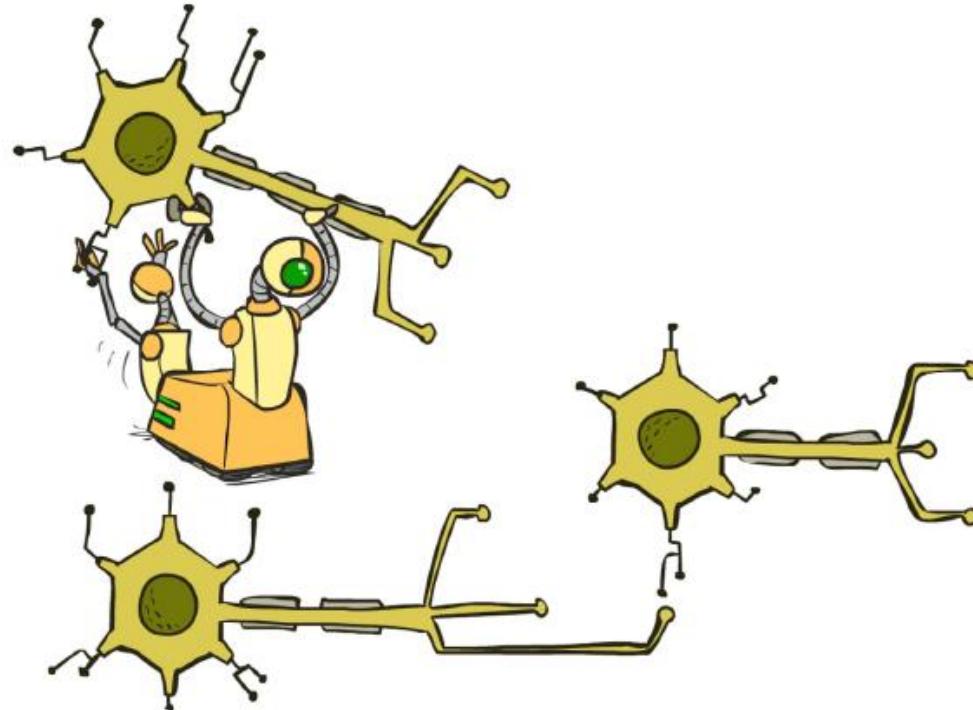


# 第十三周

- 教学计划
  - 神经网络2（上次课）
  - 项目报告：0112, 0113, 0203组
  - 神经网络3（这次课）
- 任务
  - 家作5/项目5：机器学习
  - 项目5口令：**machinelearning**

# 人工智能导论

## 神经网络3



基于CS231n课程 --- Stanford University

# Backpropagation - Recall

Backpropagation: a simple example

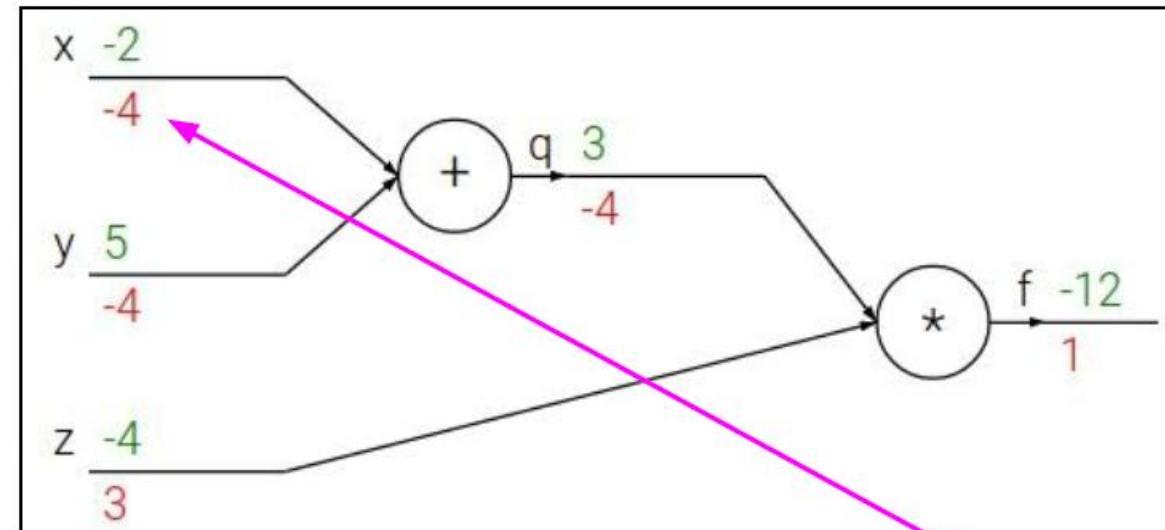
$$f(x, y, z) = (x + y)z$$

e.g.  $x = -2$ ,  $y = 5$ ,  $z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want:  $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



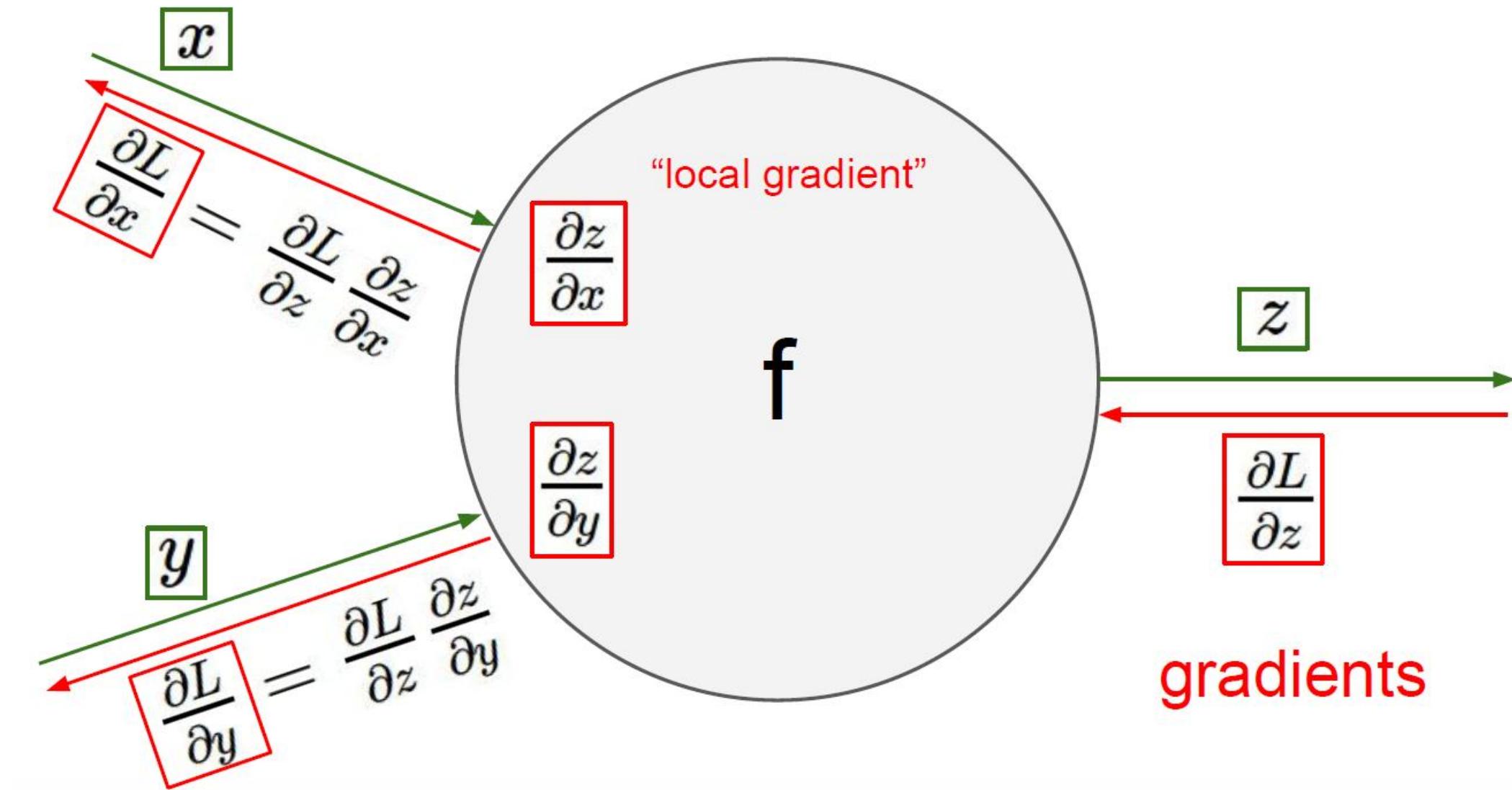
Chain rule:

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$

Upstream  
gradient

Local  
gradient

# Gradient flow



# Backpropagation

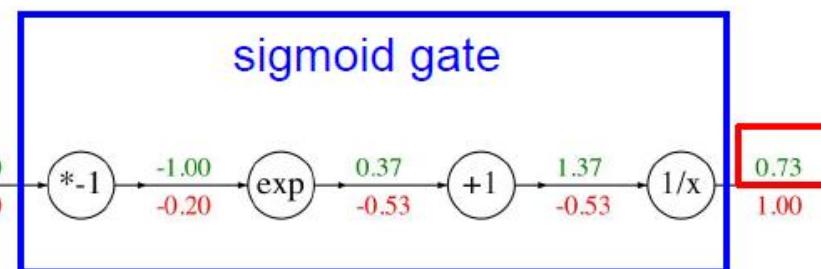
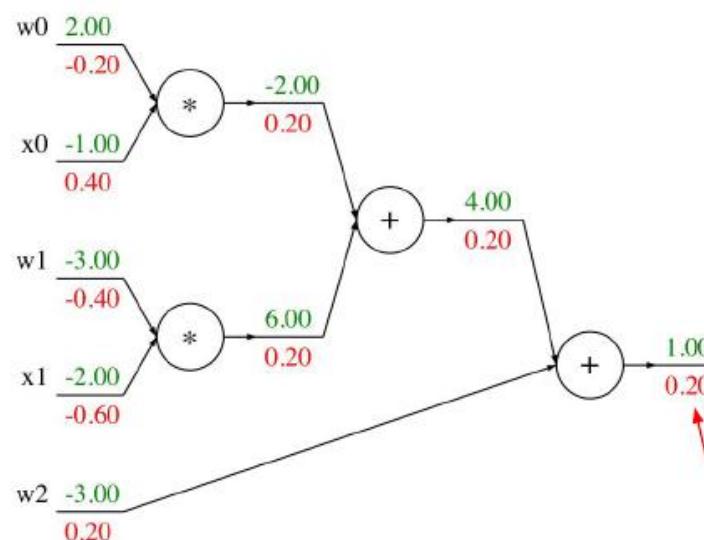
$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

Computational graph representation may not be unique. Choose one where local gradients at each node can be easily expressed!

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

sigmoid function

$$\frac{d\sigma(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2} = \left( \frac{1 + e^{-x} - 1}{1 + e^{-x}} \right) \left( \frac{1}{1 + e^{-x}} \right) = (1 - \sigma(x))\sigma(x)$$



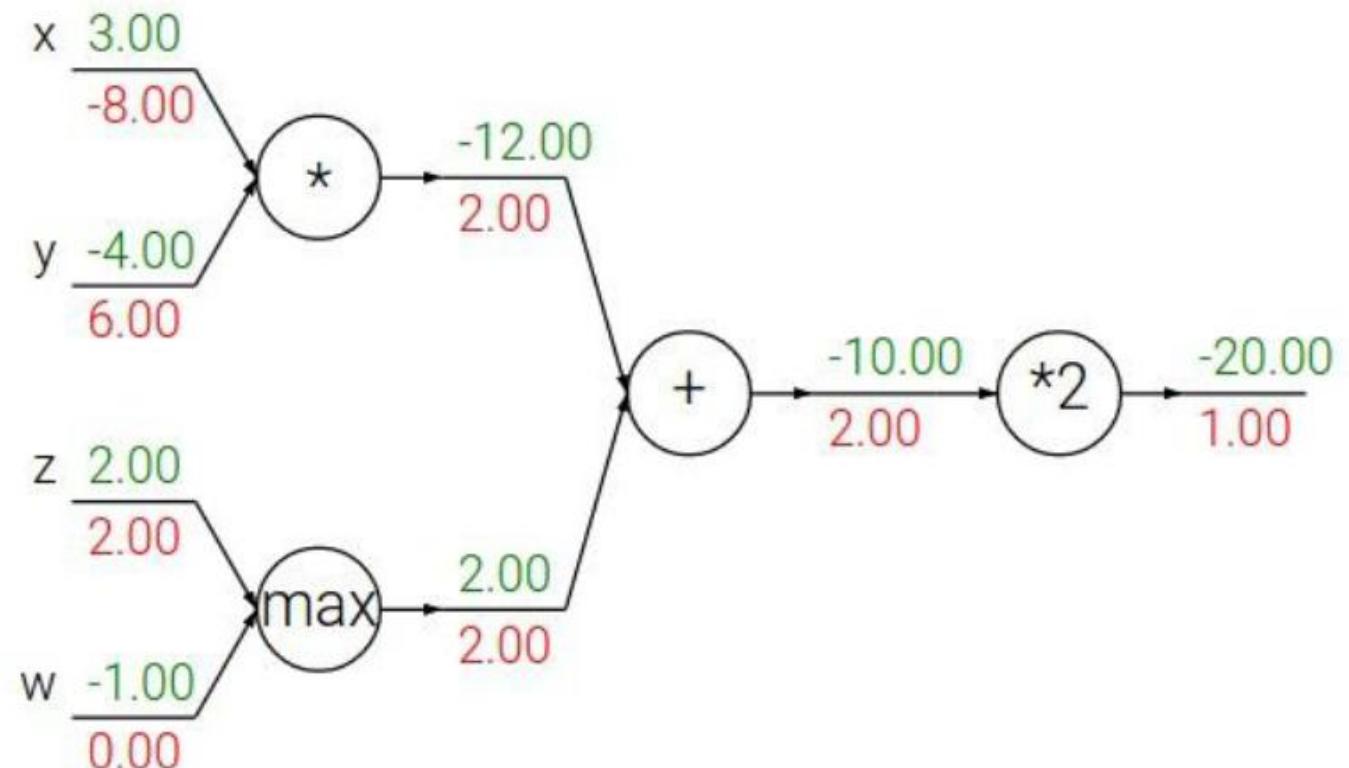
$$[\text{upstream gradient}] \times [\text{local gradient}] \\ [1.00] \times [(1 - 0.73) (0.73)] = 0.2$$

# Patterns in backward flow

**add** gate: gradient distributor

**max** gate: gradient router

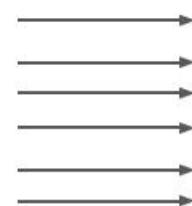
**mul** gate: gradient switcher



# Batches

## Vectorized operations

4096-d  
input vector



$f(x) = \max(0, x)$   
(elementwise)

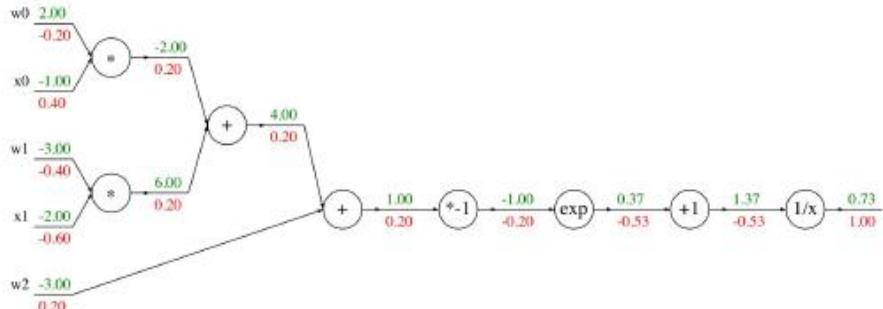
4096-d  
output vector

Q: what is the  
size of the  
Jacobian matrix?  
[4096 x 4096!]

in practice we process an  
entire minibatch (e.g. 100)  
of examples at one time:  
i.e. Jacobian would technically be a  
[409,600 x 409,600] matrix :)

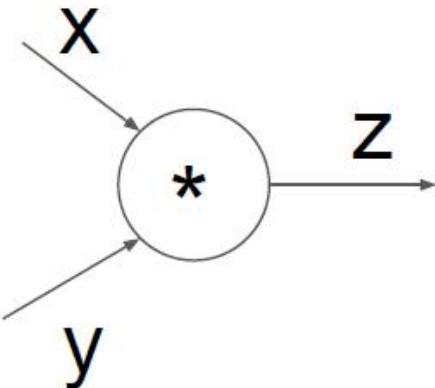
# Modular implementation: forward/backward API

Graph (or Net) object (*rough pseudo code*)



```
class ComputationalGraph(object):
    ...
    def forward(inputs):
        # 1. [pass inputs to input gates...]
        # 2. forward the computational graph:
        for gate in self.graph.nodes_topologically_sorted():
            gate.forward()
        return loss # the final gate in the graph outputs the loss
    def backward():
        for gate in reversed(self.graph.nodes_topologically_sorted()):
            gate.backward() # little piece of backprop (chain rule applied)
        return inputs_gradients
```

# Modular implementation: forward/backward API



( $x, y, z$  are scalars)

```
class MultiplyGate(object):  
    def forward(x,y):  
        z = x*y  
        self.x = x # must keep these around!  
        self.y = y  
        return z  
    def backward(dz):  
        dx = self.y * dz # [dz/dx * dL/dz]  
        dy = self.x * dz # [dz/dy * dL/dz]  
        return [dx, dy]
```

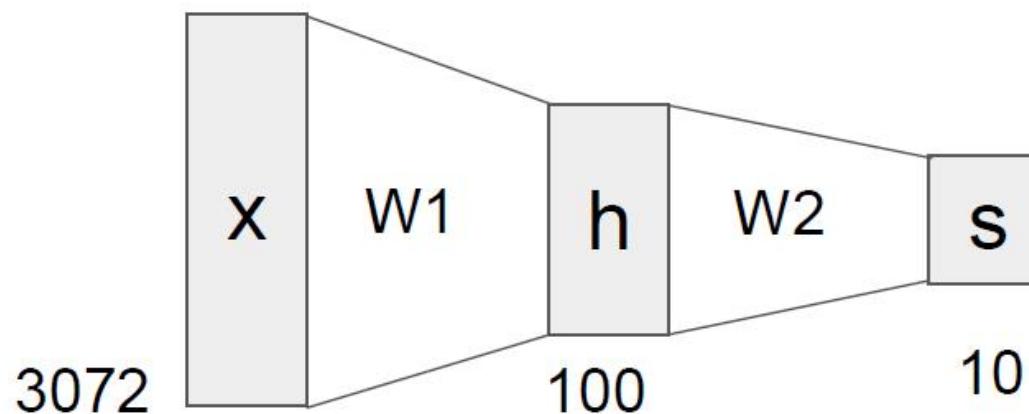
Local gradient      Upstream gradient variable

# Neural network

Neural networks: without the brain stuff

(Before) Linear score function:  $f = Wx$

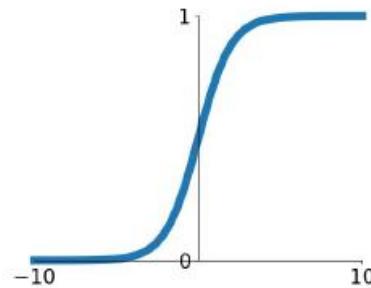
(Now) 2-layer Neural Network  $f = W_2 \max(0, W_1 x)$



# Activation function

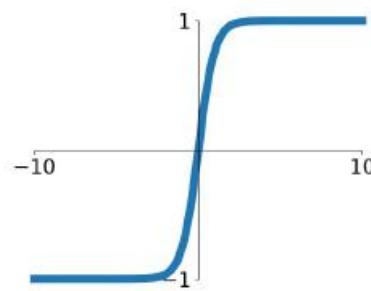
**Sigmoid**

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



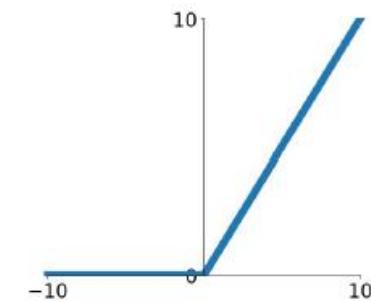
**tanh**

$$\tanh(x)$$



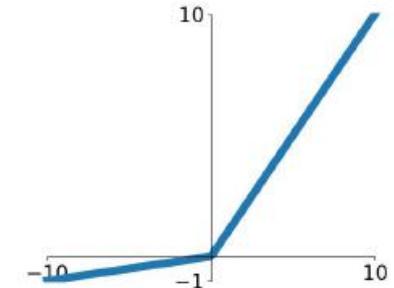
**ReLU**

$$\max(0, x)$$



**Leaky ReLU**

$$\max(0.1x, x)$$

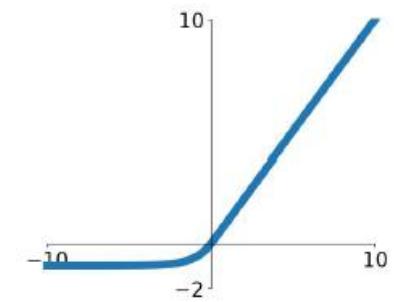


**Maxout**

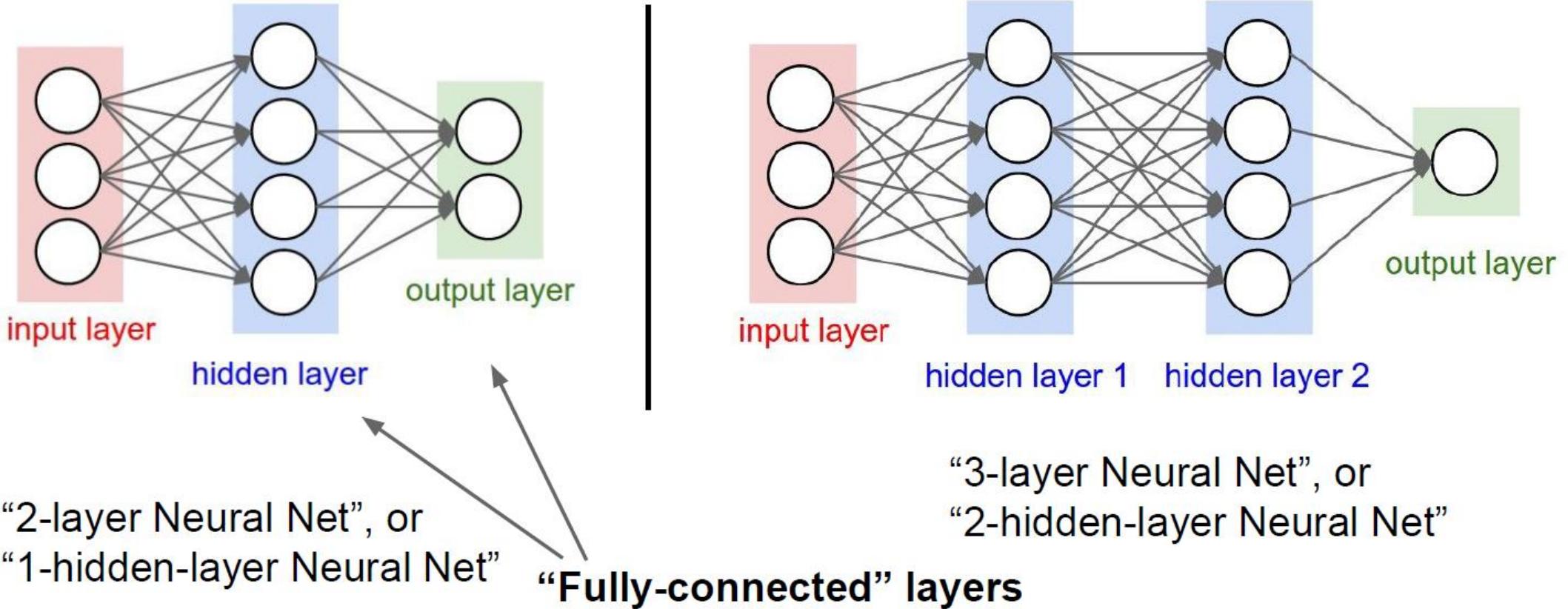
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

**ELU**

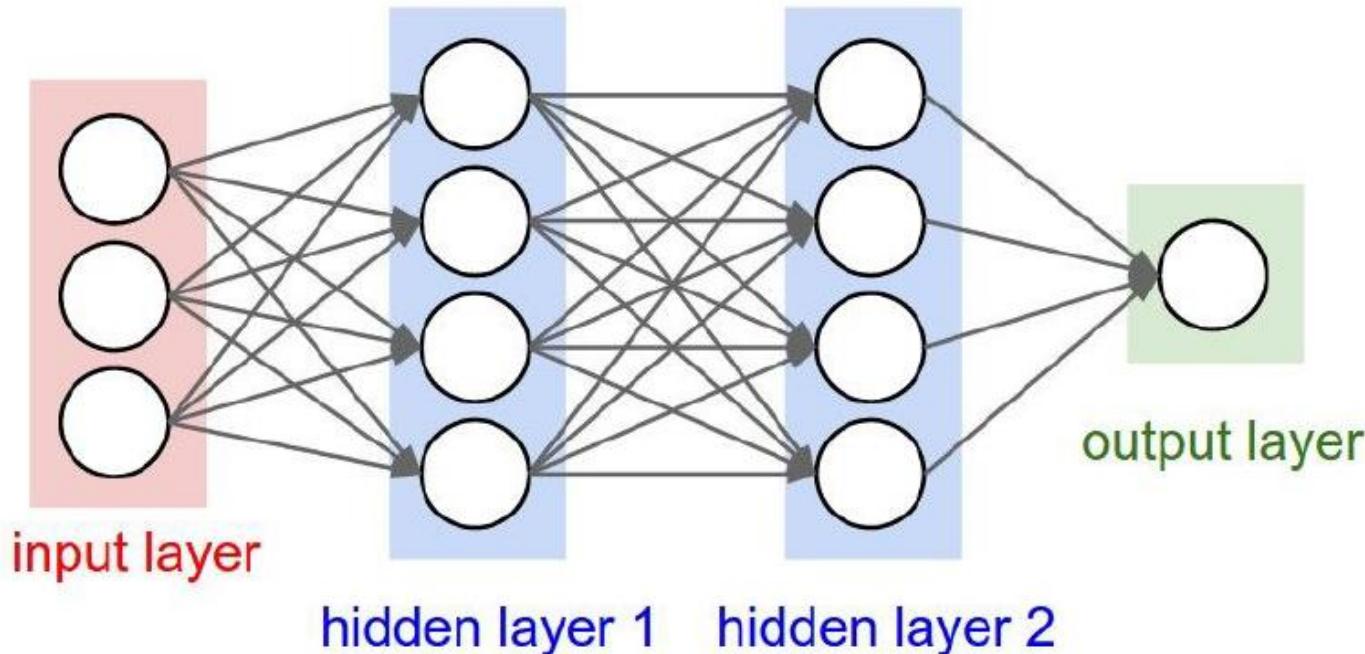
$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



# Neural network architecture



# Example feedforward computation



```
# forward-pass of a 3-layer neural network:  
f = lambda x: 1.0/(1.0 + np.exp(-x)) # activation function (use sigmoid)  
x = np.random.randn(3, 1) # random input vector of three numbers (3x1)  
h1 = f(np.dot(W1, x) + b1) # calculate first hidden layer activations (4x1)  
h2 = f(np.dot(W2, h1) + b2) # calculate second hidden layer activations (4x1)  
out = np.dot(W3, h2) + b3 # output neuron (1x1)
```

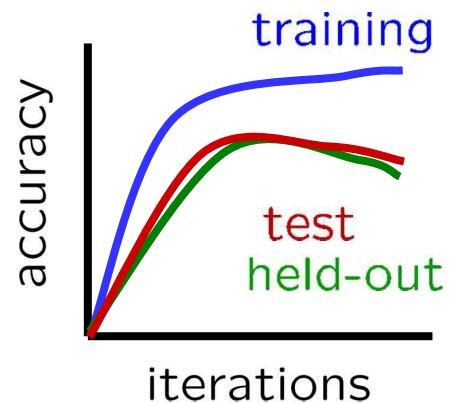
# 2-layer neural network ~20 lines

```
1 import numpy as np
2 from numpy.random import randn
3
4 N, D_in, H, D_out = 64, 1000, 100, 10
5 x, y = randn(N, D_in), randn(N, D_out)
6 w1, w2 = randn(D_in, H), randn(H, D_out)
7
8 for t in range(2000):
9     h = 1 / (1 + np.exp(-x.dot(w1)))
10    y_pred = h.dot(w2)
11    loss = np.square(y_pred - y).sum()
12    print(t, loss)
13
14    grad_y_pred = 2.0 * (y_pred - y)
15    grad_w2 = h.T.dot(grad_y_pred)
16    grad_h = grad_y_pred.dot(w2.T)
17    grad_w1 = x.T.dot(grad_h * h * (1 - h))
18
19    w1 -= 1e-4 * grad_w1
20    w2 -= 1e-4 * grad_w2
```

# Preventing Overfitting in Neural Networks

---

Early stopping:



Weight regularization

# Weight Regularization

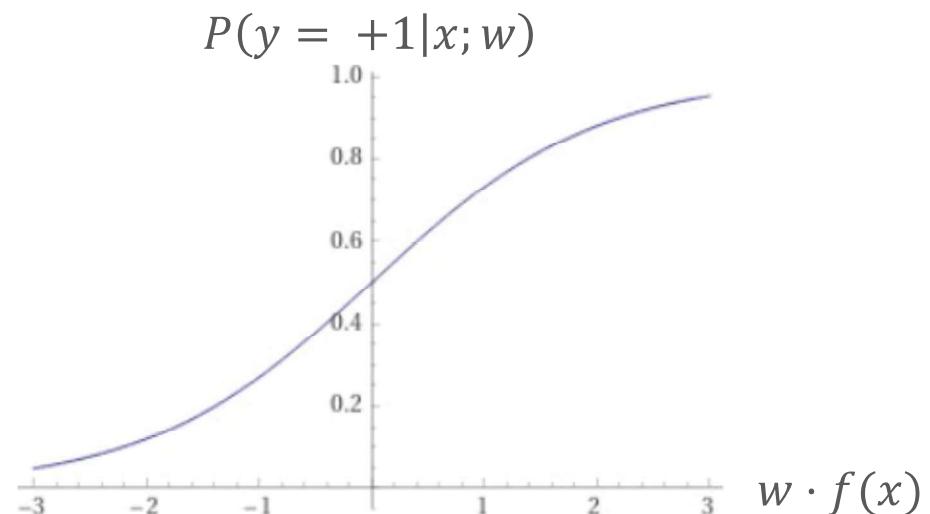
What can go wrong when we maximize log-likelihood?

Example: logistic regression

$$\max_w \sum_i \log P(y^{(i)} | x^{(i)}; w)$$

- $P(y = +1|x; w) = \frac{1}{1+e^{-w \cdot f(x)}}$
- $P(y = -1|x; w) = 1 - \frac{1}{1+e^{-w \cdot f(x)}}$

$w$  can grow very large and lead to overfitting and learning instability



# Weight Regularization

---

What can go wrong when we maximize log-likelihood?

$$\max_w \sum_i \log P(y^{(i)} | x^{(i)}; w)$$

$w$  can grow very large

Solution: add an objective term to penalize weight magnitude

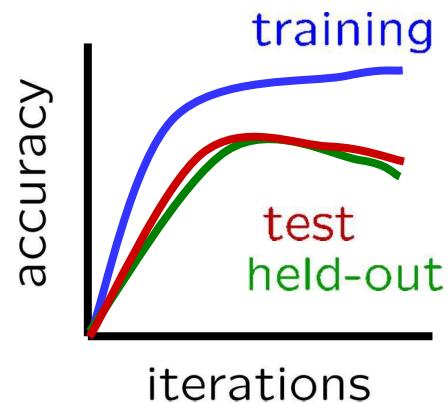
$$\max_w \sum_i \log P(y^{(i)} | x^{(i)}; w) - \frac{\lambda}{2} \sum_j w_j^2$$

$\lambda$  is a hyperparameter (typically 0.1 to 0.0001 or smaller)

# Preventing Overfitting in Neural Networks

---

Early stopping:



Weight regularization:  $\max_w \sum_i \log P(y^{(i)} | x^{(i)}; w) - \frac{\lambda}{2} \sum_j w_j^2$

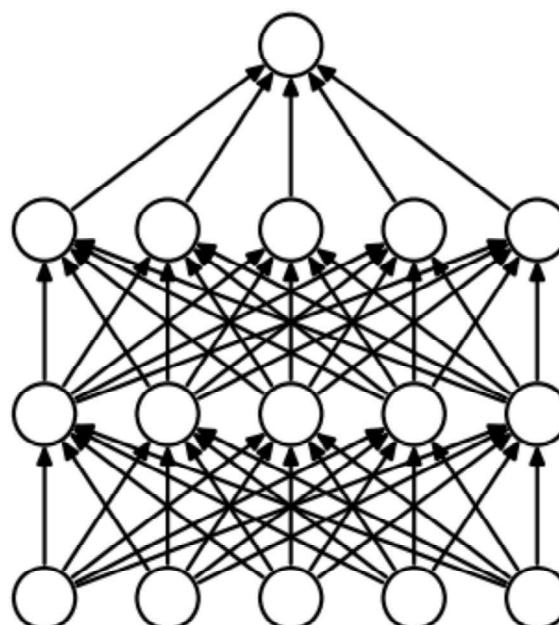
Dropout

# Dropout\*

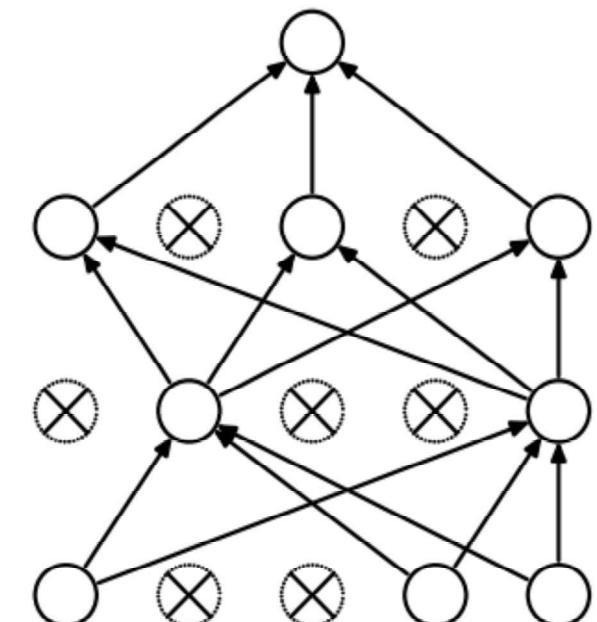
“Damage” the network during training to encourage redundancy

At each training step, with probability  $(1-p)$  set an activation to zero (drop it)

After training, don’t drop, but multiply weights by  $p$



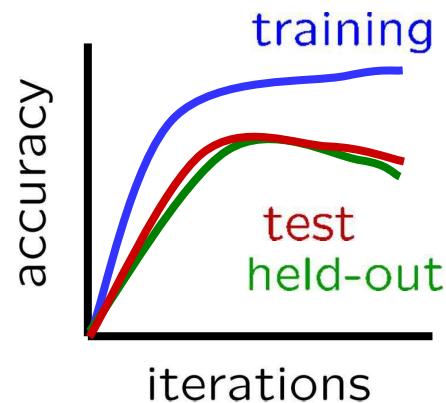
(a) Standard Neural Net



(b) After applying dropout.

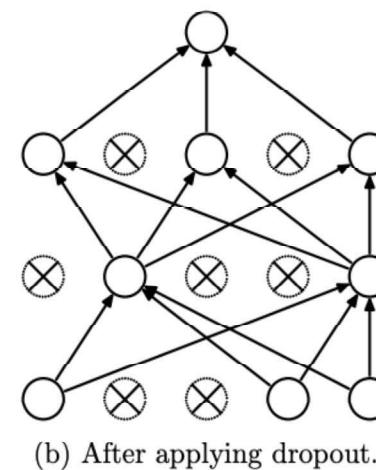
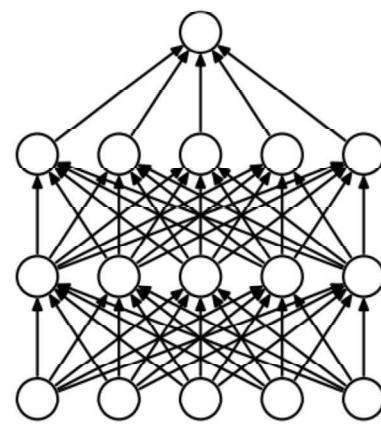
# Preventing Overfitting in Neural Networks

Early stopping:



Weight regularization:  $\max_w \sum_i \log P(y^{(i)} | x^{(i)}; w) - \frac{\lambda}{2} \sum_j w_j^2$

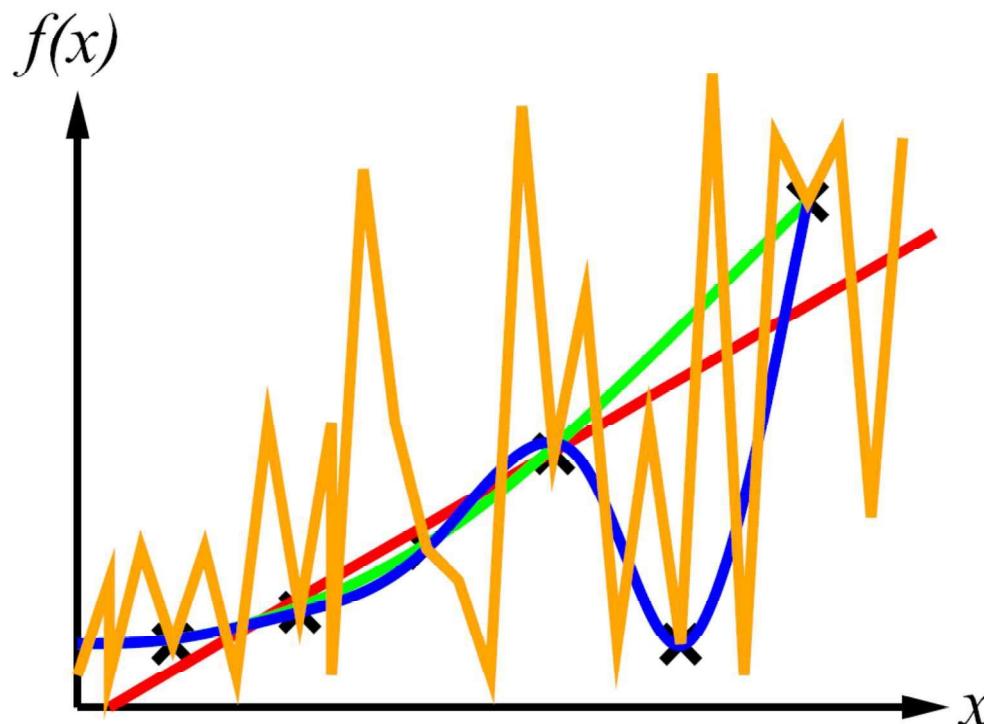
Dropout:



# Consistency vs. Simplicity

---

- Example: curve fitting (regression, function approximation)



- Consistency vs. simplicity
- Ockham's razor

# Consistency vs. Simplicity

---

- Fundamental tradeoff: bias vs. variance
- Usually algorithms prefer consistency by default (why?)
- Several ways to operationalize “simplicity”
  - Reduce the **hypothesis/model space**
    - Assume more: e.g. independence assumptions, as in naïve Bayes
    - Fewer features or neurons
    - Other limits on model structure
  - **Regularization**
    - Laplace Smoothing: cautious use of small counts
    - Small weight vectors in neural networks (stay close to zero-mean prior)
    - Hypothesis space stays big, but harder to get to the outskirts

# Fun Neural Net Demo Site

---

Demo-site:

<http://playground.tensorflow.org/>

# Summary of Key Ideas

---

Optimize probability of label given input

$$\max_w \text{ll}(w) = \max_w \sum_i \log P(y^{(i)}|x^{(i)}; w)$$

## Continuous optimization

Gradient ascent:

Compute steepest uphill direction = gradient (= just vector of partial derivatives)

Take step in the gradient direction

Repeat (until held-out data accuracy starts to drop = “early stopping”)

## Deep neural nets

Last layer = still logistic regression

Now also many more layers before this last layer

= computing the features

the features are learned rather than hand-designed

## Universal function approximation theorem

If neural net is large enough

Then neural net can represent any continuous mapping from input to output with arbitrary accuracy

But remember: need to avoid overfitting / memorizing the training data ? early stopping!

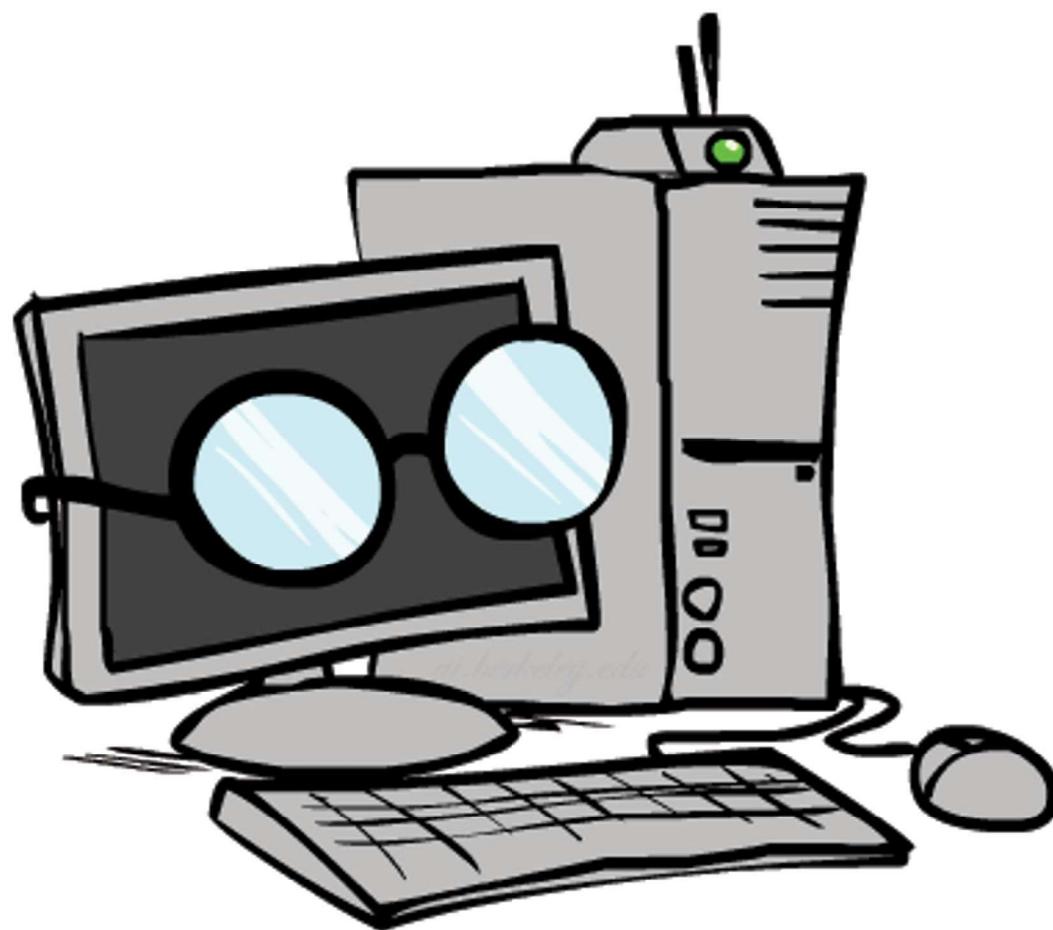
Automatic differentiation gives the derivatives efficiently (how? = outside of scope of 188)

# How well does deep learning work?

---

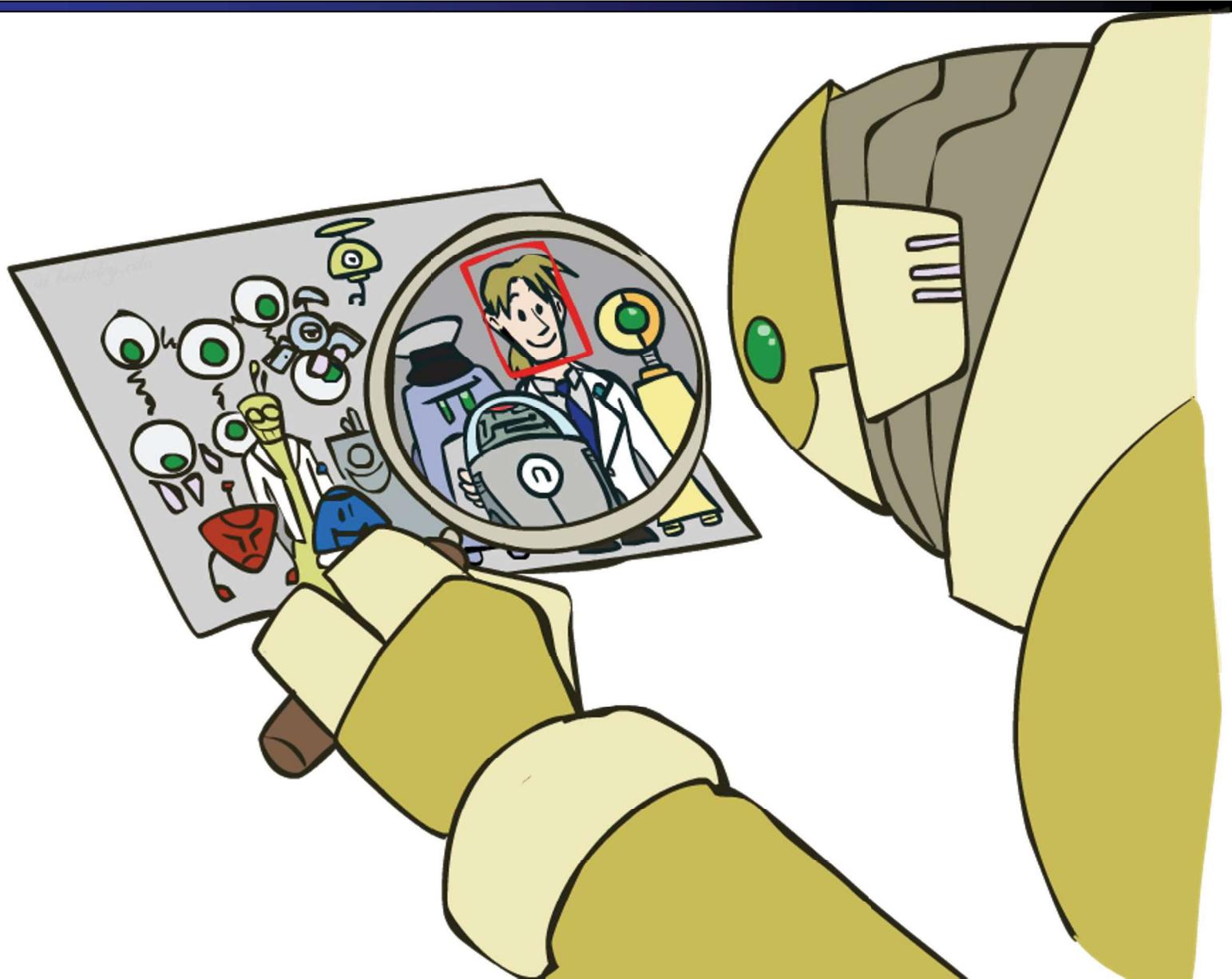
# Computer Vision

---



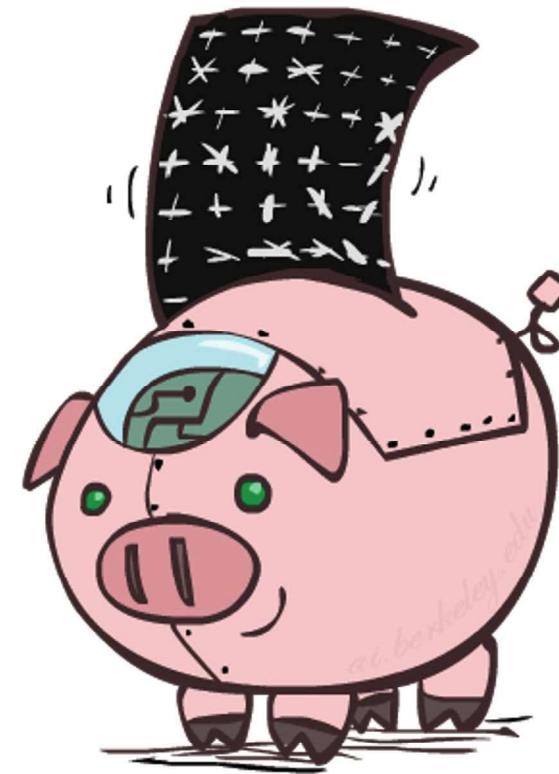
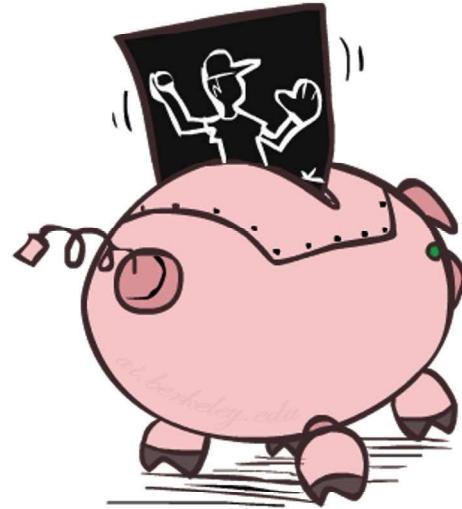
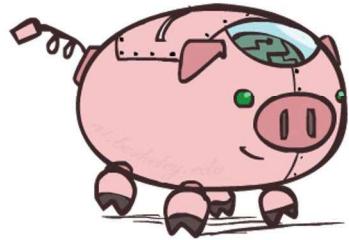
# Object Detection

---



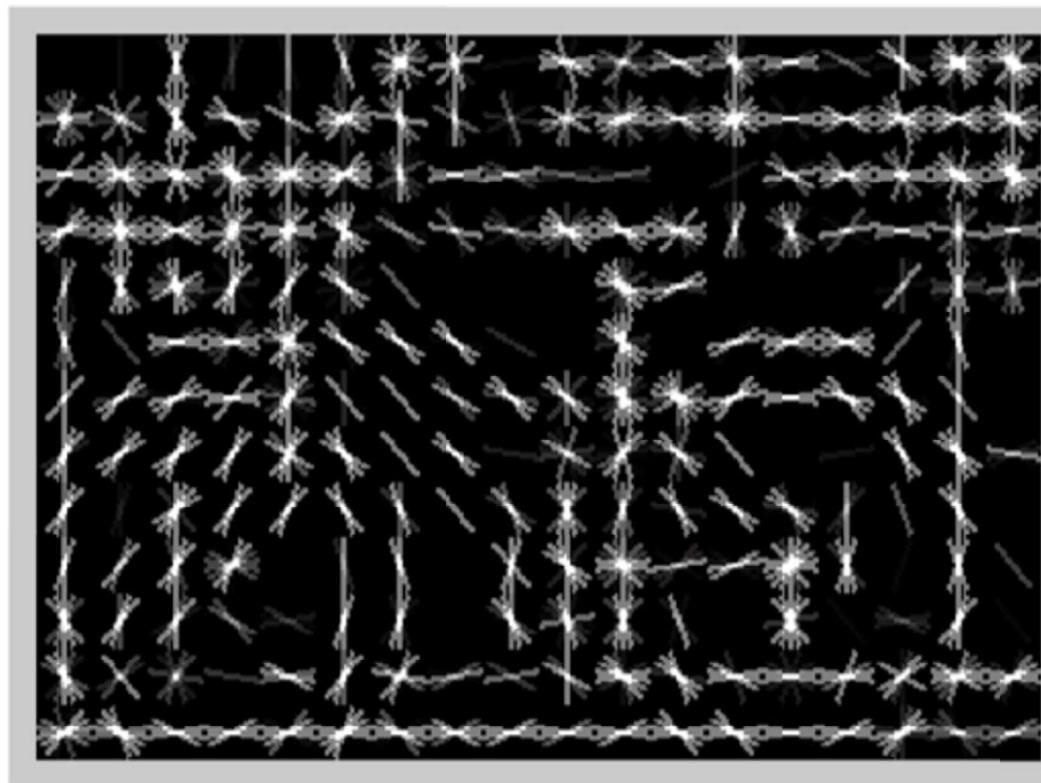
# Manual Feature Design

---



# Features and Generalization

---



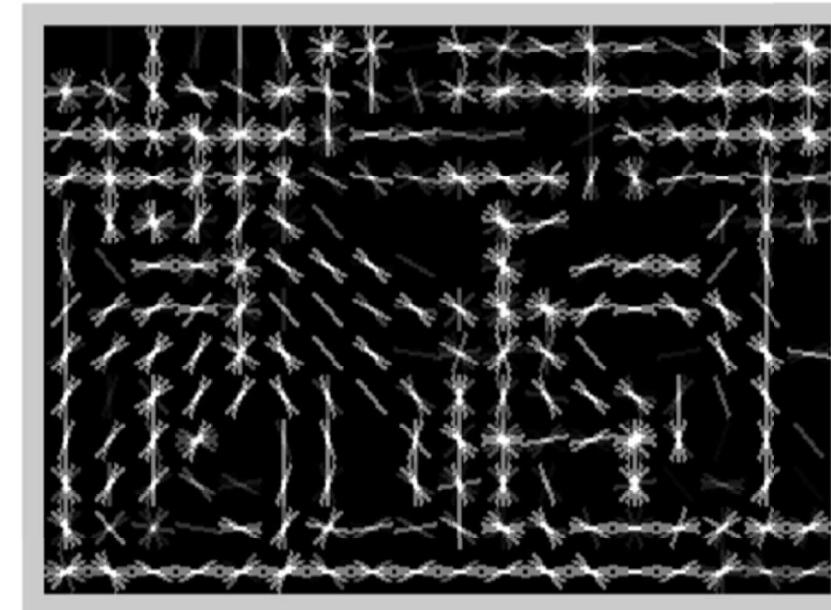
[HoG: Dalal and Triggs, 2005]

# Features and Generalization

---



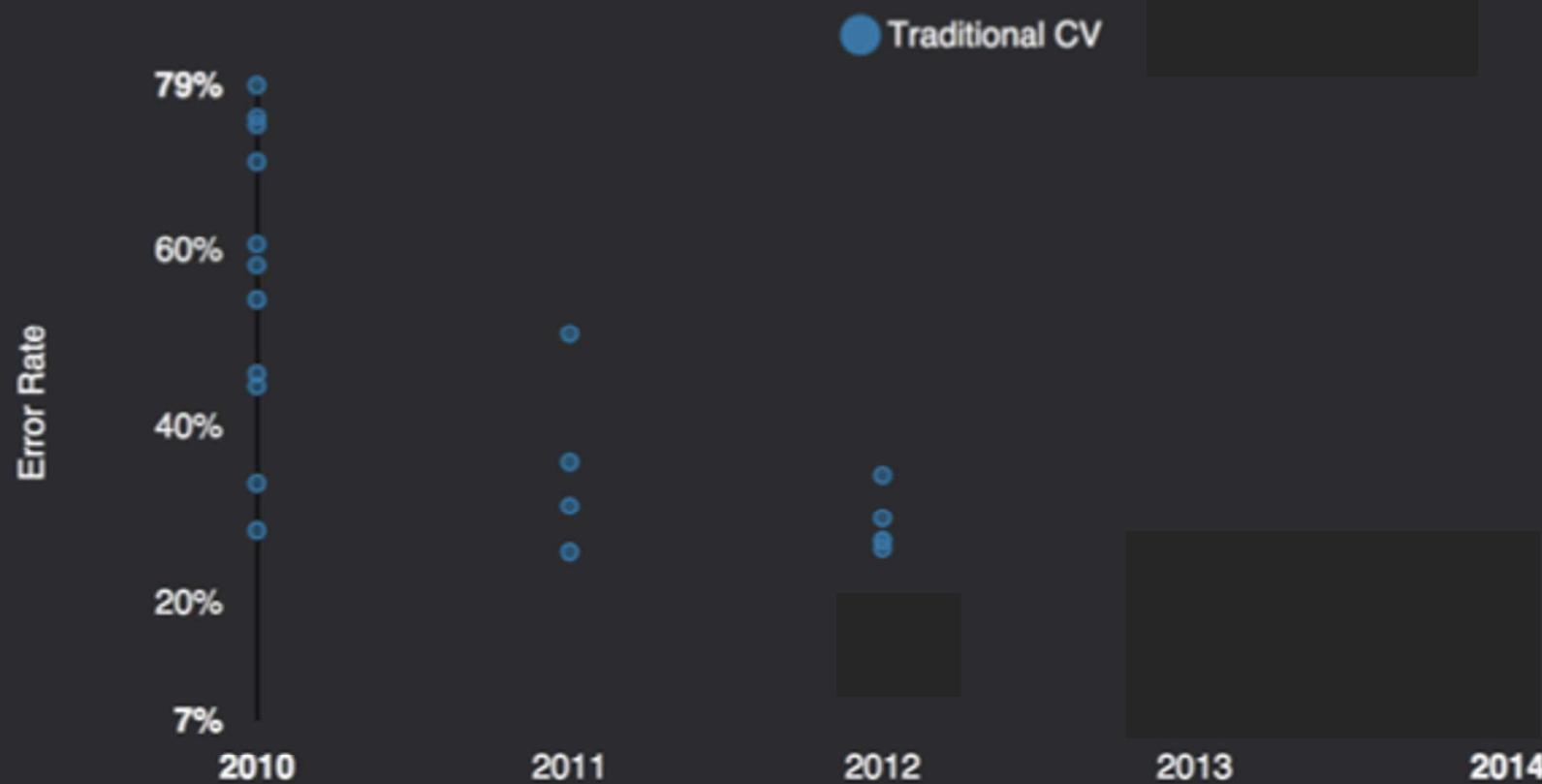
Image



HoG

# Performance

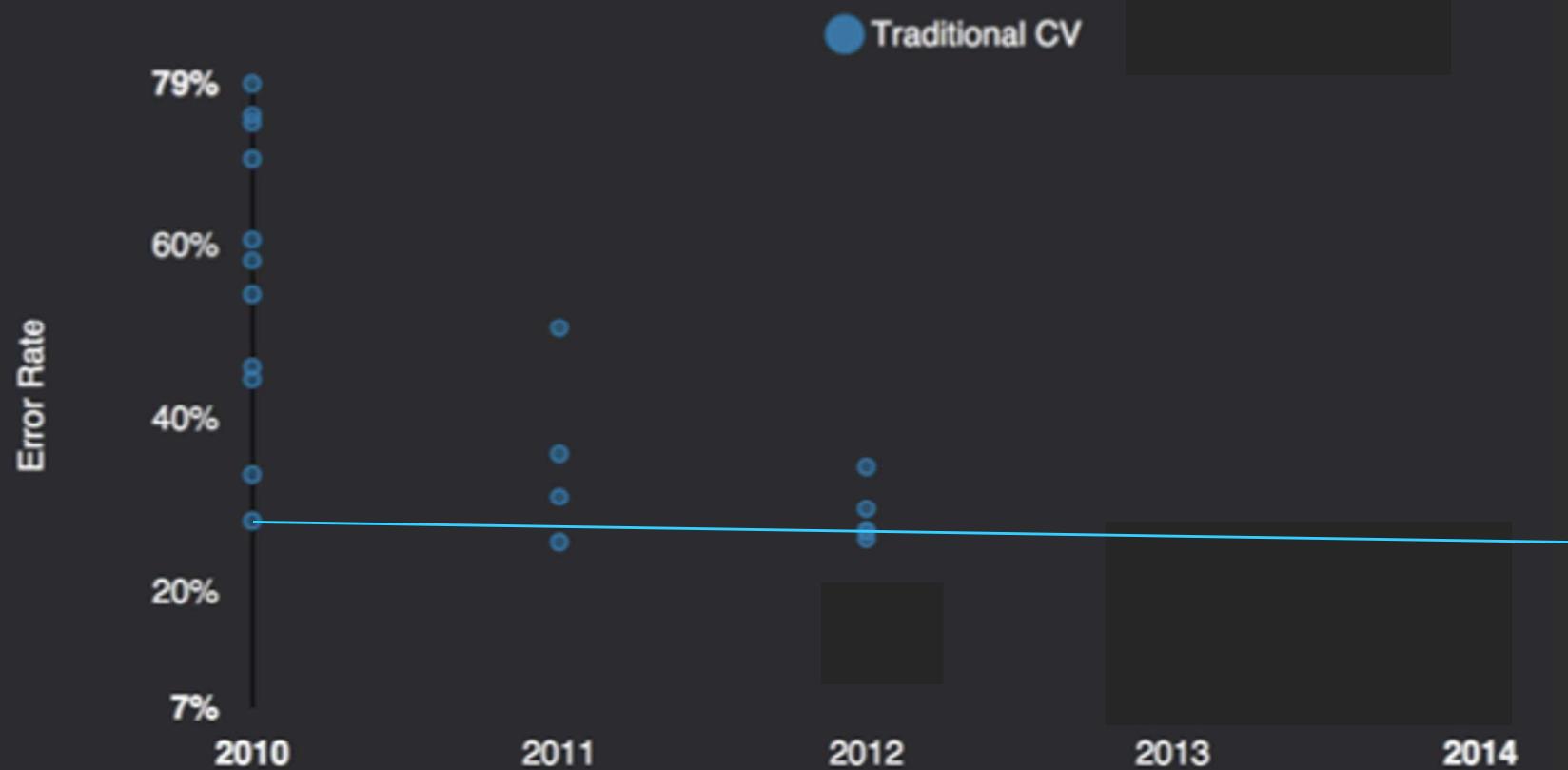
## ImageNet Error Rate 2010-2014



graph credit Matt  
Zeiler, Clarifai

# Performance

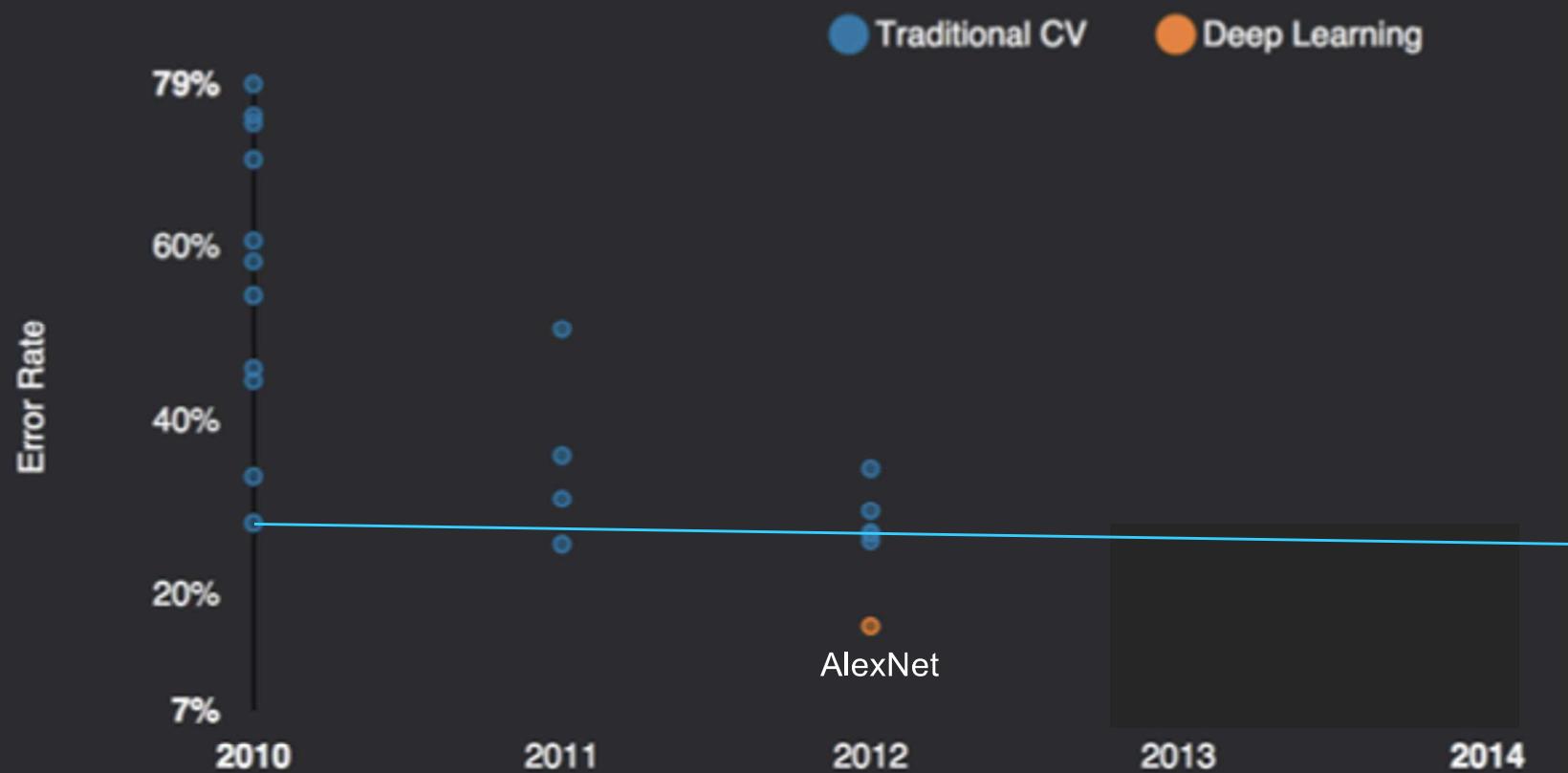
## ImageNet Error Rate 2010-2014



graph credit Matt  
Zeiler, Clarifai

# Performance

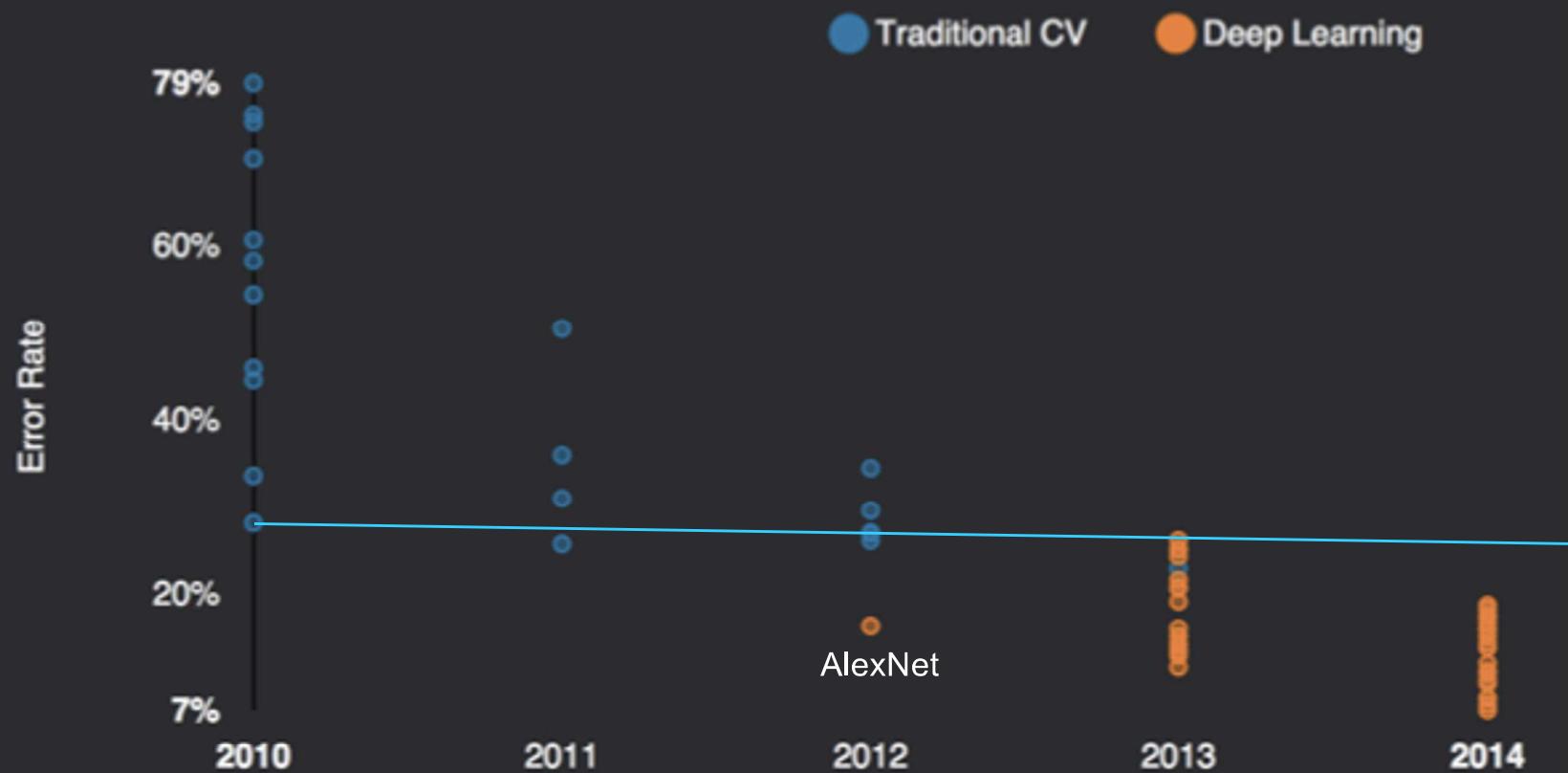
## ImageNet Error Rate 2010-2014



graph credit Matt  
Zeiler, Clarifai

# Performance

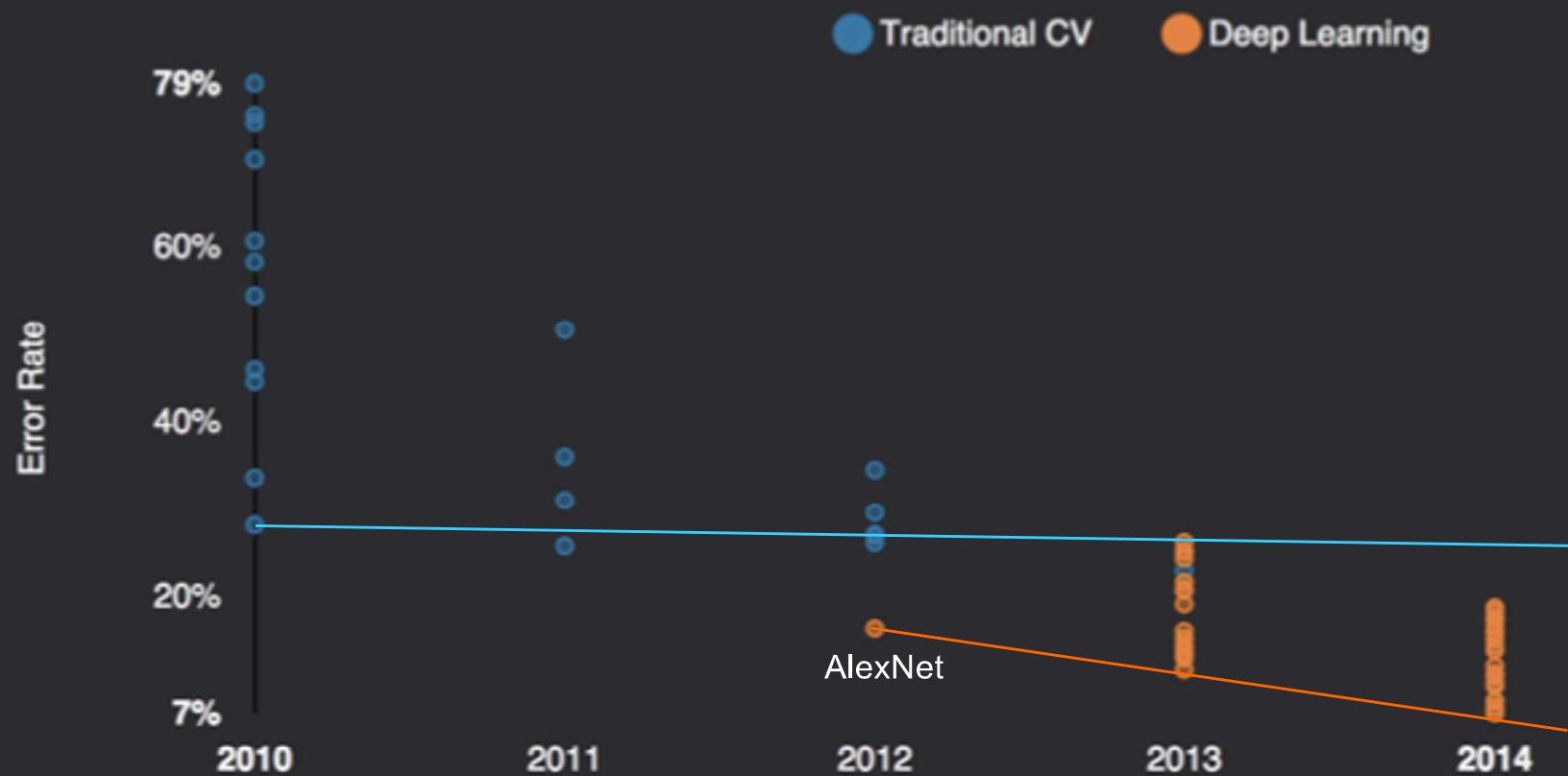
## ImageNet Error Rate 2010-2014



graph credit Matt  
Zeiler, Clarifai

# Performance

## ImageNet Error Rate 2010-2014



graph credit Matt  
Zeiler, Clarifai

# Papers With Code: ImageNet

Leaderboard

Dataset

View

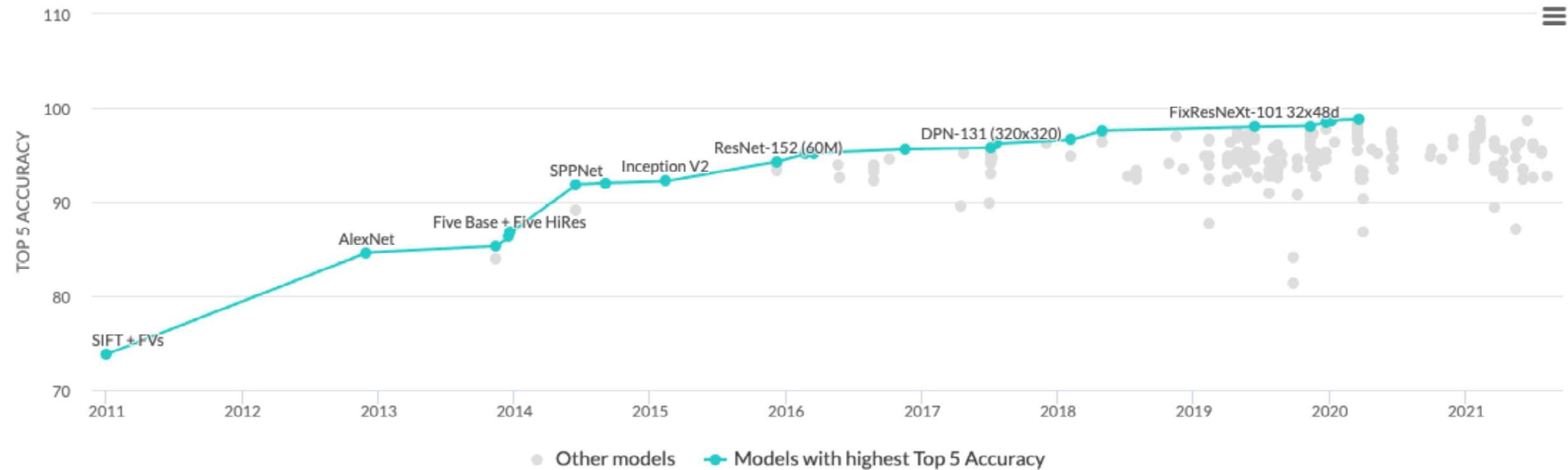
Top 5 Accuracy

by

Date

for

All models



# MS COCO Image Captioning Challenge



"man in black shirt is playing guitar."



"construction worker in orange safety vest is working on road."



"two young girls are playing with lego toy."



"boy is doing backflip on wakeboard."



"girl in pink dress is jumping in air."



"black and white dog jumps over bar."



"young girl in pink shirt is swinging on swing."



"man in blue wetsuit is surfing on wave."

Karpathy & Fei-Fei, 2015; Donahue et al., 2015; Xu et al, 2015; many more

# Visual QA Challenge

Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, Devi Parikh



What vegetable is on the plate?  
Neural Net: **broccoli**  
Ground Truth: broccoli



What color are the shoes on the person's feet ?  
Neural Net: **brown**  
Ground Truth: brown



How many school busses are there?  
Neural Net: **2**  
Ground Truth: 2



What sport is this?  
Neural Net: **baseball**  
Ground Truth: baseball



What is on top of the refrigerator?  
Neural Net: **magnets**  
Ground Truth: cereal



What uniform is she wearing?  
Neural Net: **shorts**  
Ground Truth: girl scout



What is the table number?  
Neural Net: **4**  
Ground Truth: 40



What are people sitting under in the back?  
Neural Net: **bench**  
Ground Truth: tent

# Visual Dialogue

Alayrac et al, 2022



This is a picture of two teddy bears on the moon.

What are they doing?

They are having a conversation.

What object are they using?

It looks like a computer.

Is this surprising?

Yes, it is surprising.

Why is this picture surprising to you?

I think it is surprising because teddy bears are not usually found on the moon.



What is the common thing about these three images?

They are all flamingos.

What is the difference between these three images?

The first one is a cartoon, the second one is a real flamingo, and the third one is a 3D model of a flamingo.



This is an apple with a sticker on it.

What does the sticker say?

The sticker says "iPod".

Where is the photo taken?

It looks like it's taken in a backyard.

Do you think it is printed or handwritten?

It looks like it's handwritten.

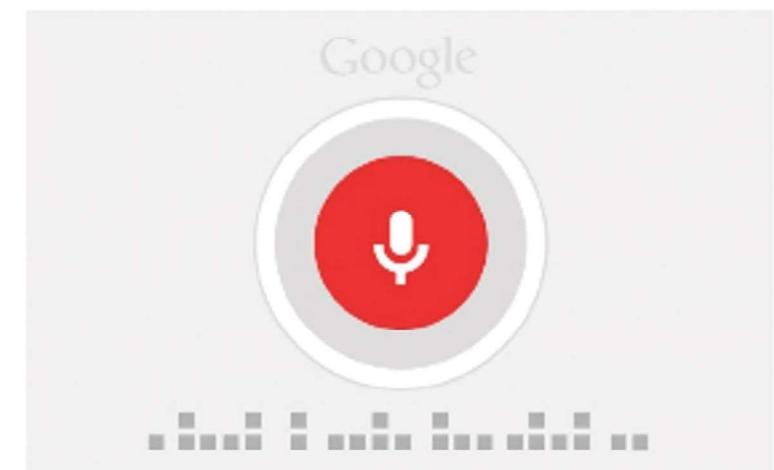
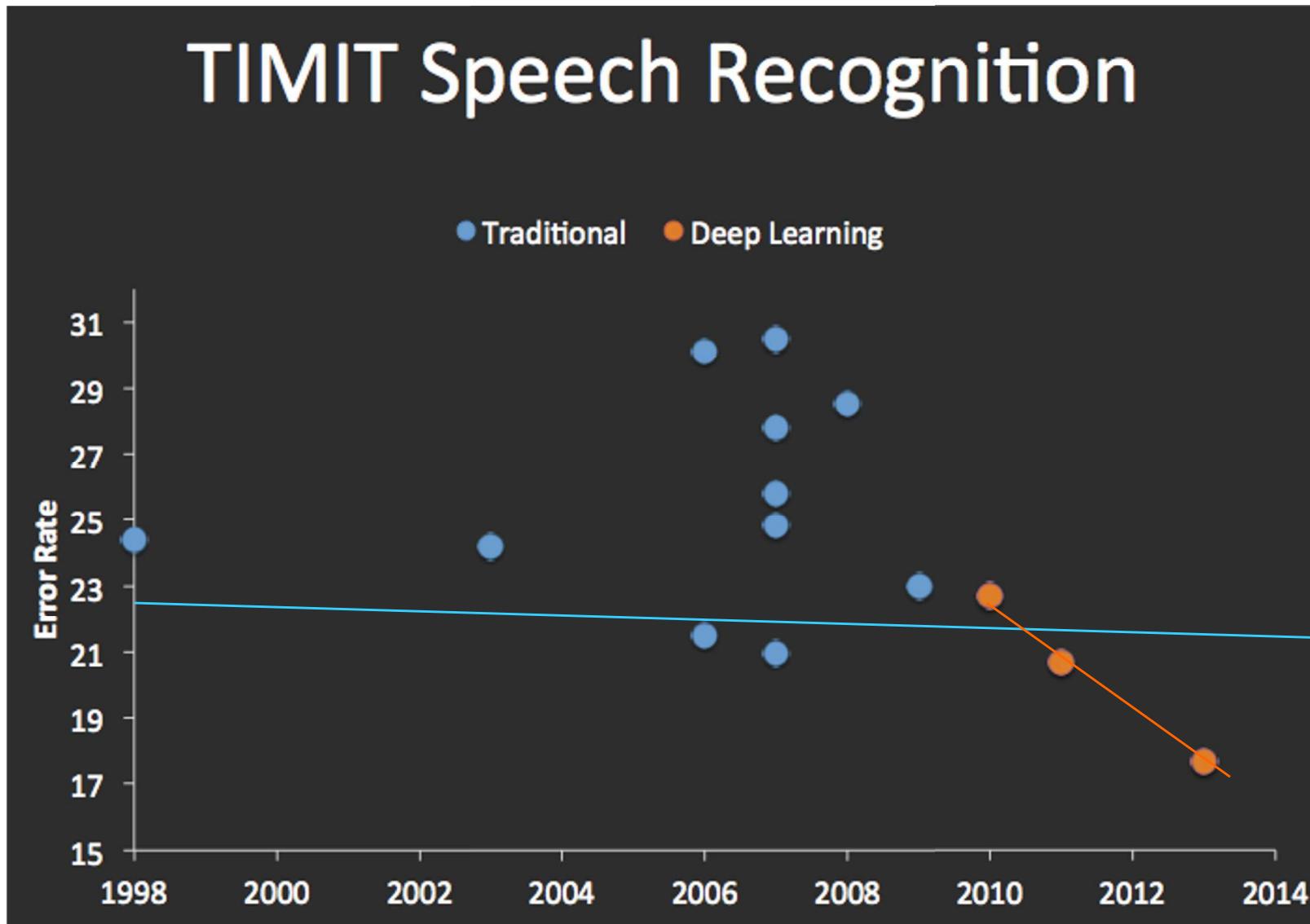
What color is the sticker?

It's white.

# Image Segmentation



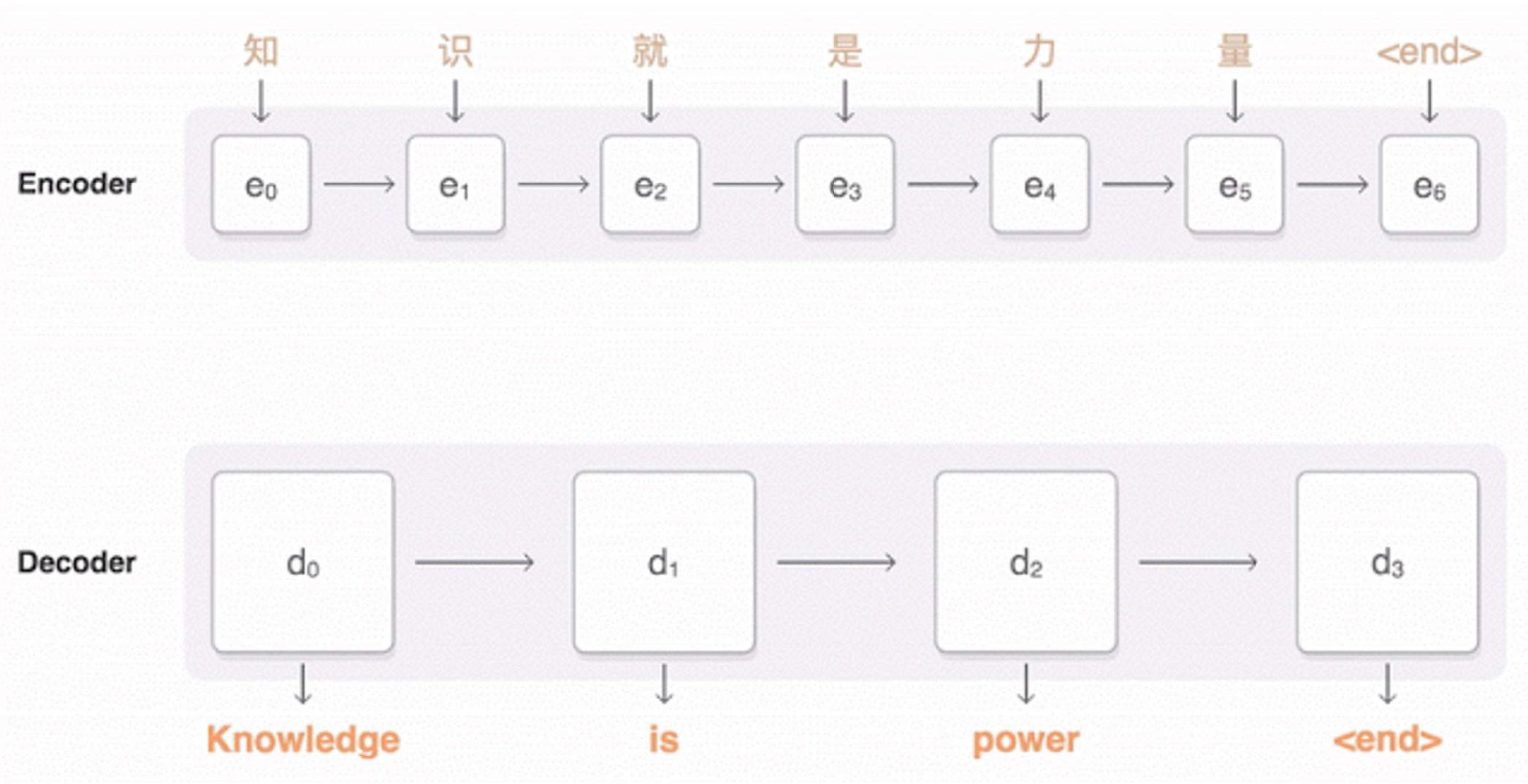
# Speech Recognition



graph credit Matt Zeiler, Clarifai

# Machine Translation

Google Neural Machine Translation (in production)



# Google and DeepMind are using AI to predict the energy output of wind farms

*To help make that energy more valuable to the power grid*

By [Nick Statt](#) | [@nickstatt](#) | Feb 26, 2019, 2:42pm EST



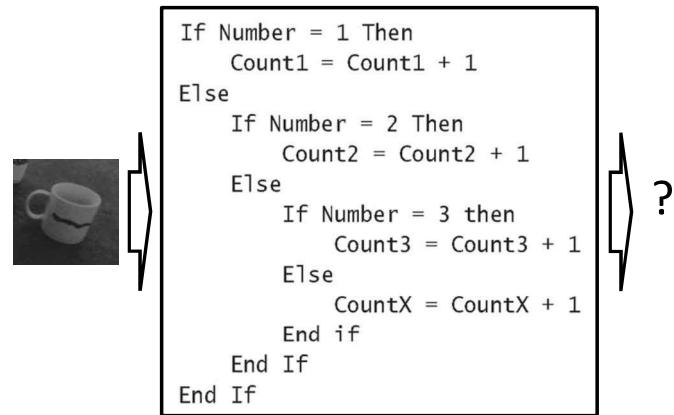
 SHARE



Google [announced today](#) that it has made energy produced by wind farms more viable using the artificial intelligence software of its London-based subsidiary DeepMind. By using DeepMind's machine learning algorithms to predict the wind output from the farms Google uses for its green energy initiatives, the company says it can now schedule set deliveries of energy output, which are more valuable to the grid than standard, non-time-based deliveries.

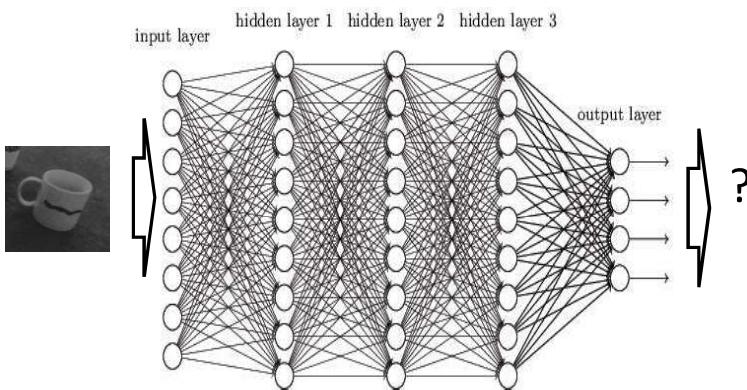
# Change in Programming Paradigm?

Traditional Programming:  
program by writing lines of code



Poor performance on AI problems

Deep Learning (“Software 2.0”):  
program by providing data



Success!