**Microsoft**
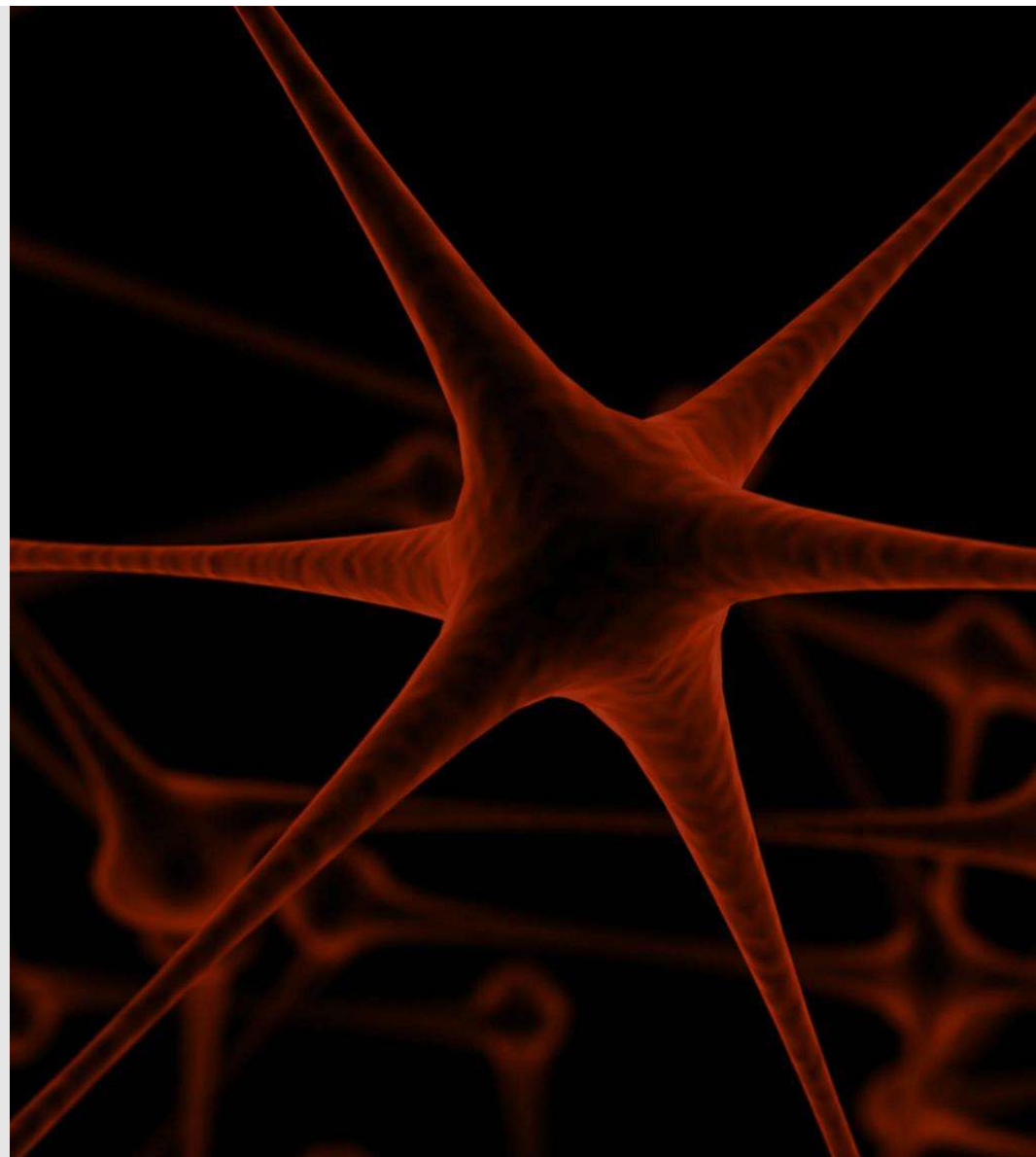
# De-mystifying Deep Learning

*Anusua Trivedi*
*Data Scientist*
*Email: antriv@microsoft.com*
*Twitter: @anurive*

# Talk Outline

❑ Deep Learning (DL)

❑ Deep Neural Networks (DNN)

❑ Types of DNNs

❑ DL Frameworks

❑ Use Cases

# Traditional ML Vs DL

Traditional ML requires manual feature extraction/engineering

Deep learning can automatically learn features in data

Feature extraction for unstructured data is very difficult

Deep learning is largely a "black box" technique, updating learned weights at each layer

# Why is DL popular?

❑ DL models has been here for a long time
- Fukushima (1980) – Neo-Cognitron
- LeCun (1989) – Convolutional Neural Network

❑ DL popularity grew recently
- With growth of Big Data
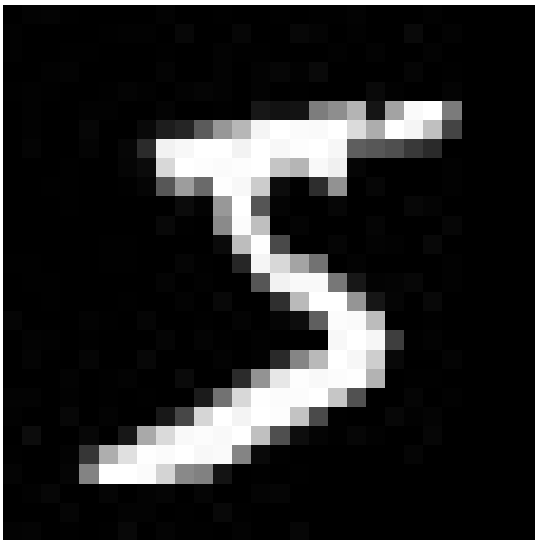- With the advent of powerful GPUs

# Deep Neural Network (DNN)

*http://neuralnetworksanddeeplearning.com/chap5.html

# Common DNNs

❑ Deep Convolutional Neural Network (DCNN)
- ▪ To extract representation from images

❑ Recurrent Neural Network (RNN)
- ▪ To extract representation from sequential data

❑ Deep Belief Neural Network (DBN)
- ▪ To extract hierarchical representation from a dataset

# Deep learning and computer vision

# Vision is hard

Vision is hard because images are big matrices of numbers.



Example from MNIST handwritten digit dataset [LeCun and Cortes, 1998].
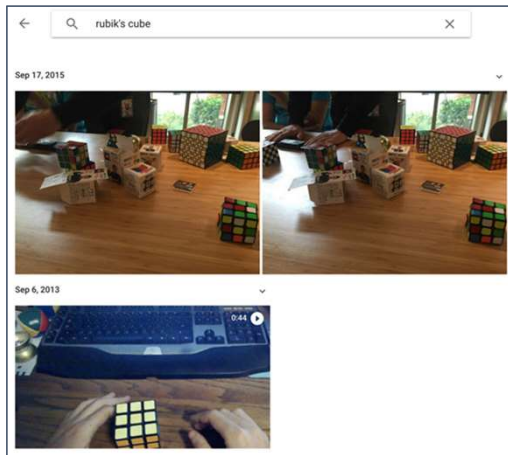
How a computer sees an image

[22, 81, 44, 88, 17, 0, ..., 45]

- Even harder for 3D objects.
- You move a bit, and everything changes.
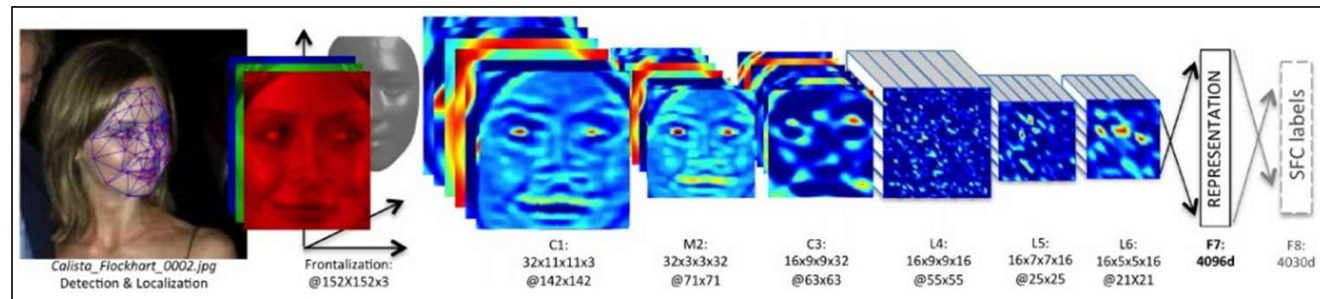
# Supervised: ConvNets are everywhere
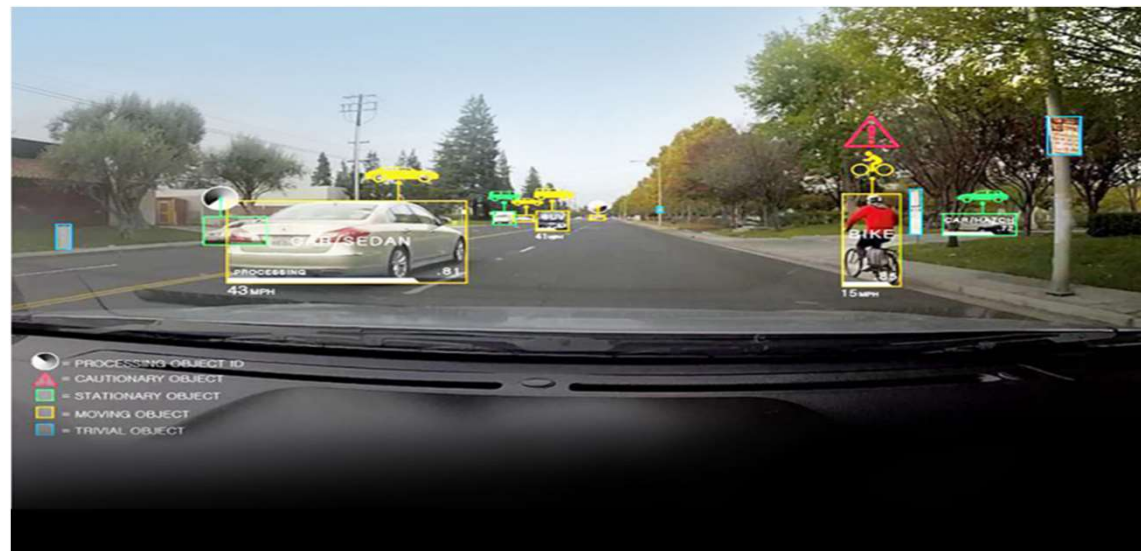


e.g. Google Photos search



Face Verification, Taigman et al. 2014 (FAIR)



[Goodfellow et al. 2014]

*Andrej Karpathy's recent presentation

Self-driving cars

# IMAGENET

*http://groups.inf.ed.ac.uk/calvin/imagenet/prototypes.html

# AlexNet

# VGGNet

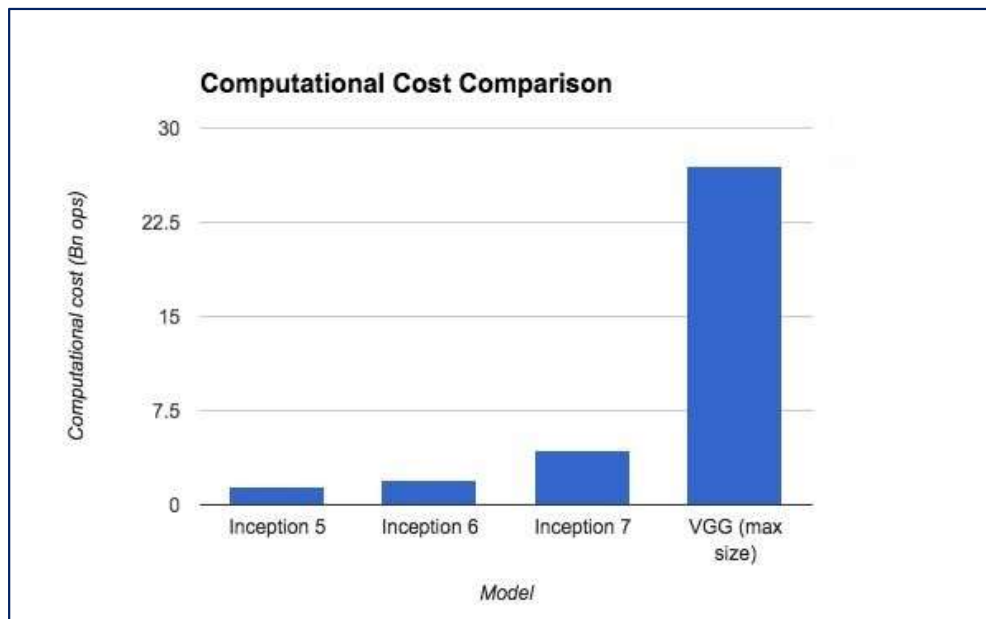| \multicolumn{6}{c}{ConvNet Configuration} | | | | | |
|---|---|---|---|---|---|
| A | A-LRN | B | C | D | E |
| 11 weight layers | 11 weight layers | 13 weight layers | 16 weight layers | 16 weight layers | 19 weight layers |
| \multicolumn{6}{c}{input ($224 \times 224$ RGB image)} | | | | | |
| conv3-64 | conv3-64 | conv3-64 | conv3-64 | conv3-64 | conv3-64 |
|  | LRN | conv3-64 | conv3-64 | conv3-64 | conv3-64 |
| \multicolumn{6}{c}{maxpool} | | | | | |
| conv3-128 | conv3-128 | conv3-128 | conv3-128 | conv3-128 | conv3-128 |
|  |  | conv3-128 | conv3-128 | conv3-128 | conv3-128 |
| \multicolumn{6}{c}{maxpool} | | | | | |
| conv3-256 | conv3-256 | conv3-256 | conv3-256 | conv3-256 | conv3-256 |
| conv3-256 | conv3-256 | conv3-256 | conv3-256 | conv3-256 | conv3-256 |
|  |  |  | conv1-256 | conv3-256 | conv3-256 |
|  |  |  |  |  | conv3-256 |
| \multicolumn{6}{c}{maxpool} | | | | | |
| conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 |
| conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 |
|  |  |  | conv1-512 | conv3-512 | conv3-512 |
|  |  |  |  |  | conv3-512 |
| \multicolumn{6}{c}{maxpool} | | | | | |
| conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 |
| conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 | conv3-512 |
|  |  |  | conv1-512 | conv3-512 | conv3-512 |
|  |  |  |  |  | conv3-512 |
| \multicolumn{6}{c}{maxpool} | | | | | |
| \multicolumn{6}{c}{FC-4096} | | | | | |
| \multicolumn{6}{c}{FC-4096} | | | | | |
| \multicolumn{6}{c}{FC-1000} | | | | | |
| \multicolumn{6}{c}{soft-max} | | | | | |

# GoogLeNet

# GoogLeNet uses Inception



Inception Module

# Inception Performance comparison

# Microsoft ResNet

# Deep learning and natural language processing

# Deep learning enables sub-symbolic processing

| | | |
|---|---|---|
| I | `<i>` | You have to remember to represent "purchased" and "automobile." |
| bought | `<bought>` | |
| a | `<a>` | What about "truck"? |
| car | `<car>` | How do you encode the meaning of the entire sentence? |
| . | `<.>` | |

# But what about a sentence?

Algorithm for generating vectors for sentences

1. Make the sentence vector be the vector for the first word.
2. For each subsequent word, combine its vector with the sentence vector.
3. The resulting vector after the last word is the sentence vector.

Can be implemented using a recurrent neural network (RNN)

# Deep learning and question answering

Bob went home.
Tim went to the junkyard.
Bob picked up the jar.
Bob went to town.
Where is the jar? A: town

Memory Networks [Weston et al., 2014]: Updates memory vectors based on a question and finds the best one to give the output.

The office is north of the yard.
The bath is north of the office.
The yard is west of the kitchen.
How do you go from the office to the kitchen? A: south, east

Neural Reasoner [Peng et al., 2015]: Encodes the question and facts in many layers, and the final layer is put through a function that gives the answer.

# Other commonly used DNNs

# Region Based CNN (RCNN)



R-CNN: *Regions with CNN features*

1. Input image
2. Extract region proposals (~2k)
3. Compute CNN features
4. Classify regions

warped region

aeroplane? no.
person? yes.
tvmonitor? no.

CNN

# Generating Image Descriptions (CNN-RNN)

# Deep Reinforcement Learning

# Increasing Re-usability of Deep Learning models

# Transfer Learning & Fine-tuning

*http://www.mathworks.com/discovery/deep-learning.html?requestedDomain=uk.mathworks.com

# GPUs in Azure