# Deep Learning on Azure

## Dr. Tim Scarfe

Data Solution Architect / Data Scientist

tim.scarfe@microsoft.com / @ecsquendor

Microsoft

"Our goal is to **democratise AI** to empower every person and every organisation to achieve more."
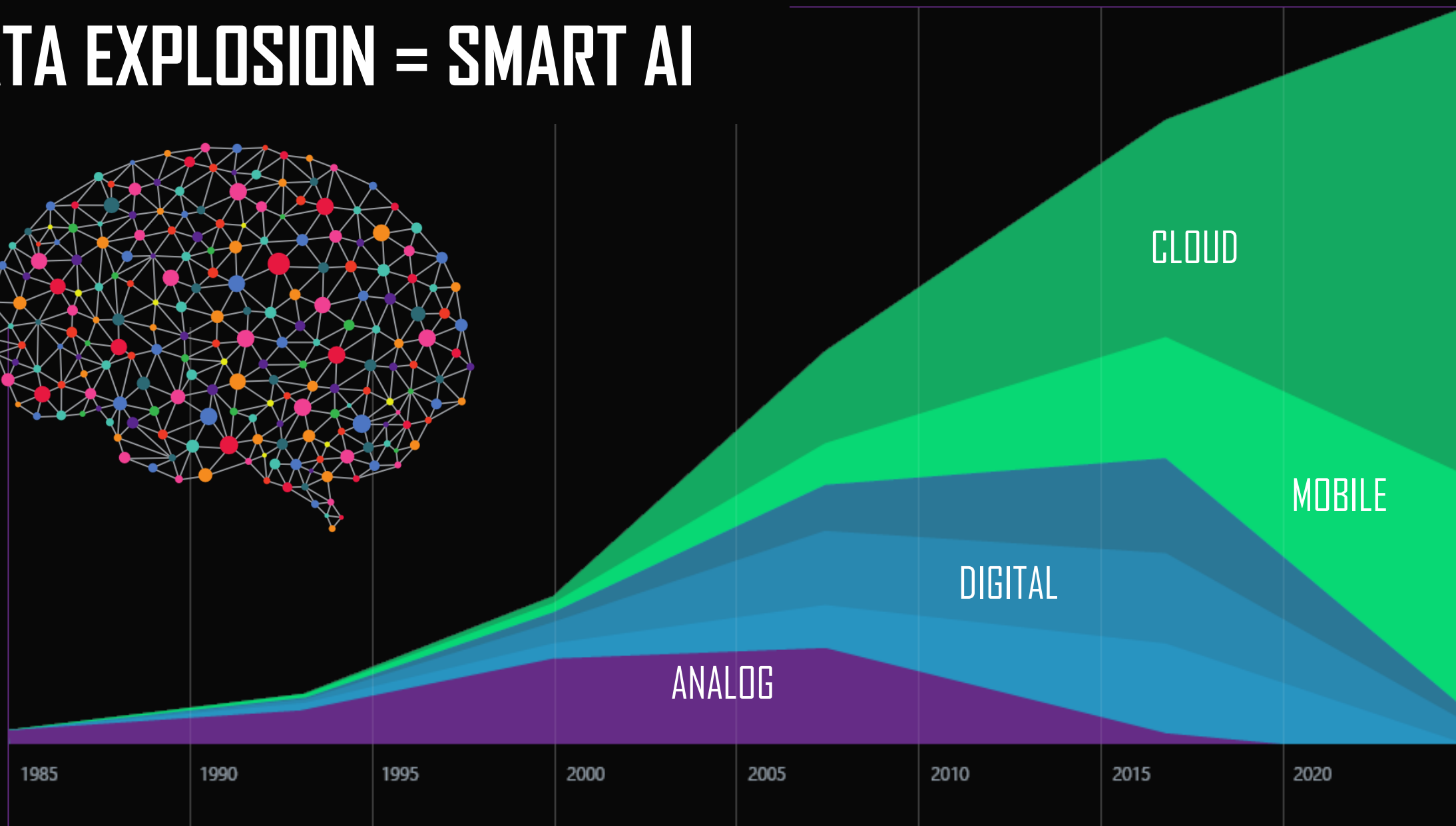
Satya Nadella

# WHAT IS MACHINE LEARNING?

Machine learning is about learning from previous experience so you can make accurate predictions about the future.

Microsoft

DATA EXPLOSION = SMART AI

CLOUD

MOBILE

DIGITAL

ANALOG

1985　1990　1995　2000　2005　2010　2015　2020

NICK BOSTROM

# SUPERINTELLIGENCE

Paths, Dangers, Strategies
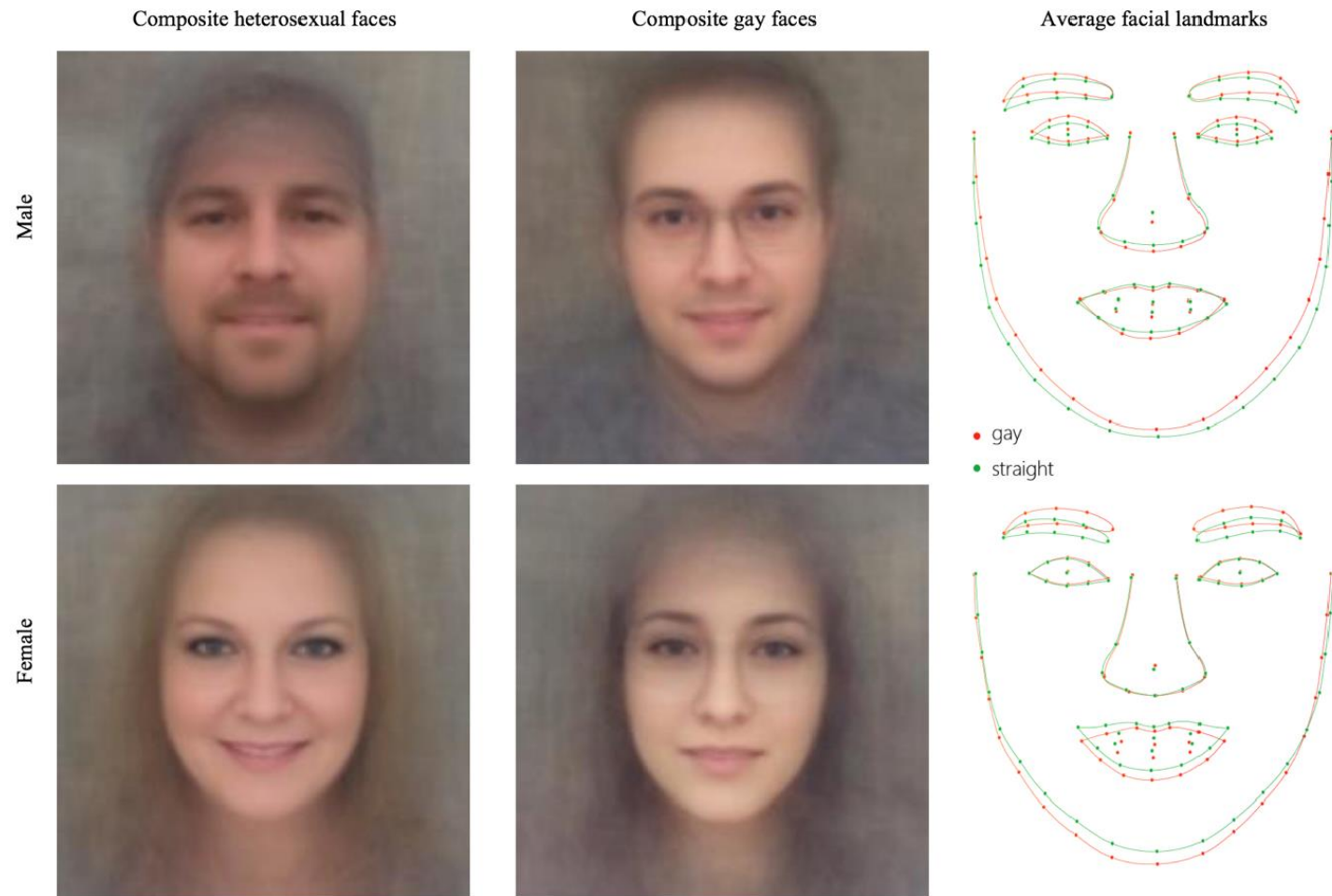
DON'T WORRY ABOUT GENERAL INTELLIGENCE

1. Privacy
2. Opaque AI
3. Data is not neutral
4. Manipulating markets and consumers/voters
5. Lack of human connection
6. Automation of labour / Socioeconomic ramifications
7. Engineers are not philosophers (moral reasoning)
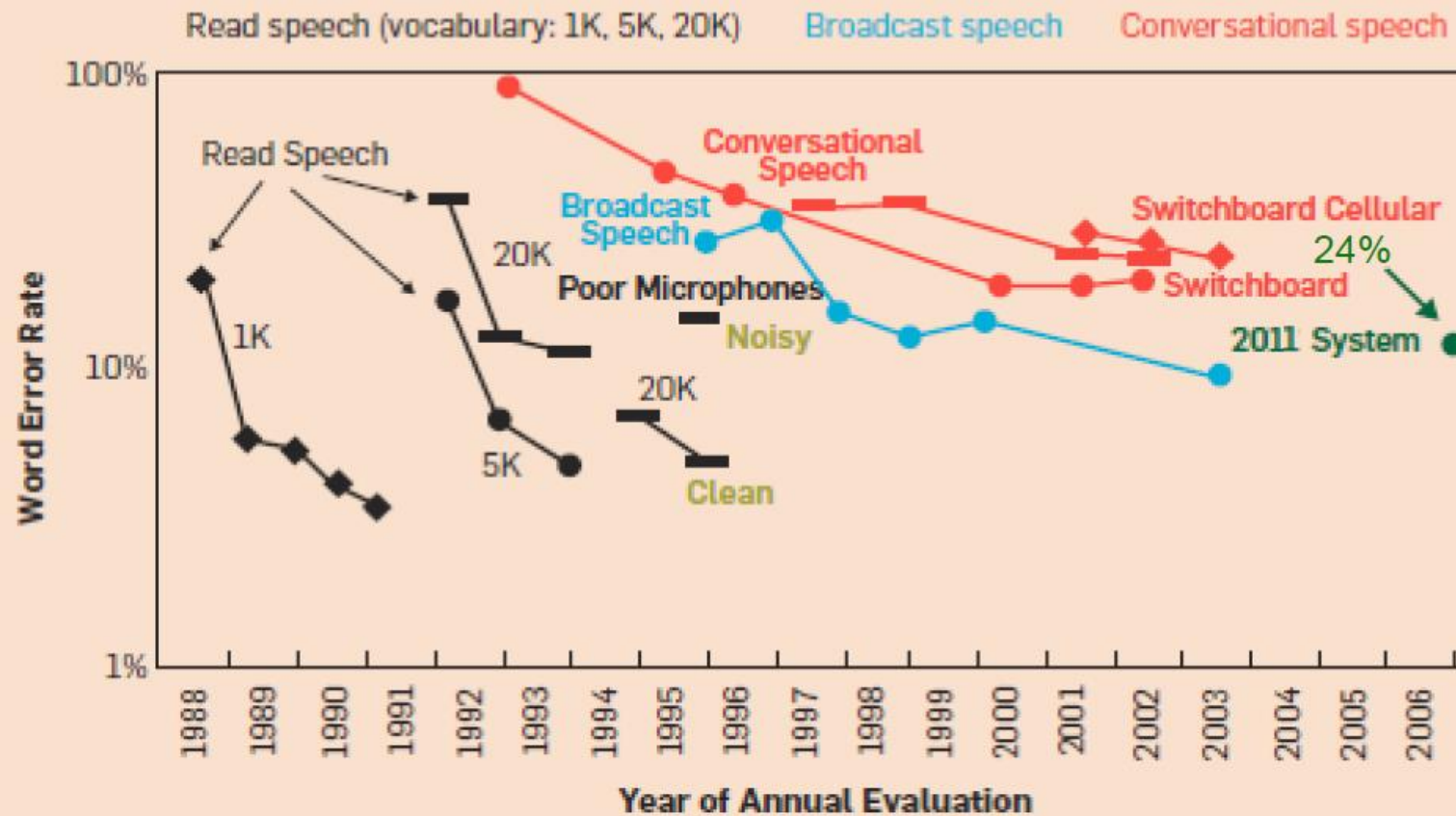
# DO WORRY ABOUT ETHICS

Composite heterosexual faces     Composite gay faces     Average facial landmarks

Male

Female

- gay
- straight

394

395 *Figure 4*. Composite faces and the average facial landmarks built by averaging faces classified as most and least likely to be gay.
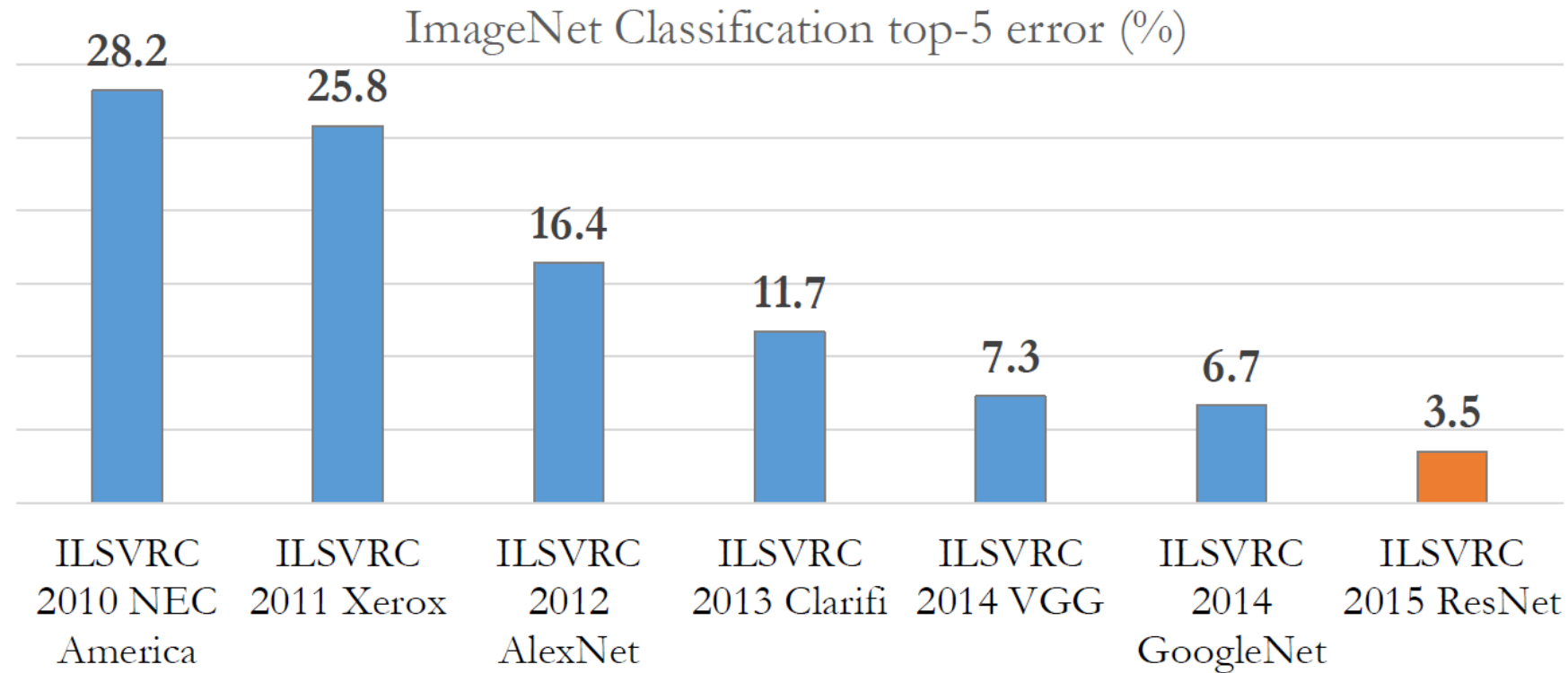
# PHYSIOGNOMY IS BACK?

Read speech (vocabulary: 1K, 5K, 20K)    Broadcast speech    Conversational speech

**Word Error Rate**

100%

Read Speech

20K

1K

Broadcast
Speech

Poor Microphones

Noisy

20K

5K

Clean

Conversational
Speech

Switchboard Cellular

24%

Switchboard

2011 System

10%

1%

1988 1989 1990 1991 1992 1993 1994 1995 1996 1997 1998 1999 2000 2001 2002 2003 2004 2005 2006

**Year of Annual Evaluation**

2017: ~5%!
Human error: ~5%

# IMPROVEMENTS IN SPEECH RECOGNITION

# IMPROVEMENTS IN COMPUTER VISION



ImageNet Classification top-5 error (%)

| Competition | Error |
|---|---|
| ILSVRC 2010 NEC America | 28.2 |
| ILSVRC 2011 Xerox | 25.8 |
| ILSVRC 2012 AlexNet | 16.4 |
| ILSVRC 2013 Clarifi | 11.7 |
| ILSVRC 2014 VGG | 7.3 |
| ILSVRC 2014 GoogleNet | 6.7 |
| ILSVRC 2015 ResNet | 3.5 |

2017: ~2.2%

Microsoft

- We are #1 contributors to open source
- Platinum member of the Linux foundation
- We support all main deep learning frameworks
- CNTK is 100% open source
- You don't have to use CNTK if you don't want to
- Project "Vienna" will support all frameworks and execution environments on-prem and cloud (cloud/containers/Spark)

# THE NEW MICROSOFT

Batch vs. on-line

Regression

Predicting sequences/images

Supervised

Classification

Anomaly Detection

Unsupervised

Agent Based Learning

Clustering

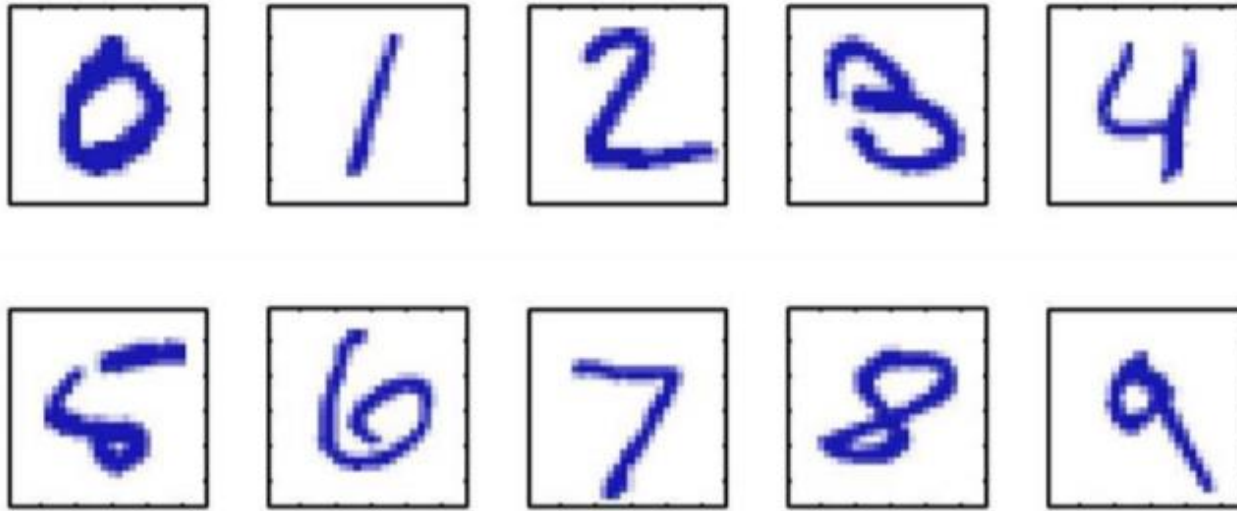Reinforcement Learning

TYPES OF MACHINE LEARNING

- Approximate a function which maps from signals (image) to labels (has-cat)
- This "decision function" can predict missing labels on new, previously unseen signals.
- Historically; different algorithms for different tasks
  - now; deep learning does everything

# MNIST Digit Classification

Images are 28 x 28 pixels

Represent input image as a vector $\mathbf{x} \in \mathbb{R}^{784}$

Learn a classifier $f(\mathbf{x})$ such that,

$$f : \mathbf{x} \rightarrow \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$$
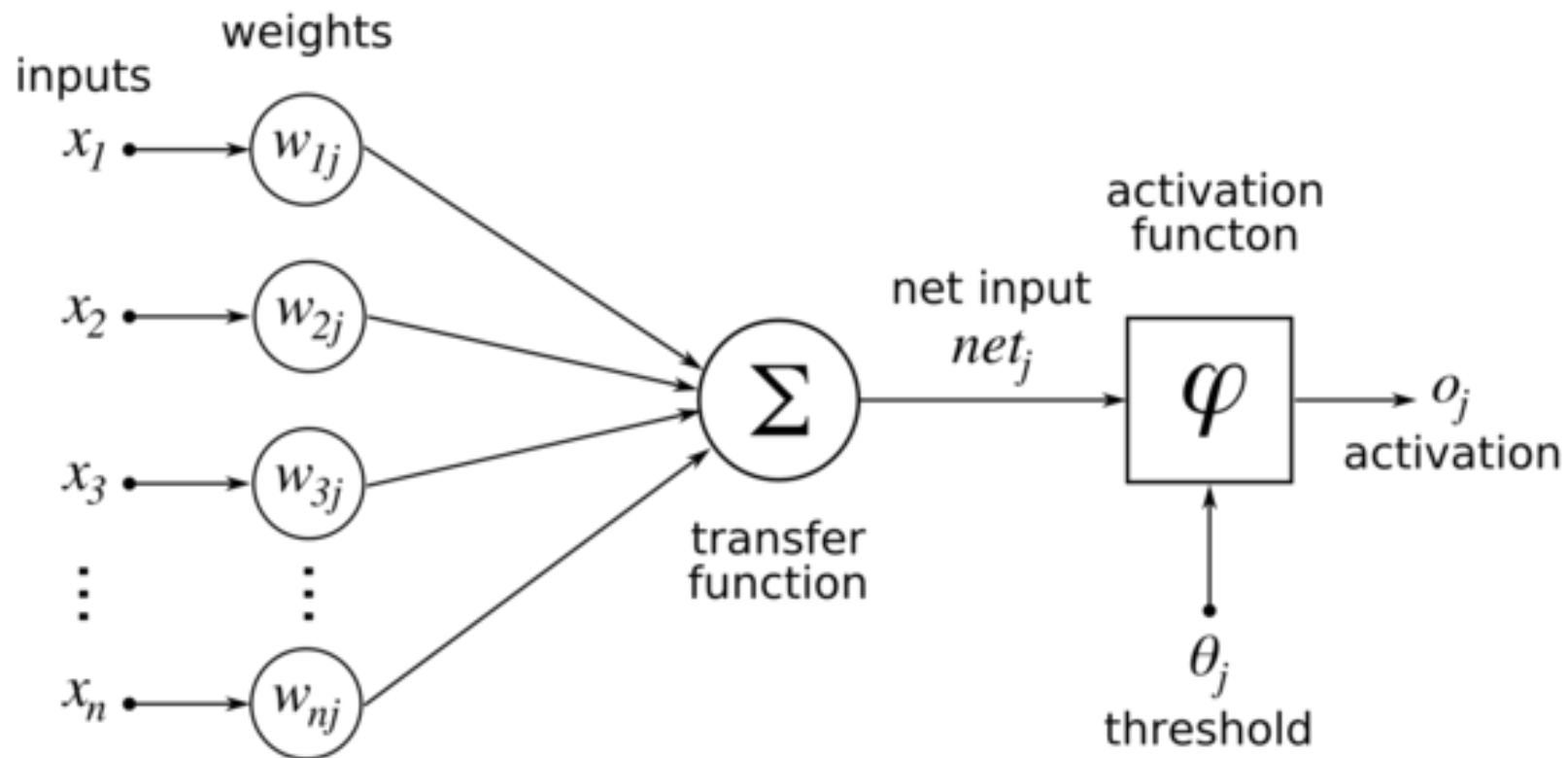
Signals

Labels

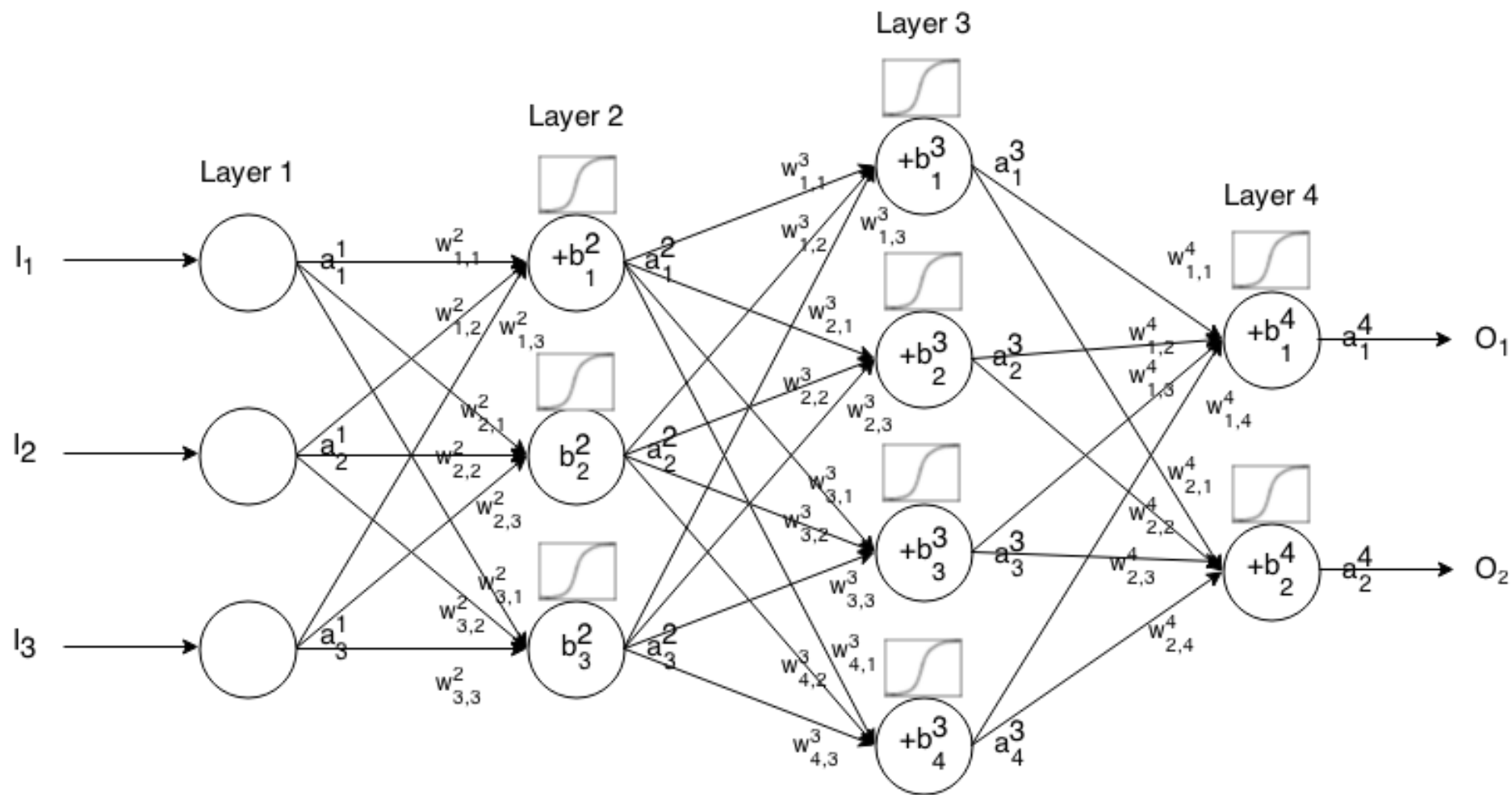# DEEP LEARNING/NEURAL NETWORK DISCUSSION

- Deep Learning = Neural Networks
- Actually, an old technology!
- Universal function approximators; extremely flexible prediction scenarios
- Less emphasis on feature extraction
- Got seriously popular after 2012 due to data+compute explosion
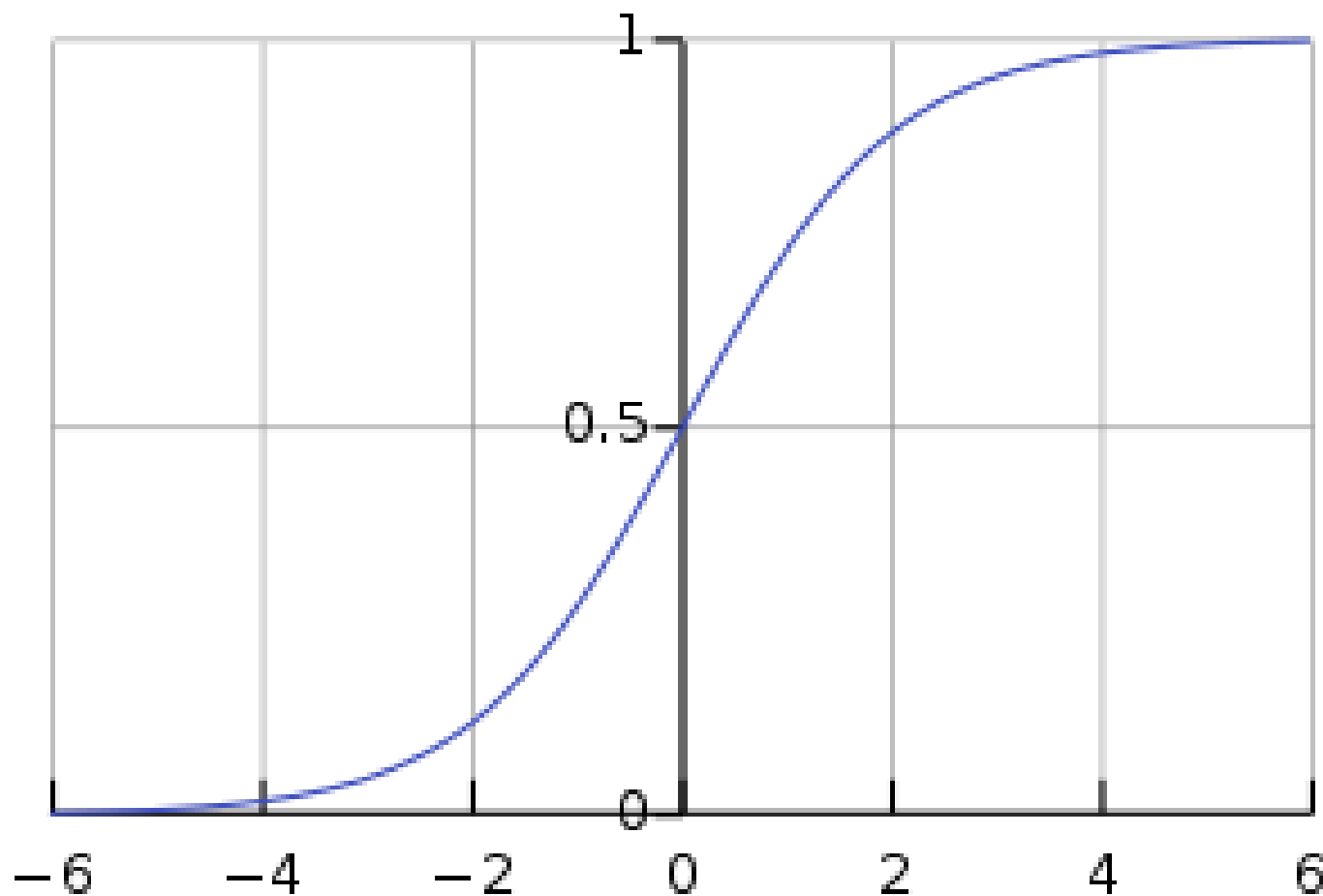- Particularly good for vision, speech, RL and NLP

Microsoft

# WHAT ARE NEURAL NETWORKS?

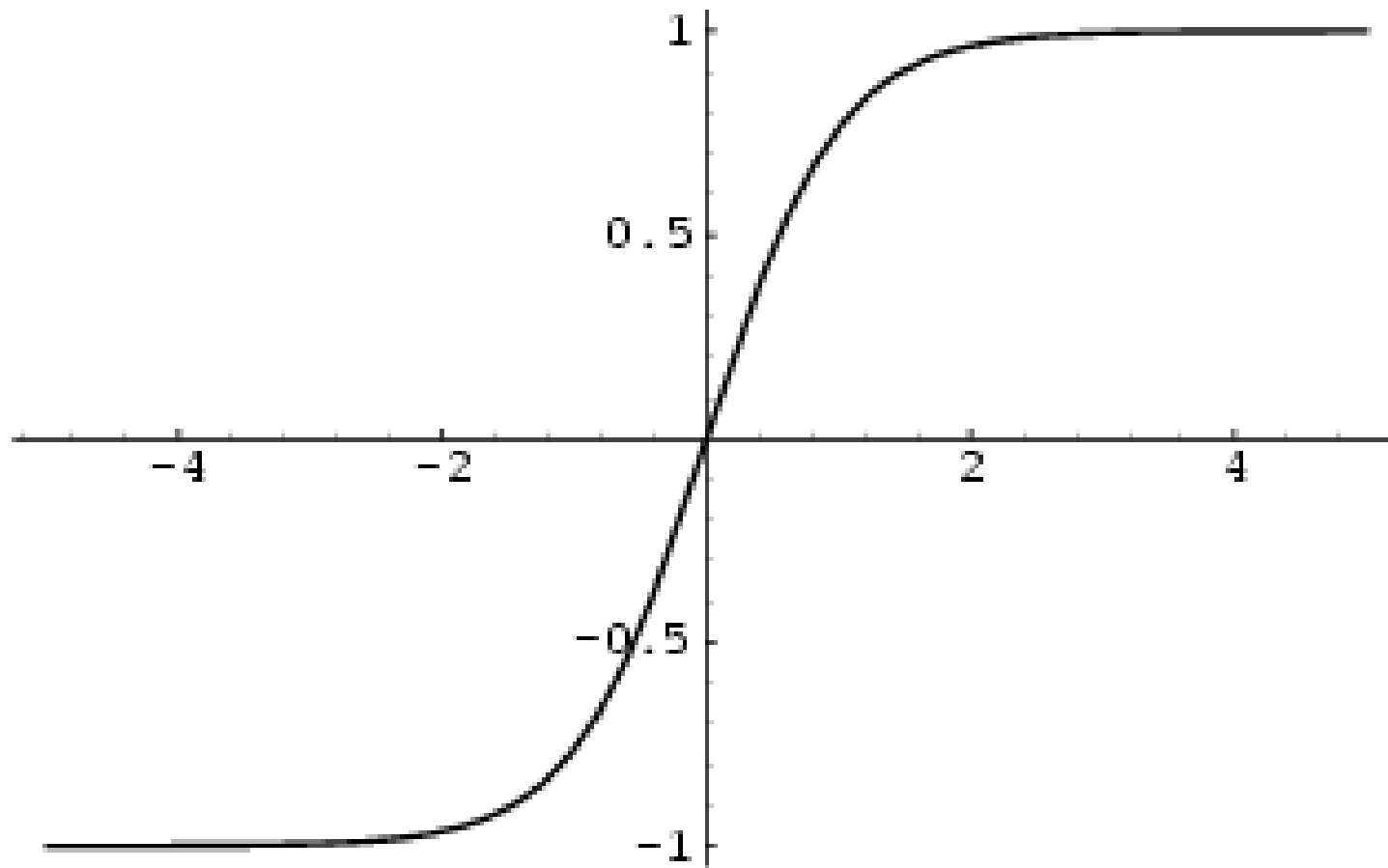# WHAT ABOUT "DEEP" NEURAL NETWORKS?
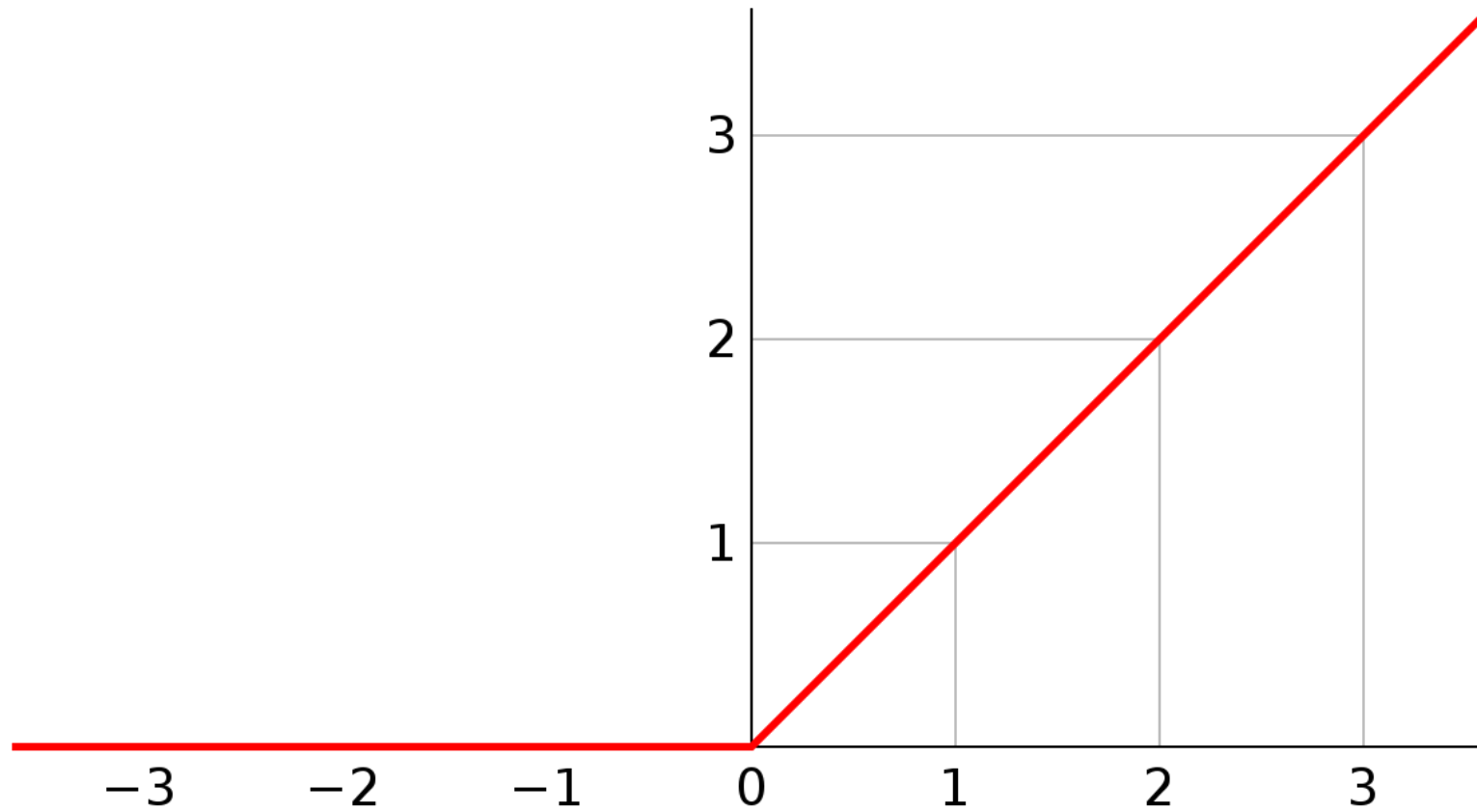
$$\frac{e^x}{e^x + 1}$$

SIGMOID SQUASHING FUNCTION

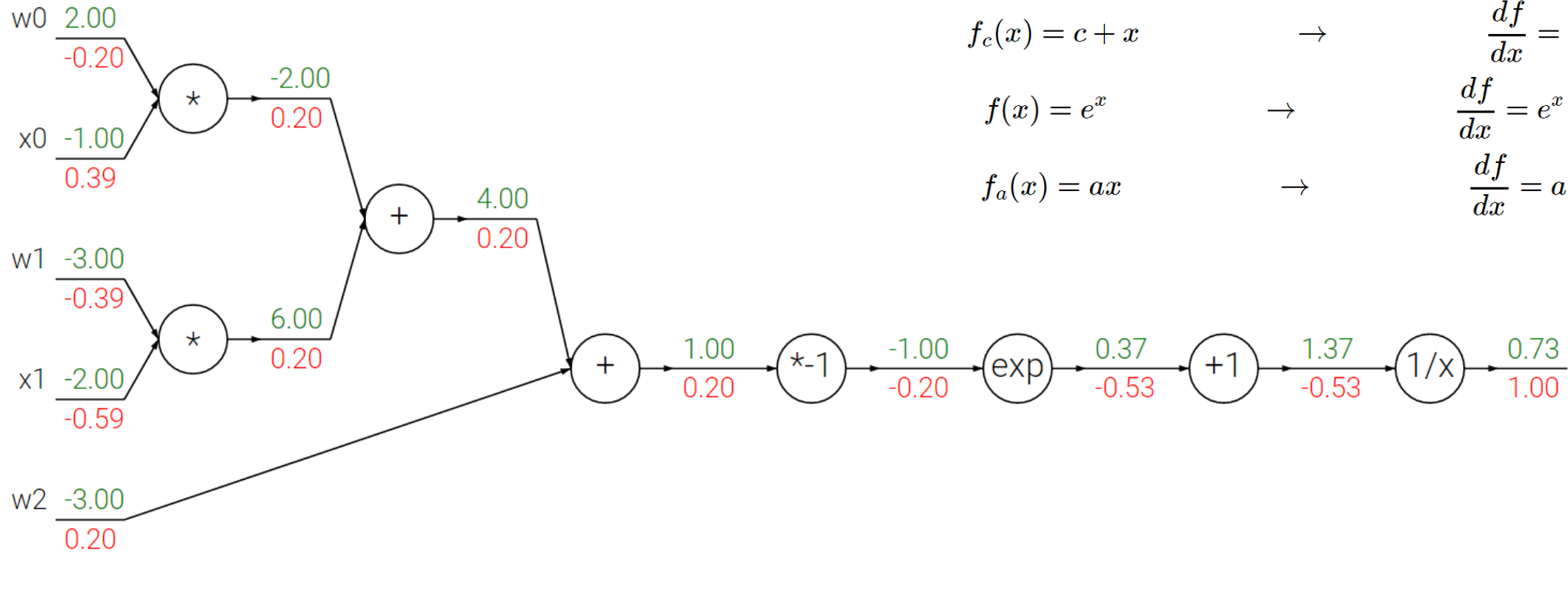$$\frac{1 - e^{-2x}}{1 + e^{-2x}}$$

TANH SQUASHING FUNCTION

$$f(x) = x^+ = \max(0, x)$$

RELU SQUASHING FUNCTION

# BACKPROPAGATION

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

w0  2.00
-0.20

x0  -1.00
0.39

-2.00
0.20

w1  -3.00
-0.39

x1  -2.00
-0.59

6.00
0.20

4.00
0.20

w2  -3.00
0.20

1.00
0.20

-1.00
-0.20

0.37
-0.53

1.37
-0.53

0.73
1.00

$*$  $+$  $+$  $*$-1  exp  $+1$  1/x

$$\frac{\partial E}{\partial w_{jk}}$$

Example circuit for a 2D neuron with a sigmoid activation function. The inputs are [x0,x1] and the (learnable) weights of the neuron are [w0,w1,w2]. As we will see later, the neuron computes a dot product with the input and then its activation,n is softly squashed by the sigmoid function to be in range from 0 to 1.

Microsoft

$$\overbrace{\Delta w_{jk} = \eta * [\; x_j * \underbrace{(o_k - t_k)}_{\text{error } e_k} * \underbrace{o_k * (1 - o_k)}_{\text{derivative of output activation } \varphi_k'}\; ]}^{\frac{\partial E}{\partial w_{jk}}}$$

learning rate

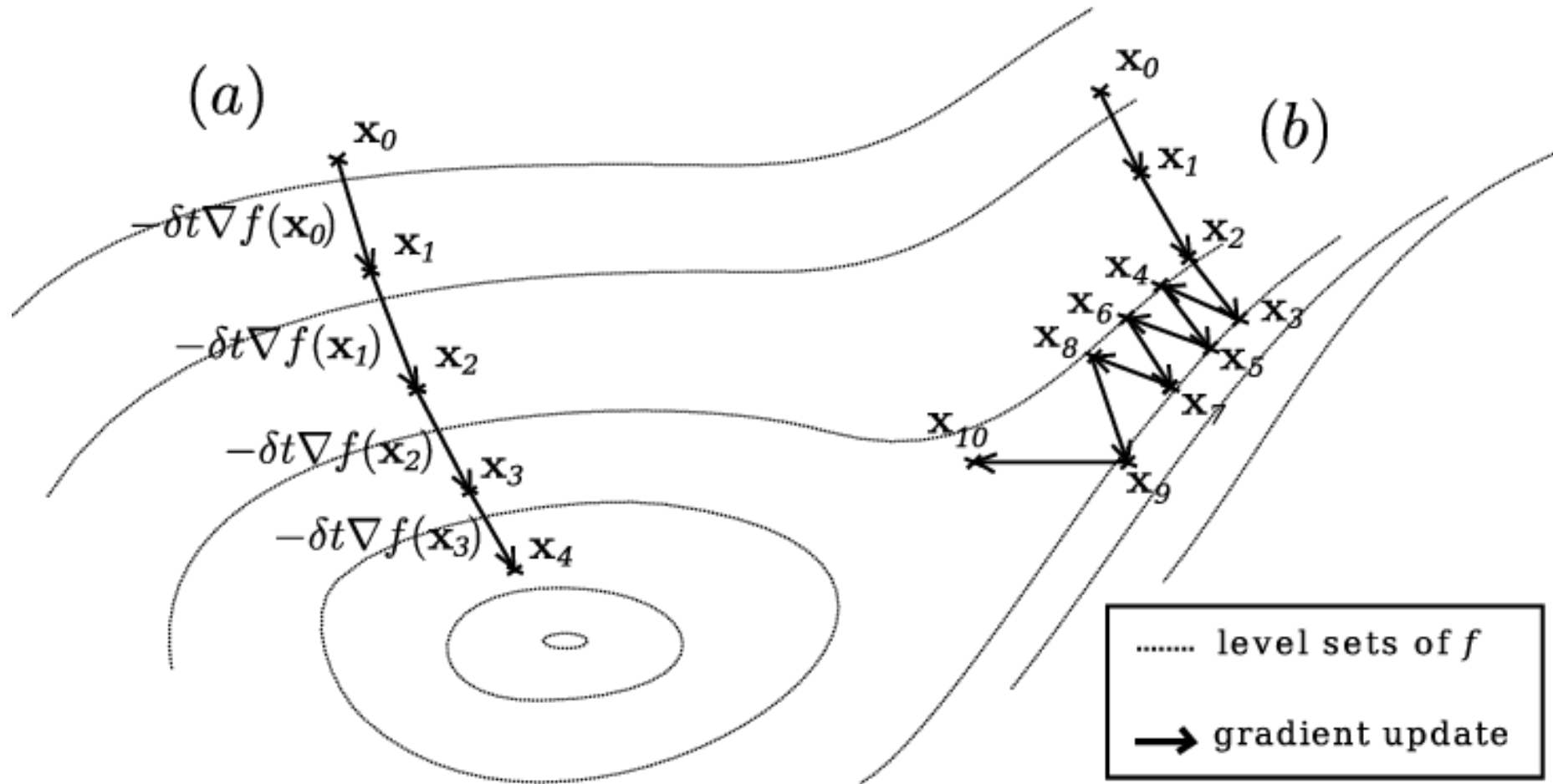signal $\delta_k$

WEIGHT UPDATE

```
loop maxEpochs times
  for-each training item
    get target values
    compute output values
    compute the gradient of each weight
    use gradient to compute delta for each weight
    update each weight using its delta
  end-for
end-loop
```

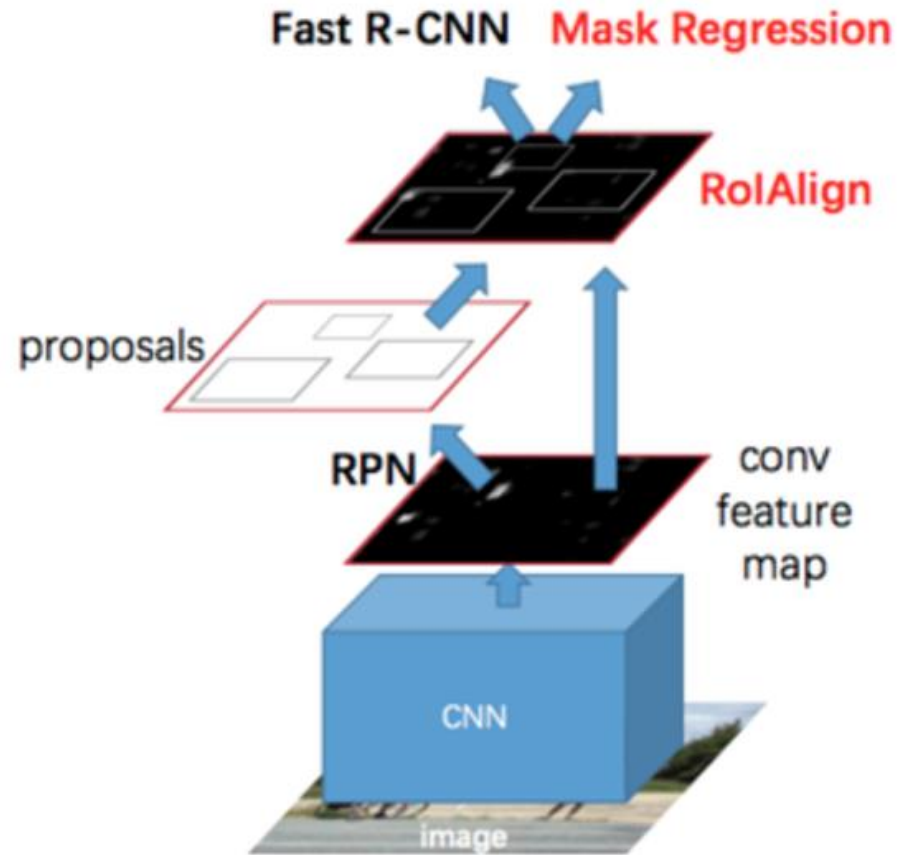# OPTIMIZATION/GRADIENT DESCENT

http://playground.tensorflow.org

# NEURAL NETWORK PLAYGROUND

Microsoft

# WHY IS DEEP LEARNING SPECIAL, IS IT A FAD?

- The way we think about neural networks now is totally different to 30 years ago
- Previous frequentist algorithms were just learning weighted combinations of hand-crafted features
- NNs learn a hierarchy of representations which work really well in many domains
- Before we used to talk about classification and regression, now we talk about *predictive architectures*

Microsoft

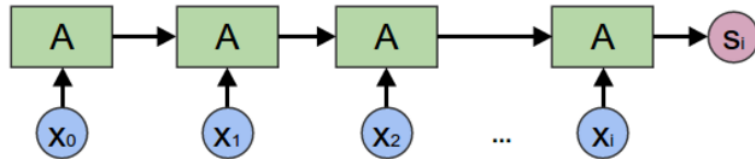# MASK R-CNN ARCHITECTURE (2017)

# PARADGM-SHIFT

- Three narratives currently exist to describe deep learning
    - Neuroscience
    - Probabilistic
    - Manifold
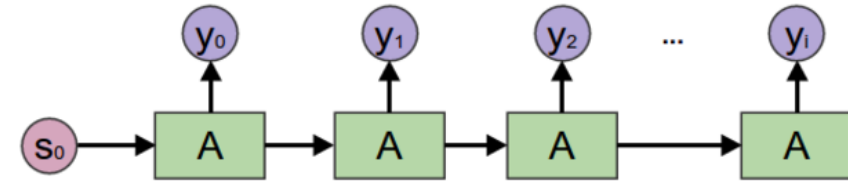- The *differentiable programming* narrative is emerging

Microsoft

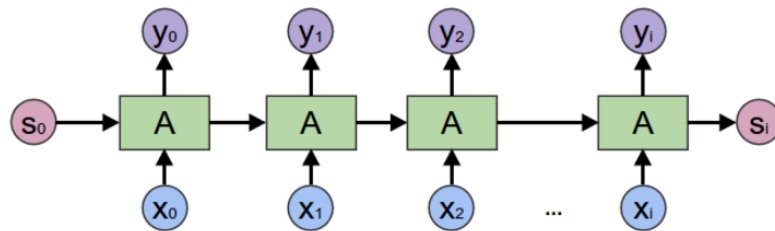# FUNCTIONAL PROGRAMMING IN DL

fold = Encoding RNN
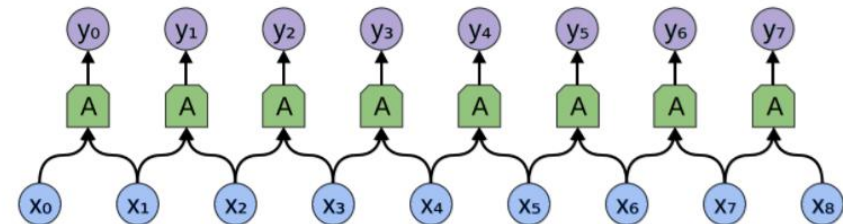
Haskell: `foldl a s`



unfold = Generating RNN

Haskell: `unfoldr a s`
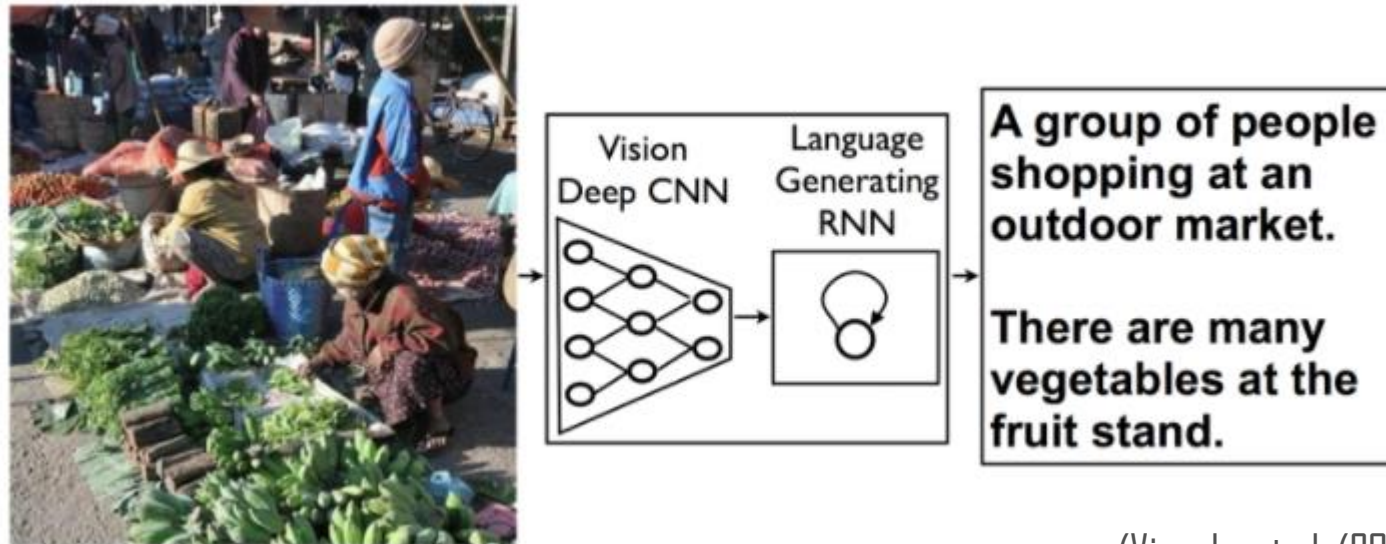


Accumulating Map = RNN

Haskell: `mapAccumR a s`



Windowed Map = Convolutional Layer

Haskell: `zipWith a xs (tail xs)`

# BUILDING PREDICTIVE ARCHITECTURES WITH COMPONENTS



(Vinyals, et al. (2014))

# WHAT IS CNTK?

DECLARITIVELY DESCRIBE AND TRAIN DEEP NEURAL NETWORKS

DOES ALL THE HARD WORK FOR YOU

80% INTERNAL MS DL WORKLOADS USE CNTK

1ST CLASS ON LINUX, WINDOWS, DOCKER

C#, PYTHON, COMMANDLINE

KERAS BINDINGS

Microsoft

Caffe: 1.0rc5(39f28e4)
CNTK: 2.0 Beta10(1ae666d)
MXNet: 0.93(32dc3a2)
TensorFlow: 1.0(4ac9c09)
Torch: 7(748f5e3)

| | Caffe | CNTK | MxNet | TensorFlow | Torch |
|---|---|---|---|---|---|
| FCN5 (1024) | 55.329ms | **51.038ms** | 60.448ms | 62.044ms | 52.154ms |
| AlexNet (256) | 36.815ms | **27.215ms** | 28.994ms | 103.960ms | 37.462ms |
| ResNet (32) | 143.987ms | **81.470ms** | 84.545ms | 181.404ms | 90.935ms |
| LSTM (256) (v7 benchmark) | - | **43.581ms** (44.917ms) | 288.142ms (284.898ms) | - (223.547ms) | 1130.606ms (906.958ms) |

# THE FASTEST TOOLKIT

Achieved with 1-bit gradient quantization algorithm

Theano only supports 1 GPU

CNTK   Theano   TensorFlow   Torch 7   Caffe

■ 1 GPU   ■ 1 x 4 GPUs   ■ 2 x 4 GPUs (8 GPUs)

MOST SCALABLE TOOLKIT (2016)

# INSTALLING CNTK

- GOOGLE "CNTK INSTALL" (WITH BING)
- USE THE "SCRIPT DRIVEN INSTALLATION"

Microsoft

# WHEN TO USE DEEP LEARNING FRAMEWORKS

- Sequence modelling (speech, language, time-series)
- Complex vision tasks (localisation, detection)
- Novel prediction architectures
- Generative models
- Reinforcement learning
- ... and many more!
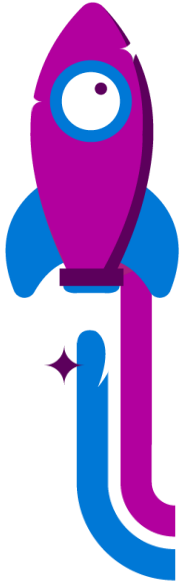
Microsoft

# DEEP LEARNING ON AZURE CLOUD

- Data Science Virtual Machine (Ubuntu and Windows)
- Batch AI Training Service
- AzureML supports some deep learning workloads
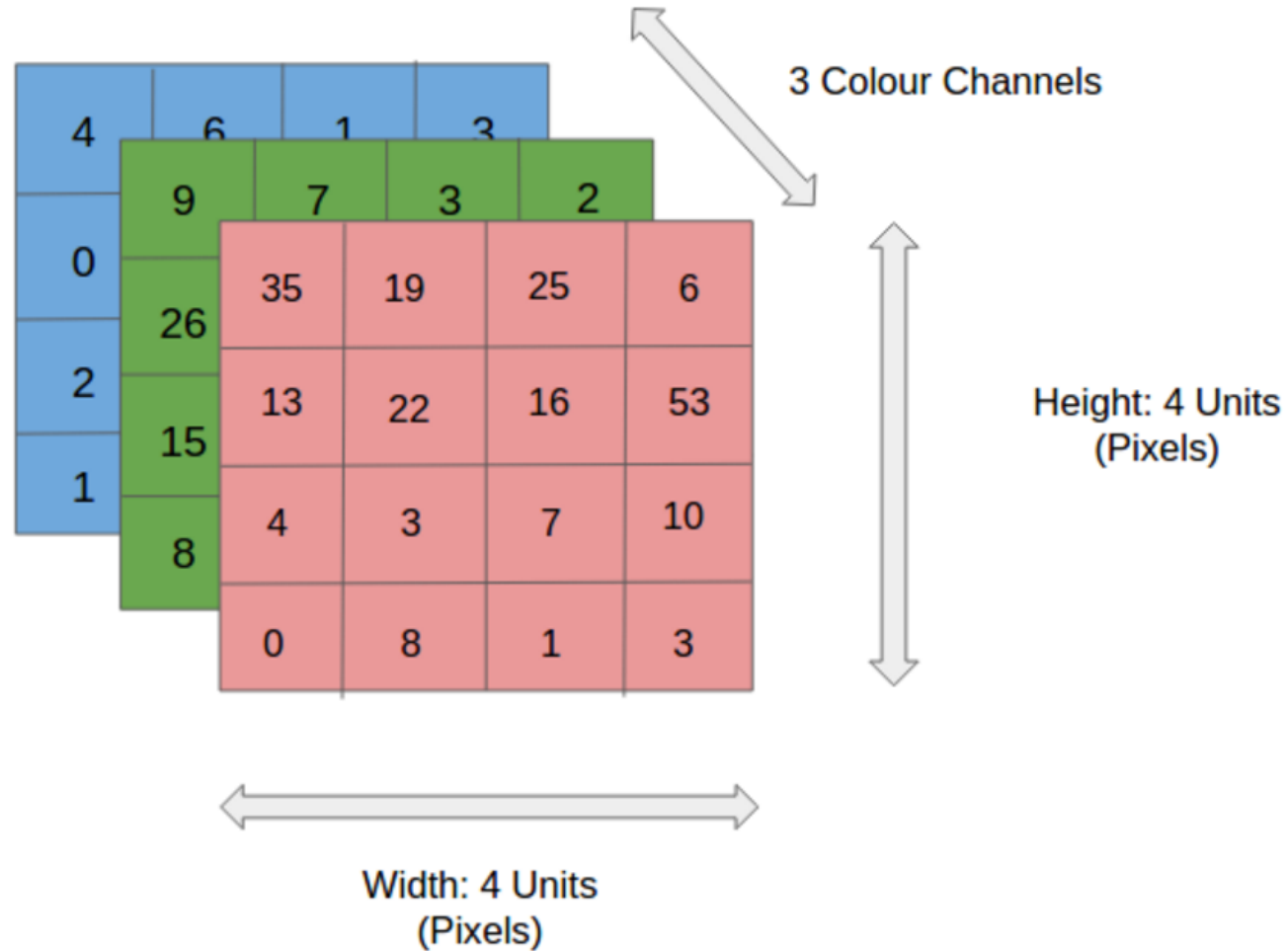- R Server supports some deep learning

Microsoft

# CNTK IRIS DEMO

Microsoft

# WHAT ABOUT VISION AND NATURAL LANGUAGE PROCESSING?

Microsoft

| 0 | 0 | 0 | 0 | 0 | 30 | 0 |
|---|---|---|---|---|----|---|
| 0 | 0 | 0 | 0 | 30 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

**Pixel representation of filter**

**Visualization of a curve detector filter**

CONVOLUTION FILTER

Visualization of the filter on the image

$$(50*30)+(50*30)+(50*30)+(20*30)+(50*30) = 6600$$

Visualization of the receptive field

Pixel representation of the receptive field

| 0 | 0 | 0 | 0 | 0 | 0 | 30 |
|---|---|---|---|---|---|----|
| 0 | 0 | 0 | 0 | 50 | 50 | 50 |
| 0 | 0 | 0 | 20 | 50 | 0 | 0 |
| 0 | 0 | 0 | 50 | 50 | 0 | 0 |
| 0 | 0 | 0 | 50 | 50 | 0 | 0 |
| 0 | 0 | 0 | 50 | 50 | 0 | 0 |
| 0 | 0 | 0 | 50 | 50 | 0 | 0 |

*

Pixel representation of filter

| 0 | 0 | 0 | 0 | 0 | 30 | 0 |
|---|---|---|---|---|----|---|
| 0 | 0 | 0 | 0 | 30 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

CONVOLUTION FILTER MATCH

MULTIPLY AND SUMMATION = 0

Visualization of the filter on the image

| 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|
| 0 | 40 | 0 | 0 | 0 | 0 | 0 |
| 40 | 0 | 40 | 0 | 0 | 0 | 0 |
| 40 | 20 | 0 | 0 | 0 | 0 | 0 |
| 0 | 50 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 50 | 0 | 0 | 0 | 0 |
| 25 | 25 | 0 | 50 | 0 | 0 | 0 |

*

| 0 | 0 | 0 | 0 | 0 | 30 | 0 |
|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 30 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 30 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Pixel representation of receptive field

Pixel representation of filter

# CONVOLUTION FILTER NO MATCH

CONVOLUTION

(i)

(ii)

(iii)

(iv)

POOLING

# Image Classification Example



INPUT     CONVOLUTION + RELU     POOLING     CONVOLUTION + RELU     POOLING     FLATTEN     FULLY CONNECTED     SOFTMAX

— CAR
— TRUCK
— VAN
— BICYCLE

**FEATURE LEARNING**        **CLASSIFICATION**

# VISUALISING THE FILTERS

# DEEP VISUALISATION TOOLBOX

# RESOLUTION OF DEPTH



ImageNet Classification top-5 error (%)

CNTK MNIST DEMO

# THANK YOU

tim.scarfe@microsoft.com