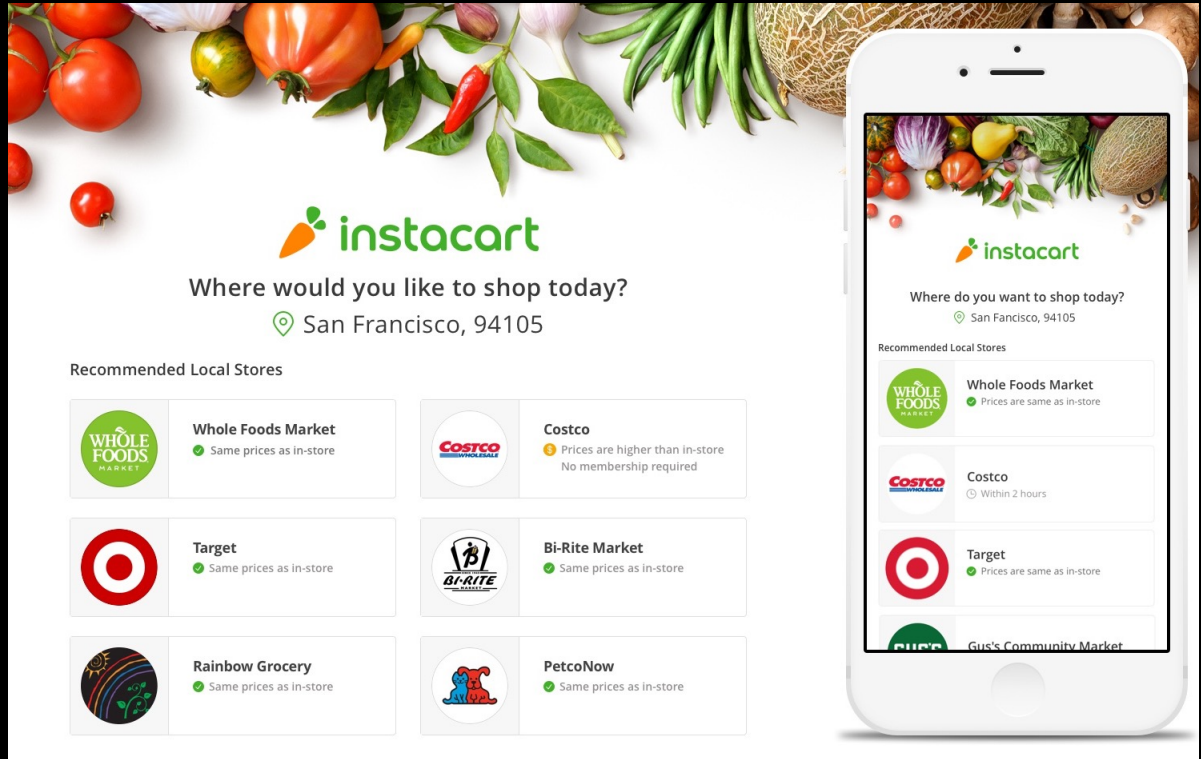
The Beach Boys

Basic information:

name	Pet Sounds
artist	The Beach Boys
release date	May 16, 1966
runtime	36 minutes 15.12 seconds

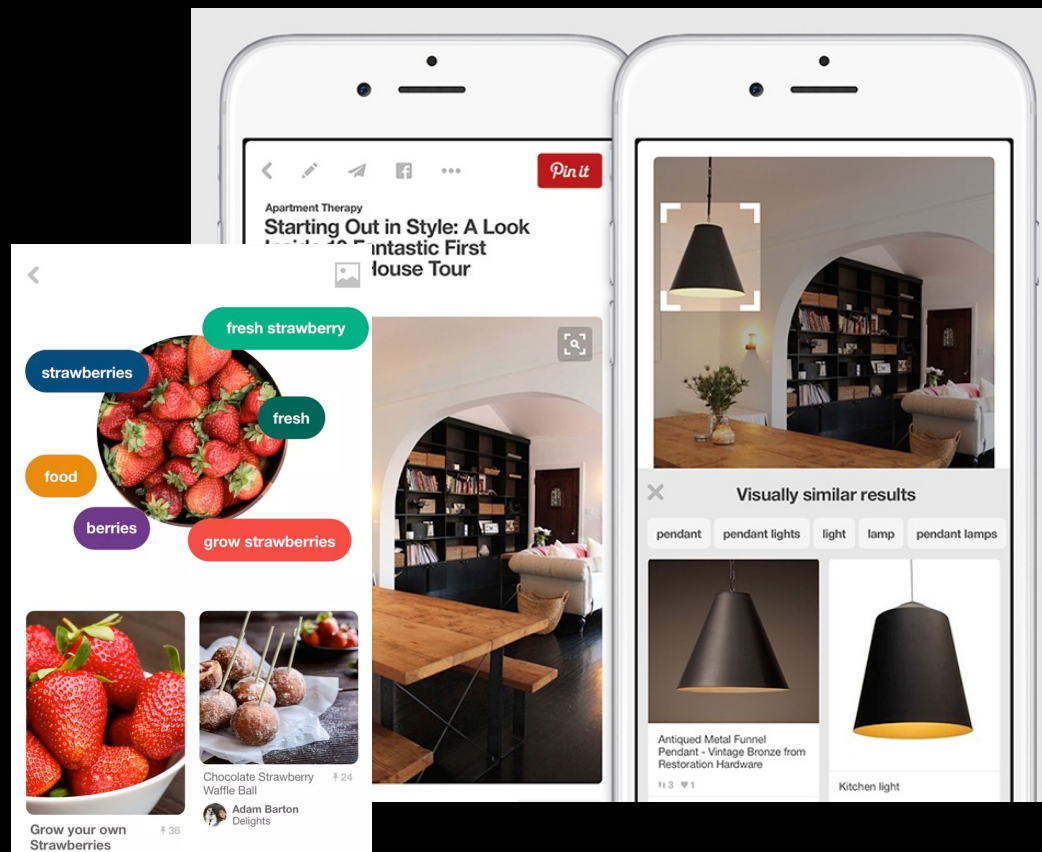


Online grocery
delivery services



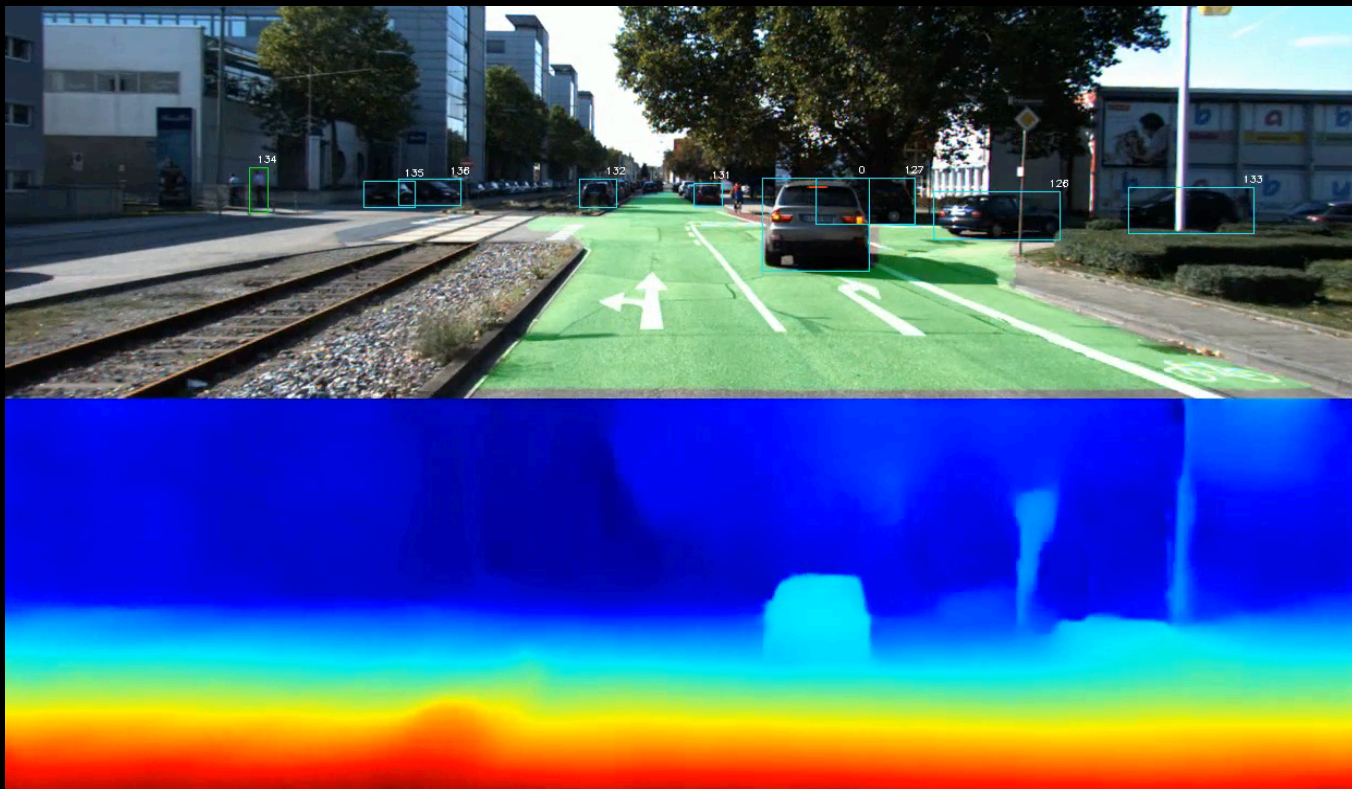


New Pinterest Lens mobile app uses TensorFlow to let users find pins based on the pictures they take









Amazon AI: Democratized Artificial Intelligence

Amazon
Rekognition

Amazon
Polly

Amazon
Lex

More to come
in 2017

Easy to Use AI
Services

Amazon
Machine Learning

Amazon Elastic
MapReduce

Spark &
SparkML

More to come
in 2017

AI Platform

**Apache
MXNet**

TensorFlow

Caffe

Torch

Theano

CNTK

Keras

**Open Source
AI Engines**

P2 16xGPU
Instance

G2 & Elastic
GPU

ECS

Lambda

AWS
Greengrass
IoT

FPGA

More to
come
in 2017

Cloud
Hardware on
Demand

Amazon Rekognition: Object And Scene Detection

Boat	99.3%
Plant	95.1%
Harbor	94.8%
Waterfront	94.8%
Yacht	78.1%
Dock	75.7%
City	72.4%
Architecture	71.8%
Urban	63.9%
Building	62.3%
Marina	60.3%
Plaza	51.1%
Spire	50.8%
Neighborhood	50.7%
Flower	50.6%



Amazon Rekognition: Facial Analysis



A man in a dark pinstripe suit, white shirt, and patterned tie is shown from the chest up. A black rectangular bounding box is drawn around his face. Yellow dots mark various facial landmarks on his face, including the eyes, nose, mouth, and eyebrows.

Age Range	38-59
Beard: False	84.3%
Emotion: Happy	86.5%
Eyeglasses: False	99.6%
Eyes Open: True	99.9%
Gender: Male	99.9%
Mouth Open: False	86.2%
Mustache: False	98.4%
Smile: True	95.9%
Sunglasses: False	99.8%

Landmarks
EyeLeft
EyeRight
Nose
MouthLeft
MouthRight
LeftPupil
RightPupil
LeftEyeBrowLeft
LeftEyeBrowRight
LeftEyeBrowUp
:

Amazon Rekognition: Facial Verification



Similarity: 98%

Amazon Rekognition: Facial Recognition



With Amazon Polly

"The temperature
in Chicago, IL is
75°F"



"The temperature in
Chicago, Illinois is 75
degrees Fahrenheit"

Amazon Polly

47 voices

24 languages

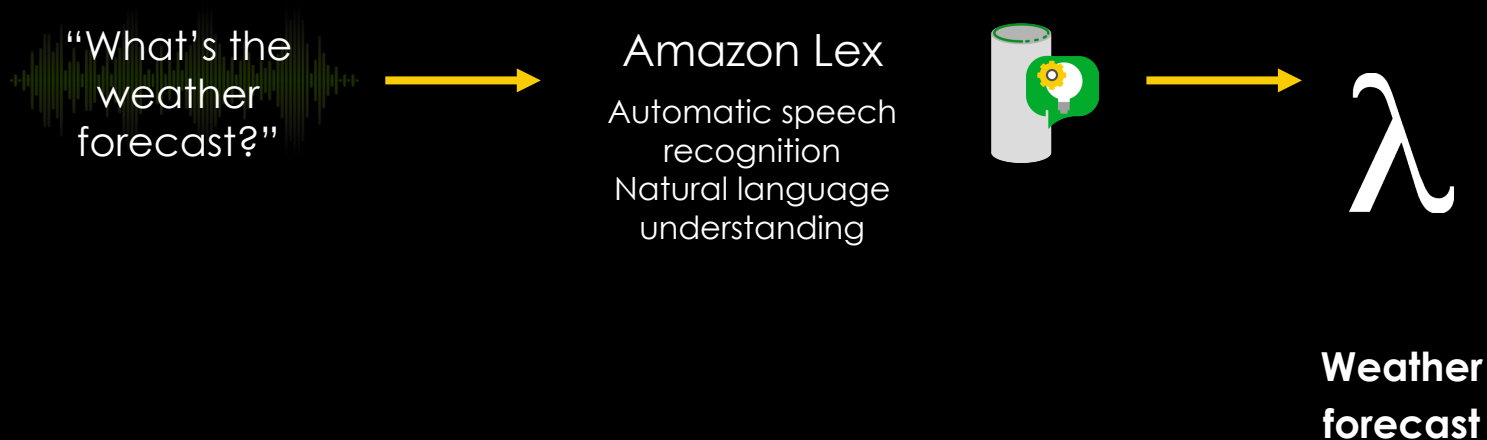
Including: Danish, Swedish, Norwegian

Amazon Polly Customers



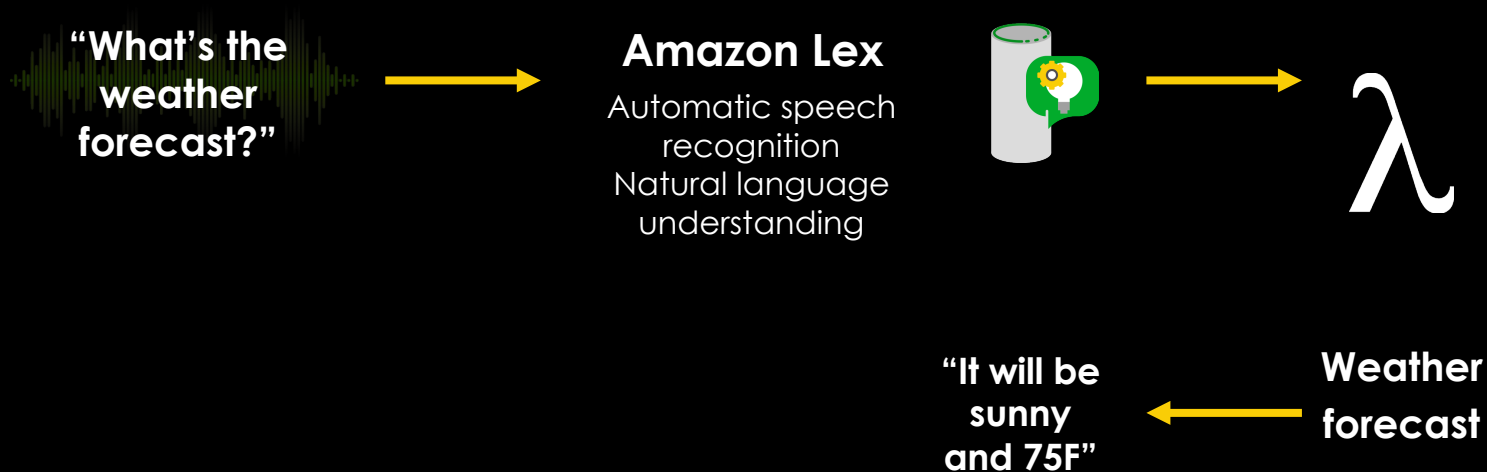
Amazon Lex

Speech recognition and natural language understanding



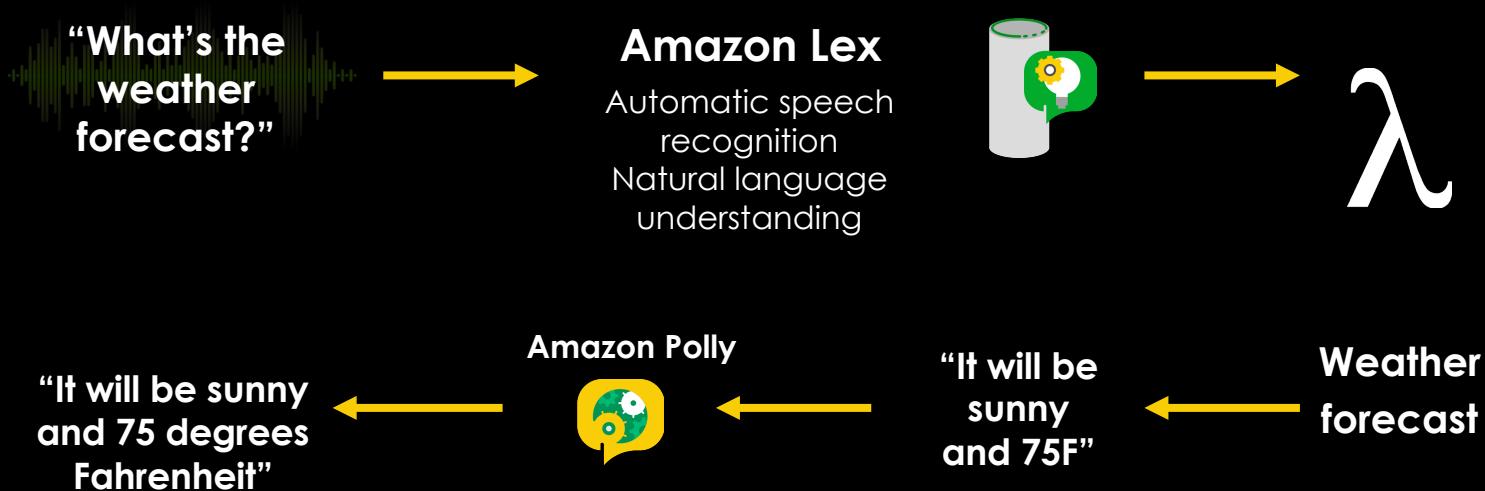
Amazon Lex

Speech recognition and natural language understanding



Amazon Lex

Speech recognition and natural language understanding



Amazon Lex Customers



Amazon Connect

Simple to use, cloud-based contact center



Easy to set
up and
manage



Scalable and
elastic



Pay as
you go

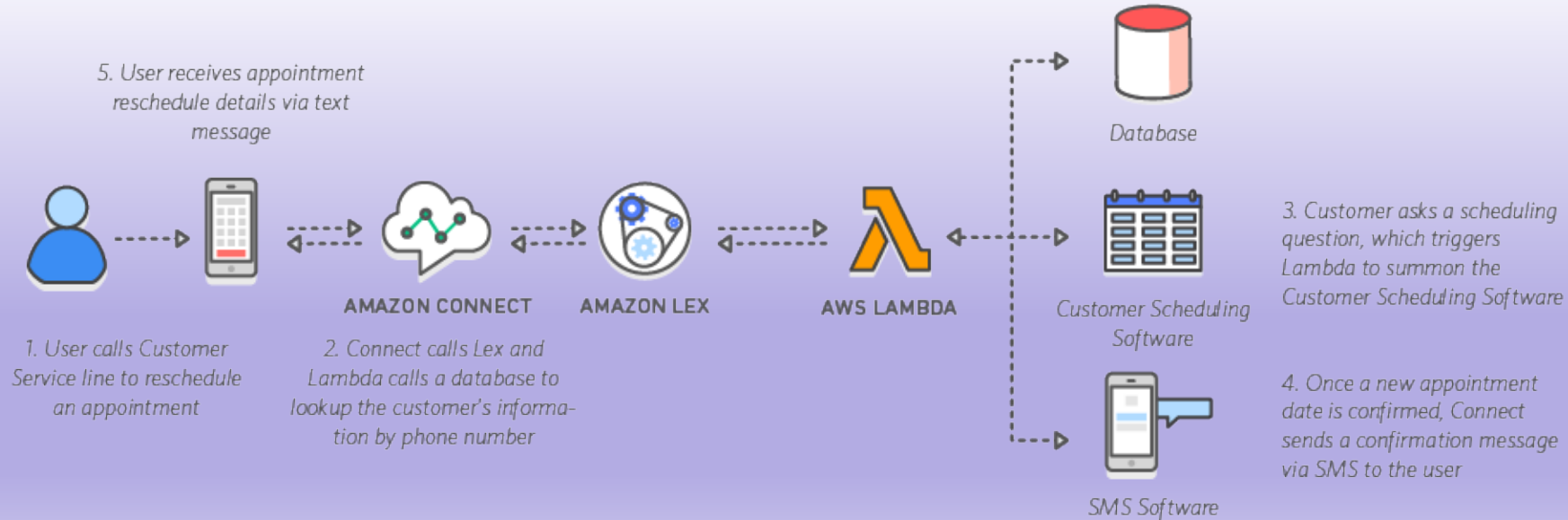


Reliable



Open
platform

Amazon Connect Contact Center Uses Amazon Lex For Natural Conversations



Securing Sensitive Data Is **Job Zero**

Identifying And Protecting Sensitive Data Can Be Challenging



I
Manual
process



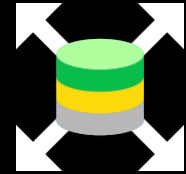
I
Time-
consuming



I
Inaccurate



I
Expensive



I
Growing
volumes of
data



AVAILABLE TODAY

Amazon Macie

Automatically discover, classify, and protect sensitive data in AWS using Machine Learning

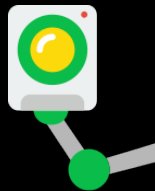
Introducing Amazon Macie



Automatically
discover and
classify your data



Understand where
sensitive data is
located and how it
is accessed

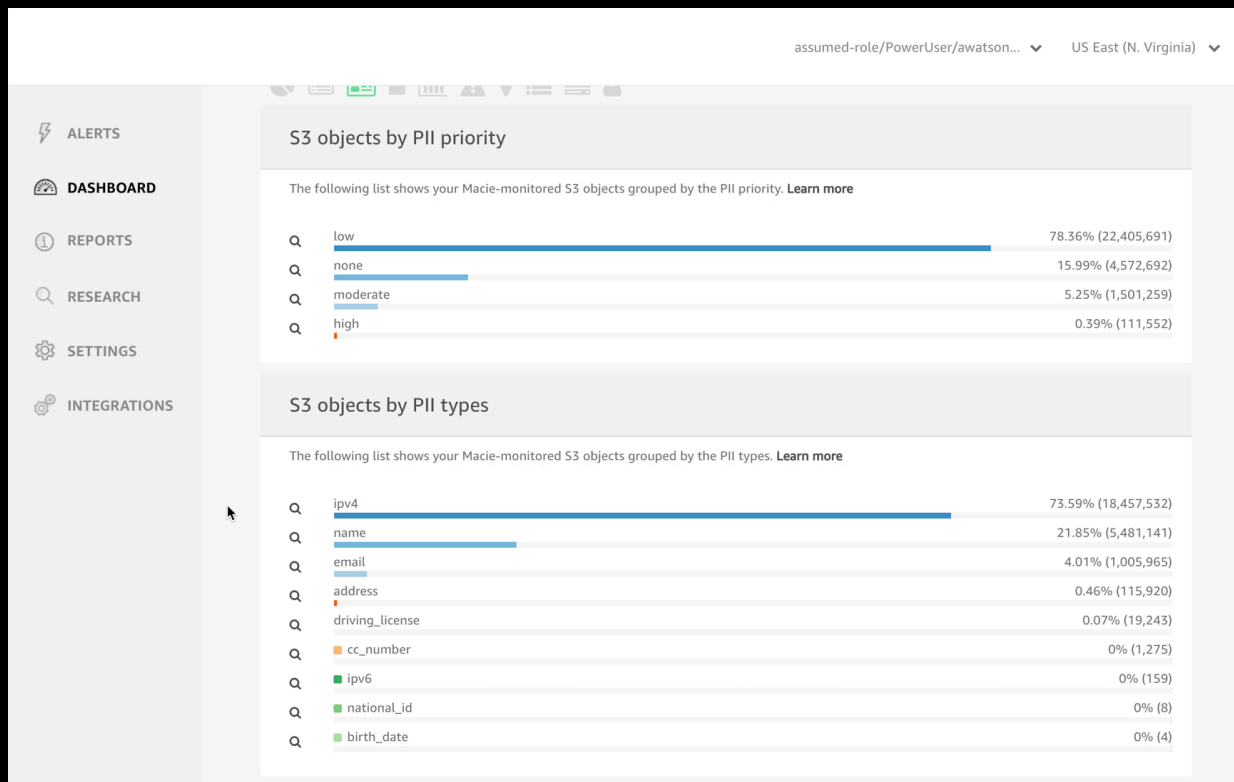


Automatically
monitor for
anomalies

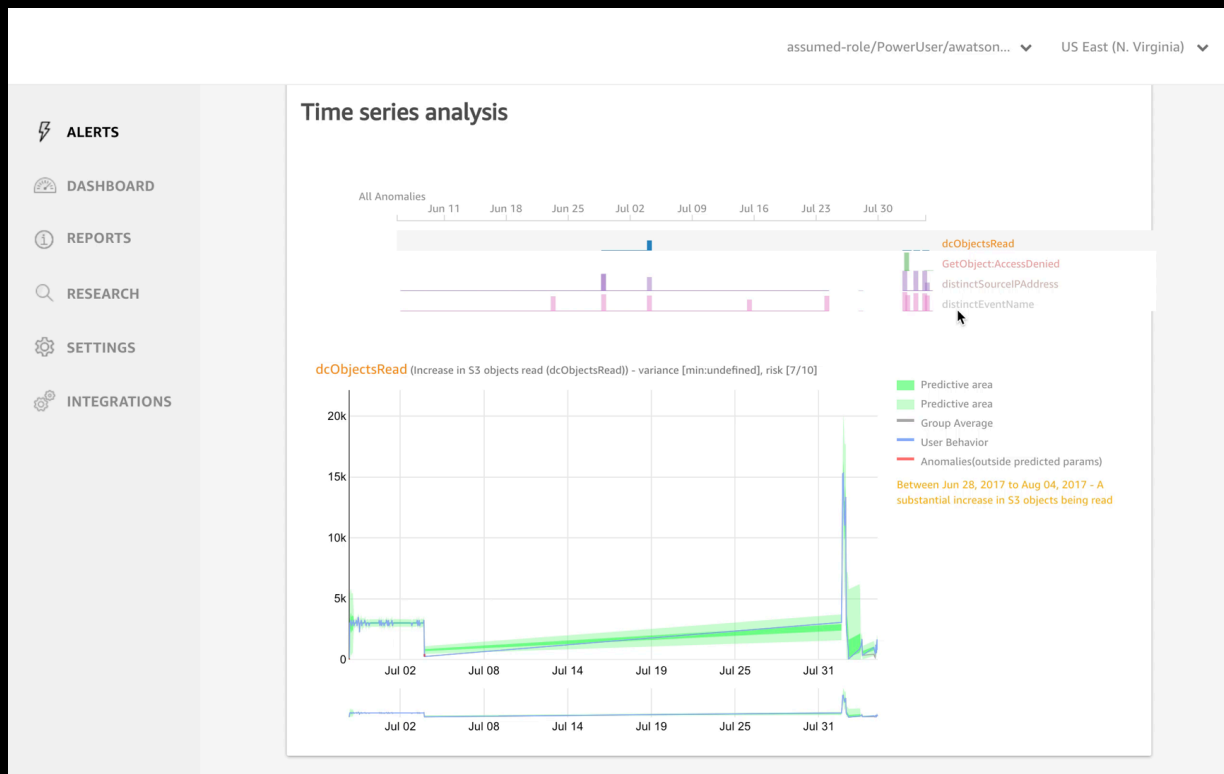


Alert your
security team
when anomalies
are detected

Gain Visibility Into Globally Shared Content

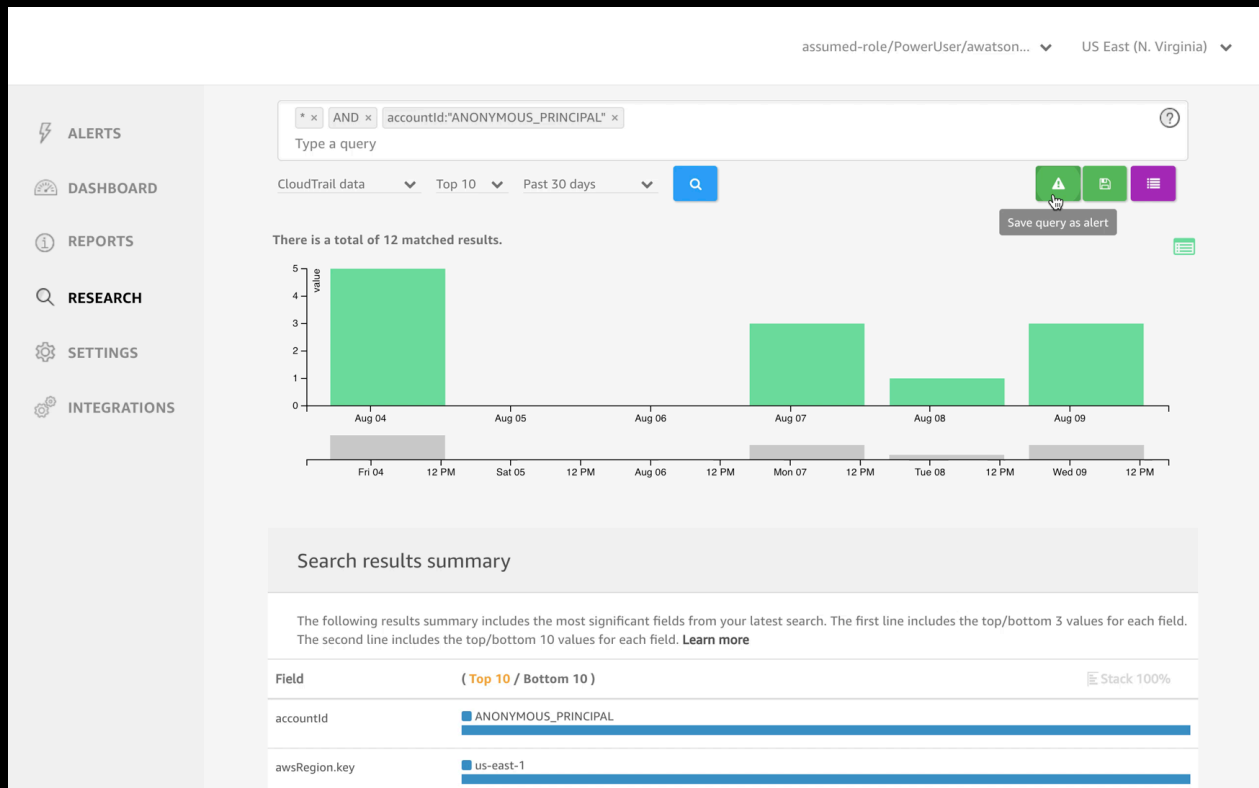


Identify anomalous accesses



Implementing Continuous Compliance

GDPR, PCI, PII



Customers Using Amazon Macie



Continuous insights into cloud infrastructure and practices



Securing PII and alerting to access anomalies



Delivering instant information and detail in the dashboard



We support, use and tune Keras,
Tensorflow and Caffe, but
Apache MXNet is the deep
learning framework
of choice for AWS

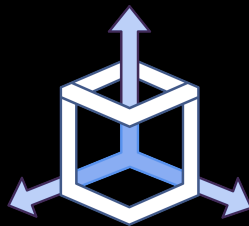
Tutorial introduction:
<http://gluon.mxnet.io/>

Apache MXNet



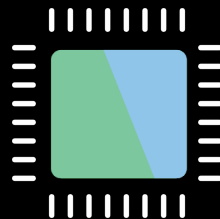
Programmable

Simple syntax,
multiple languages



Portable

Highly efficient
models for mobile
and IoT



High Performance

Near linear scaling
across hundreds of GPUs

Why Apache MXNet?



Open Model

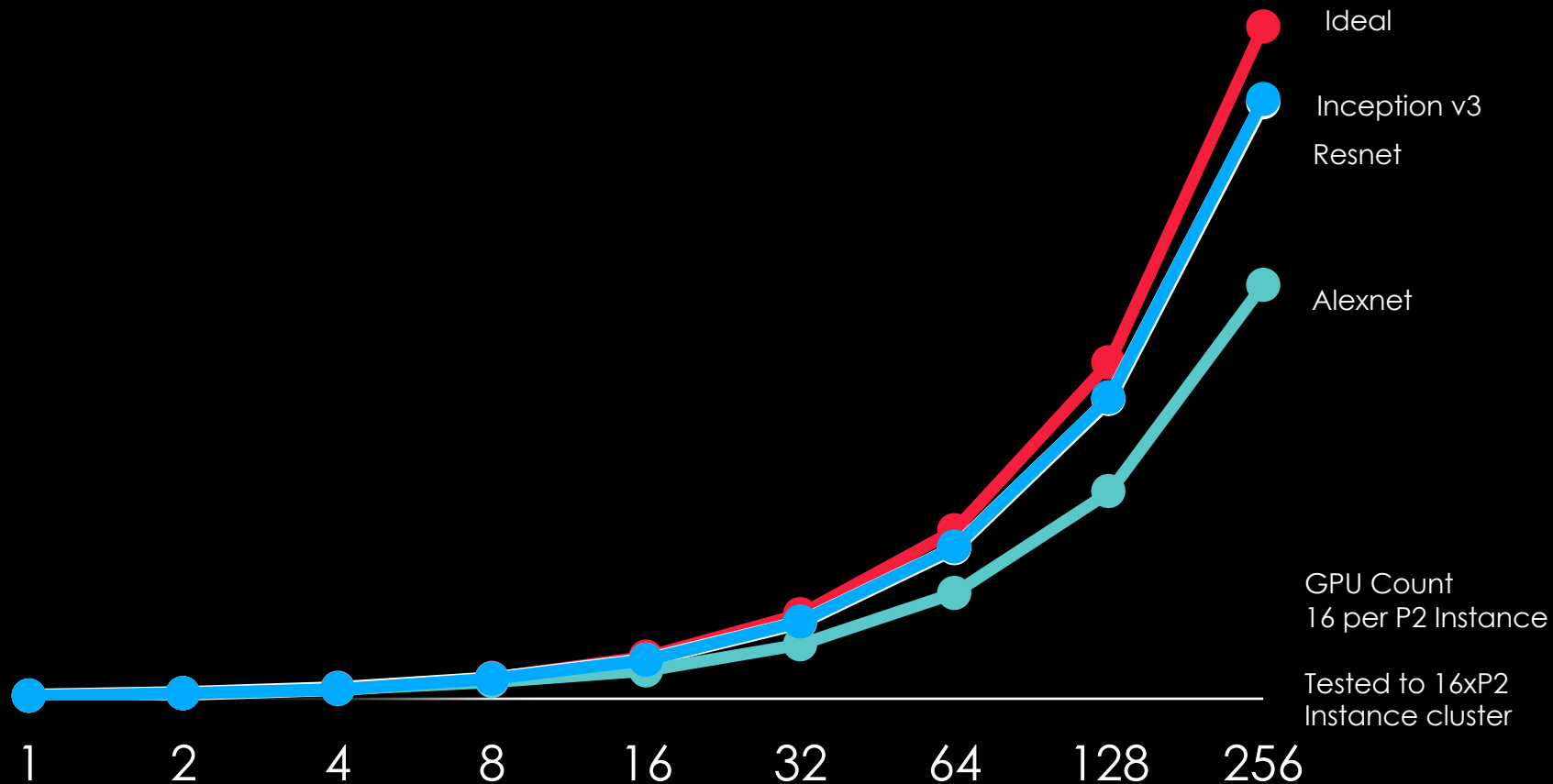
Accepted into the
Apache Incubator
Broad contributions
e.g. Apple Core ML



Better Performance

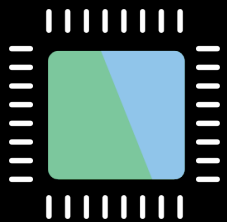
Excellent scalability
Optimized stack for
deep learning on AWS

Scaling With Apache MXNet: 88% Efficiency



Apache MXNet Introduction

AWS Deep Learning AMIs: One-Click Deep Learning



Kepler, Volta
& Skylake



Apache MXNet
TensorFlow Caffe2,
CNTK Keras,
Theano, Torch

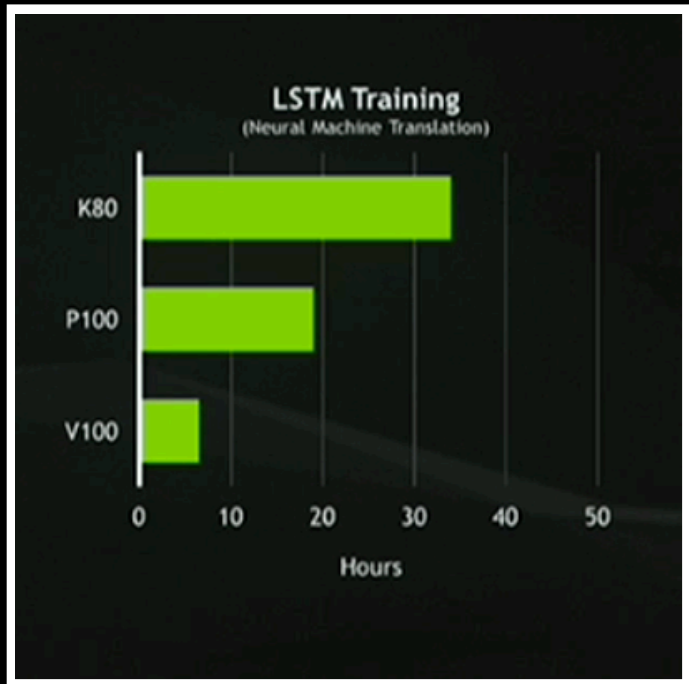


Python 3



Jupyter
Notebooks
& Examples

Training Artificial Intelligence With GPUs



NEW! NVIDIA Volta

Train in hours, not days

Custom built for artificial intelligence

Core of the next AWS GPU instance family

MXNet is already optimized for Volta

AWS Deep Learning AMI

Up to ~40k CUDA cores on P2

Apache MXNet

TensorFlow

Theano

Caffe & Caffe 2

Torch

Keras

Pre-configured CUDA drivers,

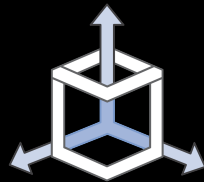
MKL

Anaconda, Python3

Ubuntu or Amazon Linux

+ CloudFormation template

+ Container Image



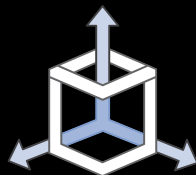
One-Click GPU & CPU
Open Source
Deep Learning
Installed, Tested, Tuned
Bootable Machine Image

Apache MXNet



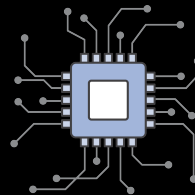
Programmable

Simple syntax,
multiple languages



Portable

Highly efficient
models for mobile
and IoT



High Performance

Near linear scaling
across hundreds of GPUs



Open Governance

Accepted into the
Apache Incubator



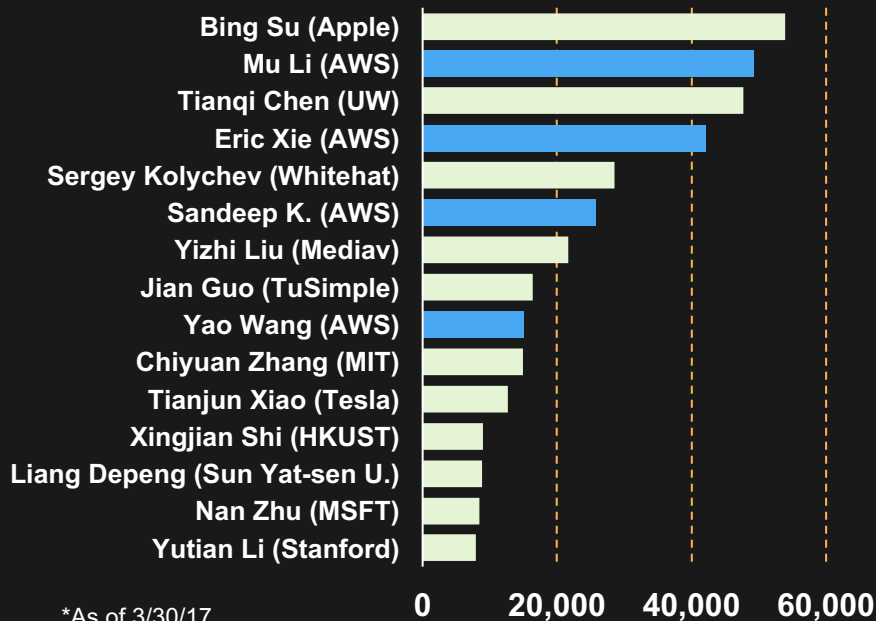
Tuned On AWS

Optimized performance and
scalability on AWS GPUs



Apache MXNet | Collaborations and Community

Diverse Community



*As of 3/30/17

**Amazon @35% of Contributions

Try out Keras models on MXNet instead of Tensorflow

Also see <https://medium.com/@julsimon/keras-shoot-out-tensorflow-vs-mxnet-51ae2b30a9c0>

Faster, smaller and more accurate...



François Chollet ✓

@fchollet

Follow

Now in beta...



Keras CNTK backend:

github.com/fchollet/keras ...



Keras MXNet backend:

Translate from Indonesian



dmlc/keras

keras - Deep Learning library for Python. Convnets, recurrent neural networks, and more. Runs on MXNet, Theano or TensorFlow.

github.com

All it takes is one line in the `~/.keras/keras.json` file.

```
1  # ~/.keras/keras.json for Tensorflow
2  {
3      "image_dim_ordering": "tf",
4      "epsilon": 1e-07,
5      "floatx": "float32",
6      "backend": "tensorflow"
7  }
8
9  # ~/.keras/keras.json for MXNet
10 {
11     "image_dim_ordering": "tf",
12     "epsilon": 1e-07,
13     "floatx": "float32",
14     "backend": "mxnet"
15 }
```

Deep Learning Framework Comparison

	Apache MXNet	TensorFlow	Cognitive Toolkit
Industry Owner	N/A – Apache Community	Google	Microsoft
Programmability	Imperative and Declarative	Declarative only	Declarative only
Language Support	R, Python, Scala, Julia, C++, Javascript, Go, Matlab, Perl...	Python, C++, Experimental Go and Java	Python, C++, Brainscript.
Code Length AlexNet (Python)	44 sloc	107 sloc using TF.Slim	214 sloc
Memory Footprint (LSTM)	2.6GB	7.2GB	N/A

Apache MXNet | Amazon Strategy



Integrate with AWS Services

Bring Scalable Deep Learning to EMR, Lambda, ECS and many more..



Amazon AI



Discovery & Search



Fulfillment & Logistics



Enhance Existing Products

Foundation for AI Services

Higher Velocity for AI Services, Research and Core AI Development



Leverage the Community

Community brings velocity and innovation with no industry ownership
Safest for long term investment

Apache MXNet API

Apache MXNet | The Basics

- ***NDArray***: Manipulate multi-dimensional arrays (tensors) in a command line paradigm (imperative).
- ***Symbol***: Symbolic expression for neural network flows (declarative).
- ***Module***: Intermediate-level and high-level interface for neural network training and inference.
- **Loading Data**: Feeding data into training/inference programs.
- **Mixed Programming**: Training algorithms developed using *NDArrays* in concert with *Symbols*.

<https://medium.com/@julsimon/an-introduction-to-the-mxnet-api-part-1-848febdcf8ab>

Imperative Programming

```
import numpy as np  
a = np.ones(10)  
b = np.ones(10) * 2  
c = b * a  
d = c + 1
```

Easy to tweak
in Python, R, Perl etc.

PROS

- Straightforward and flexible.
- Take advantage of language native features (loop, condition, debugger).
- E.g. Numpy, Matlab, Torch, ...

CONS

- Hard to optimize

Declarative Programming

```
A = Variable('A')
```

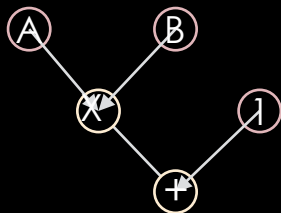
```
B = Variable('B')
```

```
C = B * A
```

```
D = C + 1
```

```
f = compile(D)
```

```
d = f(A=np.ones(10),  
      B=np.ones(10)*2)
```



C can share memory with
D because C is deleted
later

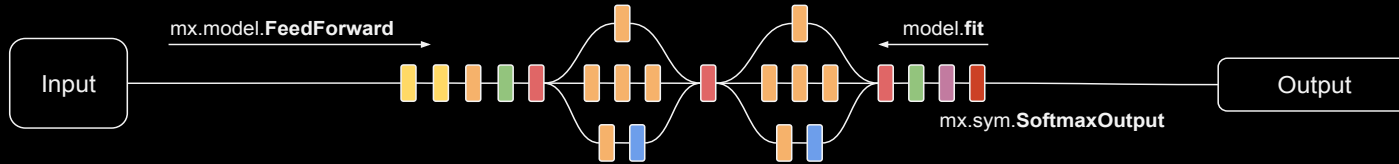
PROS

- More chances for optimization
- Separates flow structure
- E.g. TensorFlow, Theano, Caffe

CONS

- Less flexible, hard to debug

Deep Learning Models



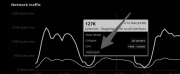
Image



Video



Speech



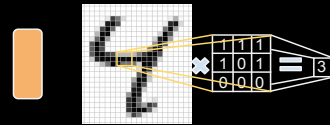
Events

"People Riding Bikes"

Text

Input $\begin{bmatrix} 1 \\ 3 \\ \dots \\ 4 \end{bmatrix} \times \text{Weights} \begin{bmatrix} 0.2 \\ -0.1 \\ \dots \\ 0.7 \end{bmatrix} = 2$

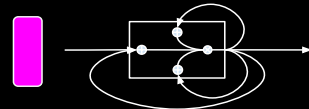
`mx.sym.FullyConnected(data, num_hidden=128)`



`mx.sym.Convolution(data, kernel=(5,5), num_filter=20)`

Max $\begin{bmatrix} 4 & 2 \\ 2 & 0 \end{bmatrix} = 4$ Avg $\begin{bmatrix} 4 & 2 \\ 2 & 0 \end{bmatrix} = 2$

`mx.sym.Pooling(data, pool_type="max", kernel=(2,2), stride=(2,2))`



`lstm.lstm_unroll(num_lstm_layer, seq_len, len, num_hidden, num_embed)`

Queen $\begin{bmatrix} 0.2 \\ -0.1 \\ \dots \\ 0.7 \end{bmatrix} \rightarrow \cos(w, \text{queen}) = \cos(w, \text{king}) - \cos(w, \text{man}) + \cos(w, \text{woman})$

`mx.symbol.Embedding(data, input_dim, output_dim = k)`

`mx.sym.Activation(data, act_type="xxxx")`

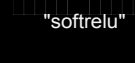
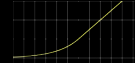
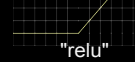
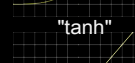
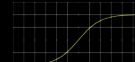
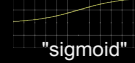
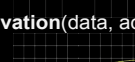
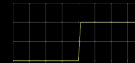


Image Segmentation



Face Search



Neural Art

"People Riding Bikes"

Image Caption

Bicycle, People, Road, Sport

Image Labels

"Οι άνθρωποι ιππασίας ποδήλατα"

Machine Translation

Do It Yourself - Robot Racing Cars



DIYrobocars Meetup

<https://www.meetup.com/diyrobocars/>

Follow @DIYrobocars for updates

Common Implementation ~\$200

1:16 scale RC Truck + Raspberry Pi3 + Camera + EC2

<https://github.com/wroscoe/donkey> software

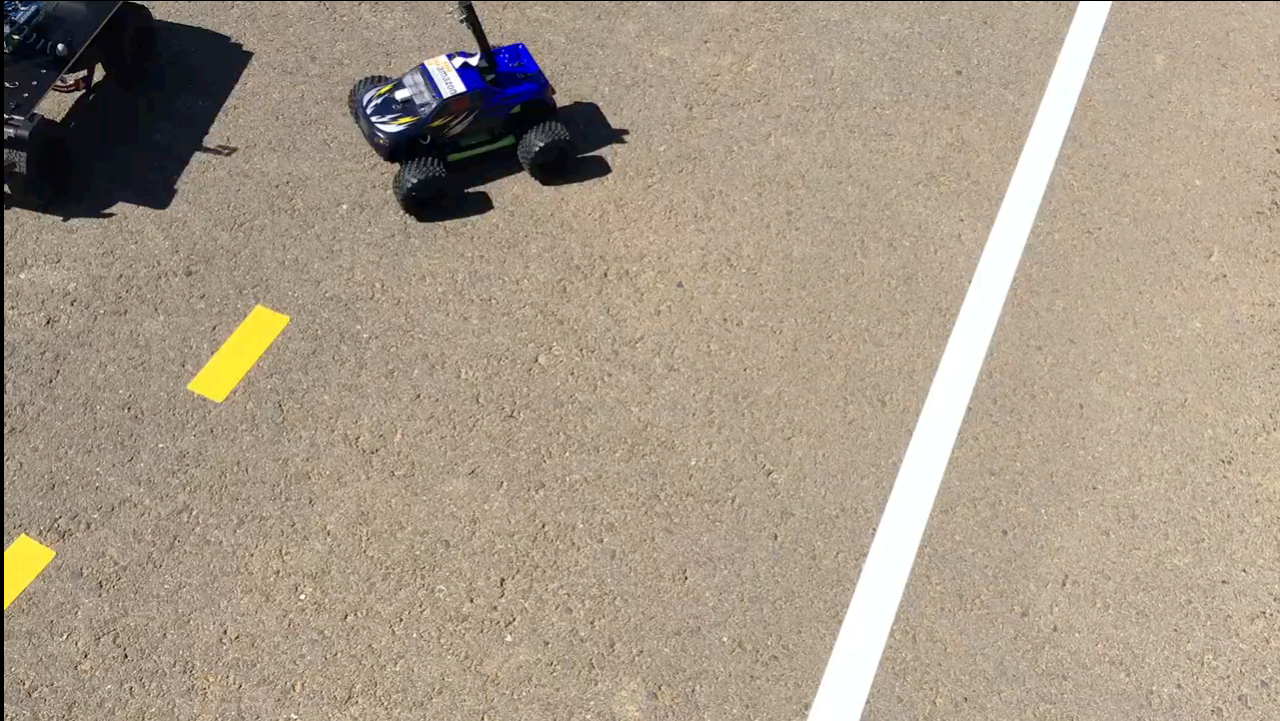
Camera feeds Donkey software with Keras or MXNet model, trained (slowly) on laptop or (quickly) on EC2 GPU instance

Will Roscoe wrote Donkey

Adrian built Donkey Hoté and instigated...

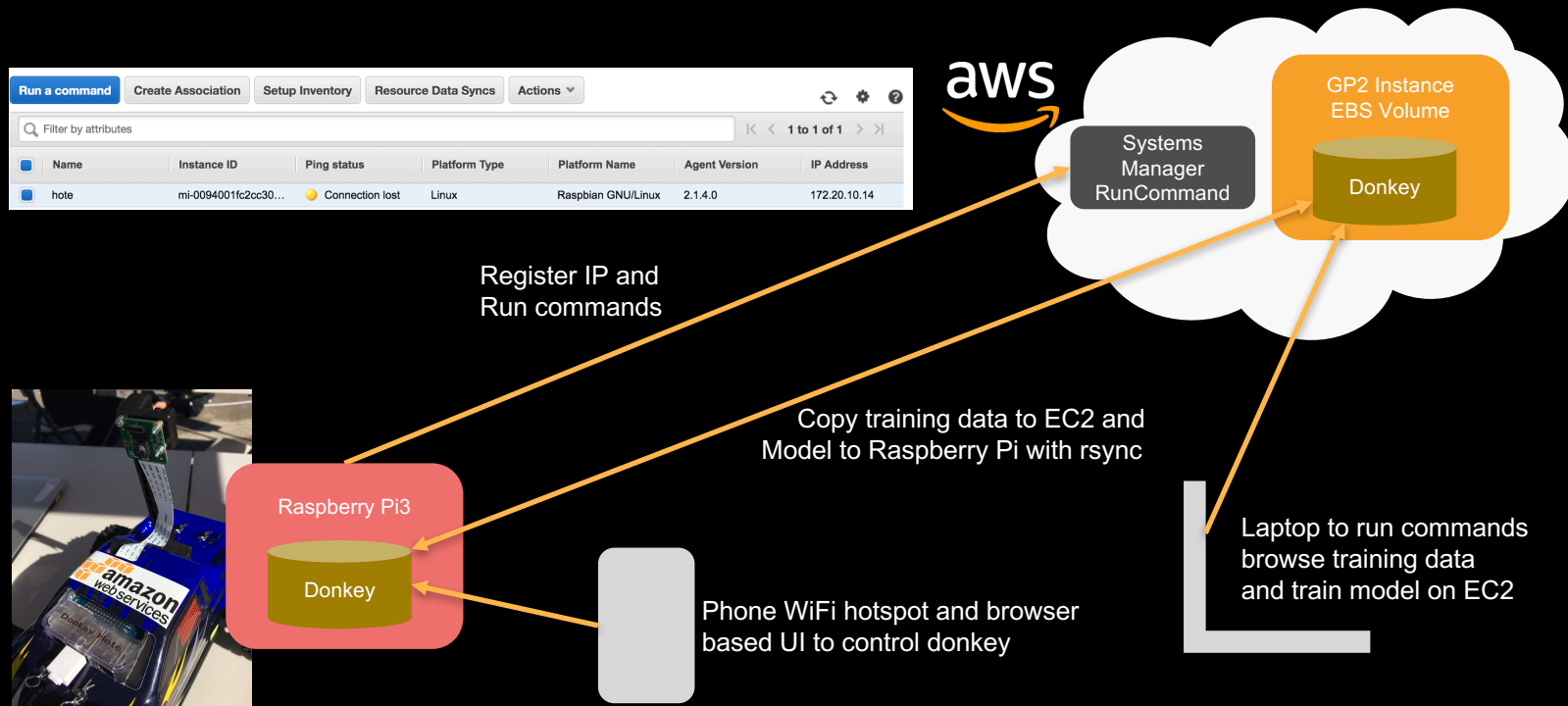
Sunil Mallya added MXNet model and got AWS officially involved

Donkey Hoté's first action...





How The Donkey Software Works



One of the Keras model options for Donkey

```
def default_categorical():
    from keras.layers import Input, Dense, merge
    from keras.models import Model
    from keras.layers import Convolution2D, MaxPooling2D, Reshape, BatchNormalization
    from keras.layers import Activation, Dropout, Flatten, Dense

    img_in = Input(shape=(120, 160, 3), name='img_in')
    x = img_in
    x = Convolution2D(24, (5,5), strides=(2,2), activation='relu')(x)
    x = Convolution2D(32, (5,5), strides=(2,2), activation='relu')(x)
    x = Convolution2D(64, (5,5), strides=(2,2), activation='relu')(x)
    x = Convolution2D(64, (3,3), strides=(2,2), activation='relu')(x)
    x = Convolution2D(64, (3,3), strides=(1,1), activation='relu')(x)
    x = Flatten(name='flattened')(x)
    x = Dense(100, activation='relu')(x)
    x = Dropout(.1)(x)
    x = Dense(50, activation='relu')(x)
    x = Dropout(.1)(x)
    # categorical output of the angle
    # Connect every input with every output, output 15 hidden units. Use Softmax to get percentage. 15 categories, find best one based off percentage 0.0-1.0
    angle_out = Dense(15, activation='softmax', name='angle_out')(x)
    # continous output of throttle
    throttle_out = Dense(1, activation='relu', name='throttle_out')(x)
    model = Model(inputs=[img_in], outputs=[angle_out, throttle_out])
    model.compile(optimizer='rmsprop',
                  loss={'angle_out': 'categorical_crossentropy',
                        'throttle_out': 'mean_absolute_error'},
                  loss_weights={'angle_out': 0.9, 'throttle_out': .001})

    return model
```

First layer, input layer, Shape comes from camera.py resolution, RGB

24 features, 5px5p kernel (convolution, feature) window, 2wx2h stride, relu

32 features, 5px5p kernel window, 2wx2h stride, relu activation

64 features, 5px5p kernal window, 2wx2h stride, relu

64 features, 3px3p kernal window, 2wx2h stride, relu

64 features, 3px3p kernal window, 1wx1h stride, relu

Flatten to 1D (Fully connected)

Classify the data into 100 features, make all negatives 0

Randomly drop out (turn off) 10% of the neurons (Prevent overfitting)

Classify the data into 50 features, make all negatives 0

Randomly drop out 10% of the neurons (Prevent overfitting)

Reduce to 1 number, Positive number only



Robocar Rally 2017

at AWS re:Invent



SUNDAY, NOV. 26 | 6:00PM – 10:00PM
MONDAY, NOV. 27 | 7:00AM – 12:00AM

Registration for Robocar Rally 2017 will launch on Oct. 19 as part of the Reserved Seating launch

<https://reinvent.awsevents.com/learn/robocar-rally/>

<https://aws.amazon.com/blogs/ai/build-an-autonomous-vehicle-on-aws-and-race-it-at-the-reinvent-robocar-rally/>

Additional Resources

MXNet Resources

- [MXNet Blog Post | AWS Endorsement](#)
 - <http://www.allthingsdistributed.com/2016/11/mxnet-default-framework-deep-learning-aws.html>
- [Read up on MXNet and Learn More:](#)
 - [See gluon.mxnet.io](#) <https://github.com/dmlc/mxnet/>
- [Re:Invent MXNet Recommender Systems Talk](#) by Leo Dirac
 - https://www.portal.reinvent.awsevents.com/connect/sessionDetail.wv?SESSION_ID=8591

AWS Resources: follow Julien Simon @julsimon, Sunil Mallya @sunilmallya

- [Deep Learning AMI](#)
 - <https://aws.amazon.com/marketplace/pp/B01M0AXXQB> | Amazon Linux
 - <https://aws.amazon.com/marketplace/pp/B06VSPXKDX> | Ubuntu
- [CloudFormation Template Instructions](#)
 - <https://github.com/dmlc/mxnet/tree/master/tools/cfn>
- [Deep Learning Benchmark](#)
 - <https://github.com/aws-labs/deeplearning-benchmark>
- [MXNet on Lambda](#)
 - <https://github.com/aws-labs/mxnet-lambda>
- [MXNet on ECS/Docker](#)
 - <https://github.com/aws-labs/ecs-deep-learning-workshop>

THANK YOU!

@adrianco