# Apache MXNet (Incubating)

Steffen Rochel
steroche@amazon.com; @srochel

AT&T **mobile**™ **lite**
@**accessibility**
for Android™

Making Android phones
accessible to the blind

**at&t**

Save

# Why Amazon's Alexa Is 'Life Changing' for the Blind

*For the blind, Amazon's Echo and Alexa is more than just neat technology; it's a lifeline.*
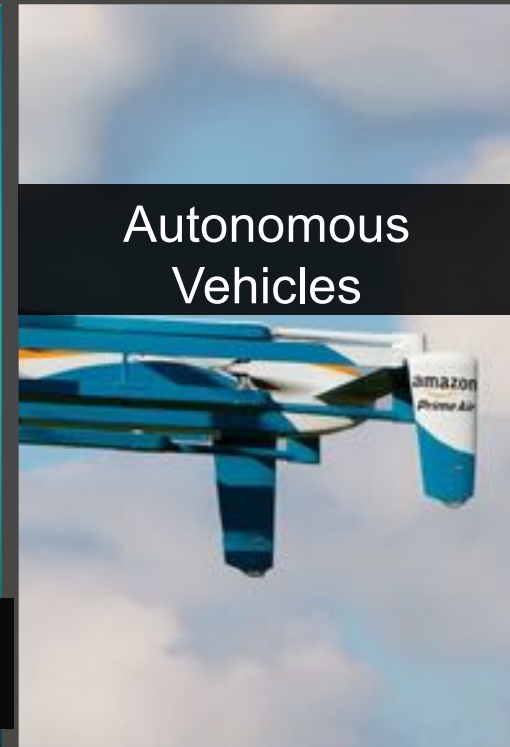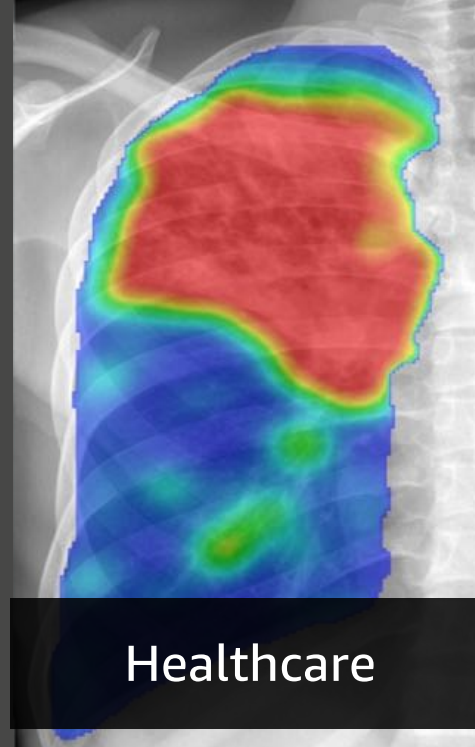
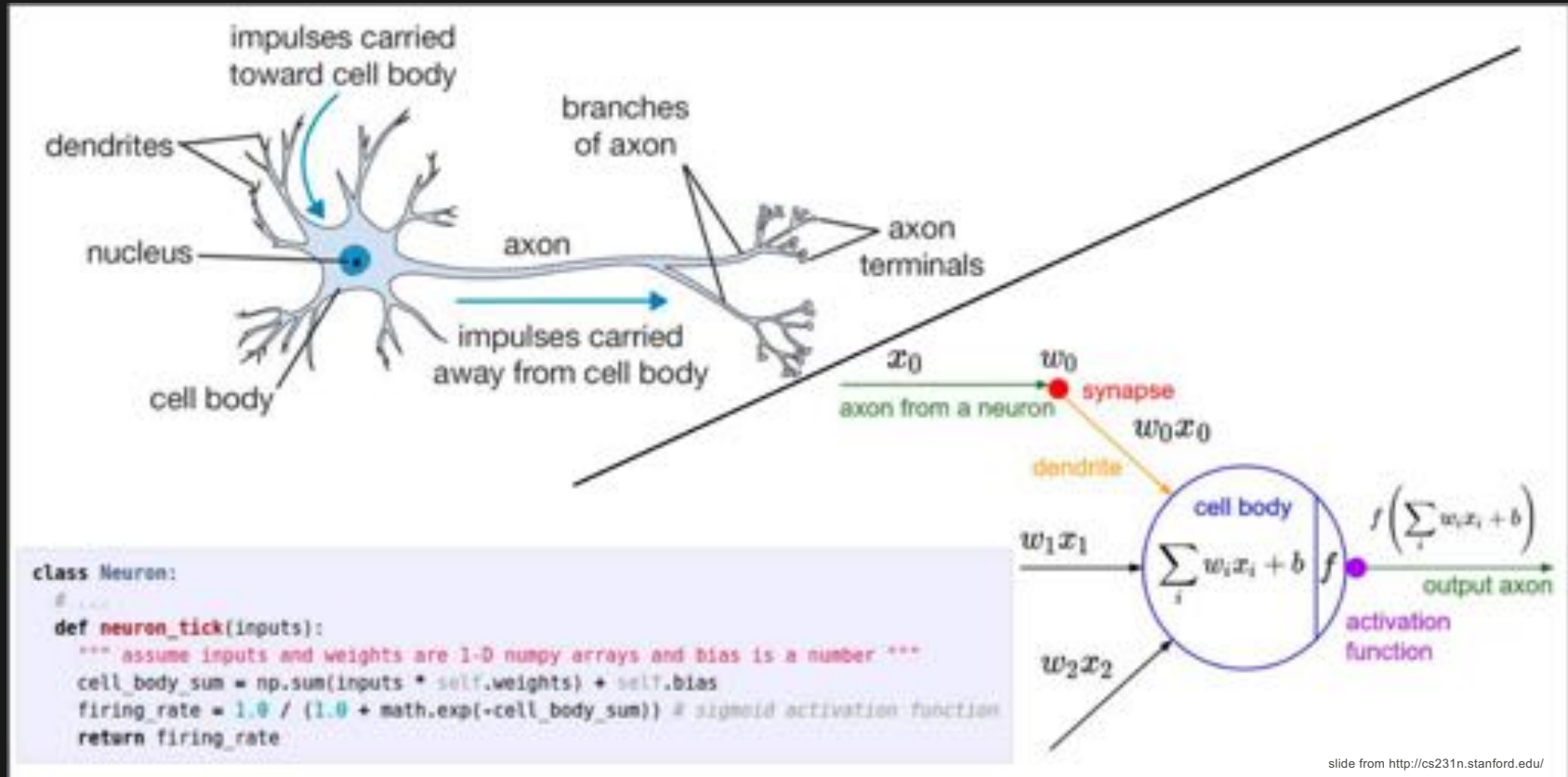By Jon Kalish   January 8, 2018 8:00AM EST

f  ☑  in  ℗  ⬡  ✉  co   **2.7K SHARES**

# Why Deep Learning?

Personalization

Logistics

Healthcare

Autonomous Vehicles

# Biological & Artificial Neuron



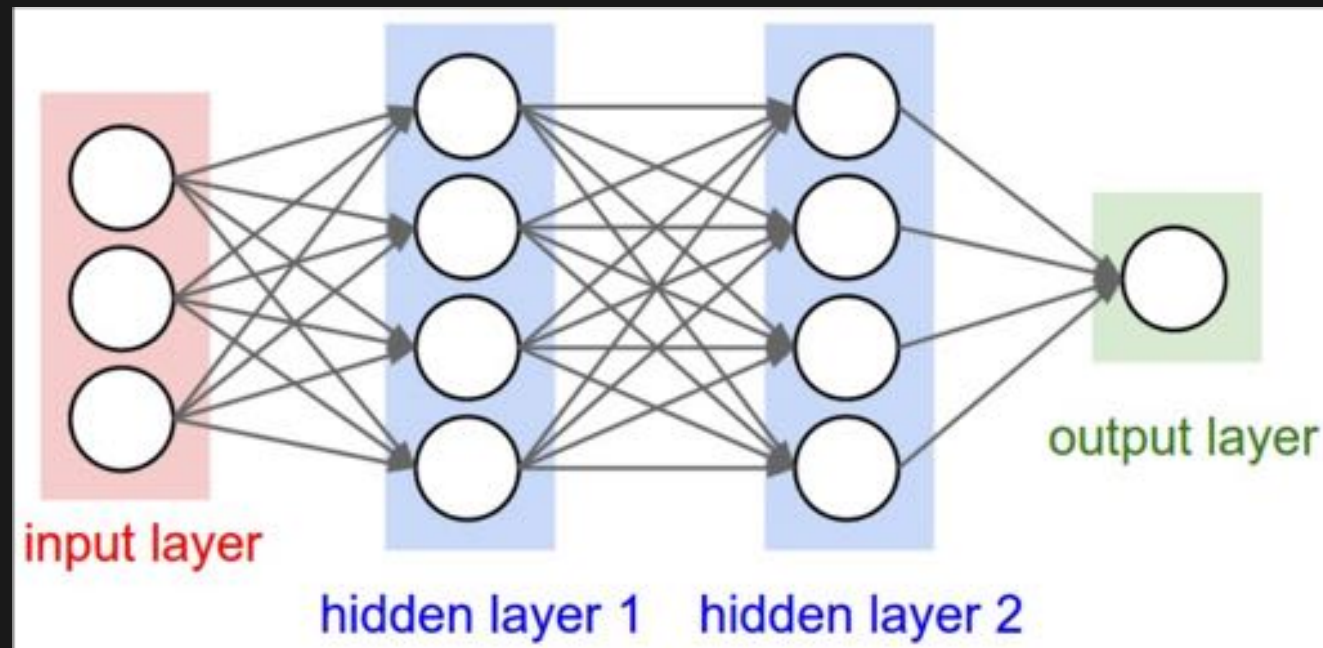slide from http://cs231n.stanford.edu/
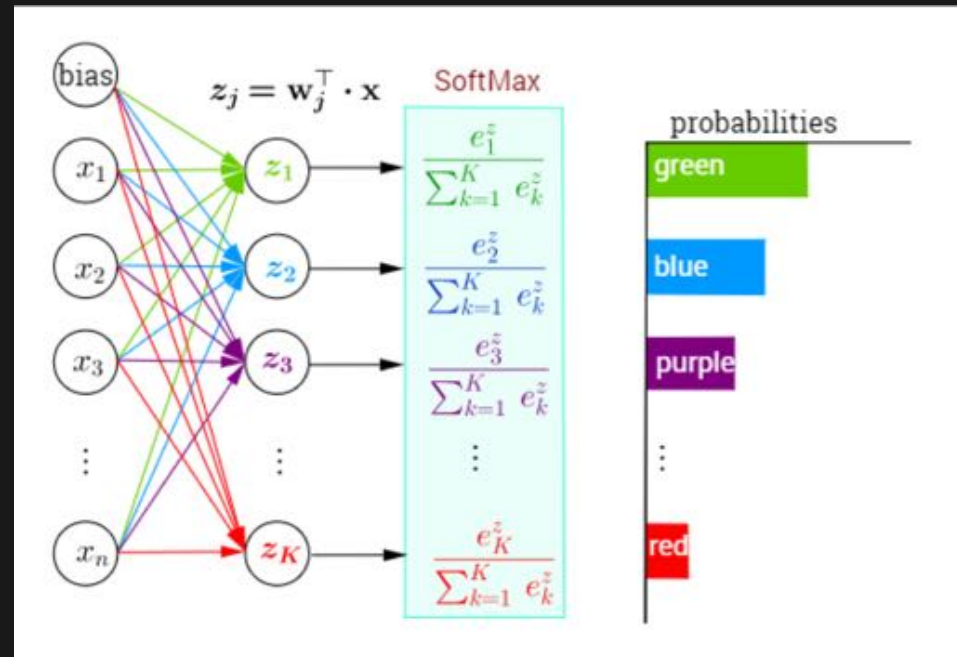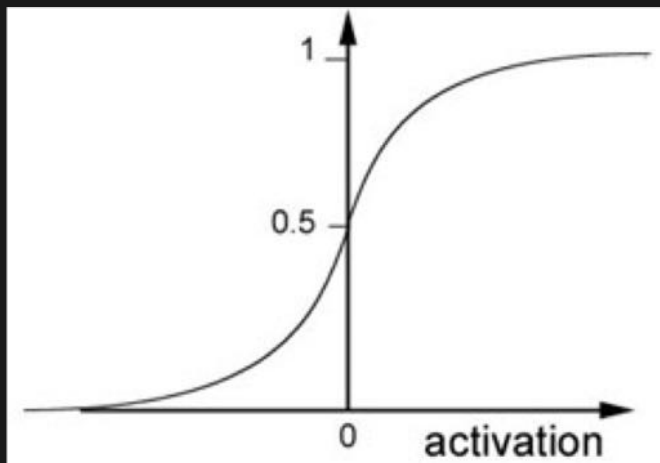
Source: http://cs231n.github.io/neural-networks-1/

# Fully Connected Layer



Each node ("neuron") in a layer is connected to every node in the previous layer
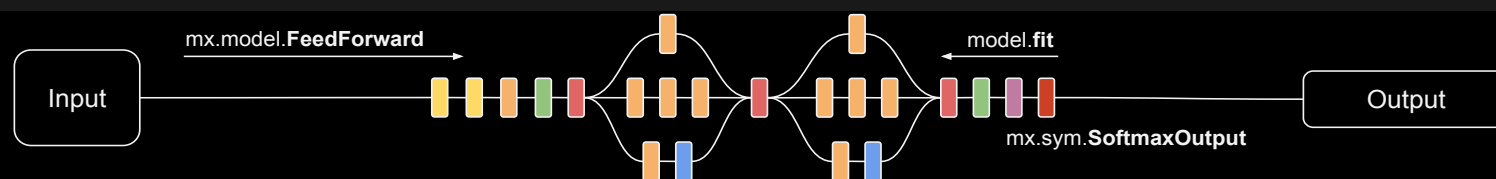
# Classification with the Softmax Function

**Softmax Function**



Softmax converts the output layer into probabilities – necessary for classification

# Deep Learning Models



mx.model.**FeedForward**

Input

model.**fit**

Output

mx.sym.**SoftmaxOutput**

Image

Video

Speech

Events

*"People Riding Bikes"*

Text

**Input**    **Weights**

| 1 |   | 0.2 |
|---|---|-----|
| 3 |   | -0.1 |
| ... | × | ... |
| 4 |   | 0.7 |

= 2

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

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

Max | 4 | 2 |    = 4    Avg | 4 | 2 |    = 2
    | 2 | 0 |              | 2 | 0 |

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 → | 0.2 |
        | -0.1 |
        | ... |
        | 0.7 |

$cos(w, \textbf{queen}) = cos(w, \textbf{king}) - cos(w, \textbf{man}) + cos(w, \textbf{woman})$

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

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

"sigmoid"

"tanh"

"relu"

"softrelu"

Image Segmentation

Face Search

Neural Art

*"People Riding Bikes"*

Image Caption

*Bicycle, People, Road, Sport*

Image Labels
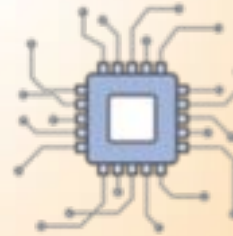
*"Οι άνθρωποι ιππασίας ποδήλατα"*

Machine Translation

# MXNet – a Deep Learning Framework

# Easily build, train, and deploy models with MXNet

### Start with high quality, pre-trained models

- Gluon CV and Gluon NLP

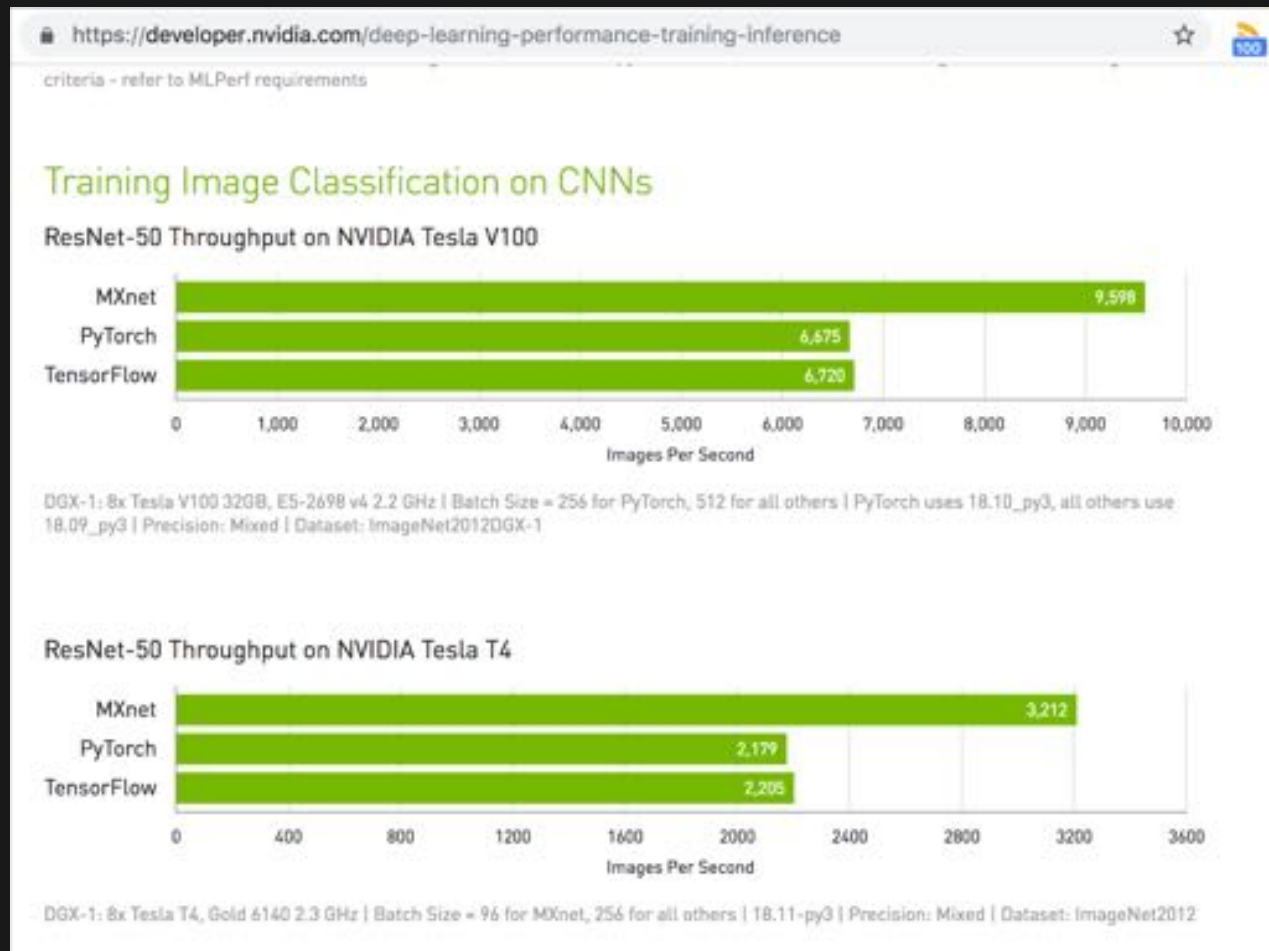### Refine with fast, scalable training

- Keras-MXNet up to 2x faster than Keras-TensorFlow
- Near-linear scalability up to 256 GPUs
- Dynamic training

### Deploy using familiar tools

- Java/Scala APIs
- MXNet Model Server

# Performance

# Why Gluon

**Simple, Easy-to-Understand Code**

**Flexible, Imperative Structure**

**Dynamic Graphs**

**High Performance**

# Gluon Code – Network Definition

```python
net = gluon.nn.HybridSequential()

with net.name_scope():

    net.add(gluon.nn.Dense(units=64, activation='relu'))

    net.add(gluon.nn.Dense(units=10))

softmax_cross_entropy = gluon.loss.SoftmaxCrossEntropyLoss()

net.initialize(mx.init.Xavier(magnitude=2.24), ctx=ctx, force_reinit=True)

trainer = gluon.Trainer(net.collect_params(), 'sgd', {'learning_rate': 0.02})
```

# Gluon Code – Training

```python
smoothing_constant = .01

for e in range(10):

    cumulative_loss = 0

    for i, (data, label) in enumerate(train_data):

        data = data.as_in_context(model_ctx).reshape((-1, 784))

        label = label.as_in_context(model_ctx)

        with autograd.record():

            output = net(data)

            loss = softmax_cross_entropy(output, label)

        loss.backward()

        trainer.step(data.shape[0])
```

aws

# MXNet EcoSystem

- Gluon Model Zoo
  https://mxnet.incubator.apache.org/api/python/gluon/model_zoo.html

- Sockeye: A Toolkit for Neural Machine Translation
  https://arxiv.org/abs/1712.05690

- GluonCV: A Deep Learning Toolkit for Computer Vision
  https://gluon-cv.mxnet.io/

- GluonNLP: a Deep Learning Toolkit for Natural Language Processing
  http://gluon-nlp.mxnet.io/

- DeepAR: Probabilistic Forecasting with Autoregressive Recurrent Networks
  https://arxiv.org/abs/1704.04110v2

- MXNet Model Server
  https://github.com/awslabs/mxnet-model-server

# MXNET.IO

Steffen Rochel

@srochel

# Visual Search

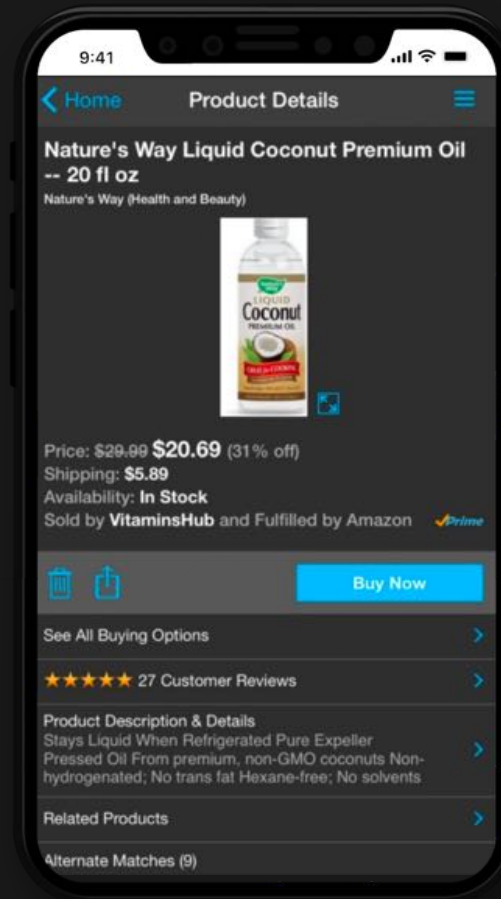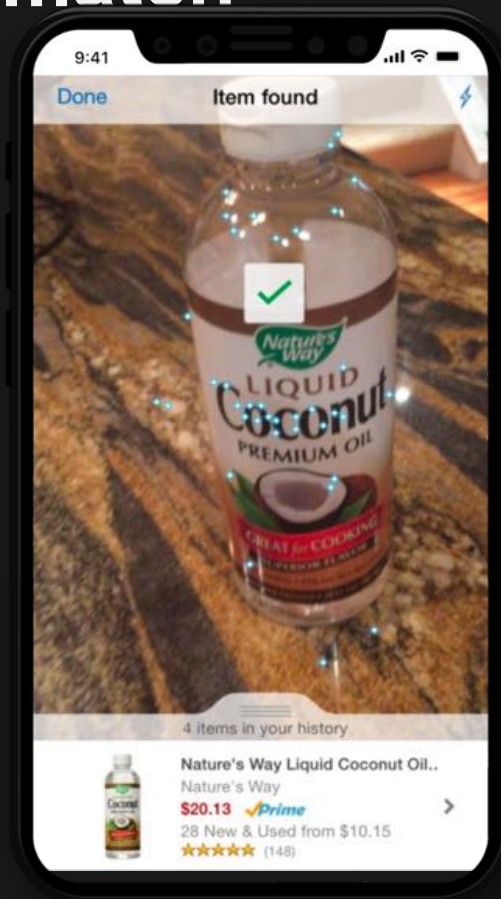# Drive content and product



Search Engines

Bing search

Google search

Social

Pinterest

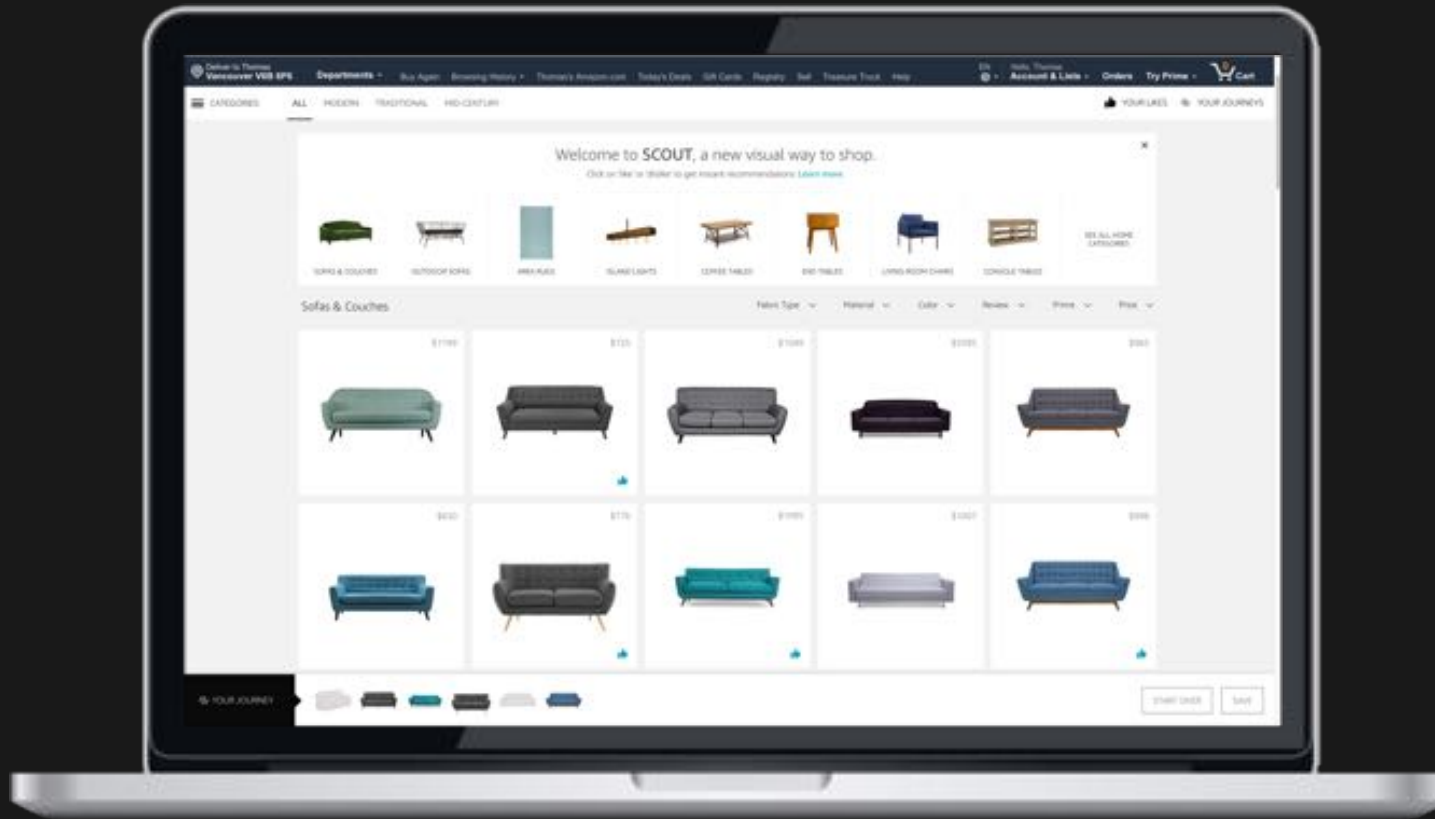(not actual screenshot)

# Exact match



Amazon  Shopping

# Visual recommendations



Amazon **Scout**

"The future of search will be about pictures rather than keywords."
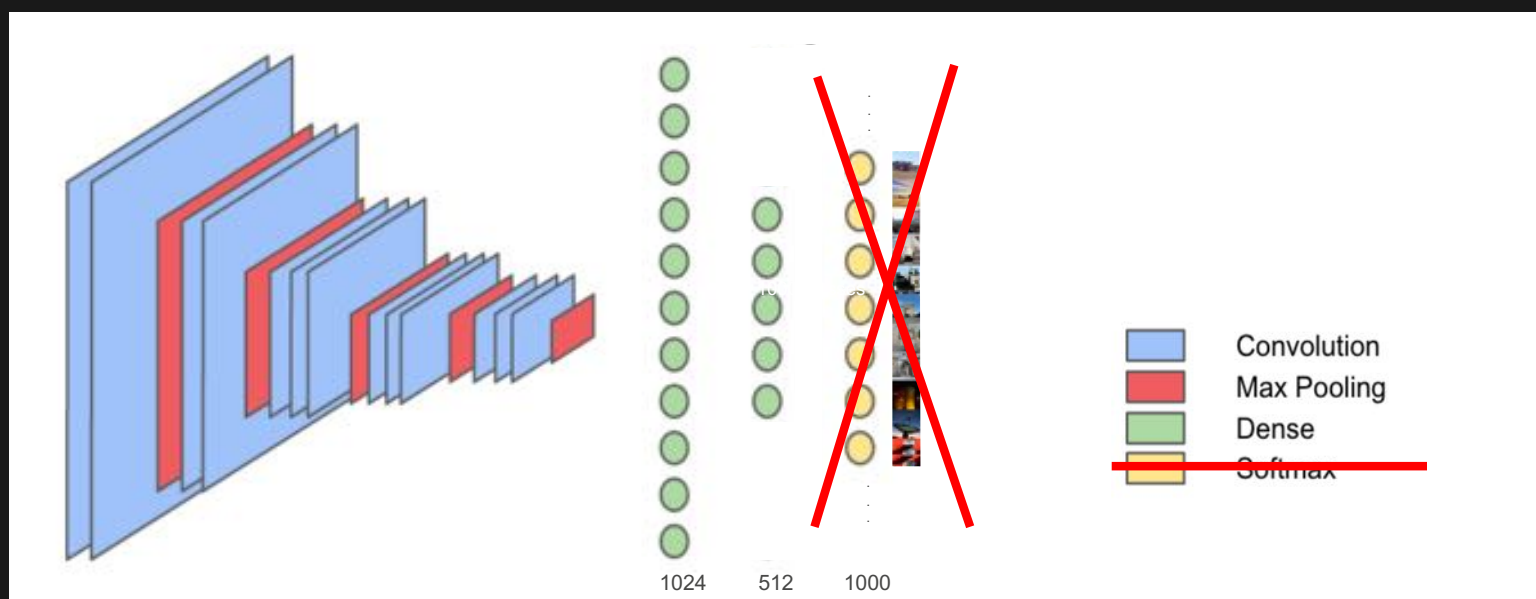
Ben Silbermann
Pinterest CEO

# Growing market

Angel.co: 76 startups

Syte.ai, Slyce.it and others

# How does it work?

aws re:Invent

aws

# Turn the network into a featurizer
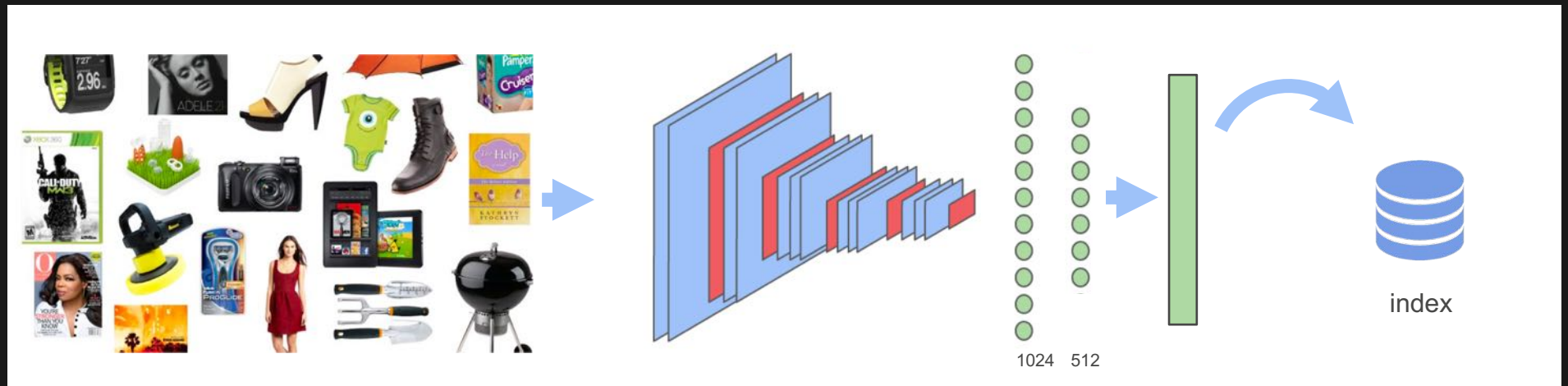
pre-trained CNN (resnet18) from the MXNet Gluon model zoo

# Index images to lower-dimensional representation

Images: 224x224x3 dimensions

pre-trained network

512 dimensions "finger-print"

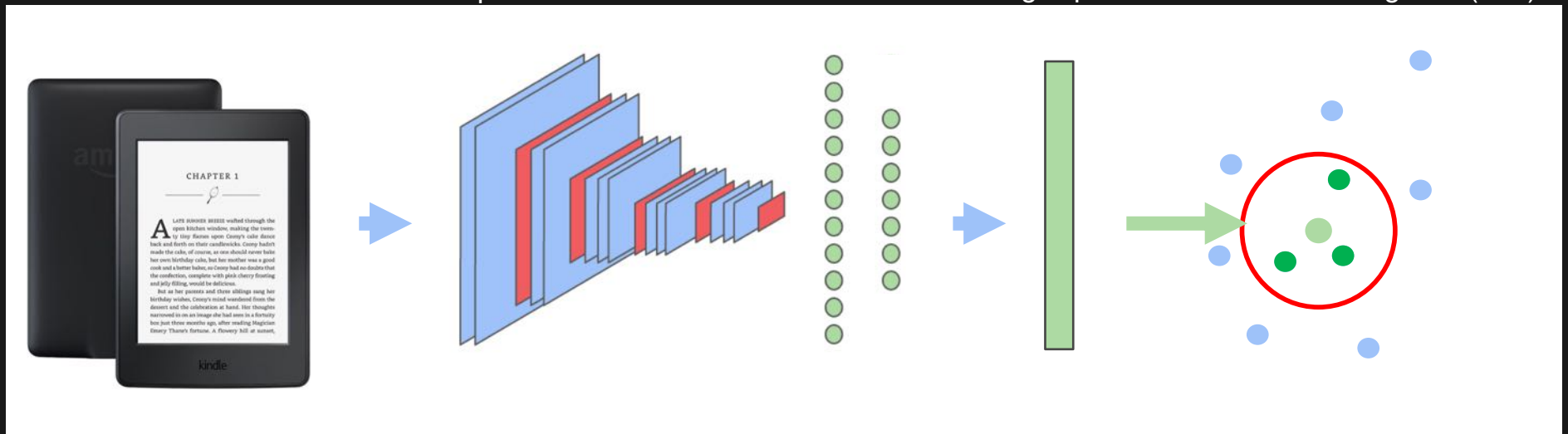

1024   512

index

# New image query
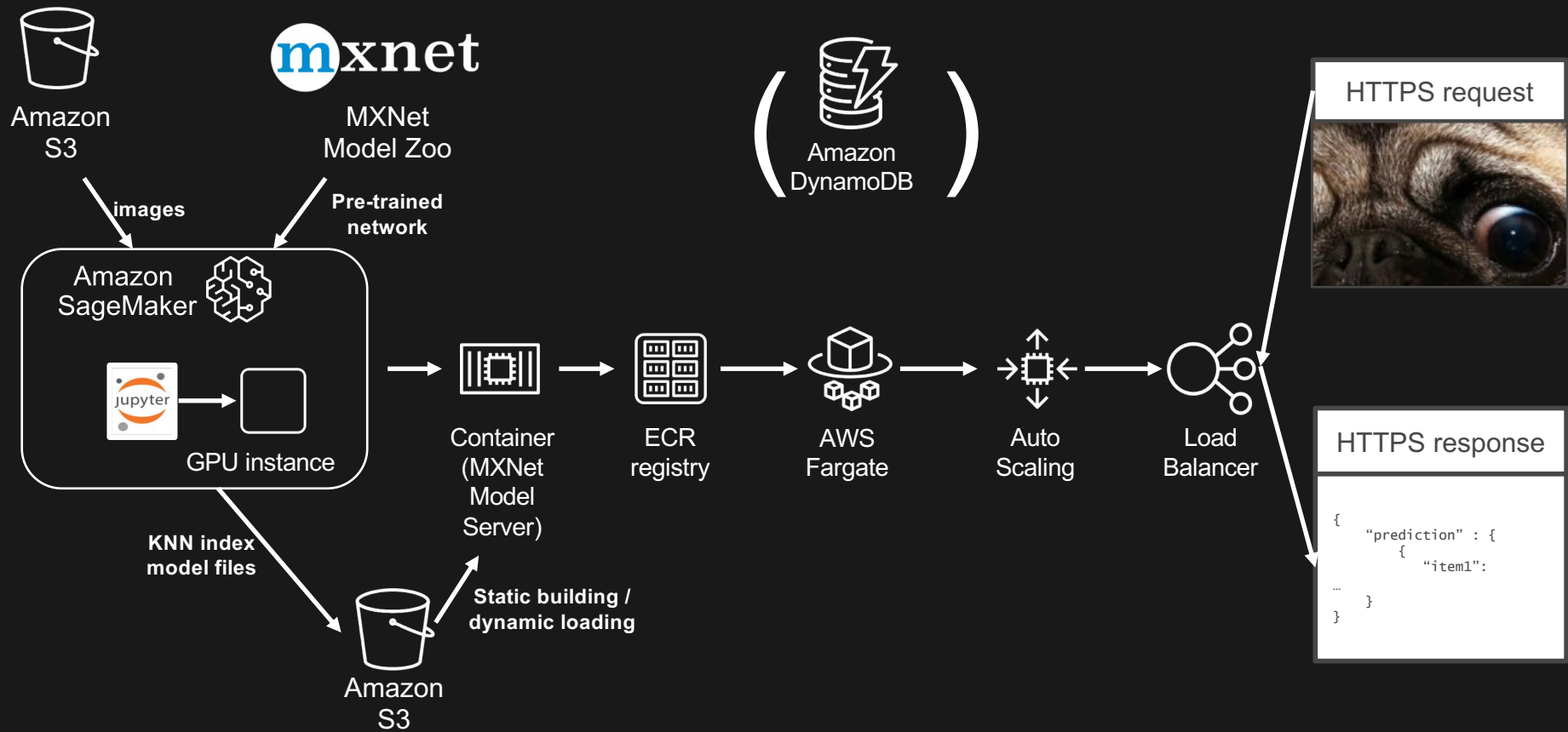


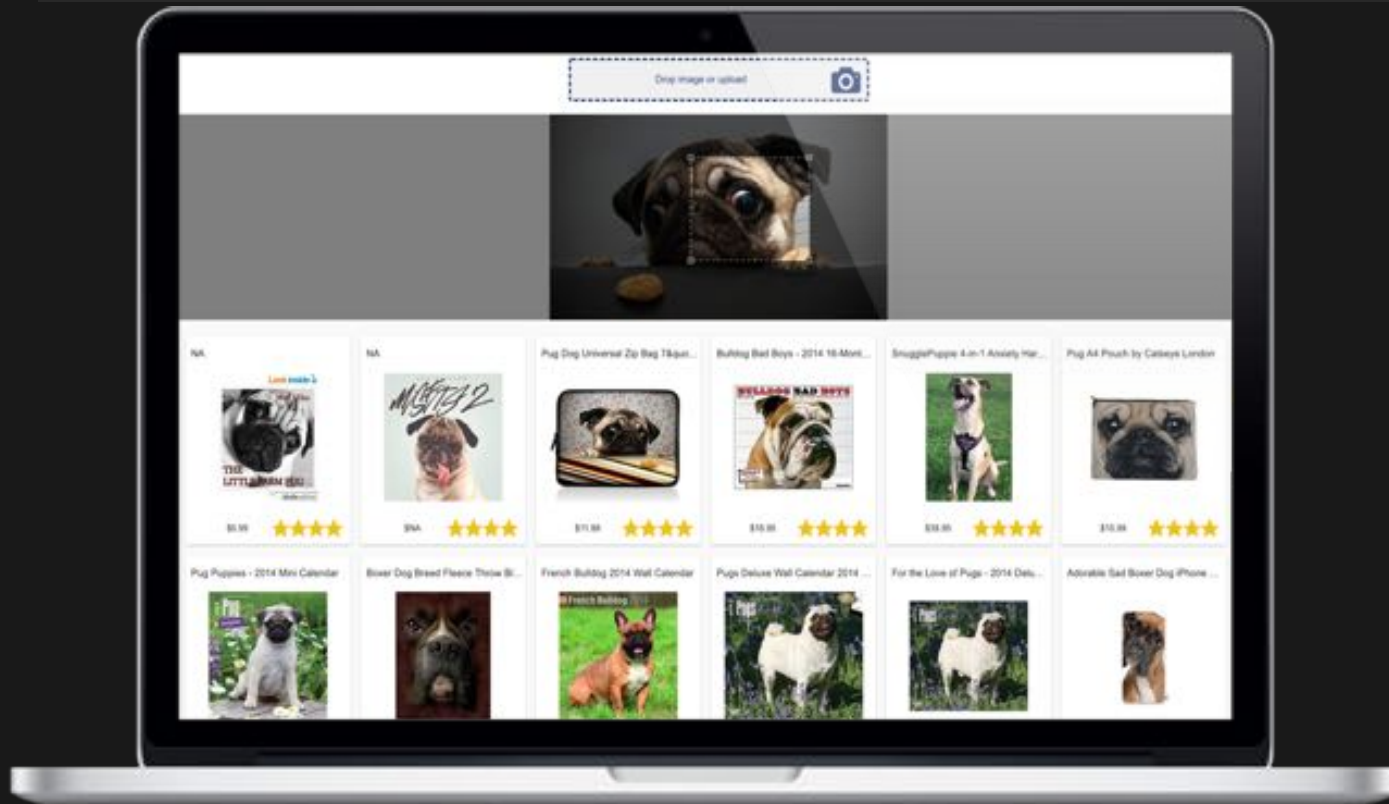224x224x3 dimensions     pre-trained CNN     512 dimensions "finger-print"     k-Nearest Neighbor (k=3)
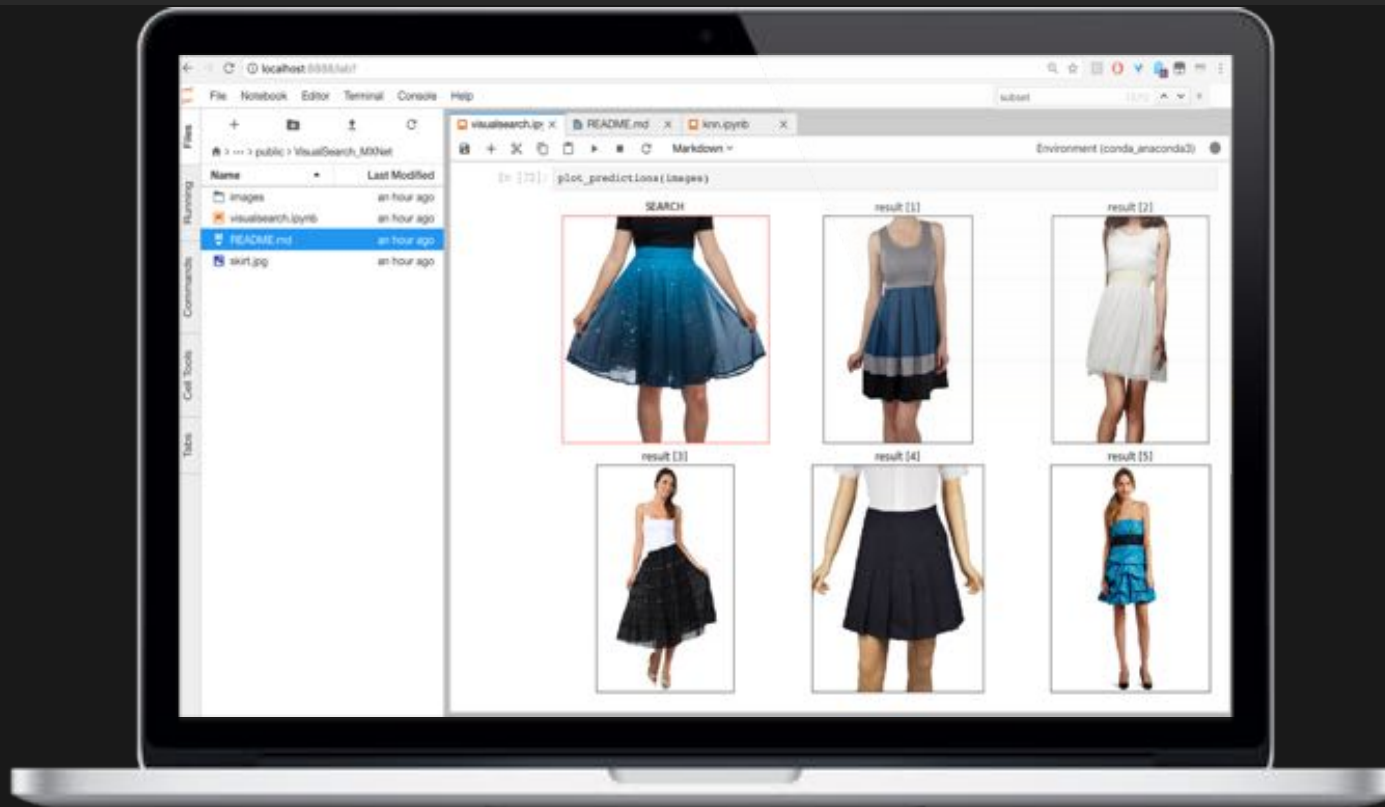
# Implementation

# Workflow and operationalization

Amazon S3

MXNet Model Zoo

Amazon DynamoDB

Amazon SageMaker

jupyter → GPU instance

images

Pre-trained network

KNN index model files

Amazon S3

Static building / dynamic loading

Container (MXNet Model Server)

ECR registry

AWS Fargate

Auto Scaling

Load Balancer

HTTPS request

HTTPS response

```
{
    "prediction" : {
        {
            "item1":
    …    }
}
```
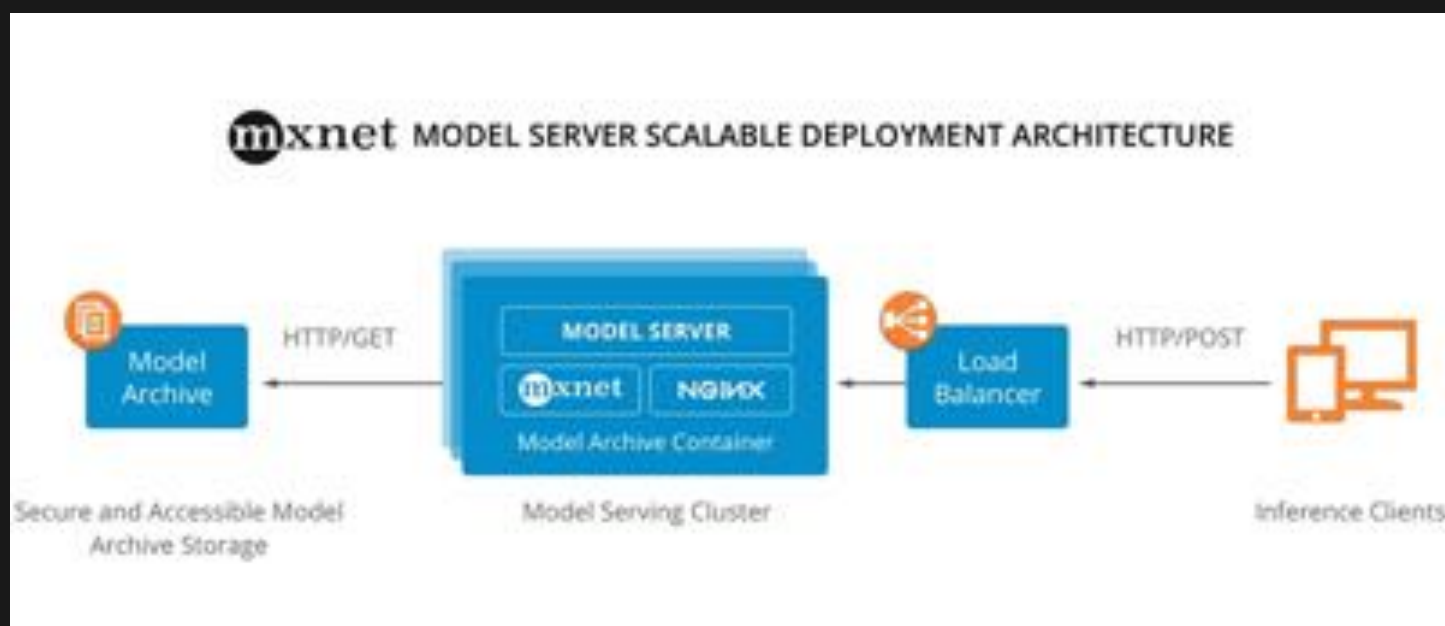
# Demo

https://thomasdelteil.github.io/VisualSearch_MXNet/

# Code

https://github.com/ThomasDelteil/VisualSearch_MXNet

# Model Serving

# MXNet Model Server Architecture

# Amazon SageMaker Neo – a path to IoT

- Model compiler supporting various hardware platforms incl. Intel, Nvidia, Arm, Cadence, Qualcomm, and Xilinx

- Supporting Tensorflow, MXNet and ONNX

- https://tvm.ai/

# Thank you!

Connect here
github.com/srochel
linkedin.com/in/steffenrochel
steroche@amazon.com

mxnet-info@amazon.com