

# QCon

全球软件开发大会【上海站】

# Topic of your presentation

Subtitle / Your Name

# AWS上的MXNet 深度学习框架

## Getting start MXNet on AWS

费良宏, AWS Technical Evangelist



[lianghon@amazon.com](mailto:lianghon@amazon.com)

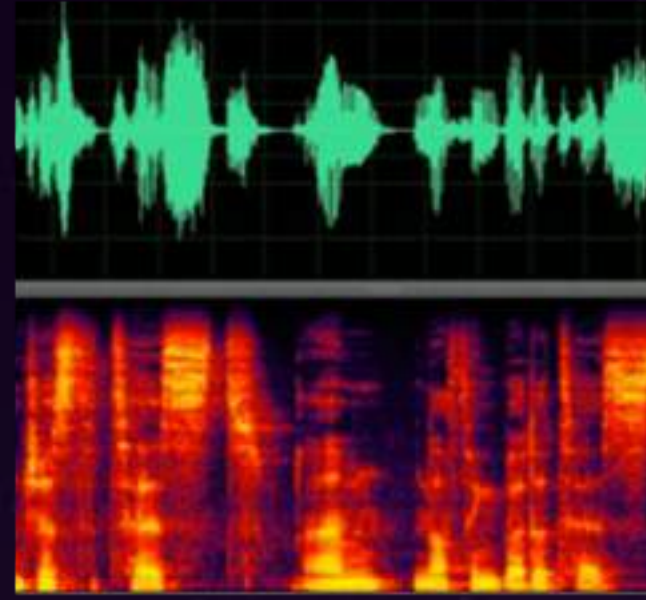


# 深度学习的发展

图像



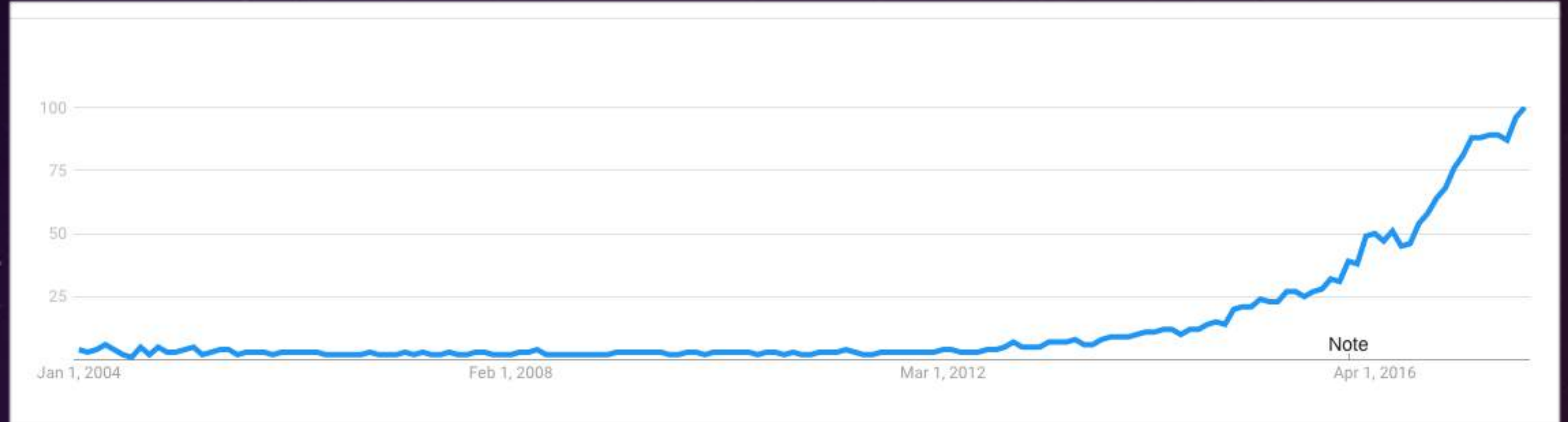
声音



自然语言处理



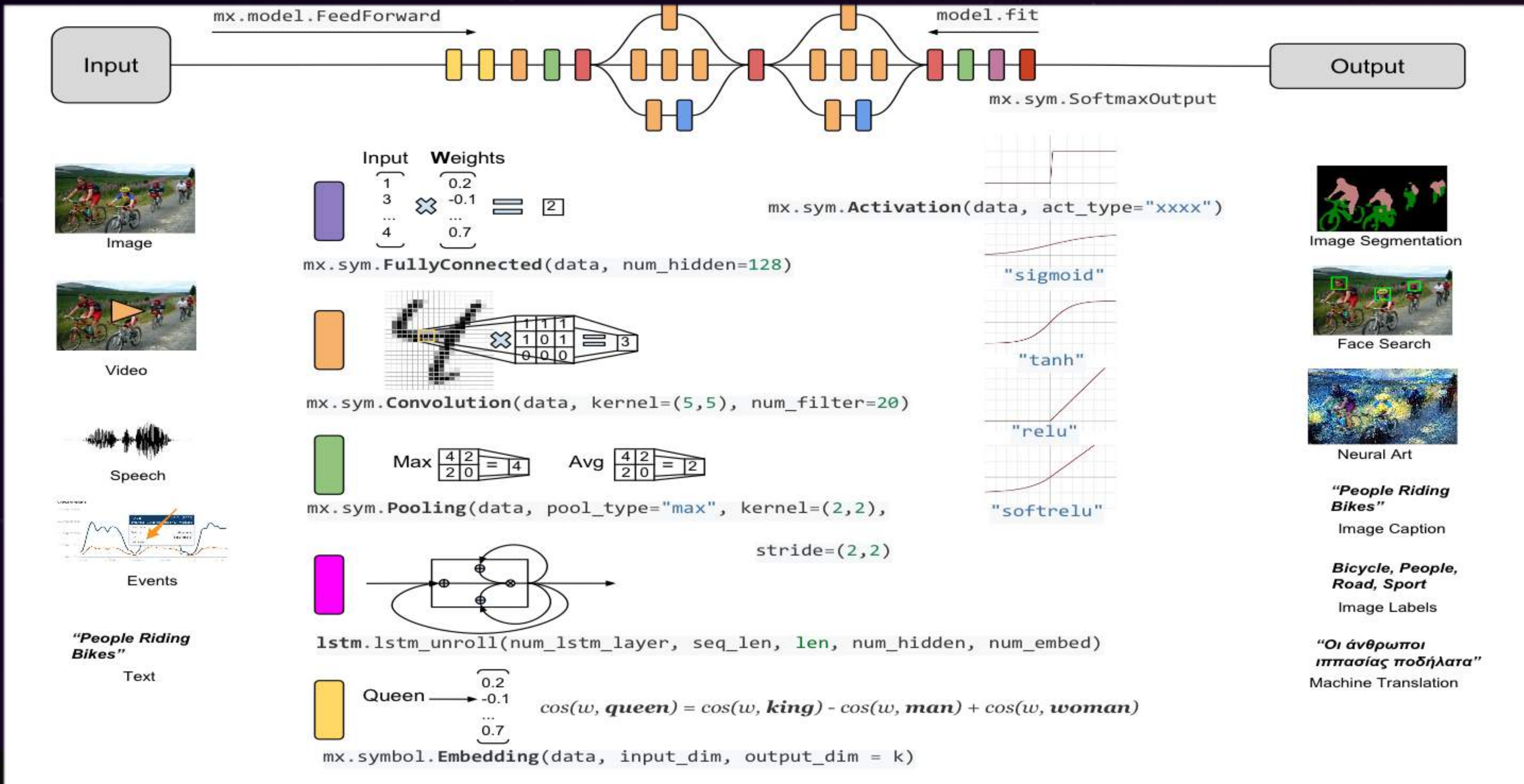
机器自主



来源：<https://trends.google.com/>

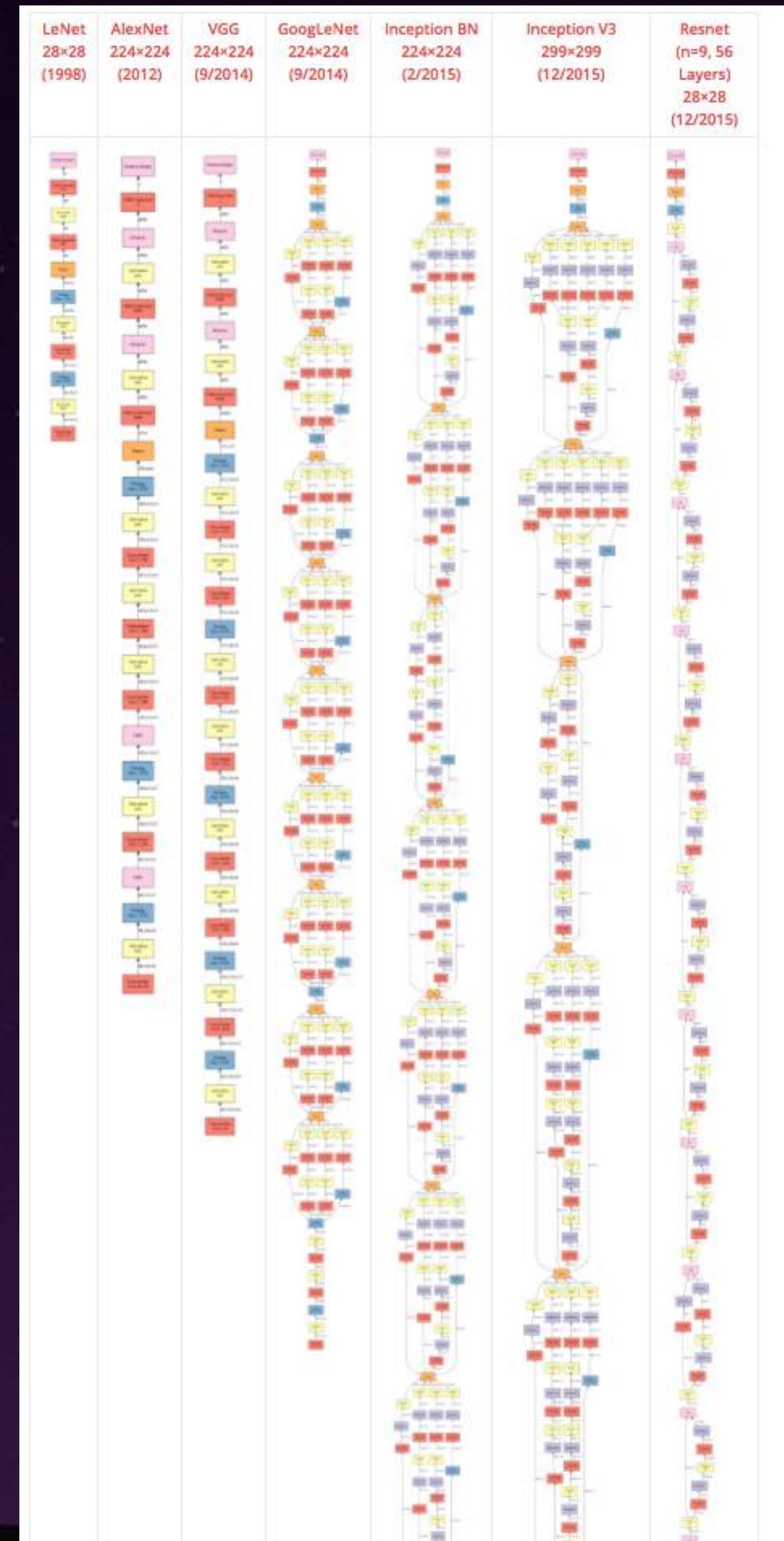


# 深度学习的模型



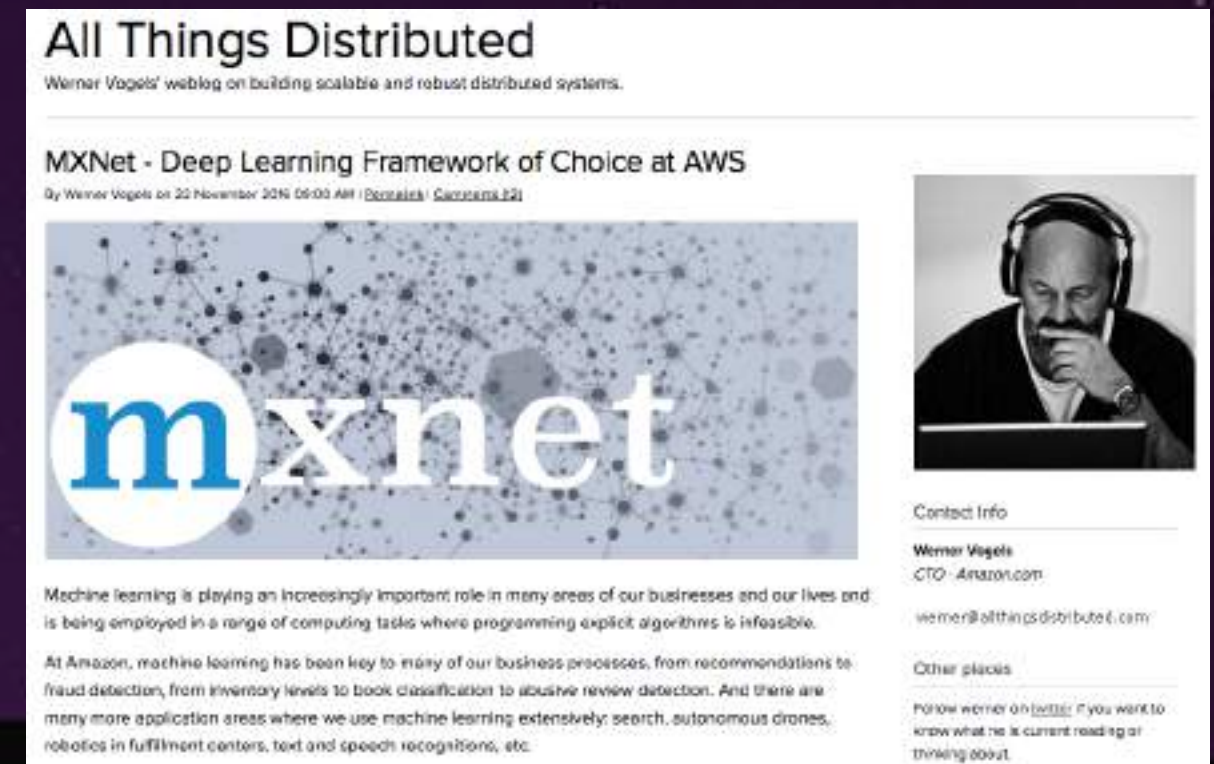
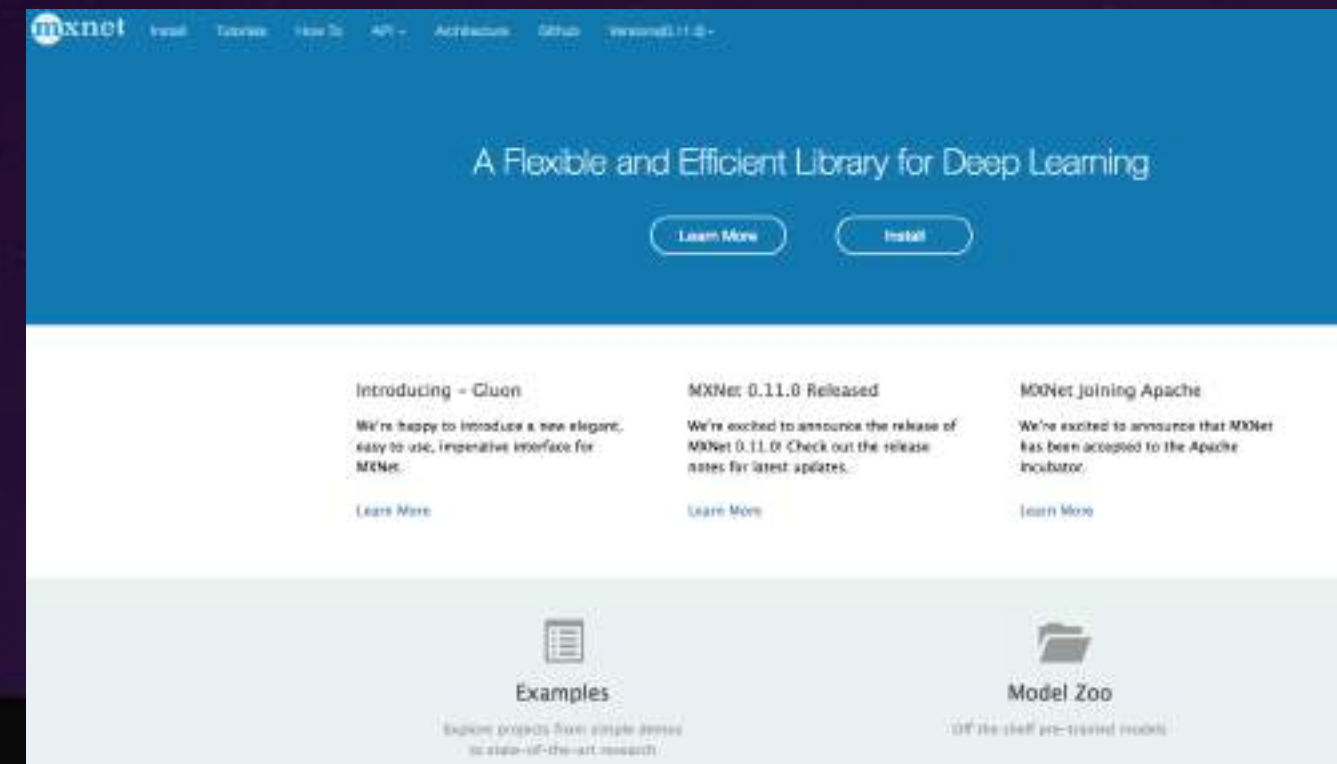
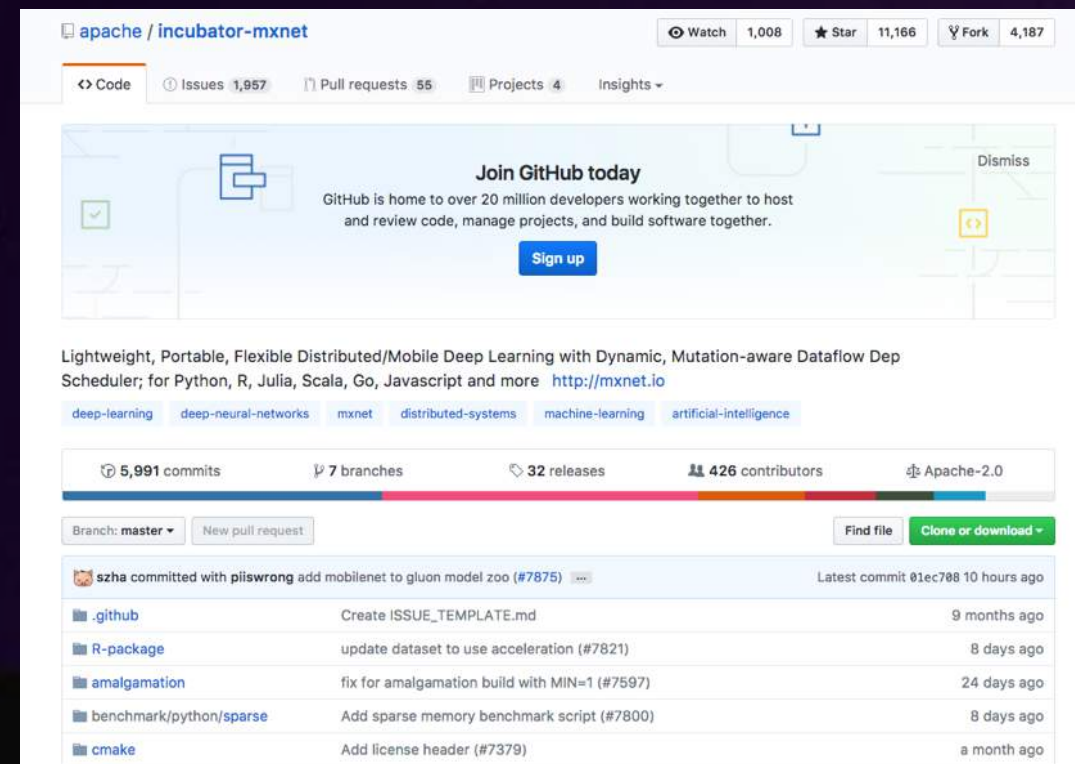
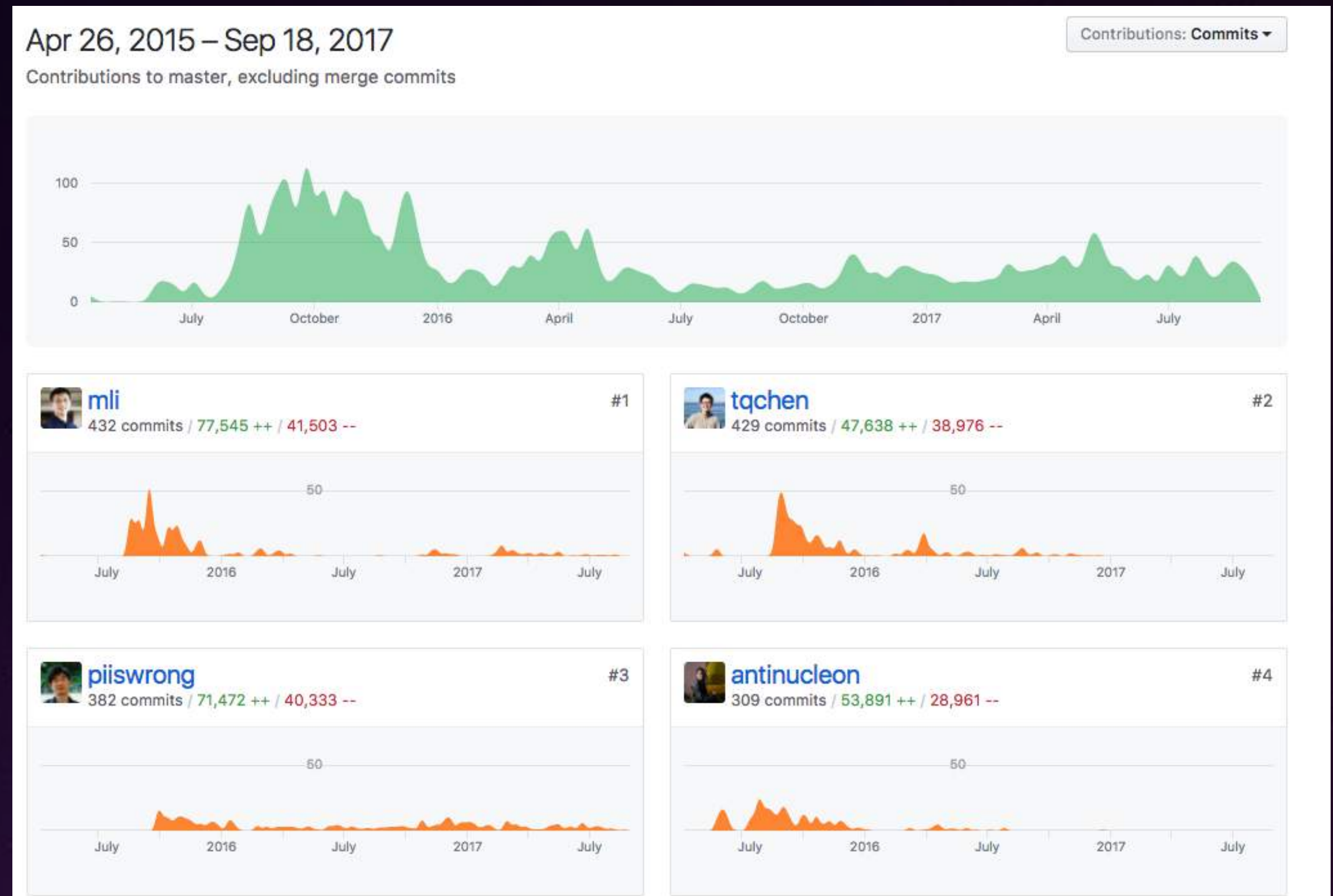
# 深度学习的模型

```
1 #!/usr/bin/env python3
2 # -*- coding:utf-8 -*-
3 # Writen by Lianghong 2017-09-18 12:17:33
4
5 import find_mxnet
6 import mxnet as mx
7 import importlib
8
9 name = "inception-v3"
10 net = importlib.import_module("symbol_" + name).get_symbol(2)
11 a = mx.viz.plot_network(
12     net, shape={"data":(1, 1, 299, 299)},
13     node_attrs={"shape":'rect',"fixedsize":'false'})
14 a.render(name)
```



# 关于 Apache MXNet

- 2015 年项目启动
- 2016 年11月 被 Amazon 选中
- 2017 年 1 月 成为 Apache 孵化项目
- 当前版本：**0.11**
- 代码贡献者：428 人
- Star 数量：11,262
- 许可证：Apache 2.0



# Apache MXNet 的新版本 0.11

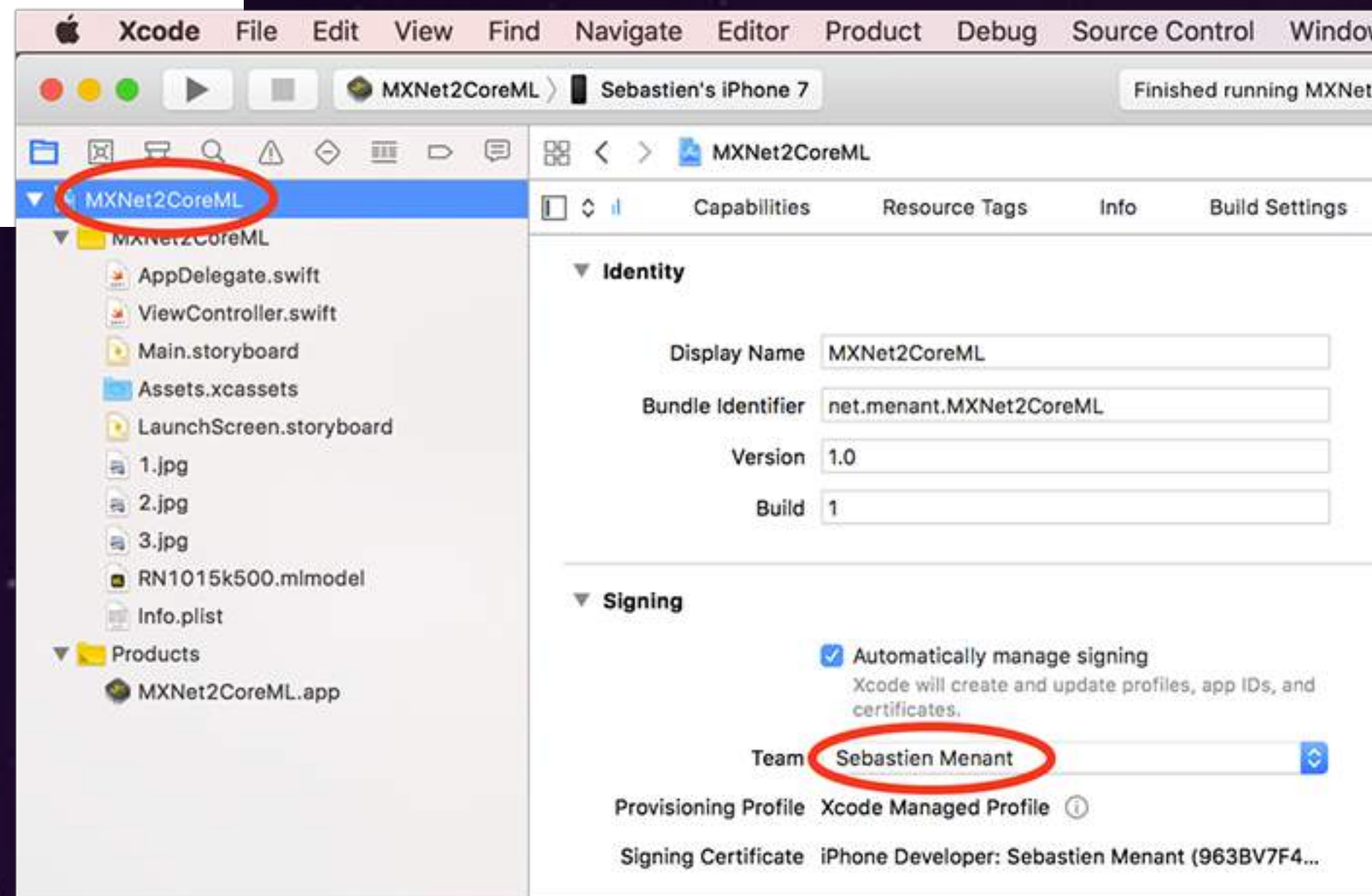
## MXNet 0.11.0

nswamy released this 20 days ago · 200 commits to master since this release

0.11.0

### Major Features

- [Apple Core ML model converter](#)
- [Support for Keras v1.2.2](#)
- [Gluon Interface \(experimental\)](#)
- Updated [LICENSE](#) and [NOTICE](#) files.
- For more information see [full release notes](#)



Holy moly! MXNet is **60% faster**: 25 seconds per epoch instead of 61. Very nice. In the same time frame, this would definitely allow us to try more things, like different model architectures or different hyper parameters. Definitely an advantage when you're experimenting.

What about memory usage? As we can see, MXNet uses over **90% less RAM** and there is plenty left for other jobs.

Julien Simon, <https://medium.com/@julsimon/keras-shoot-out-tensorflow-vs-mxnet-51ae2b30a9c0>

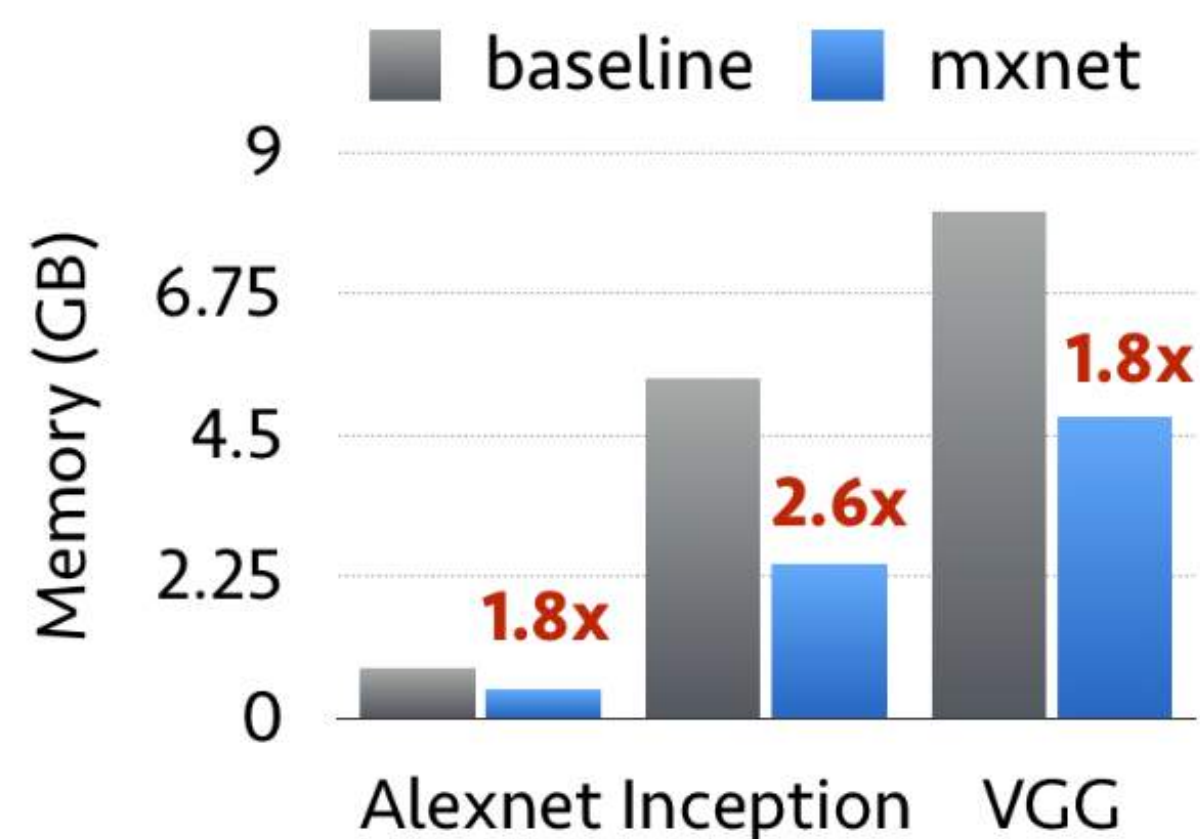


# Apache MXNet 框架特点

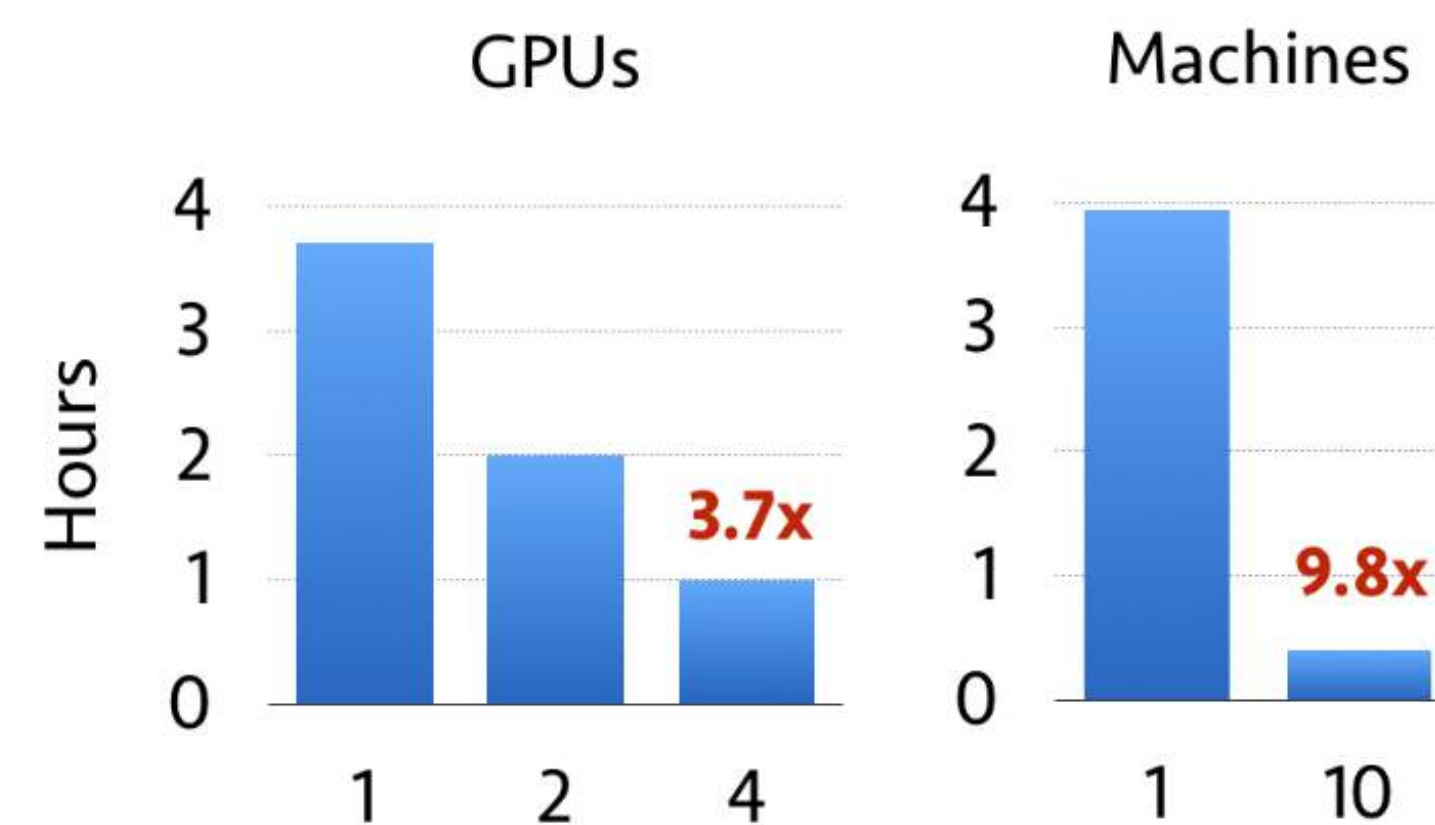
## Portable



## Efficient



## Scalable





# Apache MXNet 框架特点 - 广泛的支持

```
import picamera
from picamera.array import PiRGBArray
import inception_predict
import pollyApi as pol
import os

filename = '/tmp/capture.jpg'

def init():
    polly = pol.pollyInit()
    camera = picamera.PiCamera()
    camera.resolution=(640,480)
    return polly,camera

def capture(polly,camera):
    camera.capture(filename)
    camera.close()

    topn = inception_predict.predict_from_local_file(filename, N=5)
    print(topn)
    top1 = topn[0]
    # Convert probability to integer percentage
    prob = (str)((int)(top1[0]*100))
    # Remove category number
    item = top1[1].split(' ')
    item = ' '.join(item[1:])
    message = "I'm "+prob+"% sure that this is a "+item+"."
    os.remove(filename)
    return message

polly,cam=init()
message=capture(polly,cam)
print(message)
pol.pollySpeak(polly,message)
```

```
1 #!/usr/bin/env python3
2 # -*- coding: utf-8 -*-
3
4 import os
5 import boto3
6 import pygame
7
8 pollyRegion = 'us-east-1'
9 pollyUrl = 'https://polly.us-east-1.amazonaws.com'
10
11 def pollyInit():
12     pygame.mixer.init()
13     polly = connectToPolly()
14     return polly
15
16 def connectToPolly(regionName=pollyRegion, endpointUrl=pollyUrl):
17     return boto3.client('polly', region_name=regionName, endpoint_url=endpointUrl)
18
19 def pollySpeak(polly, text, format='mp3', voice='Amy'):
20     resp = polly.synthesize_speech(OutputFormat=format, Text=text, VoiceId=voice)
21     soundfile = open('/tmp/sound.mp3', 'wb')
22     soundBytes = resp['AudioStream'].read()
23     soundfile.write(soundBytes)
24     soundfile.close()
25     pygame.mixer.music.load('/tmp/sound.mp3')
26     pygame.mixer.music.play()
27     while pygame.mixer.music.get_busy():
28         pygame.time.Clock().tick(10)
29     os.remove('/tmp/sound.mp3')
```

```
1 import mxnet as mx
2 import numpy as np
3 import cv2, os, urllib, time
4 from collections import namedtuple
5 Batch = namedtuple('Batch', ['data'])
6
7 # Load the symbols for the networks
8 with open('synset.txt', 'r') as f:
9     synsets = [l.rstrip() for l in f]
10
11 # Load the network parameters
12 sym, arg_params, aux_params = mx.model.load_checkpoint('Inception-BN', 0)
13
14 # Load the network into an MXNet module and bind the corresponding parameters
15 mod = mx.mod.Module(symbol=sym, context=mx.cpu())
16 mod.bind(for_training=False, data_shapes=[('data', (1,3,224,224))])
17 mod.set_params(arg_params, aux_params)
18
19 '''
20 Function to predict objects by giving the model a pointer to an image file and running a forward pass through the model.
21
22 inputs:
23 filename = jpeg file of image to classify objects in
24 mod = the module object representing the loaded model
25 synsets = the list of symbols representing the model
26 N = Optional parameter denoting how many predictions to return (default is top 5)
27
28 outputs:
29 python list of top N predicted objects and corresponding probabilities
30 '''
31 def predict(filename, mod, synsets, N=5):
32     tic = time.time()
33     img = cv2.cvtColor(cv2.imread(filename), cv2.COLOR_BGR2RGB)
34     if img is None:
35         return None
36     img = cv2.resize(img, (224, 224))
37     img = np.swapaxes(img, 0, 2)
38     img = np.swapaxes(img, 1, 2)
39     img = img[np.newaxis, :]
40     print("pre-processed image in "+str(time.time()-tic))
41
42     toc = time.time()
43     mod.forward(Batch([mx.nd.array(img)]))
44     prob = mod.get_outputs()[0].asnumpy()
45     prob = np.squeeze(prob)
46     print("forward pass in "+str(time.time()-toc))
47
48     topN = []
49     a = np.argsort(prob)[::-1]
50     for i in a[0:N]:
51         print('probability=%f, class=%s' %(prob[i], synsets[i]))
52         topN.append((prob[i], synsets[i]))
53     return topN
54
55
56
57 # Code to download an image from the internet and run a prediction on it
58 def predict_from_url(url, N=5):
59     filename = url.split("/")[-1]
60     urllib.urlretrieve(url, filename)
61     img = cv2.imread(filename)
62     if img is None:
63         print("Failed to download")
64     else:
65         return predict(filename, mod, synsets, N)
66
67 # Code to predict on a local file
68 def predict_from_local_file(filename, N=5):
69     return predict(filename, mod, synsets, N)
```



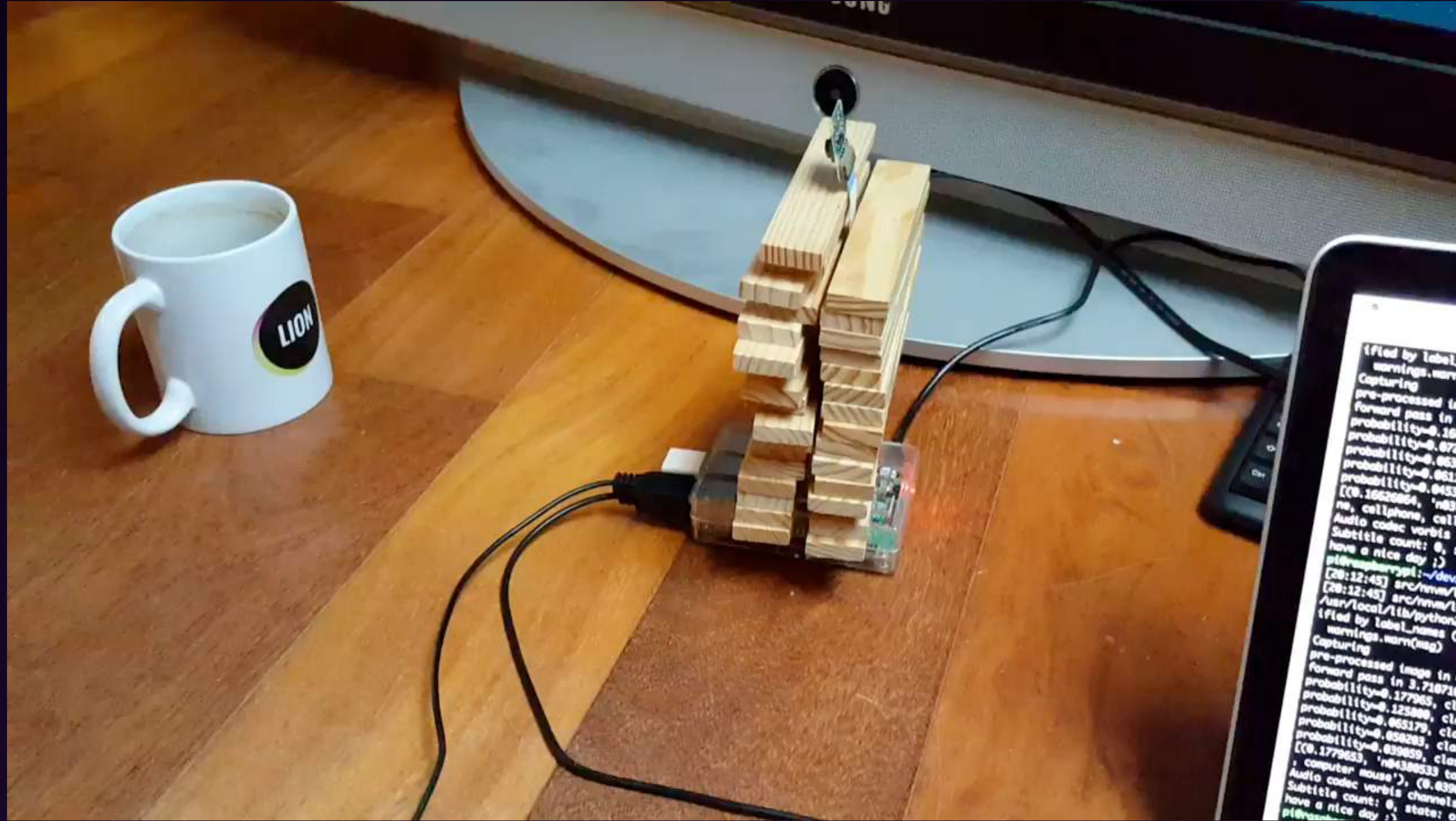
mxnet



Amazon AI

Intelligent Services Powered By Deep Learning

# 深度学习的模型



**mxnet**

 **Amazon AI**

Intelligent Services Powered By Deep Learning



# Apache MXNet 框架特点 – 融合 ( Amalgamation )

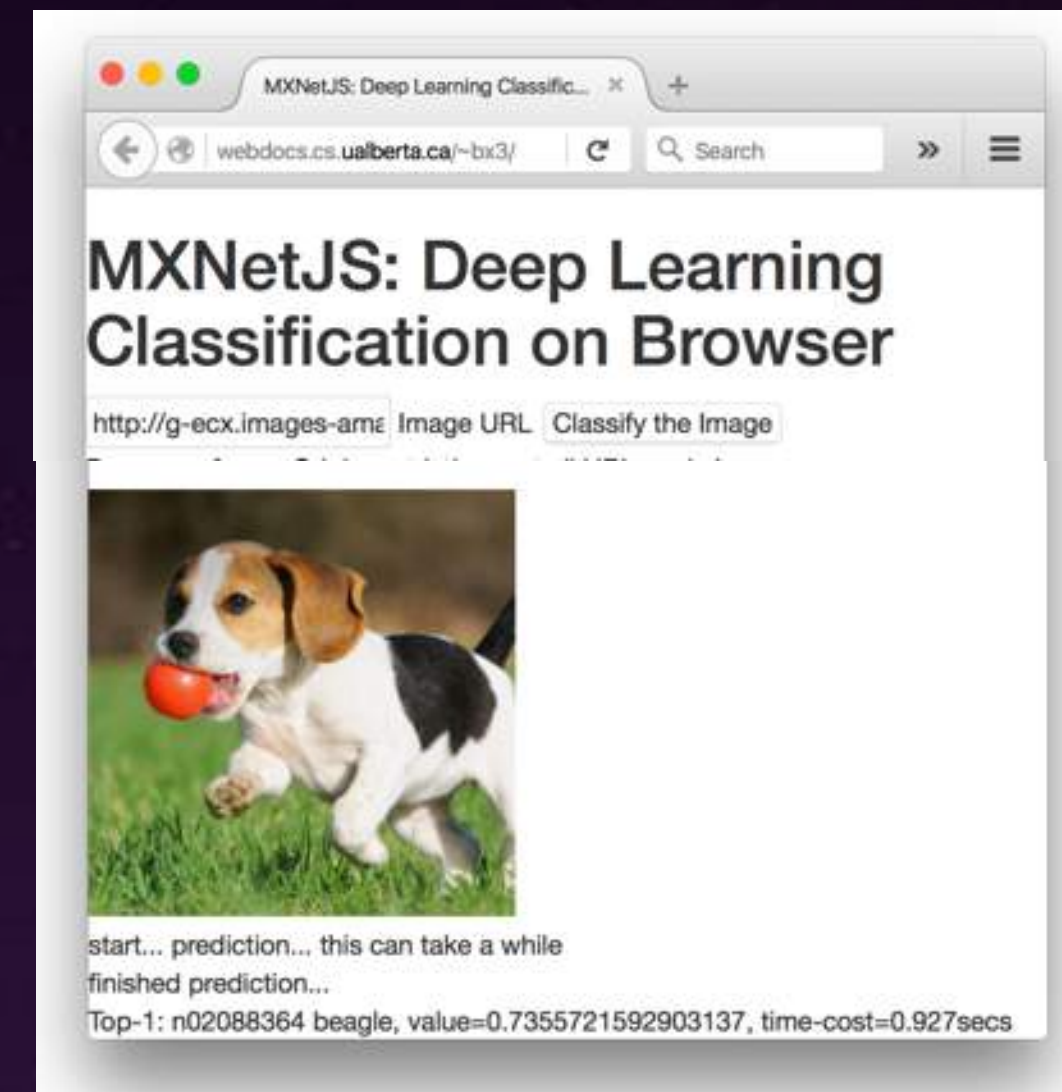
- 将所需的库合并到单独的一个C++文件中
- 可在任何平台编译



BlindTool by Joseph Paul Cohen, demo on Nexus 4



在浏览器中运行的JAVASCRIPT 代码



# Apache MXNet 框架特点 - 高效

## High-level MXNet Example

In the interest of comparison; a common (custom) data-generator (called `yield_mb(X, y, batchsize=64, shuffle=False)`) was originally used for all other frameworks - but not for MXNet. I have reproduced the MXNet example using this same generator (wrapping the results in the `mx.io.DataBatch` class) to test if MXNet is faster than other frameworks just because I was using its own data-generator. This does not appear to be the case.

```
[1]: import os
import sys
import numpy as np
import mxnet as mx
from common.params import *
from common.utils import *

/anaconda/envs/py35/lib/python3.5/site-packages/urllib3/contrib/pyopenssl.py:46: DeprecationWarning:
OpenSSL.rand is deprecated - you should use os.urandom instead
import OpenSSL.SSL
```

```
[2]: print("OS: ", sys.platform)
print("Python: ", sys.version)
print("Numpy: ", np.__version__)
print("MXNet: ", mx.__version__)
print("GPU: ", get_gpu_name())

OS: linux
Python: 3.5.2 |Anaconda custom (64-bit)| (default, Jul 2 2016, 17:53:06)
[GCC 4.4.7 20120313 (Red Hat 4.4.7-1)]
Numpy: 1.13.1
MXNet: 0.11.0
GPU: ['Tesla K80', 'Tesla K80']
```

```
[3]: def create_symbol():
    data = mx.symbol.Variable('data')
    # size = (old-size - kernel + 2*padding)/stride + 1
    # if kernel = 3, pad with 1 either side
    conv1 = mx.symbol.Convolution(data=data, num_filter=50, pad=(1,1), kernel=(3,3))
    relu1 = mx.symbol.Activation(data=conv1, act_type="relu")
    conv2 = mx.symbol.Convolution(data=relu1, num_filter=50, pad=(1,1), kernel=(3,3))
    pool1 = mx.symbol.Pooling(data=conv2, pool_type="max", kernel=(2,2), stride=(2,2))
    relu2 = mx.symbol.Activation(data=pool1, act_type="relu")
    drop1 = mx.symbol.Dropout(data=relu2, p=0.25)

    conv3 = mx.symbol.Convolution(data=drop1, num_filter=100, pad=(1,1), kernel=(3,3))
    relu3 = mx.symbol.Activation(data=conv3, act_type="relu")
    conv4 = mx.symbol.Convolution(data=relu3, num_filter=100, pad=(1,1), kernel=(3,3))
    pool2 = mx.symbol.Pooling(data=conv4, pool_type="max", kernel=(2,2), stride=(2,2))
    relu4 = mx.symbol.Activation(data=pool2, act_type="relu")
    drop2 = mx.symbol.Dropout(data=relu4, p=0.25)
```

## VGG-style CNN on CIFAR-10

DL Library	Test Accuracy (%)	Training Time (s)
Caffe2	79	149
<b>MXNet</b>	<b>77</b>	<b>149</b>
Glueon	77	157
CNTK	78	166
PyTorch	78	168
Tensorflow	78	173
Keras(CNTK)	78	200
Chainer	79	240
Keras(TF)	77	252
Lasagne(Theano)	77	253
Keras(Theano)	78	269

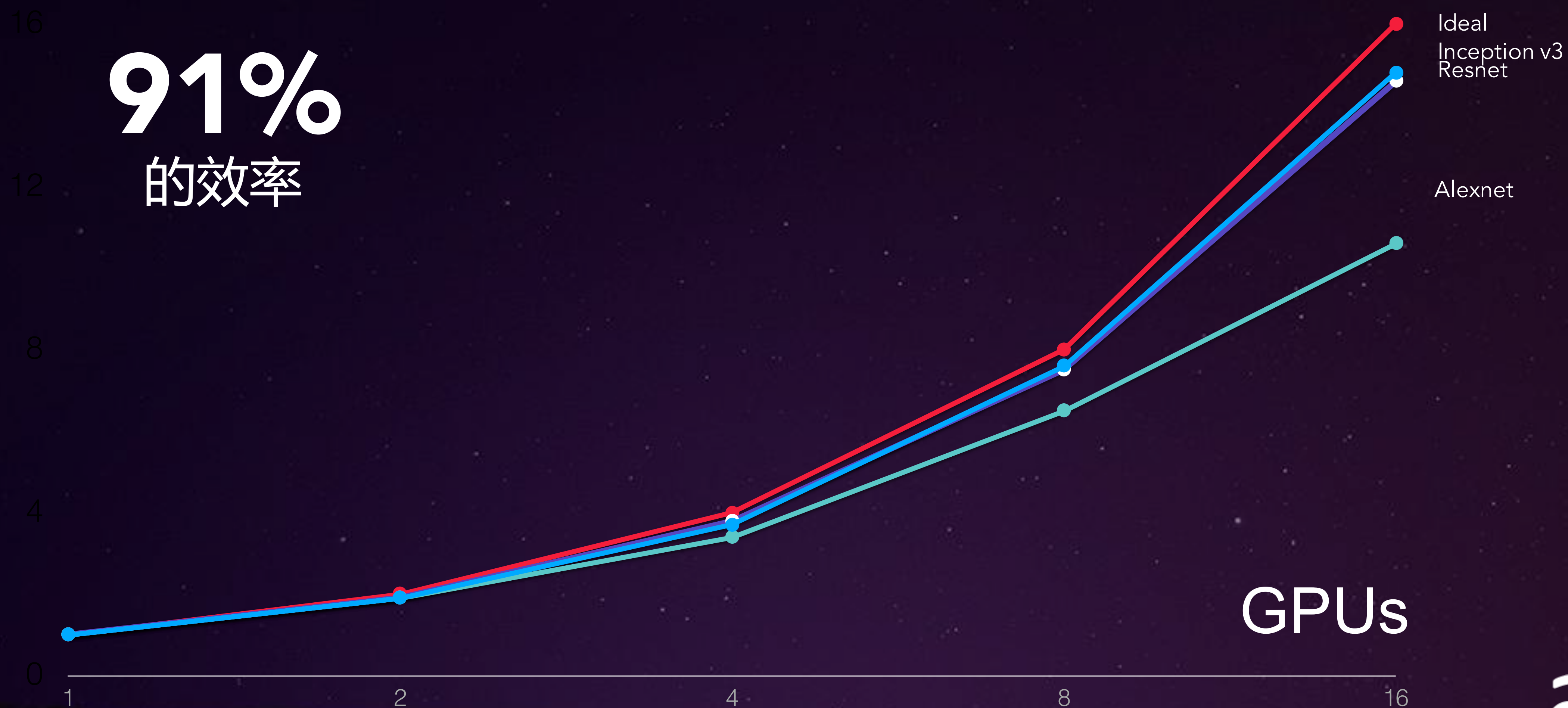
## LSTM(GRU) on IMDB

DL Library	Test Accuracy (%)	Training Time (s)
Keras(CNTK)	86	223
Tensorflow	86	79
Pytorch	87	36
<b>MXNet</b>	<b>88</b>	<b>12</b>



# Apache MXNet 框架特点 - 扩展性

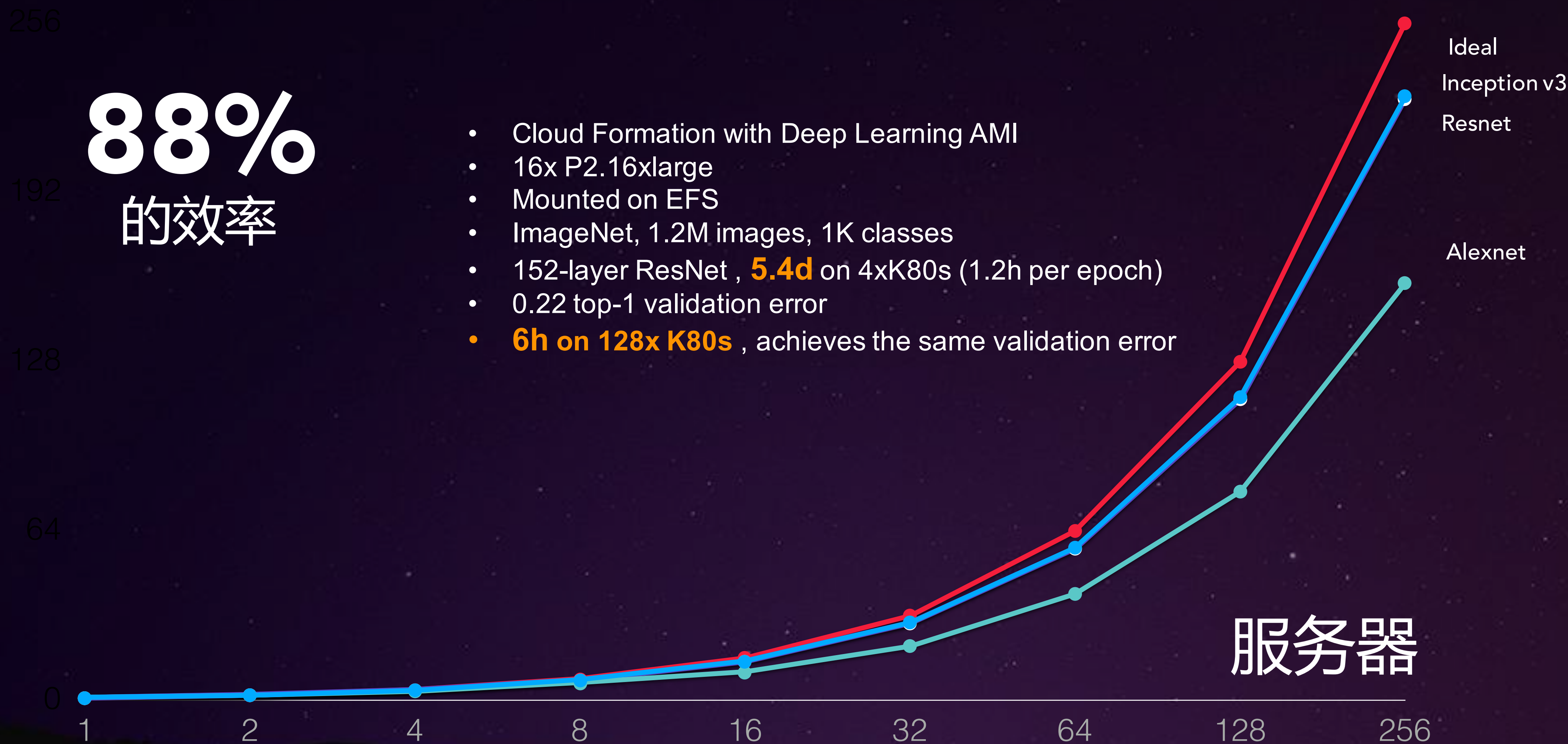
91%  
的效率



# Apache MXNet 框架特点 – 扩展性

88%  
的效率

- Cloud Formation with Deep Learning AMI
- 16x P2.16xlarge
- Mounted on EFS
- ImageNet, 1.2M images, 1K classes
- 152-layer ResNet , **5.4d** on 4xK80s (1.2h per epoch)
- 0.22 top-1 validation error
- **6h on 128x K80s** , achieves the same validation error



服务器



# Apache MXNet 框架基础

- **NDArray**: 多维数组运算 ( **imperative** )
- **Symbol**: 神经网络符号表达式 ( **declarative** )
- **Module**: 神经网络模型训练及接口
- **Loading Data**: 提供数据给训练 / 推理过程
- **混合编程**: 使用NDArray 和 Symbolic 进行模型训练



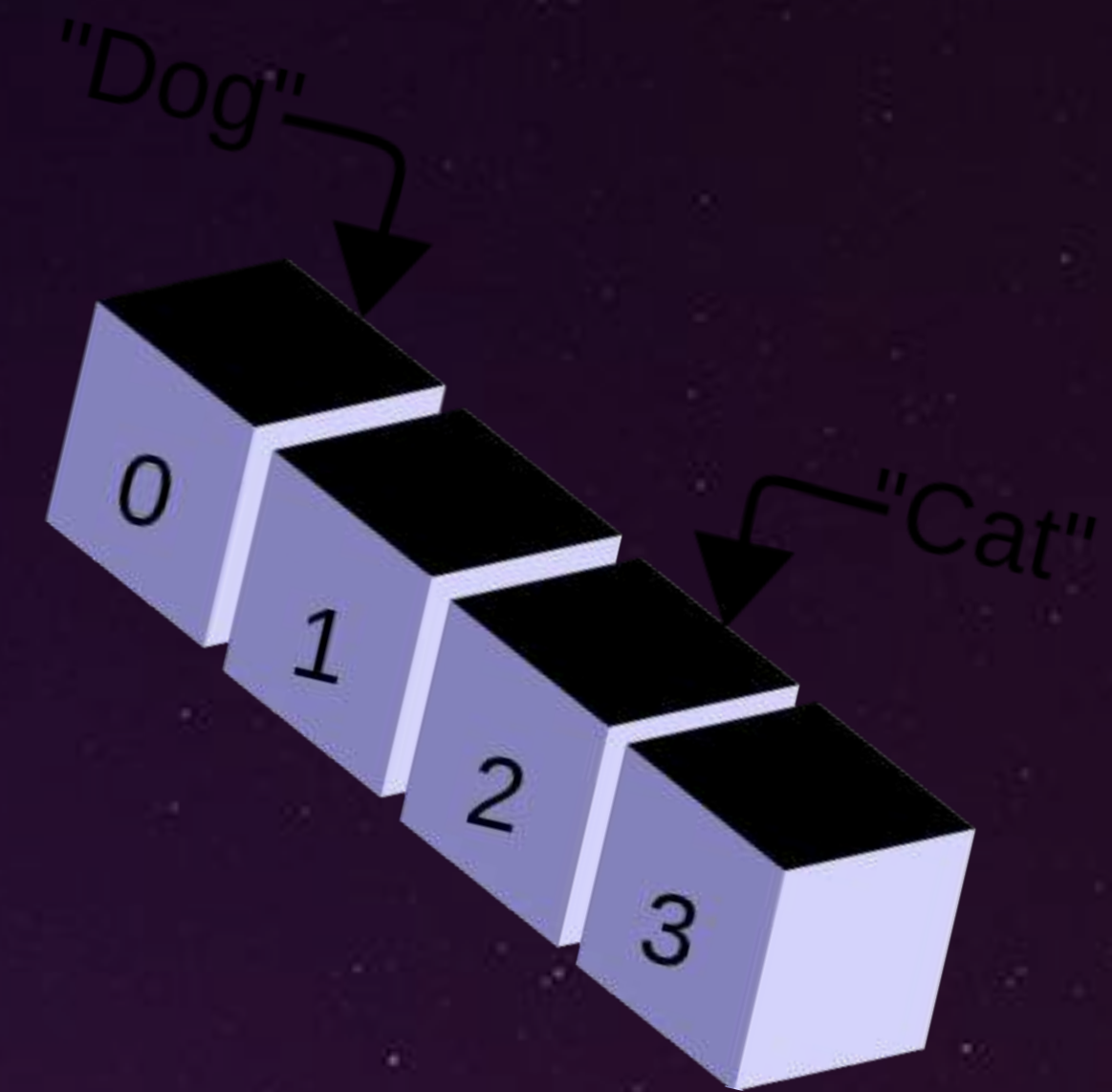
# Apache MXNet 框架基础 – NDArray

## NumPy

- 仅限于使用 CPU
- 没有自动微分 ( automatic differentiation )
- **阻塞式 ( 计算完成返回结果 )**

## NDArray

- 通过设备上下文的多个 CPUs/GPUs
- 云计算环境下的分布式环境
- 自动微分 ( automatic differentiation )
- **非阻塞式的 Python 前端 ( 无 GIL ) & 懒惰评估 ( lazy evaluation )**





# Apache MXNet 框架基础 – NDArray

context = mx.cpu()



**注意！  
NDArray 在GPU上 是不可变的！**

context = mx.gpu(0)  
context = mx.gpu(1)



g = copyto(c)  
g = c.as\_in\_context(mx.gpu(0))



# 命令式编程 ( Imperative )

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

易于针对Python代码进行调整

## 优点

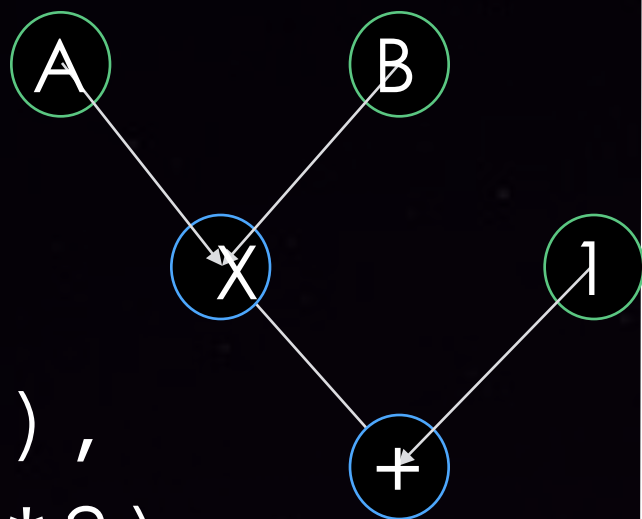
- 简单直接并且灵活
- 可利用程序语言自身的特点 ( 循环、条件判断、调试器 )
- 例如 : Numpy、Matlab,以及 Torch

## 缺点

- 难于优化

# 声明式编程 (Declarative)

```
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可以与D共享内存，因为C稍后会被删除

## 优点

- 更多优化的机会
- 可以适用多种程序语言
- 例如：TensorFlow、Theano 以及 Caffe

## 缺点

- 灵活性不足
- 难于调试
- 无法使用原生的程序代码
- 无动态图



# 混合编程的范式

## 命令式 NDARRAY API

```
>>> import mxnet as mx
>>> a = mx.nd.zeros((100, 50))
>>> b = mx.nd.ones((100, 50))
>>> c = a + b
>>> c += 1
>>> print(c)
```

## 声明式 SYMBOLIC EXECUTOR

```
>>> import mxnet as mx
>>> net = mx.symbol.Variable('data')
>>> net = mx.symbol.FullyConnected(data=net, num_hidden=128)
>>> net = mx.symbol.SoftmaxOutput(data=net)
>>> texec = mx.module.Module(net)
>>> texec.forward(data=c)
>>> texec.backward()
```

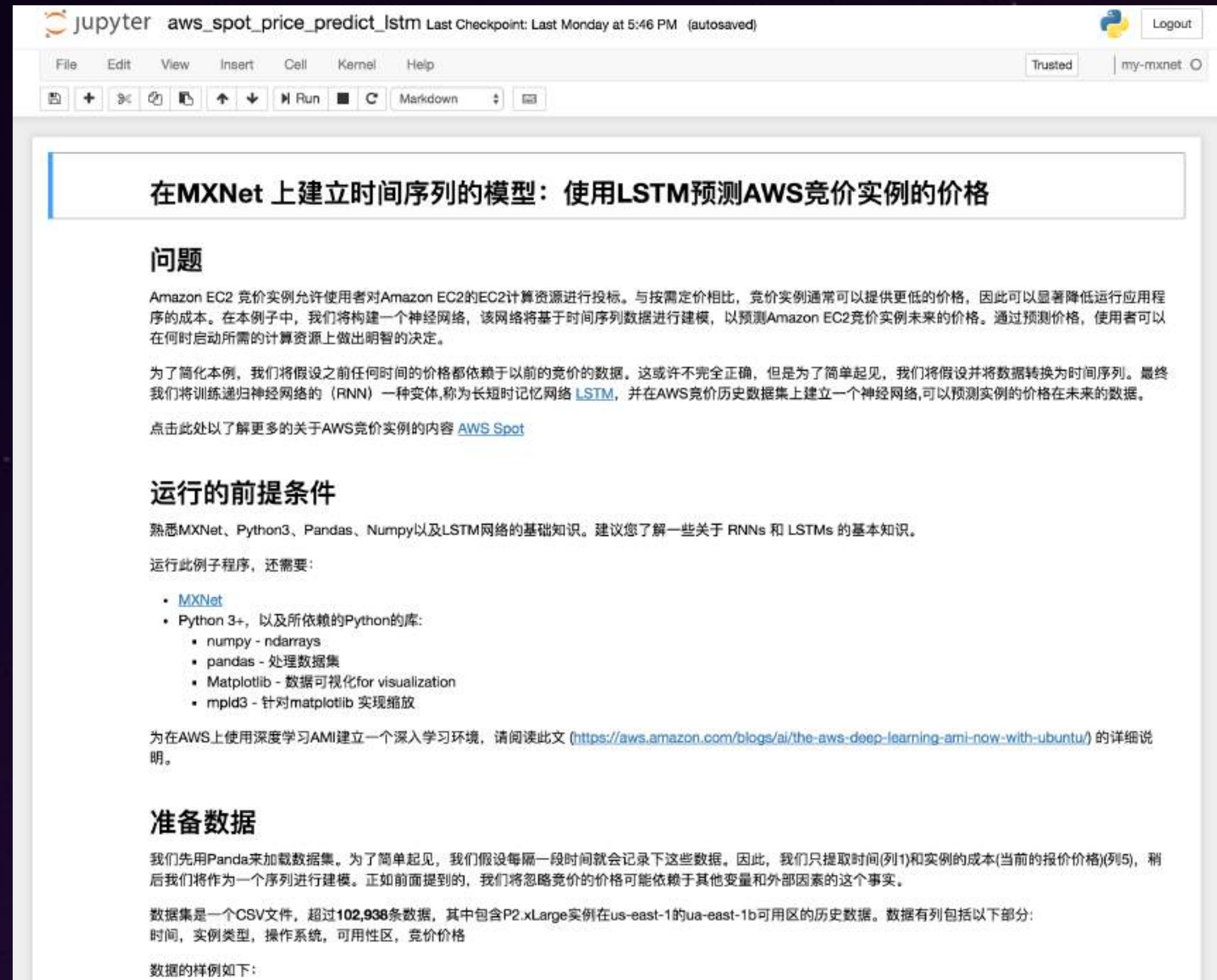
NDArray 可以用来作为 Graph 的输入

# Apache MXNet 应用演示 – AWS EC2 竞价实例价格预测

**AWS EC2竞价型实例** - 客户可以利用竞价型实例对未使用的 EC2 容量进行竞价，并可在竞价高于当前现货价格的期间内运行此类实例。现货价格根据供需情况定期变化，出价达到或超过竞价的客户可获得相应的竞价型实例。



# Apache MXNet 应用演示 – AWS EC2 竞价实例价格预测



jupyter aws\_spot\_price\_predict\_lstm Last Checkpoint: Last Monday at 5:46 PM (autosaved) Logout

File Edit View Insert Cell Kernel Help Trusted my-mxnet

## 在MXNet 上建立时间序列的模型：使用LSTM预测AWS竞价实例的价格

### 问题

Amazon EC2 竞价实例允许使用者对Amazon EC2的EC2计算资源进行投标。与按需定价相比，竞价实例通常可以提供更低的价格，因此可以显著降低运行应用程序的成本。在本例子中，我们将构建一个神经网络，该网络将基于时间序列数据进行建模，以预测Amazon EC2竞价实例未来的价格。通过预测价格，使用者可以在何时启动所需的计算资源上做出明智的决定。

为了简化本例，我们将假设之前任何时间的价格都依赖于以前的竞价的数据。这或许不完全正确，但是为了简单起见，我们将假设并将数据转换为时间序列。最终我们将训练递归神经网络的（RNN）一种变体，称为长短期记忆网络 [LSTM](#)，并在AWS竞价历史数据集上建立一个神经网络，可以预测实例的价格在未来的数据。

[点击此处](#)以了解更多的关于AWS竞价实例的内容 [AWS Spot](#)

### 运行的前提条件

熟悉MXNet、Python3、Pandas、Numpy以及LSTM网络的基础知识。建议您了解一些关于 RNNs 和 LSTMs 的基本知识。

运行此例子程序，还需要：

- [MXNet](#)
- Python 3+，以及所依赖的Python的库：
  - numpy - ndarrays
  - pandas - 处理数据集
  - Matplotlib - 数据可视化for visualization
  - mpld3 - 针对matplotlib 实现缩放

为在AWS上使用深度学习AMI建立一个深入学习环境，请阅读此文 (<https://aws.amazon.com/blogs/ai/the-aws-deep-learning-ami-now-with-ubuntu/>) 的详细说明。

### 准备数据

我们先用Panda来加载数据集。为了简单起见，我们假设每隔一段时间就会记录下这些数据。因此，我们只提取时间(列1)和实例的成本(当前的报价价格)(列5)，稍后我们将作为一个序列进行建模。正如前面提到的，我们将忽略竞价的价格可能依赖于其他变量和外部因素的这个事实。

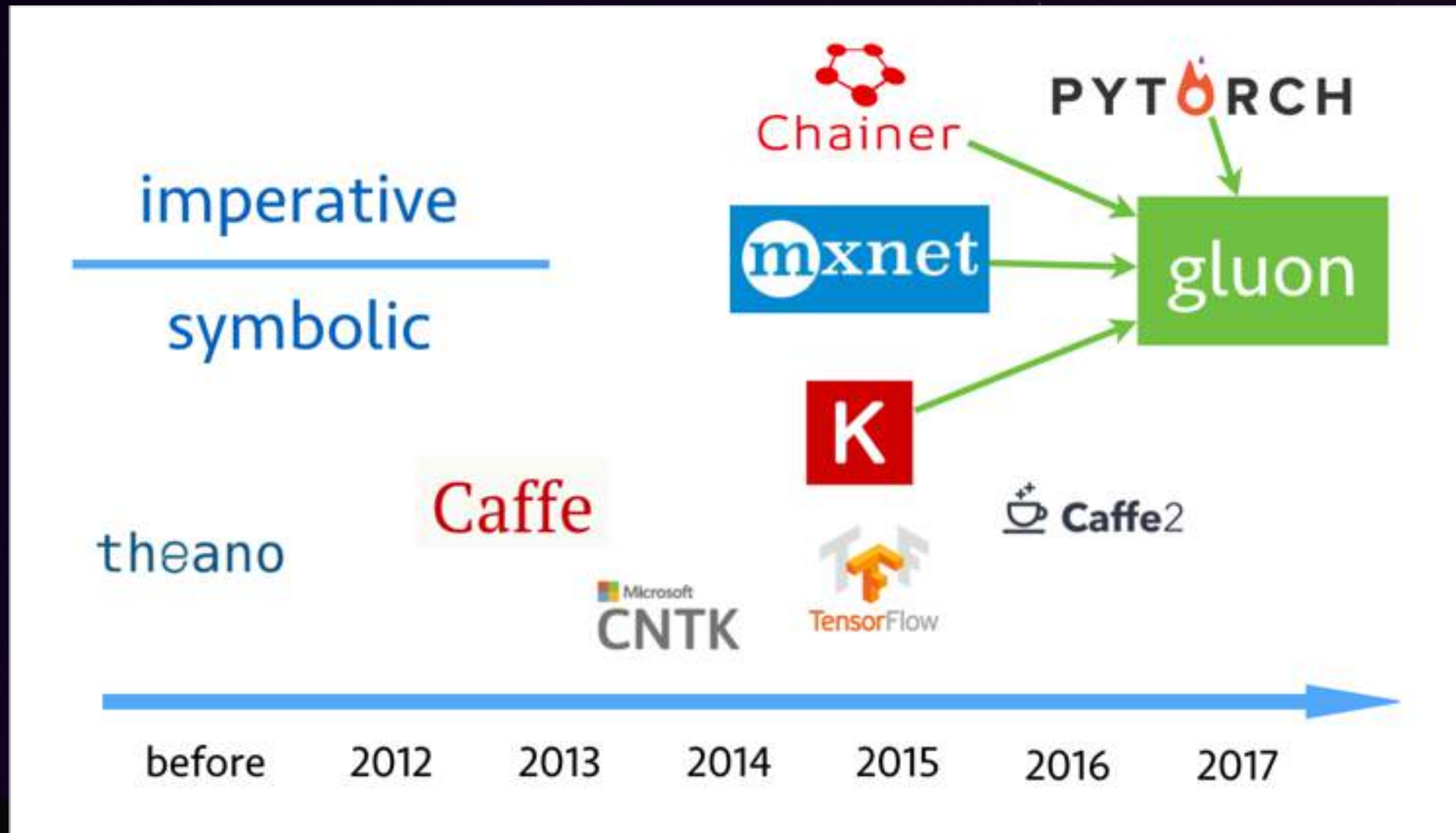
数据集是一个CSV文件，超过**102,938**条数据，其中包含P2.xLarge实例在us-east-1的ua-east-1b可用区的历史数据。数据有列包括以下部分：  
时间，实例类型，操作系统，可用性区，竞价价格

数据的样例如下：

链接：[https://lianghong.info/ipython/notebooks/spotprice\\_predict/aws\\_spot\\_price\\_predict\\_lstm.ipynb](https://lianghong.info/ipython/notebooks/spotprice_predict/aws_spot_price_predict_lstm.ipynb)



# Apache MXNet 框架的新接口 – Gluon



# Apache MXNet 框架的新接口 – Gluon

```
net = gluon.nn.Sequential()
with net.name_scope():
    net.add(gluon.nn.Dense(128, activation='relu'))
    net.add(gluon.nn.Dense(64, activation='relu'))
    net.add(gluon.nn.Dense(10))

loss = gluon.loss.SoftmaxCrossEntropyLoss()

for data, label in get_batch():
    with autograd.record():
        l = loss(net(data), label)
    l.backward()
    trainer.step(batch_size=data.shape[0])
```

net.hybridize()

从声明式转换到符号执行