

# CS5670: Computer Vision

## Convolutional neural networks, Part II

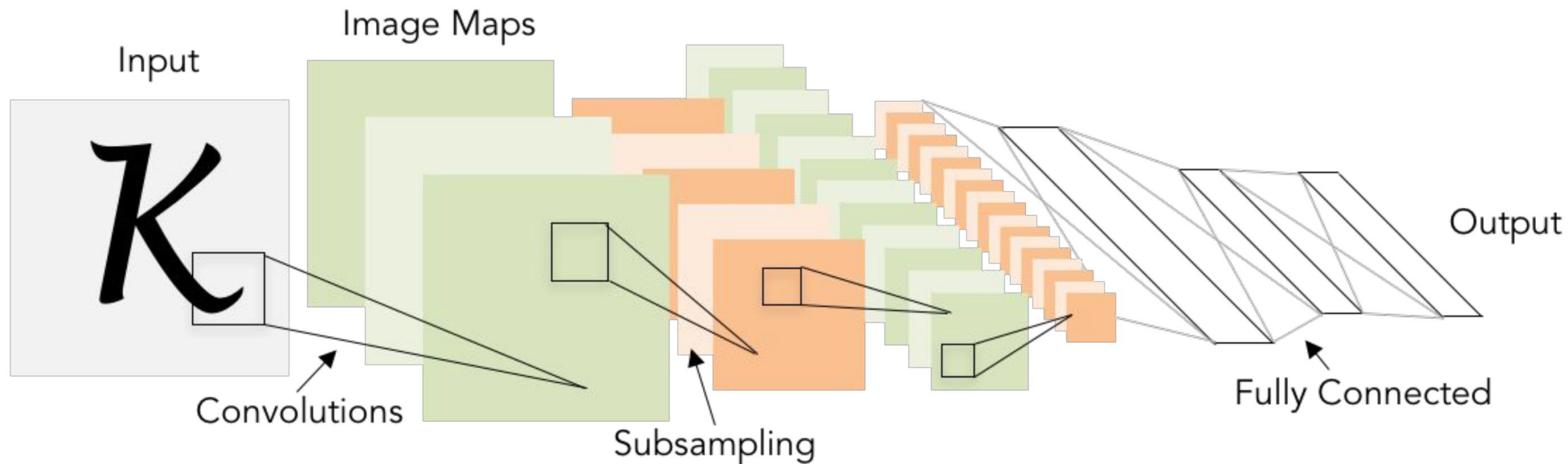


Illustration of LeCun et al. 1998 from CS231n 2017 Lecture 1

# Announcements

- Project 5 (Convolutional Neural Networks) released today
  - Due Tuesday, May 11, 2021 (7:00 pm)
- Take-home final exam to be released Wednesday, May 12, 2021; due Monday, May 17, 2021
- Sample final available on Ed Stem

# Readings

- Convolutional neural networks
  - <http://cs231n.github.io/convolutional-networks/>
- Stochastic Gradient Descent & Backpropagation
  - <http://cs231n.github.io/optimization-1/>
  - <http://cs231n.github.io/optimization-2/>
- Best practices for training CNNs
  - <http://cs231n.github.io/neural-networks-2/>
  - <http://cs231n.github.io/neural-networks-3/>

# **Project 5 Demo (Ruojin)**

# Last time

- Neural networks
- Convolutional neural networks

# Today

- Convolutional neural networks (continued)
- Training neural networks with backpropagation
- Stochastic gradient descent
- Data processing and augmentation
- CNN architectures
- Transfer learning

# Image Classification: a core task in computer vision

- Assume given set of discrete labels, e.g.  
{cat, dog, cow, apple, tomato, truck, ... }

$f(\text{apple image}) = \text{"apple"}$

$f(\text{tomato image}) = \text{"tomato"}$

$f(\text{cow image}) = \text{"cow"}$

# Recap: Neural networks

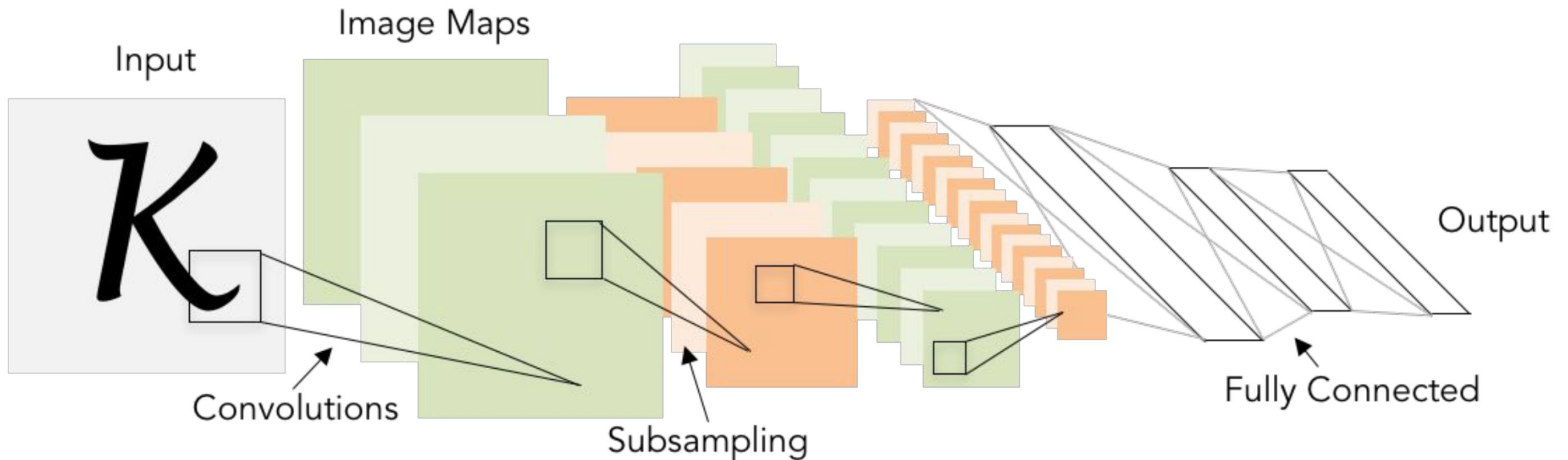
- Very coarse generalization:
  - Linear functions chained together and separated by nonlinearities (*activation functions*), e.g. "max"

$$f = W_3 \max(0, W_2 \max(0, W_1 x))$$



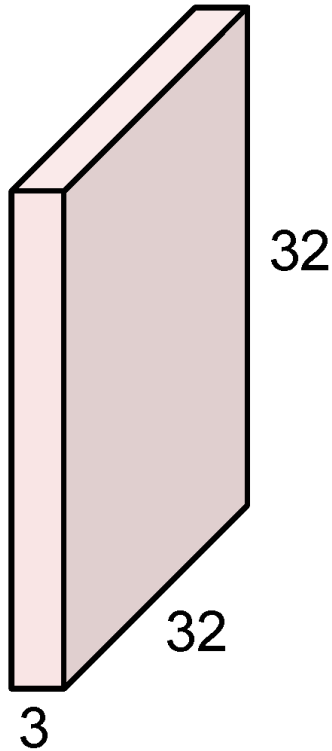
# Convolutional neural networks

- Made up of many layers of a few different types (mainly *convolution layers*)



# Convolutions as network layers

32x32x3 image



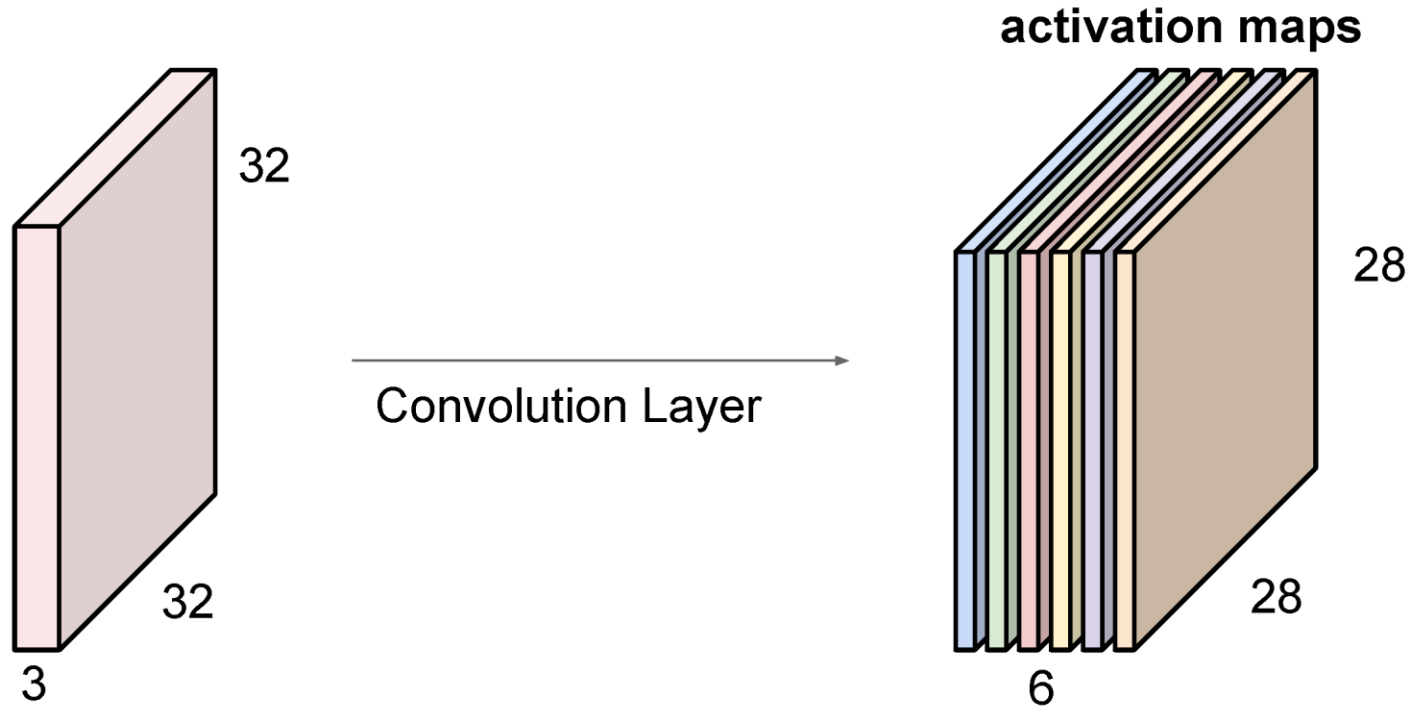
5x5x3 filter (weights are learned)



**Convolve** the filter with the image  
i.e. “slide over the image spatially,  
computing dot products”

# Convolutional layer

For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

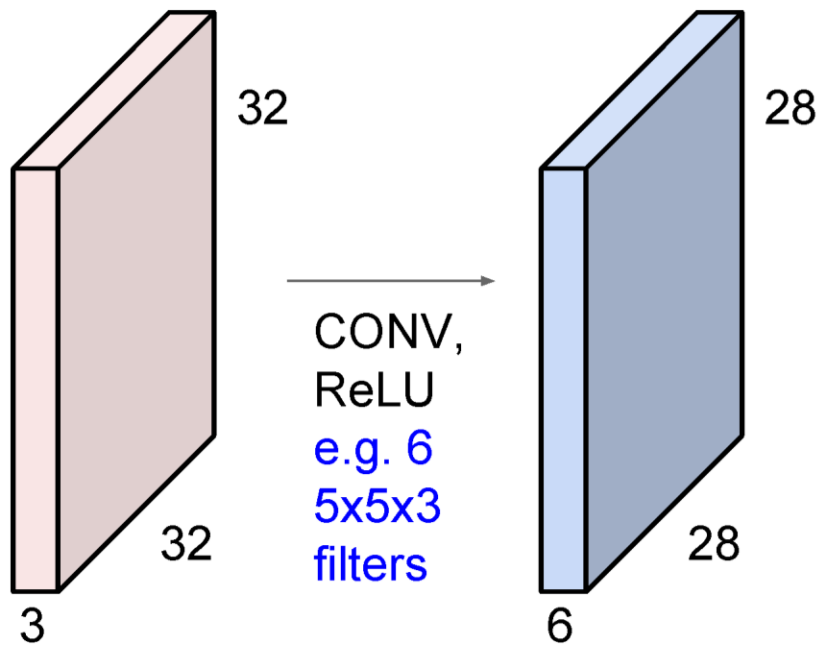


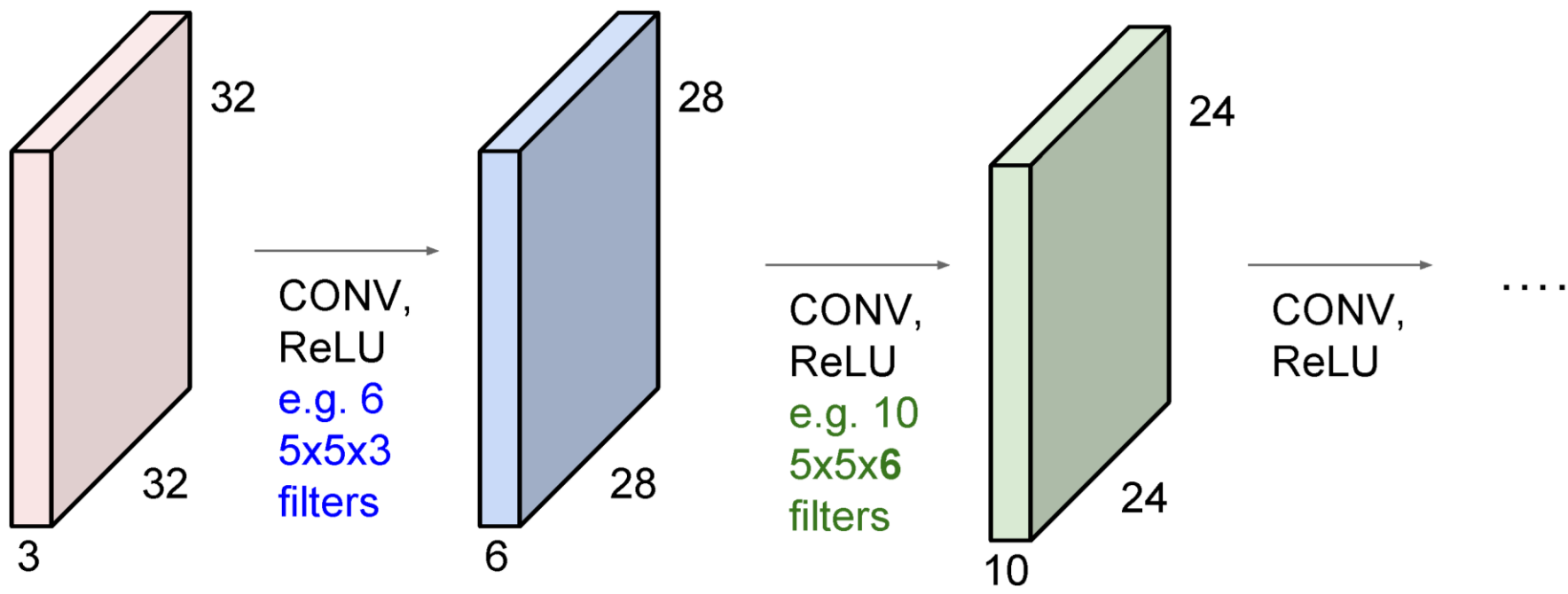
We stack these up to get a “new image” of size 28x28x6!

(total number of parameters:  $6 \times (75 + 1) = \mathbf{456}$ )

# Convolution layer parameters

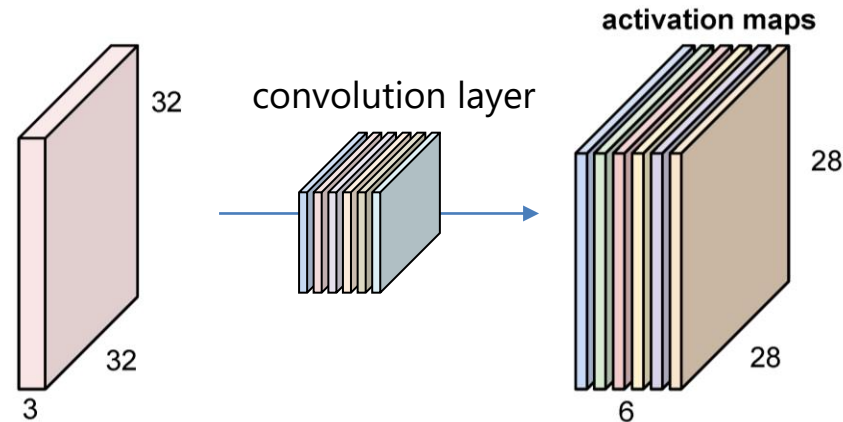
- Kernel size
- Number of kernels
- Stride



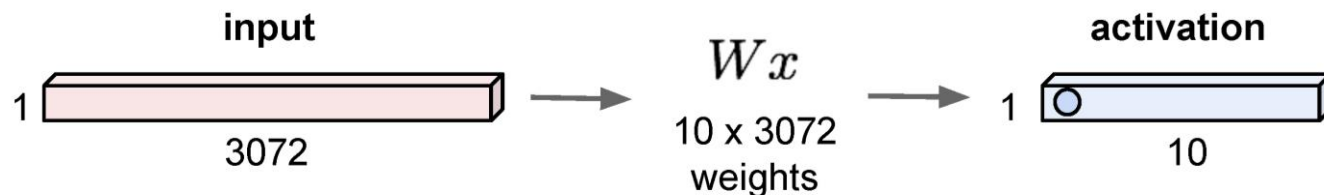


# Some convolutional network layer types

- **Convolution layers** (some parameters)

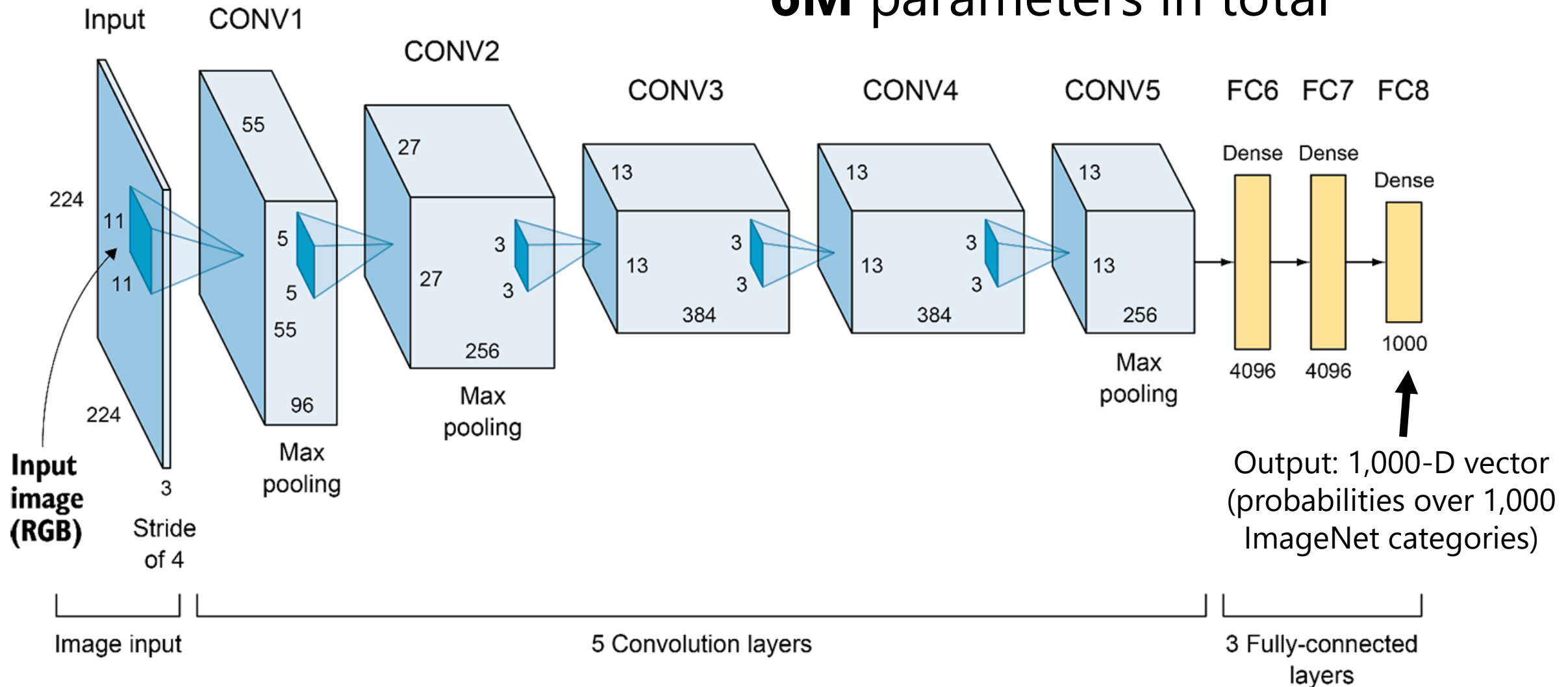


- Pooling layers (no parameters)
- Fully connected layers (many, many parameters)



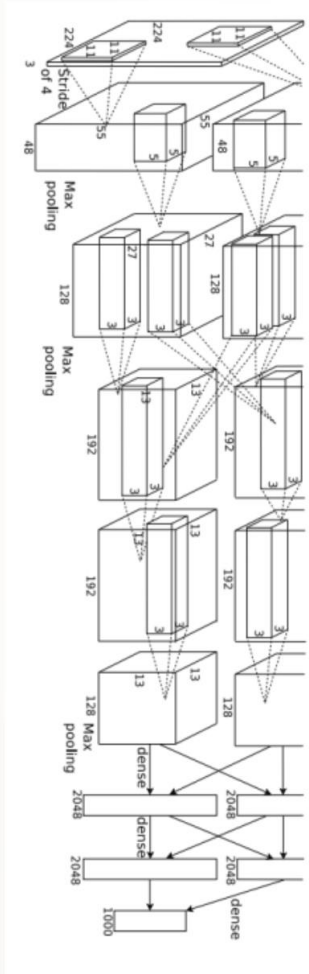
# AlexNet (2012)

**6M** parameters in total



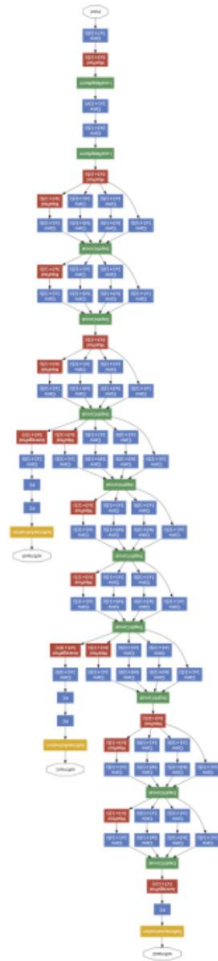


# “AlexNet”



[Krizhevsky et al. NIPS 2012]

# “GoogLeNet”



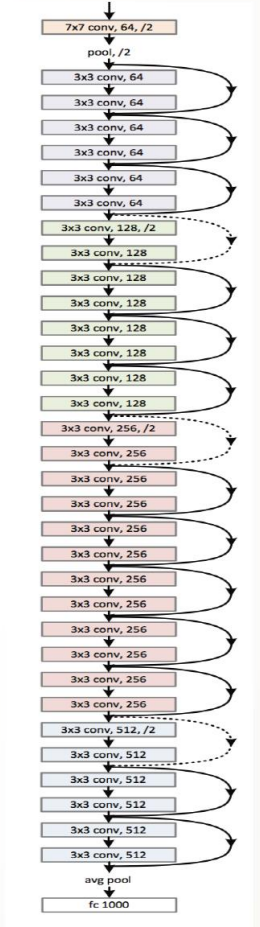
[Szegedy et al. CVPR 2015]

# “VGG Net”



[Simonyan & Zisserman, ICLR 2015]

# “ResNet”



[He et al. CVPR 2016]

# Big picture

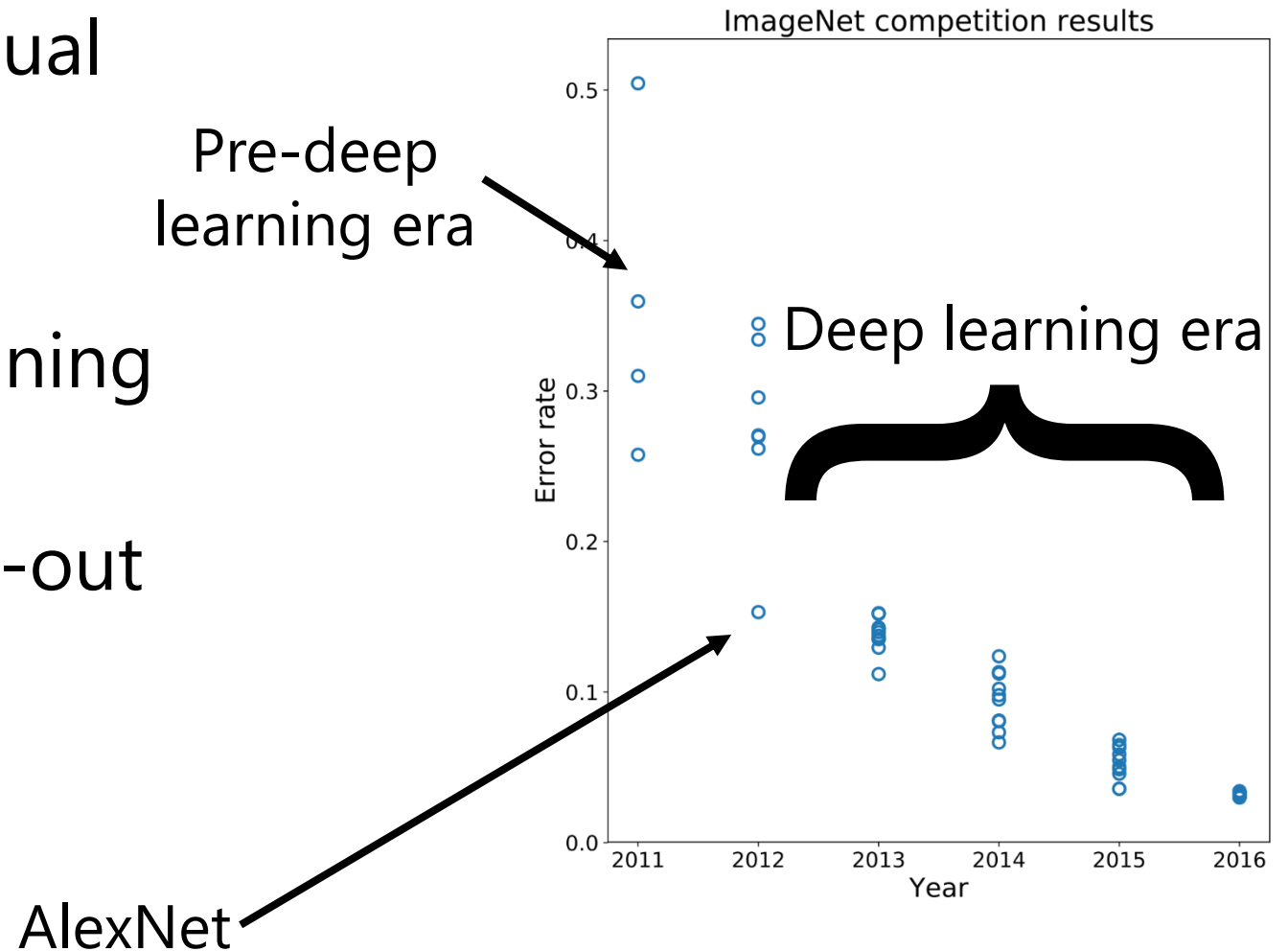
- A convolutional neural network can be thought of as a function from images to class scores
  - With millions of adjustable weights...
  - ... leading to a very non-linear mapping from images to features / class scores.
  - We will set these weights based on classification accuracy on training data...
  - ... and hopefully our network will generalize to new images at test time

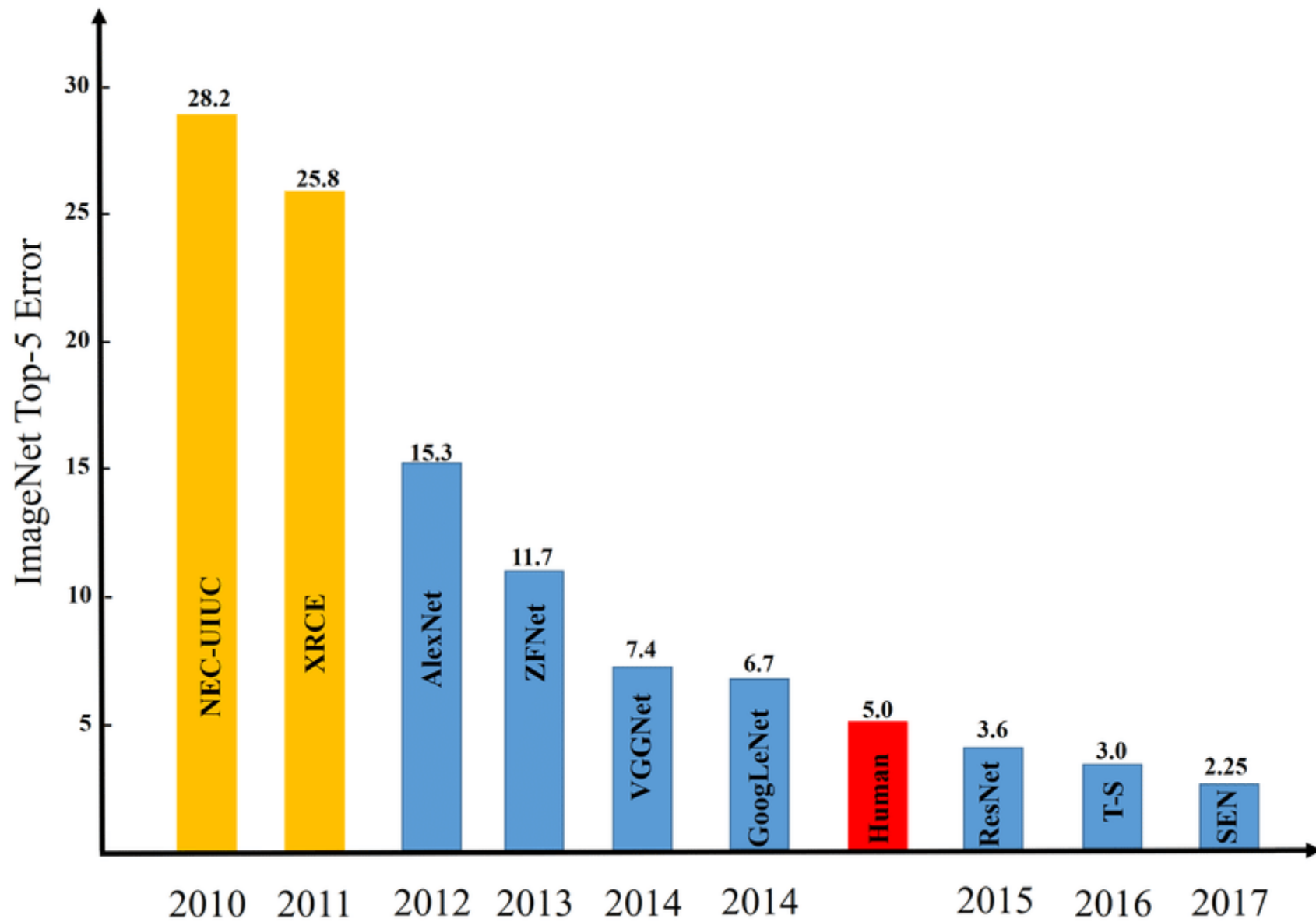
# Data is key—enter ImageNet

- ImageNet (and the ImageNet Large-Scale Visual Recognition Challenge, aka **ILSVRC**) has been key to training deep learning methods
  - J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li and L. Fei-Fei, **ImageNet: A Large-Scale Hierarchical Image Database**. *CVPR*, 2009.
- **ILSVRC**: 1,000 object categories, each with ~700-1300 training images. Test set has 100 images per category (100,000 total).
- Standard ILSVRC error metric: top-5 error
  - if the correct answer for a given test image is in the top 5 categories, your answer is judged to be correct

# Performance improvements on ILSVRC

- ImageNet Large-Scale Visual Recognition Challenge
- Held from 2011-2017
- 1000 categories, 1000 training images per category
- Test performance on held-out test set of images





**Questions?**

# Training the network

- Now we know what the structure of our function from images -> class scores is (a CNN)
- How do we set the weights given training data?

# How do we set the weights?

- Need to solve an optimization problem:
  - Find weights  $W$  that minimize training loss  $L$  over a training set
- In general this is a non-linear, non-convex problem
  - Closed-form solvers do not generally exist, unlike with e.g. least squares problems
  - Might not find the globally optimal weights
- (Side note: some learning problems, such as linear SVMs, do have convex loss functions)



# (Bad) idea #1: Random search

```
# assume X_train is the data where each column is an example (e.g. 3073 x 50,000)
# assume Y_train are the labels (e.g. 1D array of 50,000)
# assume the function L evaluates the loss function

bestloss = float("inf") # Python assigns the highest possible float value
for num in xrange(1000):
    W = np.random.randn(10, 3073) * 0.0001 # generate random parameters
    loss = L(X_train, Y_train, W) # get the loss over the entire training set
    if loss < bestloss: # keep track of the best solution
        bestloss = loss
        bestW = W
    print 'in attempt %d the loss was %f, best %f' % (num, loss, bestloss)

# prints:
# in attempt 0 the loss was 9.401632, best 9.401632
# in attempt 1 the loss was 8.959668, best 8.959668
# in attempt 2 the loss was 9.044034, best 8.959668
# in attempt 3 the loss was 9.278948, best 8.959668
# in attempt 4 the loss was 8.857370, best 8.857370
# in attempt 5 the loss was 8.943151, best 8.857370
# in attempt 6 the loss was 8.605604, best 8.605604
# ... (truncated: continues for 1000 lines)
```

Lets see how well this works on the test set...

```
# Assume X_test is [3073 x 10000], Y_test [10000 x 1]  
scores = Wbest.dot(Xte_cols) # 10 x 10000, the class scores for all test examples  
# find the index with max score in each column (the predicted class)  
Yte_predict = np.argmax(scores, axis = 0)  
# and calculate accuracy (fraction of predictions that are correct)  
np.mean(Yte_predict == Yte)  
# returns 0.1555
```

15.5% accuracy! not bad!  
(SOTA is ~95%)

# (Good) Idea #2: Gradient descent



# (Good) Idea #2: Gradient descent

In 1-dimension, the derivative of a function:

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

In multiple dimensions, the **gradient** is the vector of (partial derivatives) along each dimension

The slope in any direction is the **dot product** of the direction with the gradient  
The direction of steepest descent is the **negative gradient**

# Scores, losses, and gradients

- Function  $f$  maps images to class scores

$$s = f(x; W) = \cancel{Wx} \quad f \text{ is a deep CNN}$$

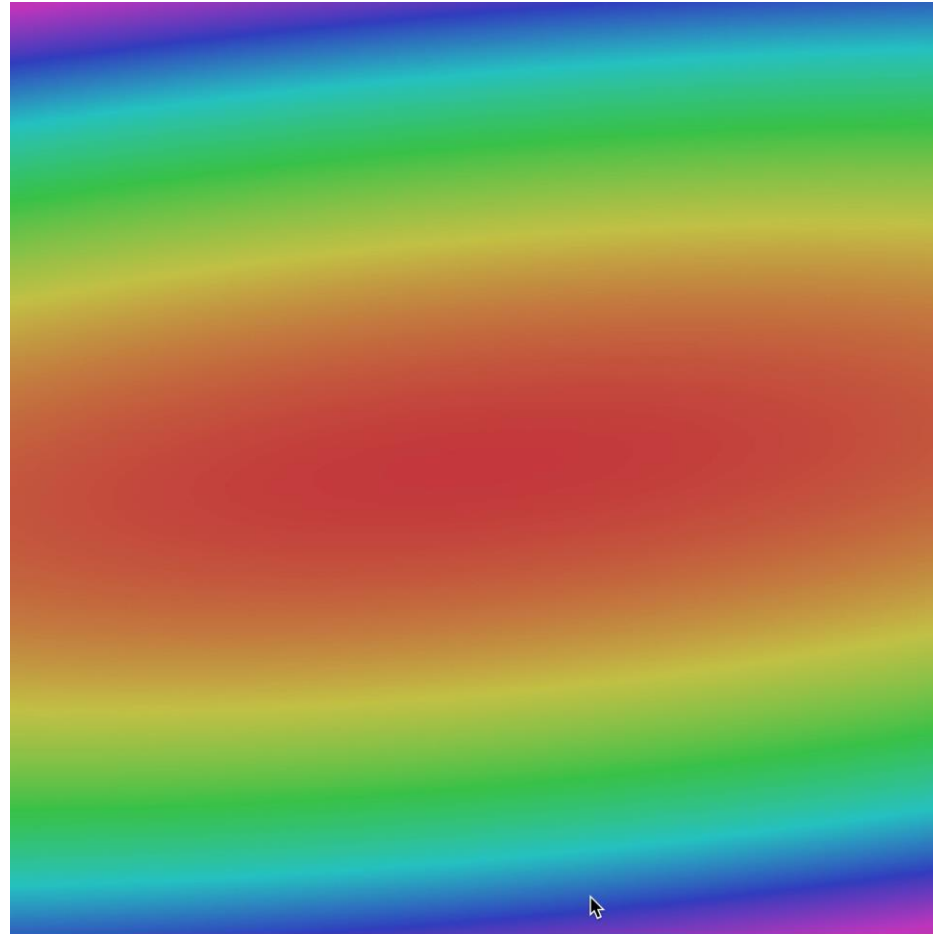
- Loss function maps class scores to “badness”

$$L_i = -\log \left( \frac{e^{f_{y_i}}}{\sum_j e^{f_j}} \right) \quad \text{Cross-entropy loss}$$

$$L = \frac{1}{N} \sum_{i=1}^N L_i + \sum_k W_k^2 \quad \text{Data loss + regularization}$$

want  $\nabla_W L$  (gradient of  $L$  w.r.t.  $W$ , computed analytically)

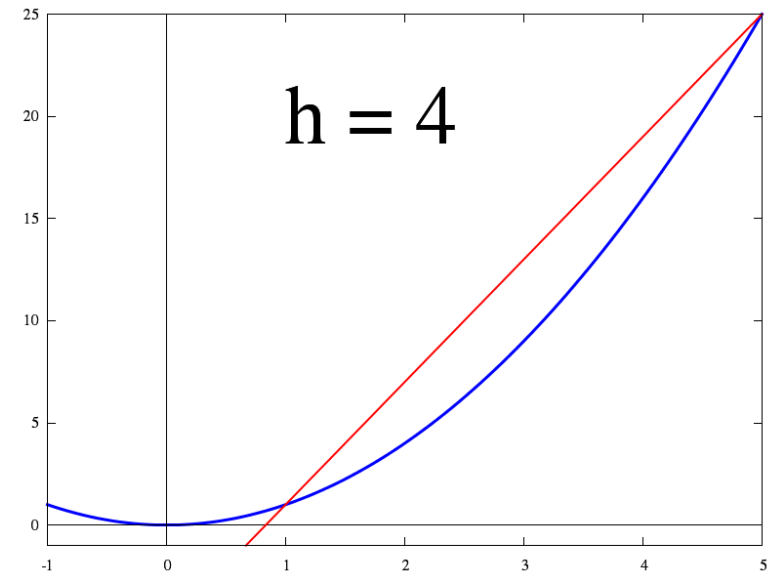
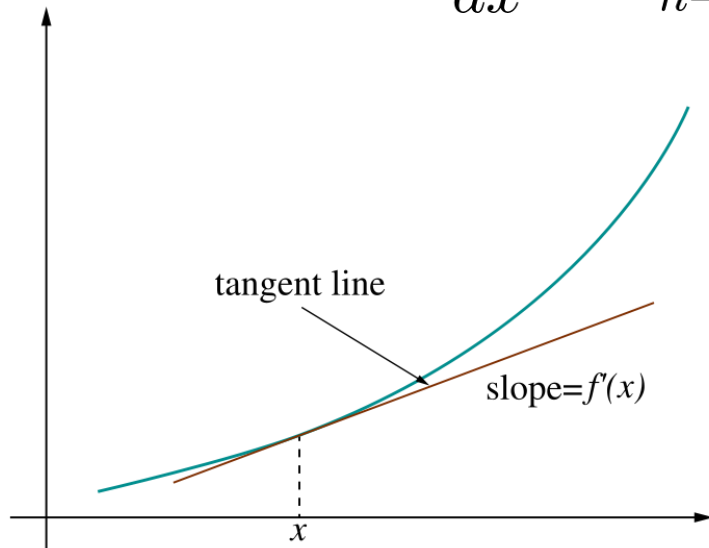
# Gradient descent: iteratively follow the slope



# How do we compute gradients for CNNs?

- Recall: a function with a single with N parameters
- Our loss function involves millions of parameters
- **Idea 1:** Numerically compute derivatives (finite differences)

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$



**current W:**

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25347**

**gradient dW:**

[?,  
?,  
?,  
?,  
?,  
?,  
?,  
?,  
?,...]



**current W:**

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25347**

**W + h (first dim):**

[0.34 + **0.0001**,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25322**

**gradient dW:**

[?,  
?,  
?,  
?,  
?,  
?,  
?,  
?,  
?,...]

**current W:**

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25347**

**W + h (first dim):**

[0.34 + **0.0001**,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25322**

**gradient dW:**

[-2.5,  
?,  
?,

$$(1.25322 - 1.25347)/0.0001 = -2.5$$

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

?,  
?,...]

**current W:**

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25347**

**W + h (second dim):**

[0.34,  
-1.11 + **0.0001**,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25353**

**gradient dW:**

[-2.5,  
?,  
?,  
?,  
?,  
?,  
?,  
?,  
?,...]

**current W:**

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25347**

**W + h (second dim):**

[0.34,  
-1.11 + **0.0001**,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25353**

**gradient dW:**

[-2.5,  
**0.6**,  
?,  
?,

$$(1.25353 - 1.25347)/0.0001 = 0.6$$

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

?,...]

**current W:**

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25347**

**W + h (third dim):**

[0.34,  
-1.11,  
0.78 + **0.0001**,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25347**

**gradient dW:**

[-2.5,  
0.6,  
?,  
?,  
?,  
?,  
?,  
?,  
?,...]

**current W:**

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25347**

**W + h (third dim):**

[0.34,  
-1.11,  
0.78 + **0.0001**,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25347**

**gradient dW:**

[-2.5,  
0.6,  
**0**,  
?,  
?

$$(1.25347 - 1.25347)/0.0001 = 0$$

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

?, ...]

**current W:**

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25347**

**W + h (third dim):**

[0.34,  
-1.11,  
0.78 + **0.0001**,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25347**

**gradient dW:**

[-2.5,  
0.6,  
**0**,  
?,  
?

**Numeric Gradient**

- Slow! Need to loop over all dimensions
- Approximate

?,...]

# But the loss is just a function of $W$ !

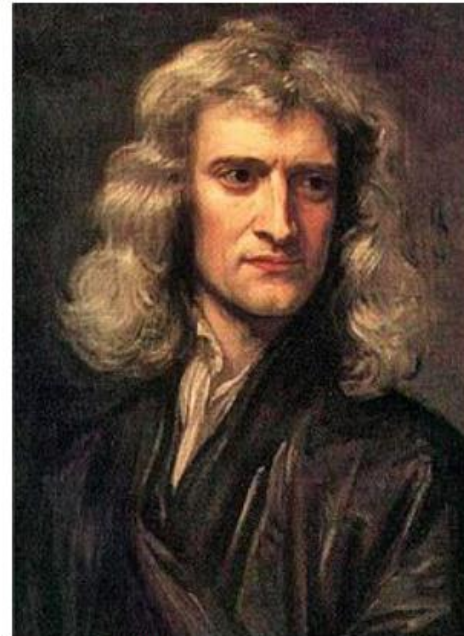
$$L = \frac{1}{N} \sum_{i=1}^N L_i + \sum_k W_k^2$$

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

$$s = f(x; W) = Wx$$

want  $\nabla_W L$

Use calculus to compute an  
**analytic gradient**



[This image](#) is in the public domain



[This image](#) is in the public domain



**current W:**

[0.34,  
-1.11,  
0.78,  
0.12,  
0.55,  
2.81,  
-3.1,  
-1.5,  
0.33,...]

**loss 1.25347**

$dW = \dots$   
(some function  
data and W)



**gradient dW:**

[-2.5,  
0.6,  
0,  
0.2,  
0.7,  
-0.5,  
1.1,  
1.3,  
-2.1,...]

# Idea #2: Calculating gradients analytically

$$s = f(x; W) = Wx$$

$$\begin{aligned} L_i &= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \\ &= \sum_{j \neq y_i} \max(0, W_{j,:} \cdot x + W_{y_i,:} \cdot x + 1) \end{aligned}$$

$$\begin{aligned} L &= \frac{1}{N} \sum_{i=1}^N L_i + \lambda \sum_k W_k^2 \\ &= \frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, W_{j,:} \cdot x + W_{y_i,:} \cdot x + 1) + \lambda \sum_k W_k^2 \end{aligned}$$

$$\nabla_W L = \nabla_W \left( \frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, W_{j,:} \cdot x + W_{y_i,:} \cdot x + 1) + \lambda \sum_k W_k^2 \right)$$

# Idea #2: Calculating gradients analytically

$$s = f(x; W) = Wx$$

$$L_i = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$$

$$= \sum_{j \neq y_i} \max(0, W_{j,:} \cdot x + W_{y_i,:} \cdot x + 1)$$

$$L = \frac{1}{N} \sum_{i=1}^N L_i + \lambda \sum_k W_k^2$$

$$= \frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, W_{j,:} \cdot x + W_{y_i,:} \cdot x + 1) + \lambda \sum_k W_k^2$$

$$\nabla_W L = \nabla_W \left( \frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, W_{j,:} \cdot x + W_{y_i,:} \cdot x + 1) + \lambda \sum_k W_k^2 \right)$$

**Problem:** Very tedious: Lots of matrix calculus, need lots of paper

**Problem:** What if we want to change loss? E.g. use softmax instead of SVM? Need to re-derive from scratch =

**Problem:** Not feasible for very complex models!

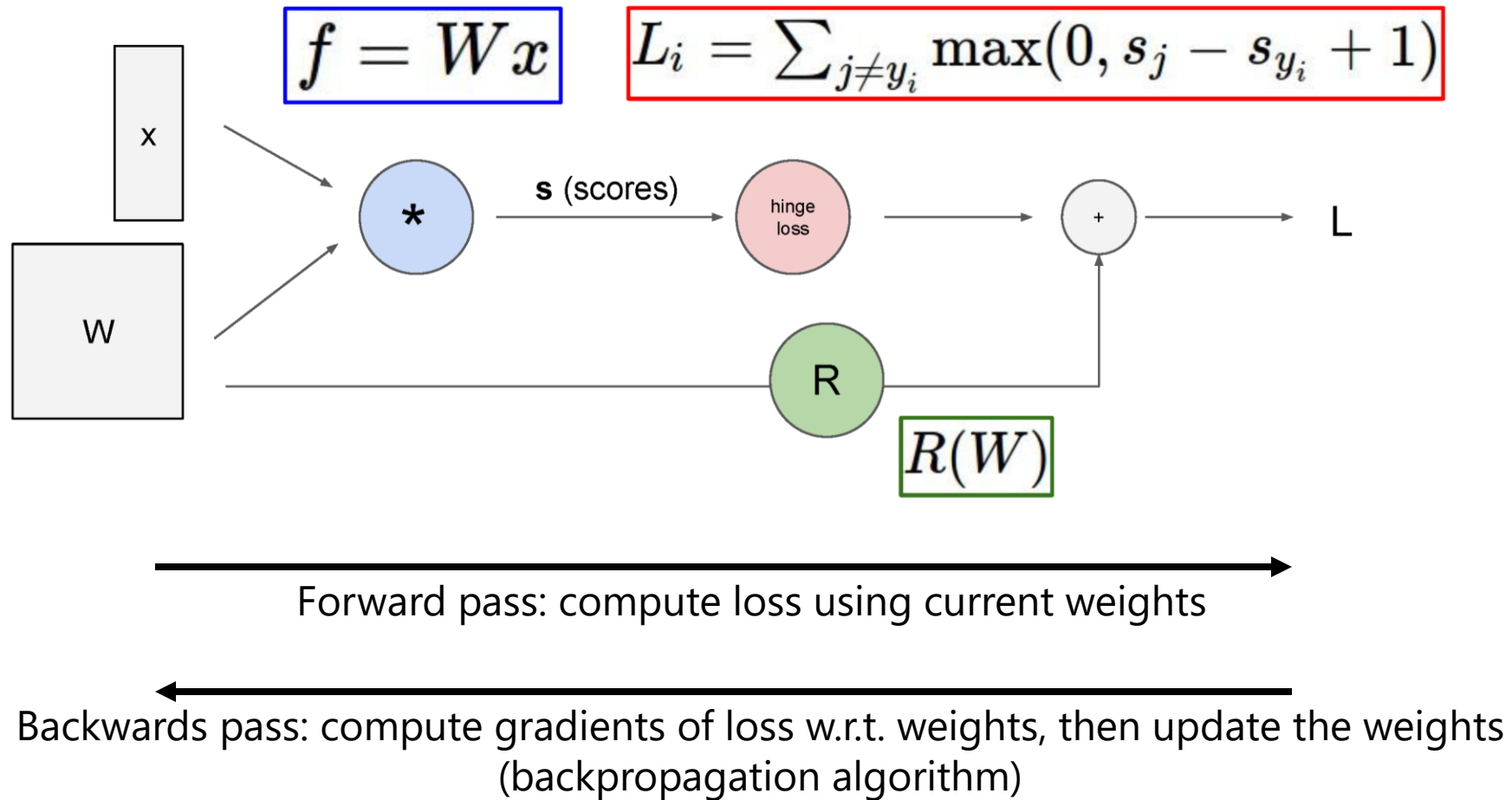
## In summary:

- Numerical gradient: approximate, slow, easy to write
- Analytic gradient: exact, fast, error-prone

=>

In practice: Always use analytic gradient, but check implementation with numerical gradient. This is called a **gradient check**.

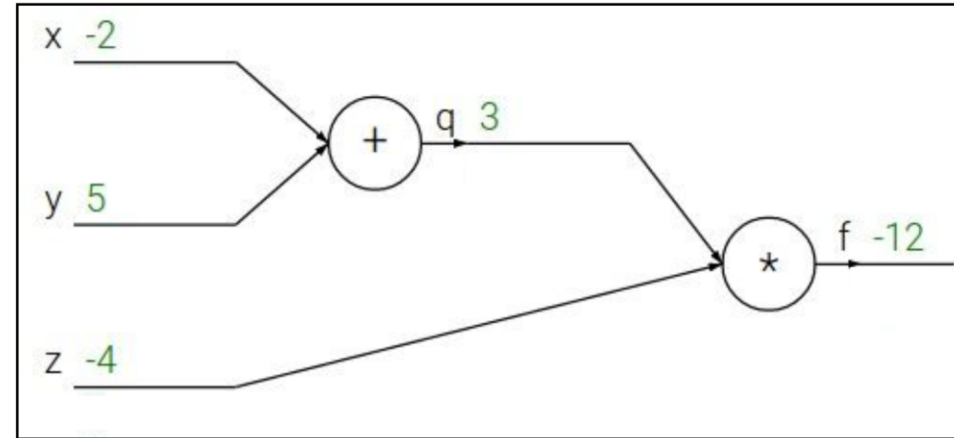
# Better idea: computation graphs + backpropagation



Backpropagation: a simple example

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

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



Backpropagation: a simple example

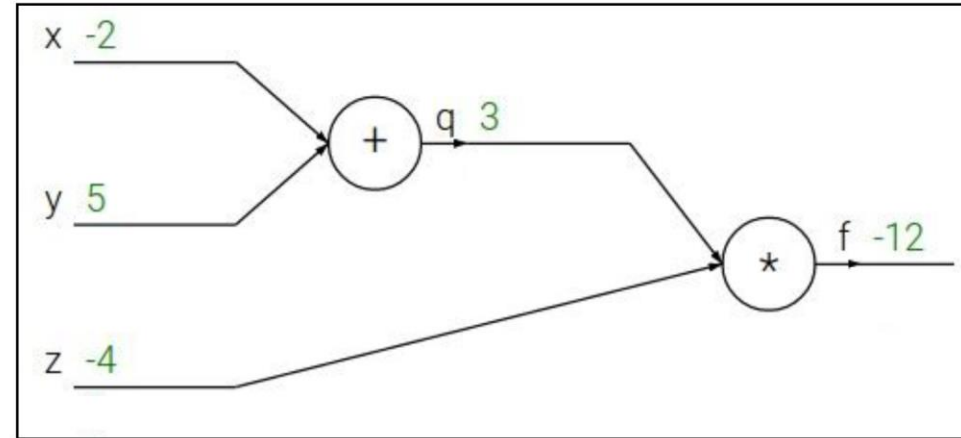
$$f(x, y, z) = (x + y)z$$

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



Backpropagation: a simple example

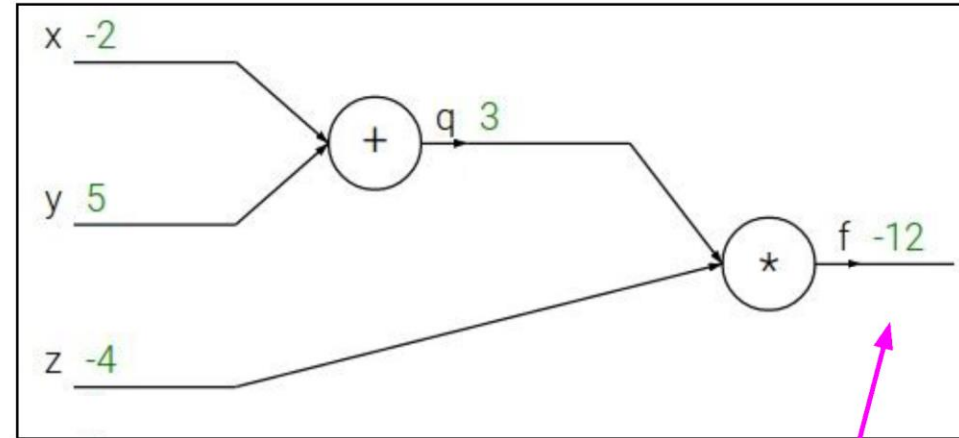
$$f(x, y, z) = (x + y)z$$

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



$$\frac{\partial f}{\partial f}$$



Backpropagation: a simple example

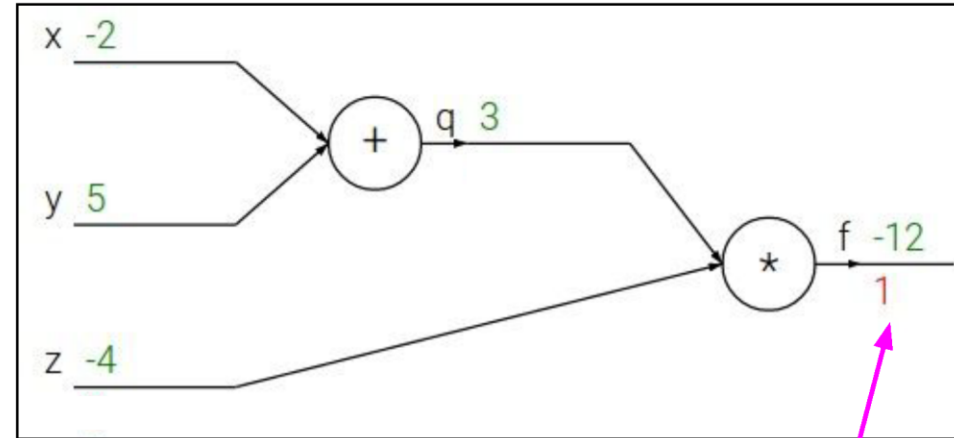
$$f(x, y, z) = (x + y)z$$

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



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

Backpropagation: a simple example

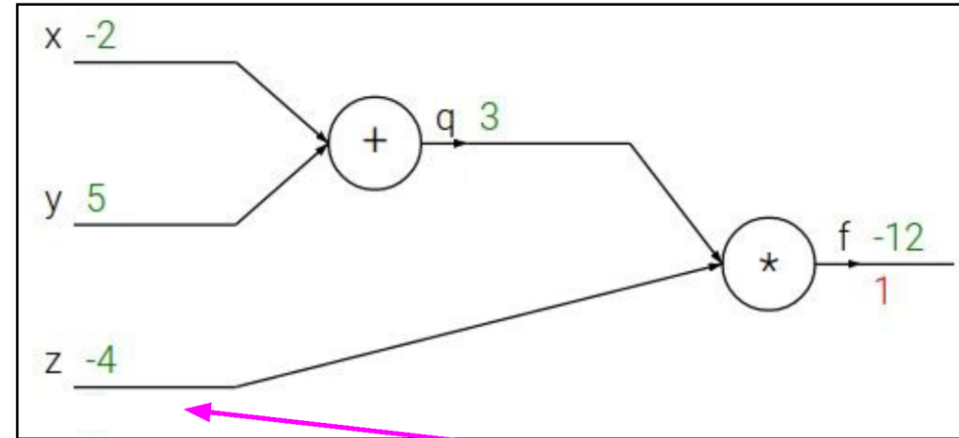
$$f(x, y, z) = (x + y)z$$

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



$$\frac{\partial f}{\partial z}$$

Backpropagation: a simple example

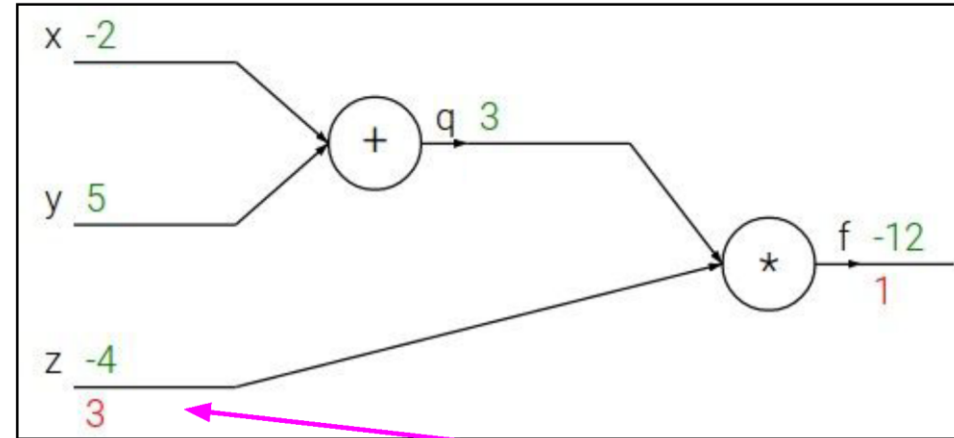
$$f(x, y, z) = (x + y)z$$

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



$$\frac{\partial f}{\partial z}$$

Backpropagation: a simple example

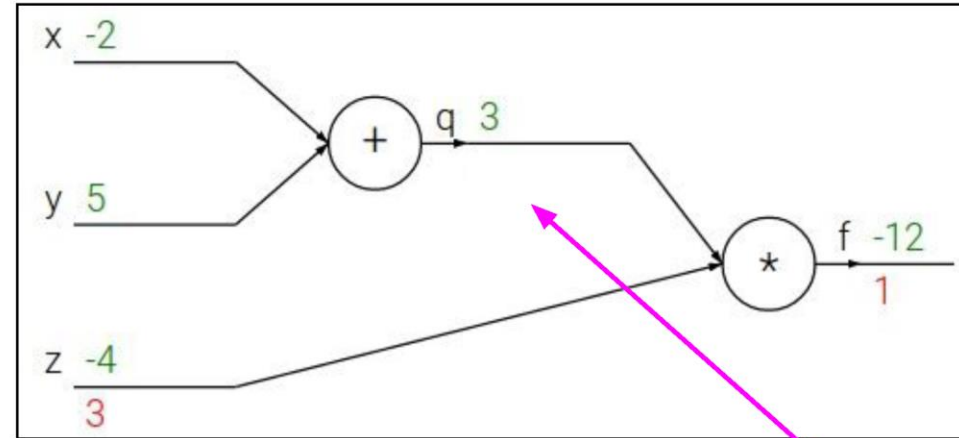
$$f(x, y, z) = (x + y)z$$

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



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

Backpropagation: a simple example

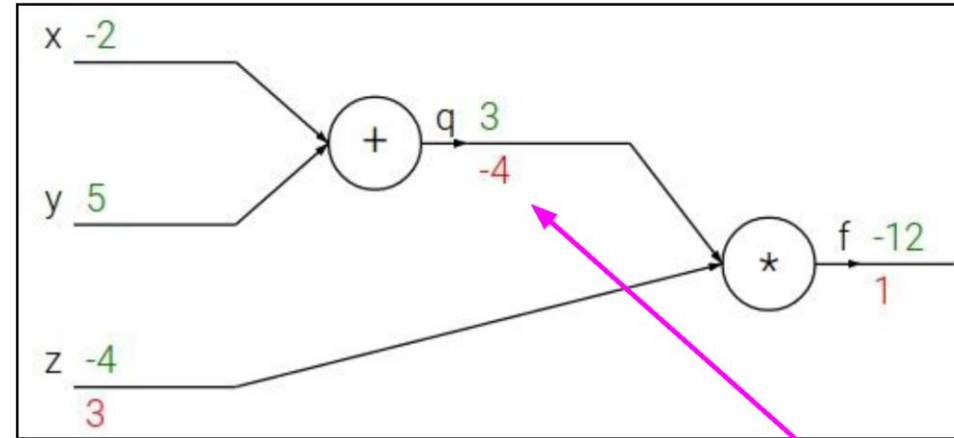
$$f(x, y, z) = (x + y)z$$

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



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

Backpropagation: a simple example

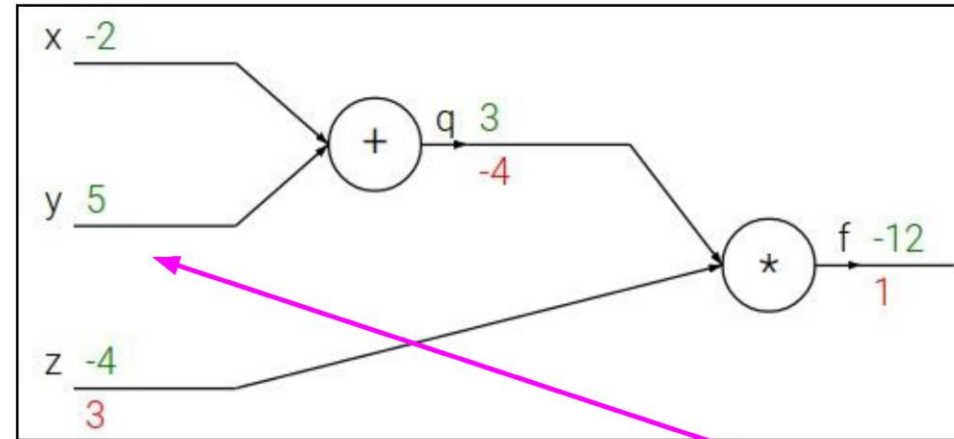
$$f(x, y, z) = (x + y)z$$

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:

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

Upstream  
gradient

Local  
gradient

$$\frac{\partial f}{\partial y}$$

Backpropagation: a simple example

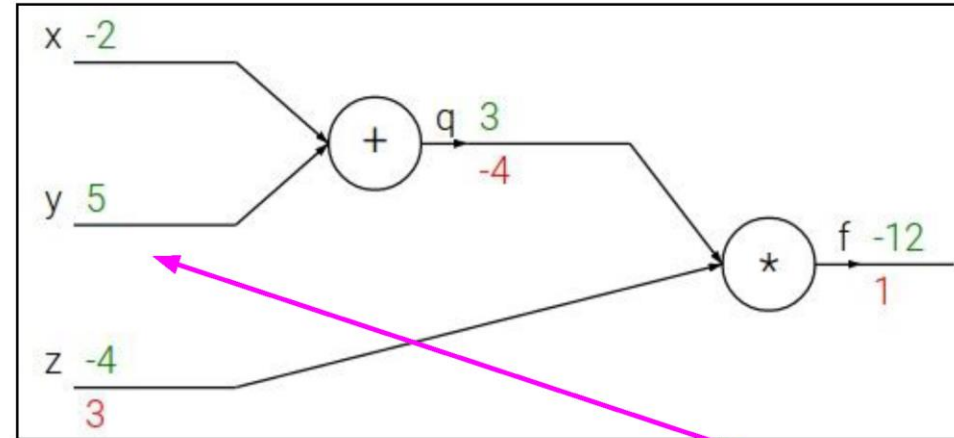
$$f(x, y, z) = (x + y)z$$

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



$$\frac{\partial f}{\partial y}$$

Chain rule:

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

Upstream  
gradient

Local  
gradient

Backpropagation: a simple example

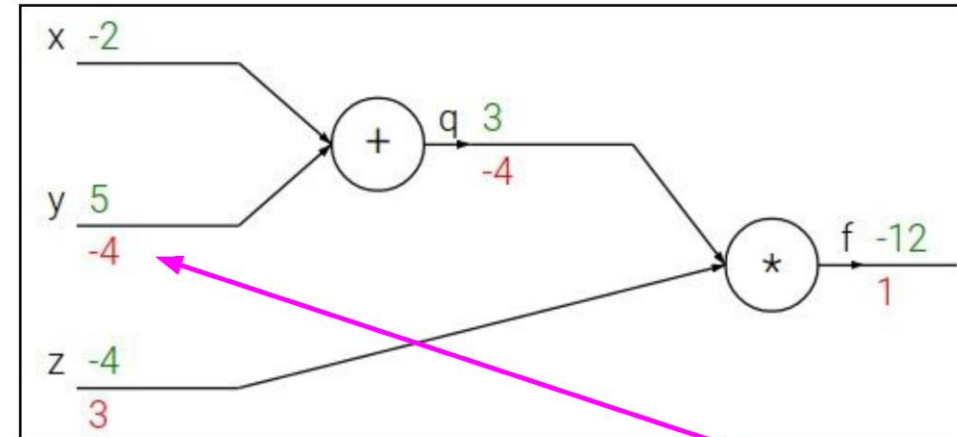
$$f(x, y, z) = (x + y)z$$

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:

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

Upstream  
gradient

Local  
gradient

$$\frac{\partial f}{\partial y}$$



Backpropagation: a simple example

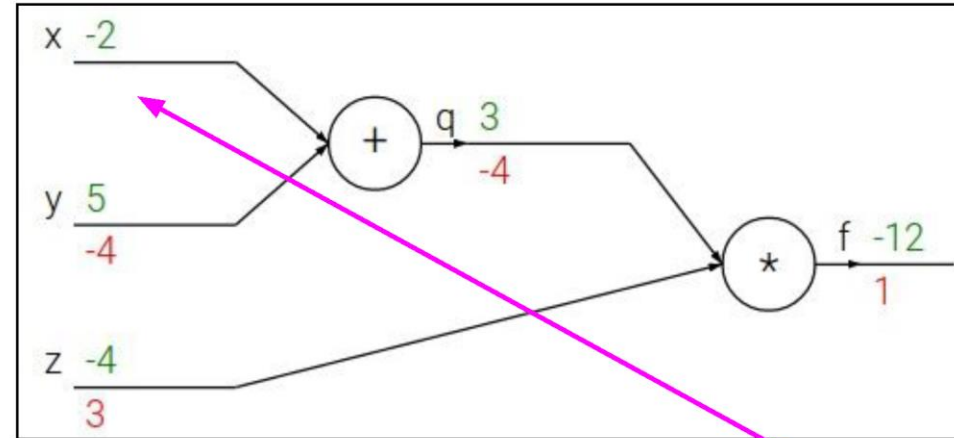
$$f(x, y, z) = (x + y)z$$

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



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

Chain rule:

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

Upstream  
gradient

Local  
gradient

Backpropagation: a simple example

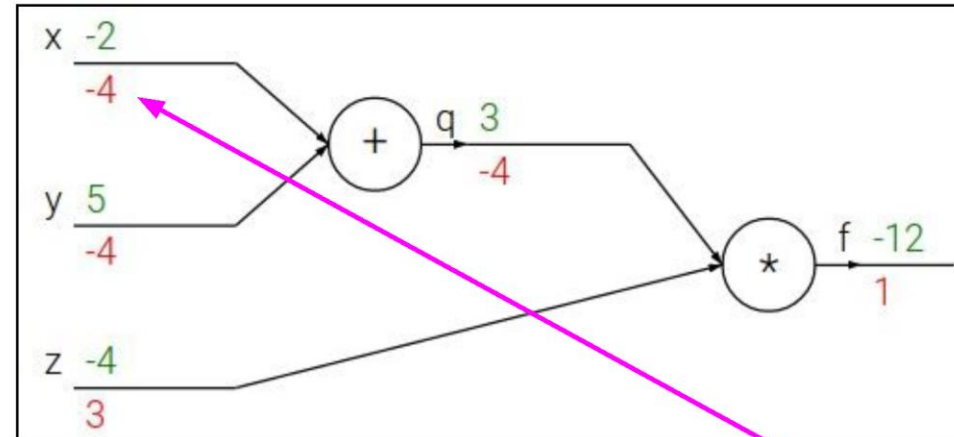
$$f(x, y, z) = (x + y)z$$

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



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

Chain rule:

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

Upstream  
gradient

Local  
gradient

# Backpropagation

- General idea: Recursive application of the chain rule (calculus 101) backwards through a computation graph
- Can reuse intermediate calculations computed during the forwards pass during the backwards pass
- Natural extensions from scalar computations to vector computations
- Deep learning frameworks like Pytorch / TensorFlow support efficient automated backpropagation via automatic differentiation and GPU acceleration

**Questions?**

# What if the training data is very large?

- Recall that ImageNet has > 1.2M training images

$$L = \frac{1}{N} \sum_{i=1}^N L_i \quad \leftarrow \text{Loss function is summed over all } N \text{ training images}$$

$$\nabla L = \frac{1}{N} \sum_{i=1}^N \nabla L_i \quad \leftarrow \text{Gradient is also summed over all } N \text{ training images}$$

- Computing the value of the loss and its gradient over the entire training set is **very** expensive in terms of computation

# Alternative: stochastic gradient descent

- Approximate the sum using a **minibatch** of examples
  - e.g., 32, 64, or 128 examples

$$L = \frac{1}{B} \sum_{i=1}^B L_i$$



Where B (e.g. 32) is the minibatch size

$$\nabla L = \frac{1}{B} \sum_{i=1}^B \nabla L_i$$

- For each step of gradient descent, choose a different batch

# Stochastic gradient descent (SGD)

- A full pass through the dataset (i.e., using batches that cover the training data) is called an **epoch**
- Usually need to train for multiple epochs, i.e., multiple full passes through the dataset to converge
- Stochastic gradient descent approximates the true gradient, but works remarkably well in practice

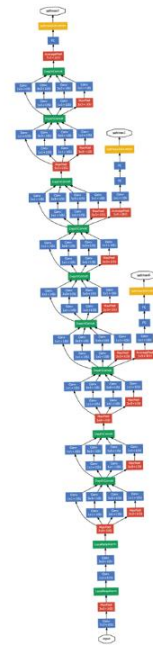
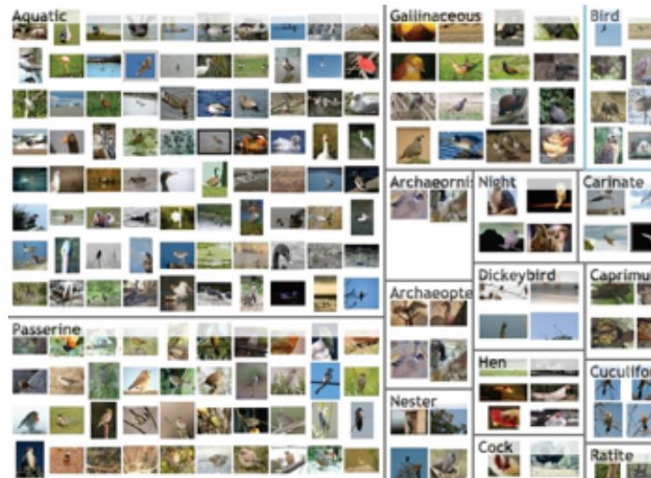
# How do you actually train these things?

## Roughly speaking:

Gather  
labeled data

Find a ConvNet  
architecture

Minimize  
the loss





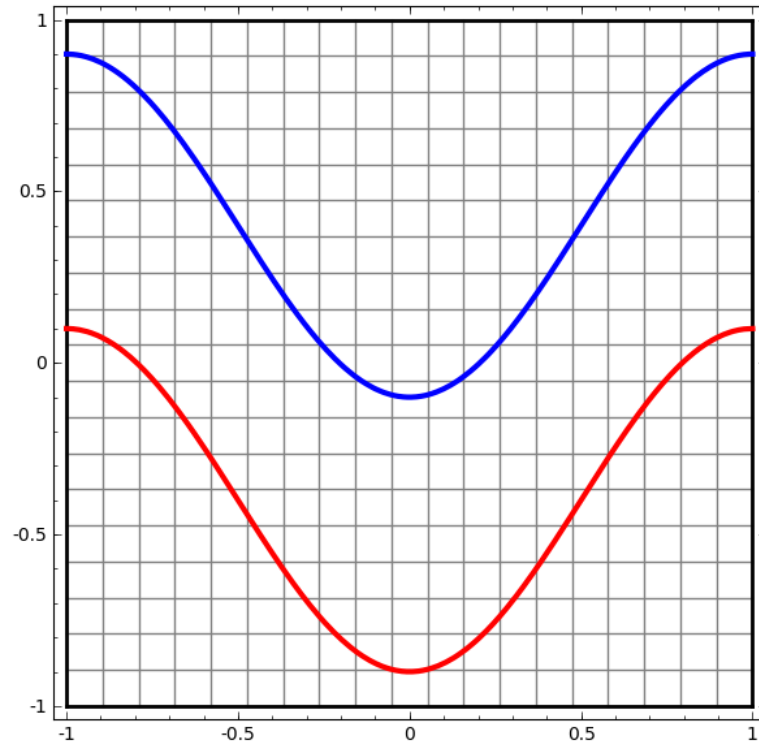
# Training a convolutional neural network

- Split and preprocess your data
- Choose your network architecture
- Initialize the weights
- Find a learning rate and regularization strength
- Minimize the loss and monitor progress
- Fiddle with knobs

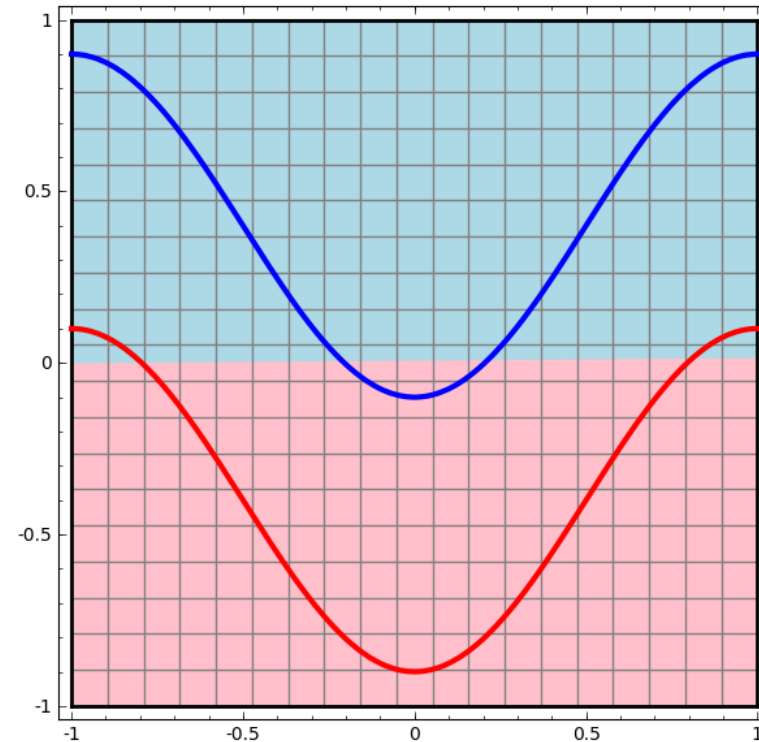
# Why so complicated?

- Training deep networks can be finicky – lots of parameters to learn, complex, non-linear optimization function

# Visualizing Linear Classification (slides from Abe Davis)

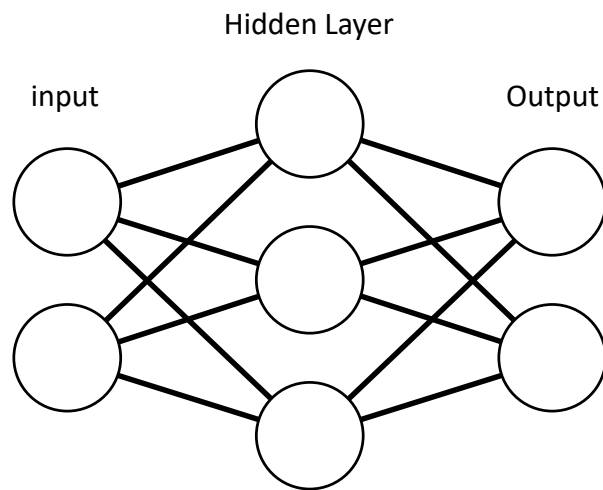


**Classification Problem:  
Separate Red & Blue**

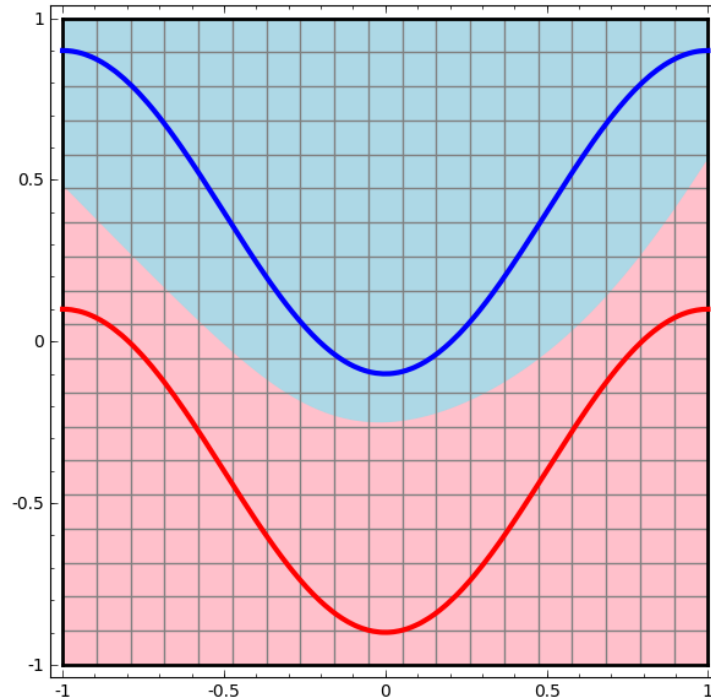


**Linear Solution**

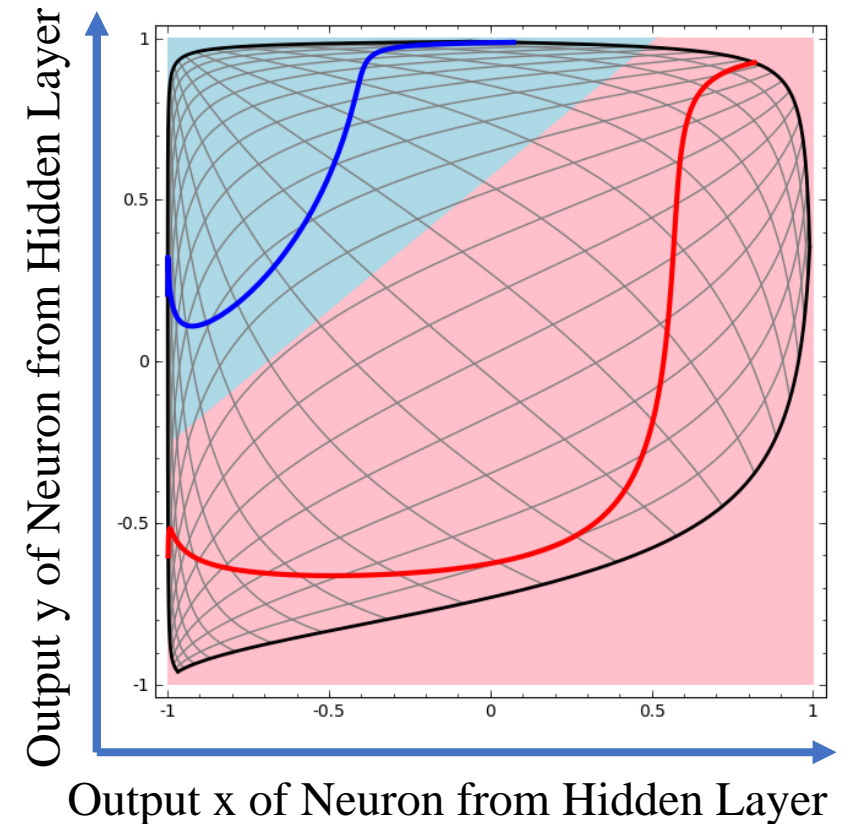
# Visualizing Classification With a Neural Network



Example Network



Classification Results for Every Point in Original Space



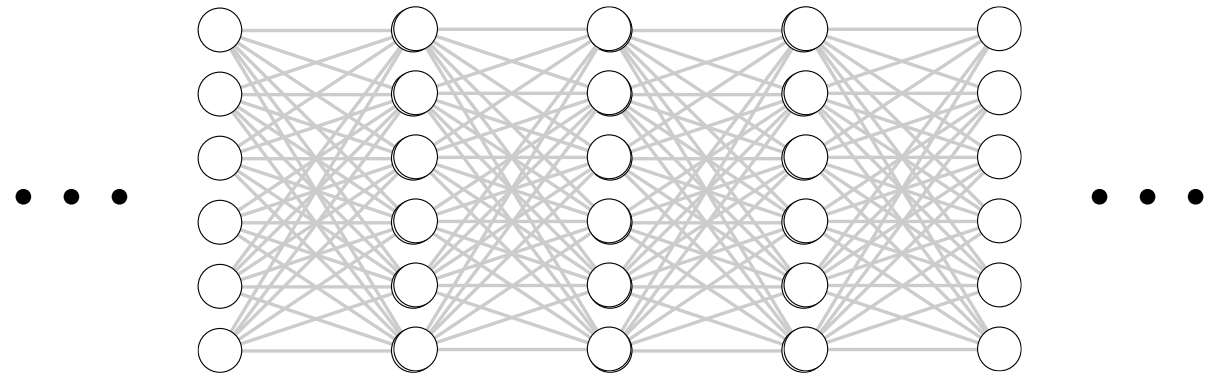
Classification Results for Every Point in Transformed Feature Space

# Demo

<https://cs.stanford.edu/people/karpathy/convnetjs/demo/classify2d.html>

# What Makes Training Deep Nets Hard?

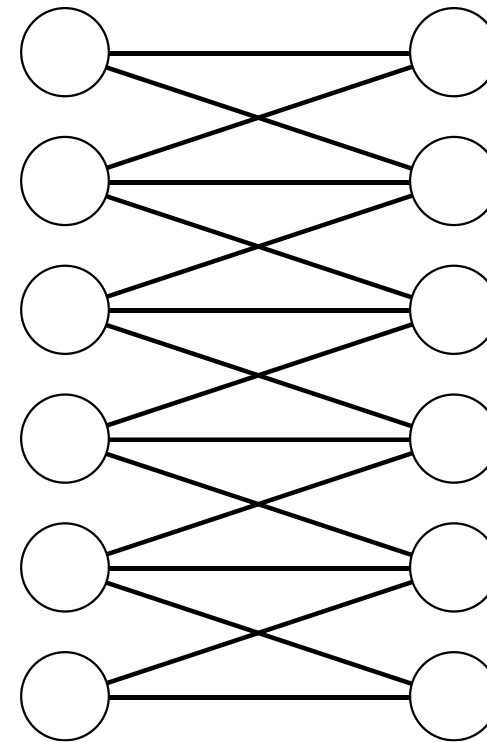
- It's easy to get high training accuracy:
  - Use a huge, fully connected network with tons of layers
  - Let it memorize your training data
- It's hard to get high *test* accuracy



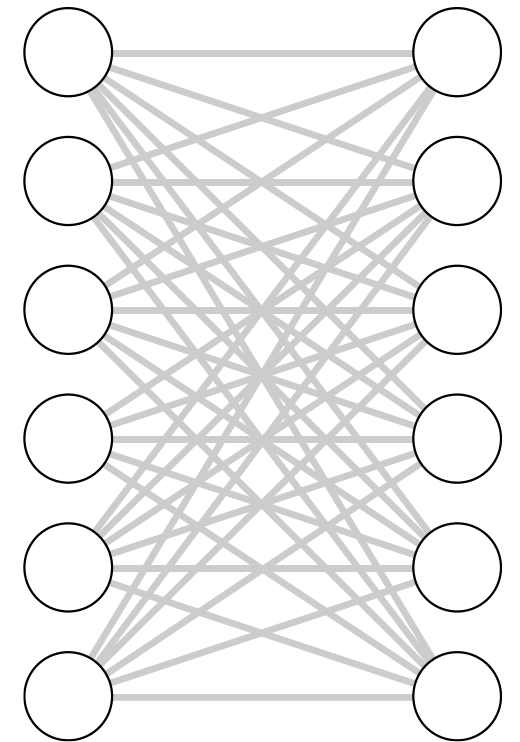
This would be an example of *overfitting*

# Related Question: Why Convolutional Layers?

- A fully connected layer can generally represent the same functions as a convolutional one
  - Think of the convolutional layer as a version of the FC layer with constraints on parameters
- What is the advantage of CNNs?



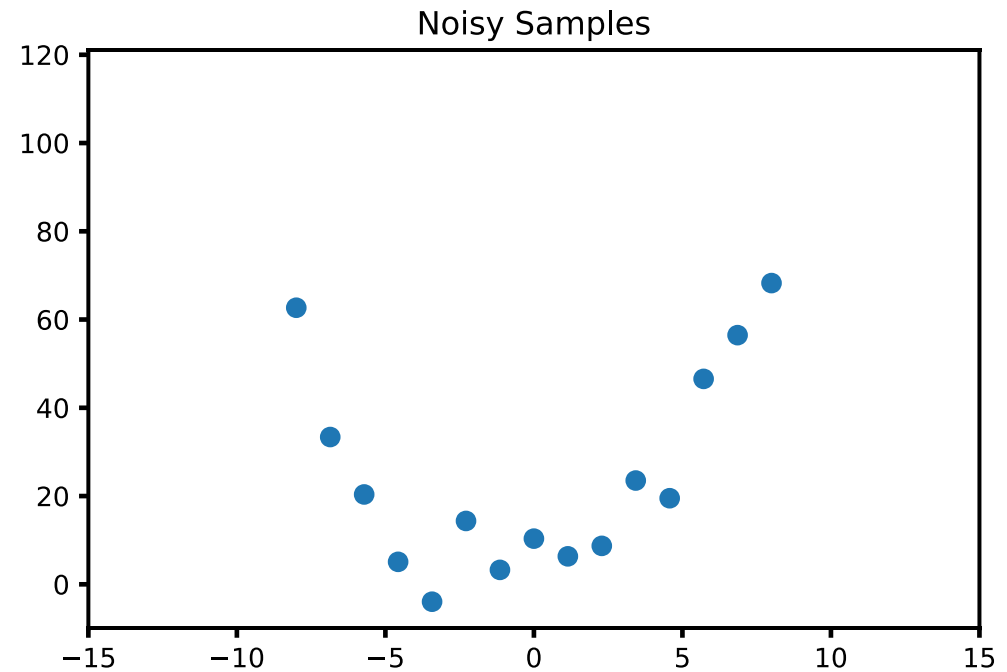
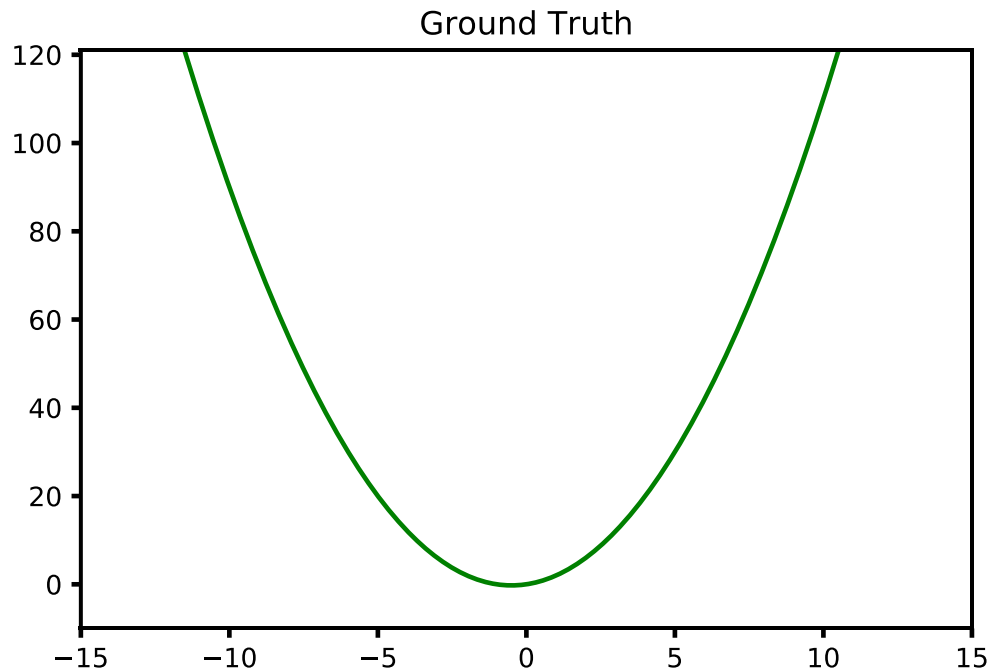
**Convolutional Layer**



**Fully Connected Layer**

# Overfitting: More Parameters, More Problems

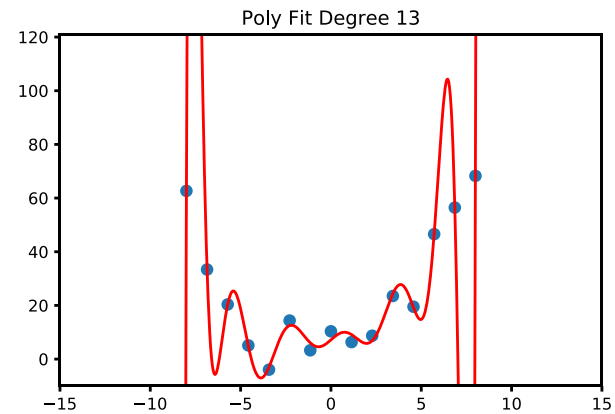
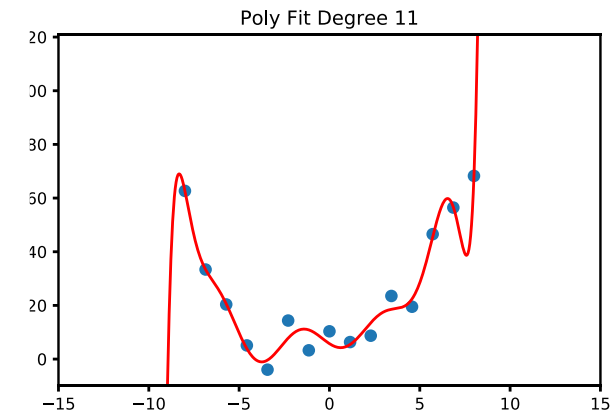
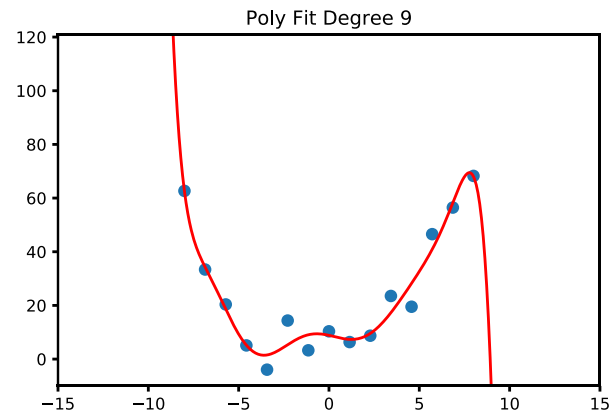
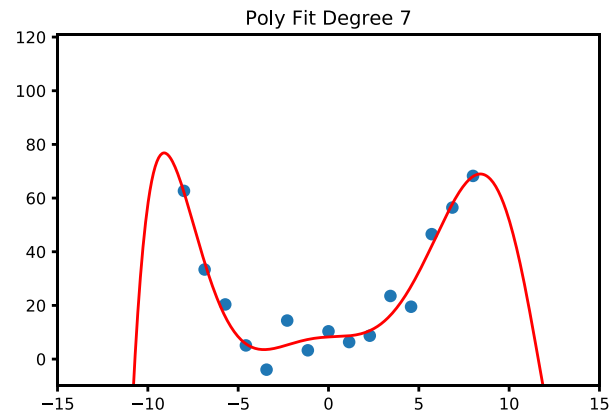
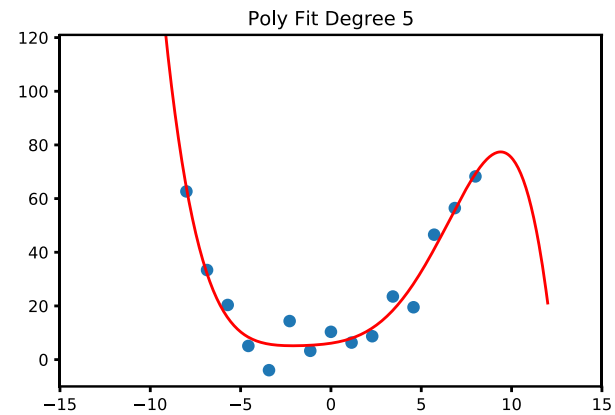
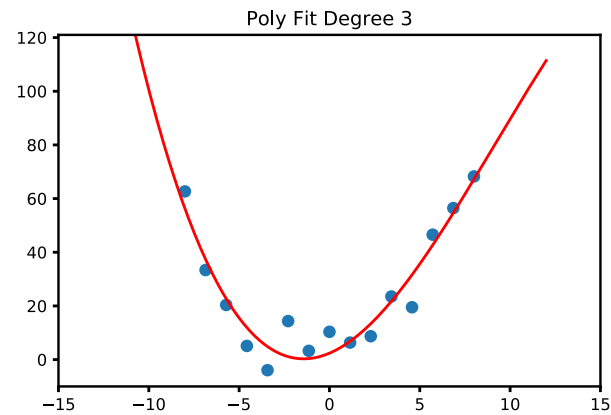
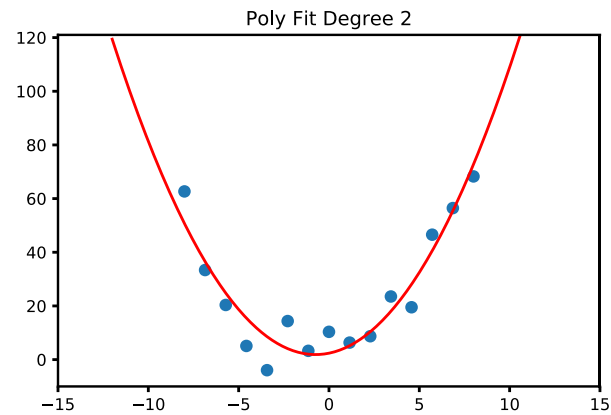
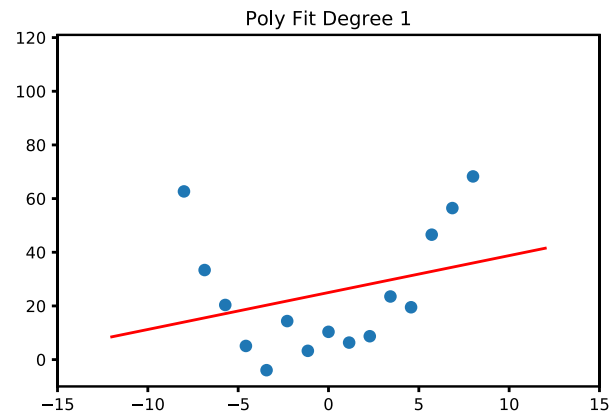
- Non-Deep Example: consider the function  $x^2 + x$
- Let's take some noisy samples of the function...





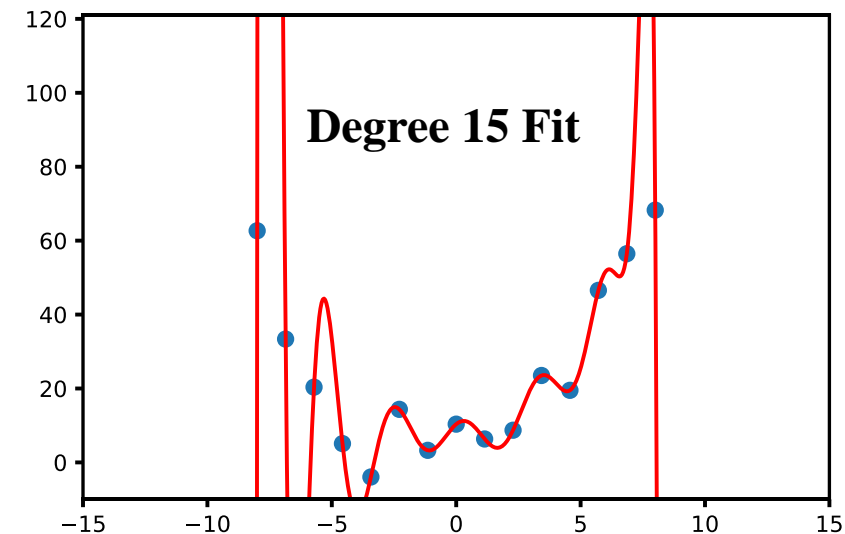
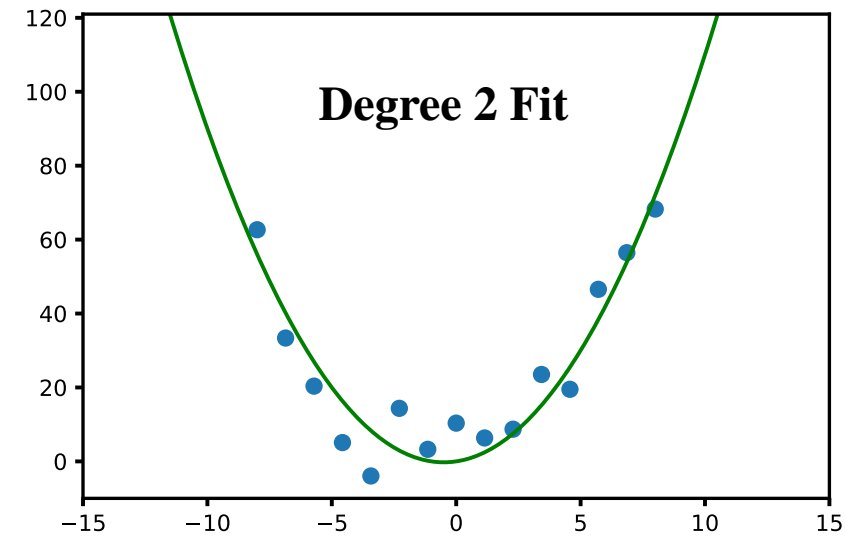
# Overfitting: More Parameters, More Problems

- Now let's fit a polynomial to our samples of the form  $P_N(x) = \sum_{k=0}^N x^k p_k$



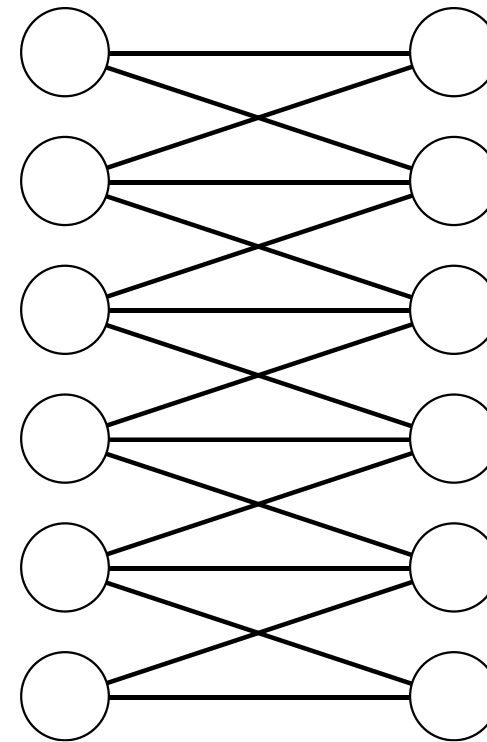
# Overfitting: More Parameters, More Problems

- A model with more parameters can represent more functions
- E.g.,: if  $P_N(x) = \sum_{k=0}^N x^k p_k$  then  $P_2 \in P_{15}$
- More parameters will often **reduce training error** but **increase testing error**. This is *overfitting*.
- When overfitting happens, models do not generalize well.

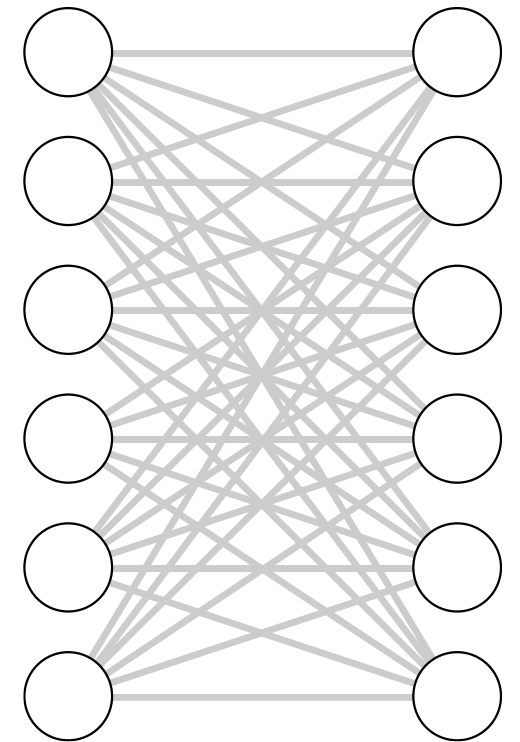


# Deep Learning: More Parameters, More Problems?

- More parameters let us represent a larger space of functions
- The larger that space is, the harder our optimization becomes
- This means we need:
  - More data
  - More compute resources
  - Etc.



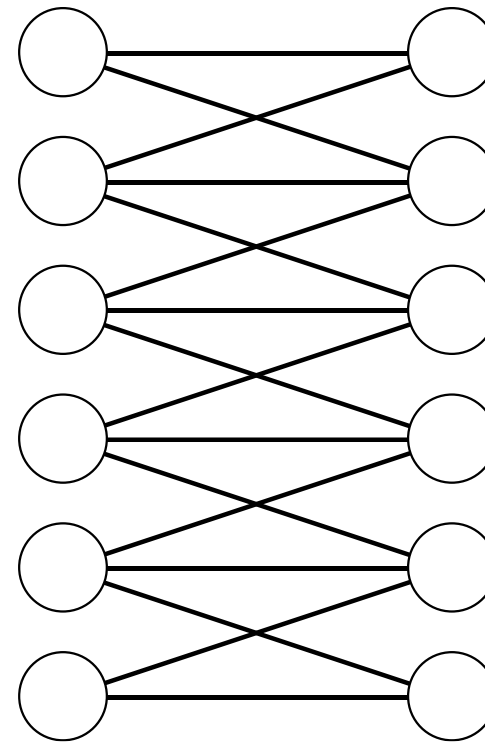
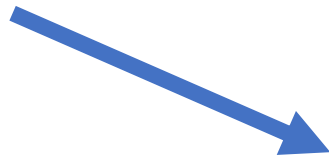
**Convolutional Layer**



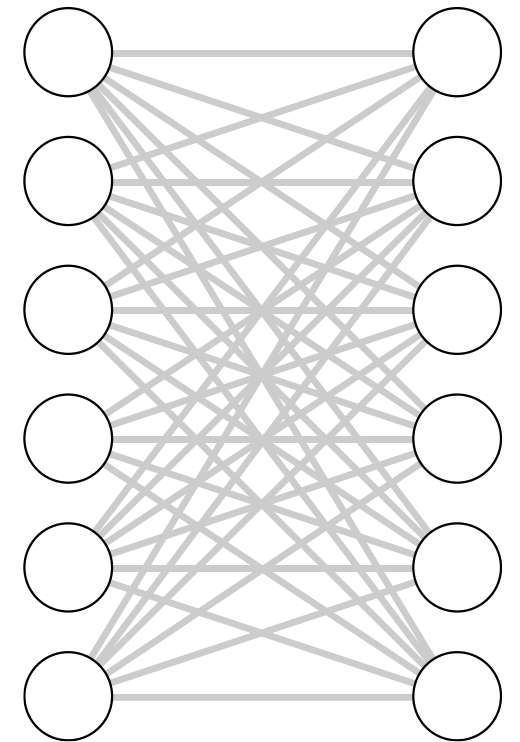
**Fully Connected Layer**

# Deep Learning: More Parameters, More Problems?

A convolutional layer looks for components of a function that are spatially-invariant



Convolutional Layer



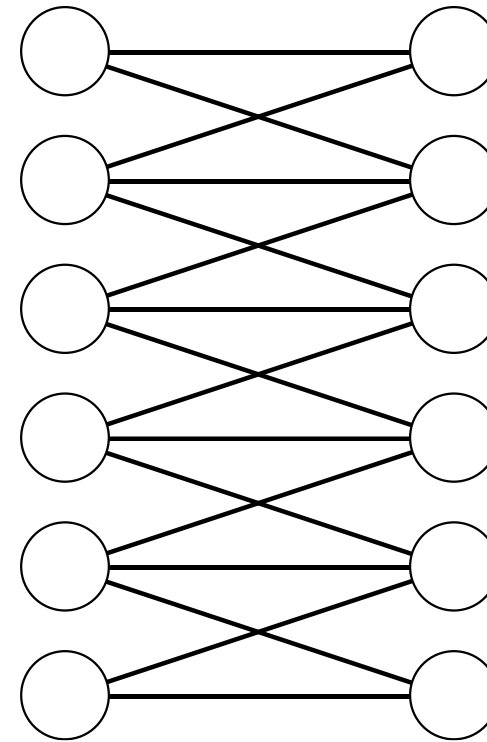
Fully Connected Layer

# How to Avoid Overfitting: Regularization

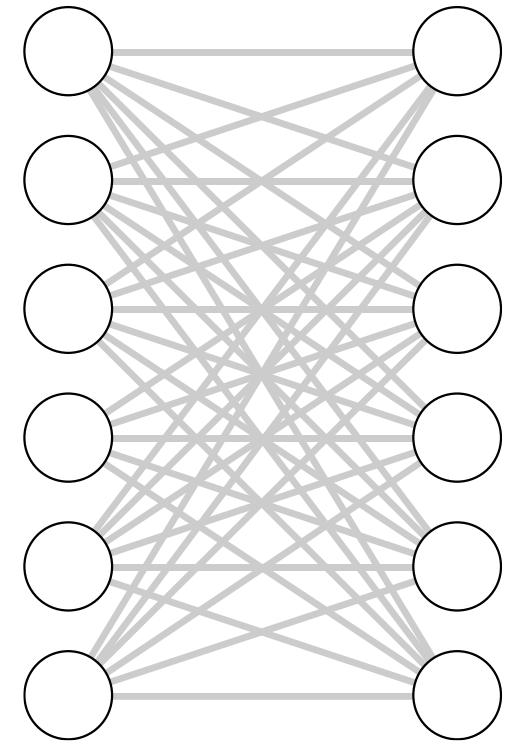
- In general:
  - More parameters means higher risk of overfitting
  - More constraints/conditions on parameters can help
- If a model is overfitting, we can
  - Collect more data to train on
  - *Regularize*: add some additional information or assumptions to better constrain learning
- Regularization can be done through:
  - the design of architecture
  - the choice of loss function
  - the preparation of data
  - ...

# Regularization: Architecture Choice

- “Bigger” architectures (typically, those with more parameters) tend to be more at risk of overfitting.



**Convolutional Layer**

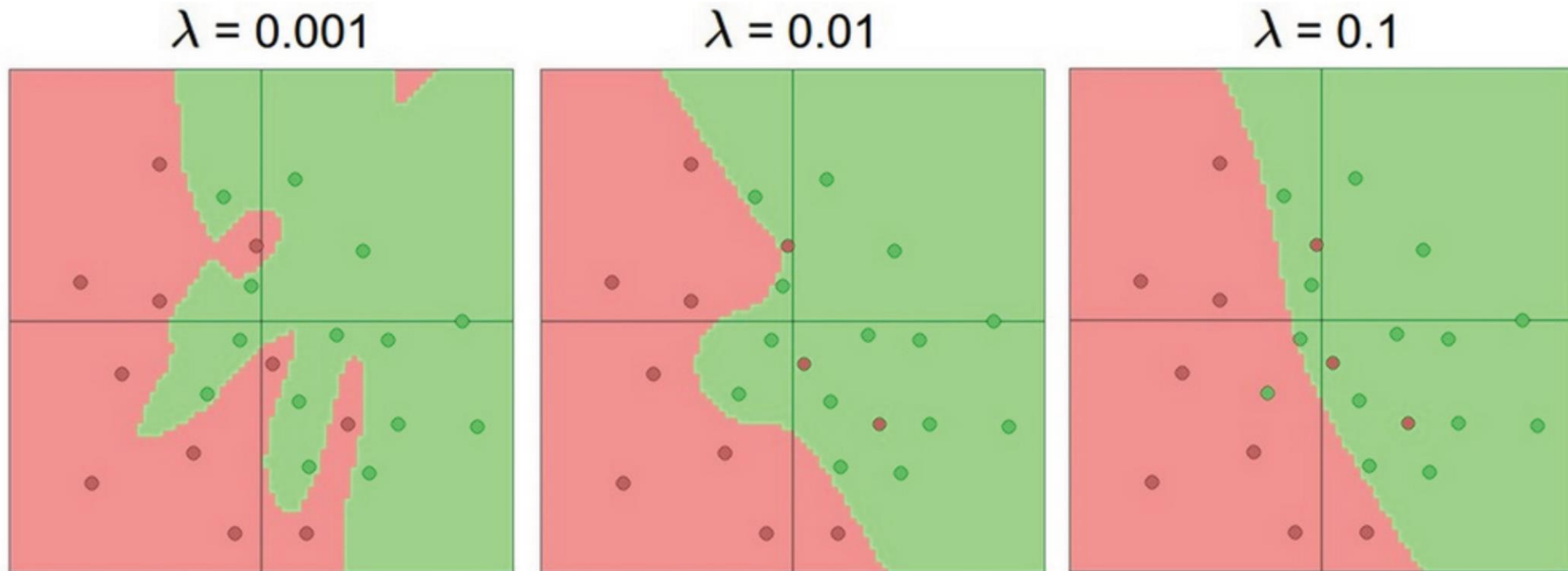


**Fully Connected Layer**

# Regularization

**Regularization reduces overfitting:**

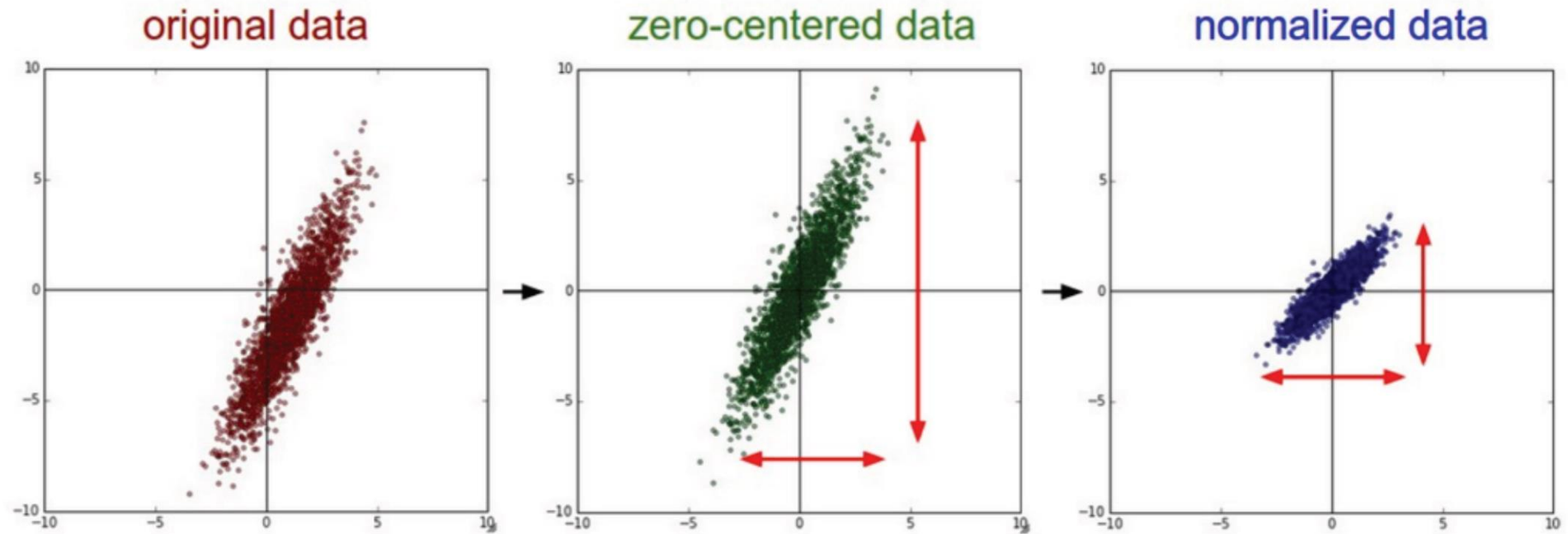
$$L = L_{\text{data}} + L_{\text{reg}} \quad L_{\text{reg}} = \lambda \frac{1}{2} \|W\|_2^2$$



[Andrej Karpathy <http://cs.stanford.edu/people/karpathy/convnetjs/demo/classify2d.html>]

# (1) Data preprocessing

Preprocess the data so that learning is better conditioned:



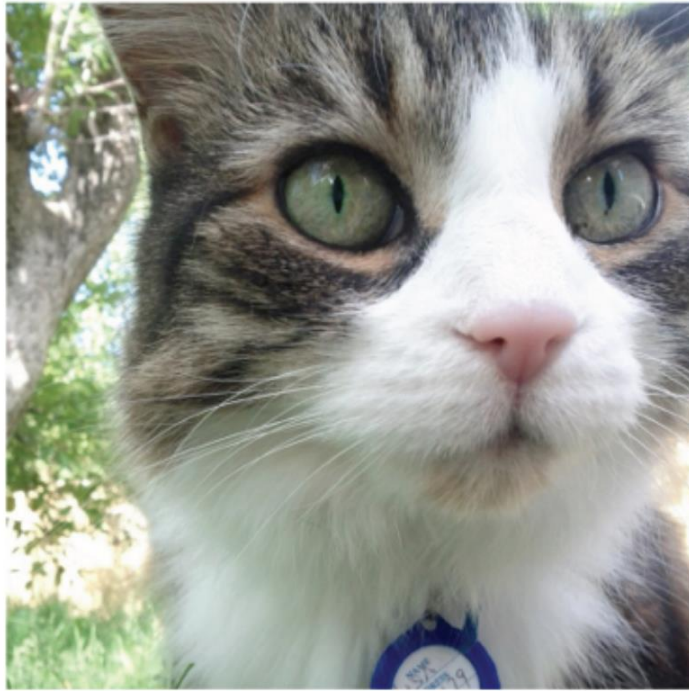
```
X -= np.mean(axis=0, keepdims=True)
```

```
X /= np.std(axis=0, keepdims=True)
```



# (1) Data preprocessing

For ConvNets, typically only the mean is subtracted.



An input image (256x256)



Minus sign



The mean input image

A per-channel mean also works (one value per R,G,B).

# Batch normalization

- Side note – can also perform normalization after each layer of the network to stabilize network training ("*batch normalization*")

# (1) Data preprocessing

**Augment the data** — extract random crops from the input, with slightly jittered offsets. Without this, typical ConvNets (e.g. [Krizhevsky 2012]) overfit the data.



**E.g.** 224x224 patches  
extracted from 256x256 images

Randomly reflect horizontally

Perform the augmentation live  
during training

*Figure: Alex Krizhevsky*

# (2) Choose your architecture

The screenshot shows the TensorFlow Playground interface with the following settings and components:

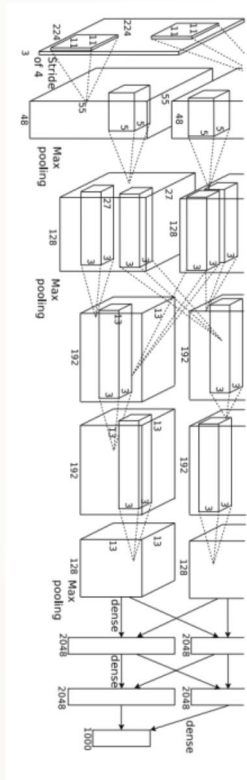
- Browser:** <https://playground.tensorflow.org/#activation=tanh&batchSize=10&dataset=circle&regDataset=reg-plane&learningRate=0...>
- Controls:** Play button, Epoch: 000,000, Learning rate: 0.03, Activation: Tanh, Regularization: None, Regularization rate: 0, Problem type: Classification.
- DATA:** Which dataset do you want to use? (Circle dataset selected), Ratio of training to test data: 50%, Noise: 0, Batch size: 10, REGENERATE button.
- FEATURES:** Which properties do you want to feed in? (Selected features:  $X_1$ ,  $X_2$ ,  $X_1^2$ ,  $X_2^2$ ,  $X_1X_2$ ,  $\sin(X_1)$ ,  $\sin(X_2)$ ).
- Architecture:** 2 HIDDEN LAYERS. Layer 1: 4 neurons. Layer 2: 2 neurons. Weights are shown by line thickness.
- OUTPUT:** Test loss 0.507, Training loss 0.504. A scatter plot shows data points (orange and blue) and decision boundaries (shaded regions).
- Legend:** Colors show data, neuron, and weight values. A color scale from -1 (blue) to 1 (orange) is provided.

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

# (2) Choose your architecture

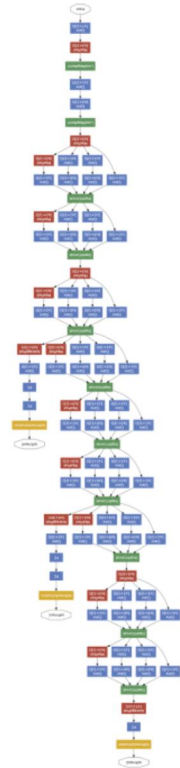
Very common modern choice

*“AlexNet”*



[Krizhevsky et al. NIPS 2012]

*“GoogLeNet”*



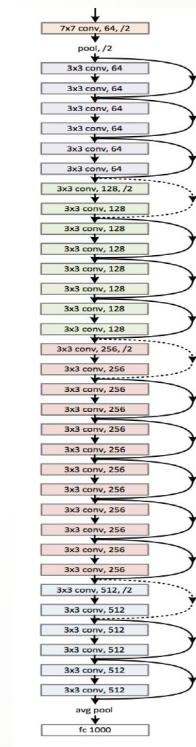
[Szegedy et al. CVPR 2015]

*“VGG Net”*

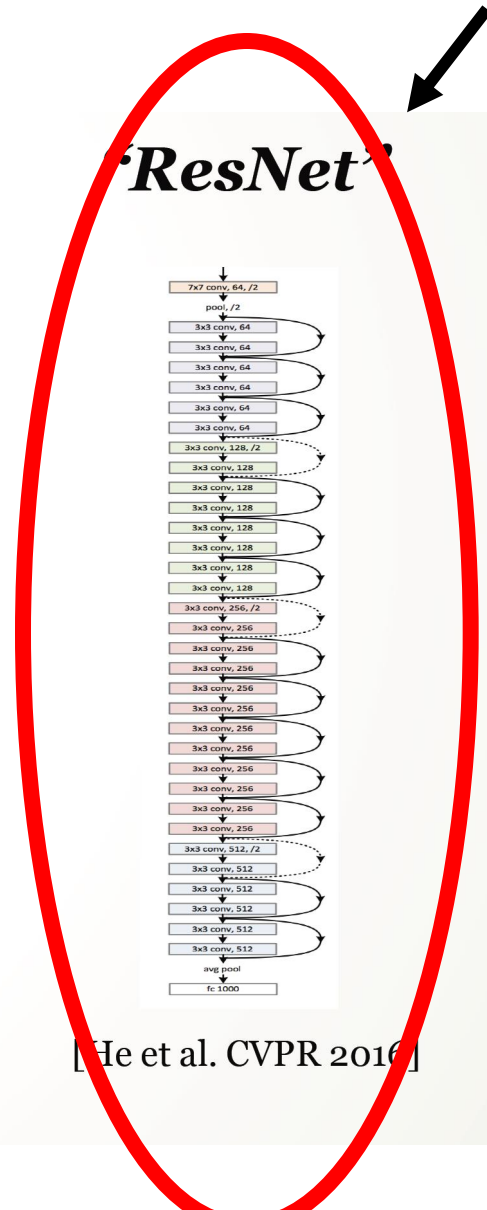


[Simonyan & Zisserman, ICLR 2015]

*“ResNet”*



[He et al. CVPR 2016]



# (3) Initialize your weights

**Set the weights to small random numbers:**

```
W = np.random.randn(D, H) * 0.001
```

(matrix of small random numbers drawn from a Gaussian distribution)

**Set the bias to zero (or small nonzero):**

```
b = np.zeros(H)
```

(if you use ReLU activations, folks tend to initialize bias to small positive number)

## (4) Overfit a small portion of the data

```
model = init_two_layer_model(32*32*3, 50, 10) # input size, hidden size, number of classes
trainer = ClassifierTrainer()
X_tiny = X_train[:20] # take 20 examples ←
y_tiny = y_train[:20]
best_model, stats = trainer.train(X_tiny, y_tiny, X_tiny, y_tiny,
                                  model, two_layer_net,
                                  num_epochs=200, reg=0.0,
                                  update='sgd', learning_rate_decay=1,
                                  sample_batches = False,
                                  learning_rate=1e-3, verbose=True)
```

The above code:

- take the first 20 examples from CIFAR-10
- turn off regularization (reg = 0.0)
- use simple vanilla 'sgd'

## (4) Overfit a small portion of the data

```
model = init_two_layer_model(32*32*3, 50, 10) # input size, hidden size, number of classes
trainer = ClassifierTrainer()
X_tiny = X_train[:20] # take 20 examples ←
y_tiny = y_train[:20]
best_model, stats = trainer.train(X_tiny, y_tiny, X_tiny, y_tiny,
                                  model, two_layer_net,
                                  num_epochs=200, reg=0.0,
                                  update='sgd', learning_rate_decay=1,
                                  sample_batches = False,
                                  learning_rate=1e-3, verbose=True)
```

### Details:

'sgd': vanilla gradient descent (no momentum etc)

learning\_rate\_decay = 1: constant learning rate

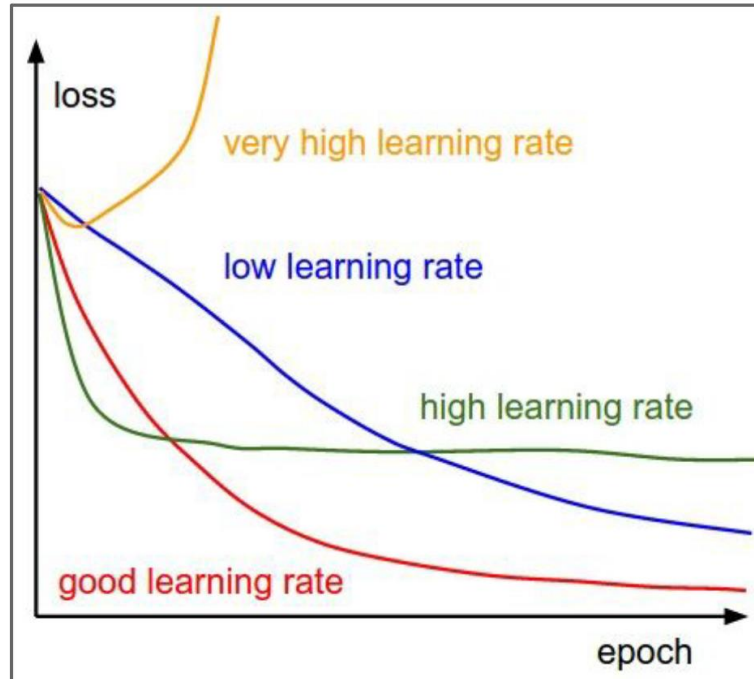
sample\_batches = False (full gradient descent, no batches)

epochs = 200: number of passes through the data





# (4) Find a learning rate



Q: Which one of these learning rates is best to use?

# Learning rate schedule

**How do we change the learning rate over time?**

**Various choices:**

- Step down by a factor of 0.1 every 50,000 mini-batches (used by SuperVision [Krizhevsky 2012])
- Decrease by a factor of 0.97 every epoch (used by GoogLeNet [Szegedy 2014])
- Scale by  $\sqrt{1-t/\max\_t}$  (used by BVLC to re-implement GoogLeNet)
- Scale by  $1/t$
- Scale by  $\exp(-t)$

# Summary of things to fiddle

- Network architecture
- Learning rate, decay schedule, update type
- Regularization (L2, L1, maxnorm, dropout, ...)
- Loss function (softmax, SVM, ...)
- Weight initialization

Neural network  
parameters



# Summary of things to fiddle

- Network architecture
- Learning rate, decay schedule, update type (+batch size)
- Regularization (L2, L1, maxnorm, dropout, ...)
- Loss function (softmax, SVM, ...)
- Weight initialization

Neural network  
parameters



**Questions?**

# Transfer Learning

“You need a lot of a data if you want to train/use CNNs”

# Transfer Learning

“You need a lot of data if you want to train/use CNNs”

**BUSTED**



# Transfer Learning with CNNs

## 1. Train on Imagenet



Donahue et al, "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014  
Razavian et al, "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014

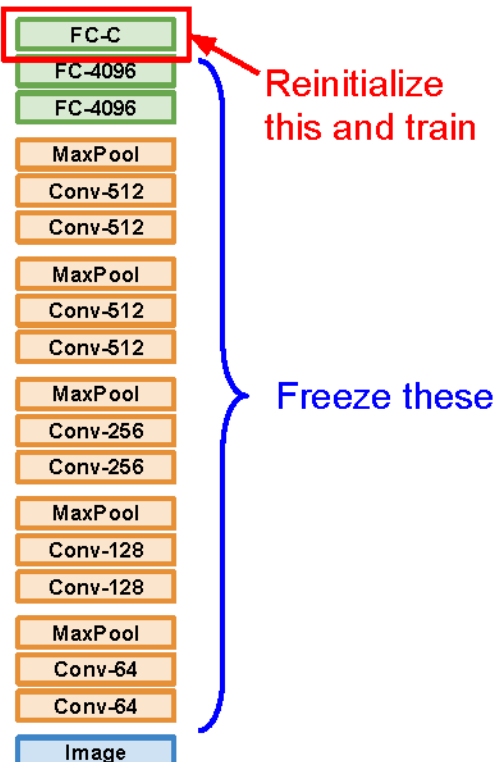
# Transfer Learning with CNNs

Donahue et al, "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014  
Razavian et al, "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014

## 1. Train on Imagenet



## 2. Small Dataset (C classes)



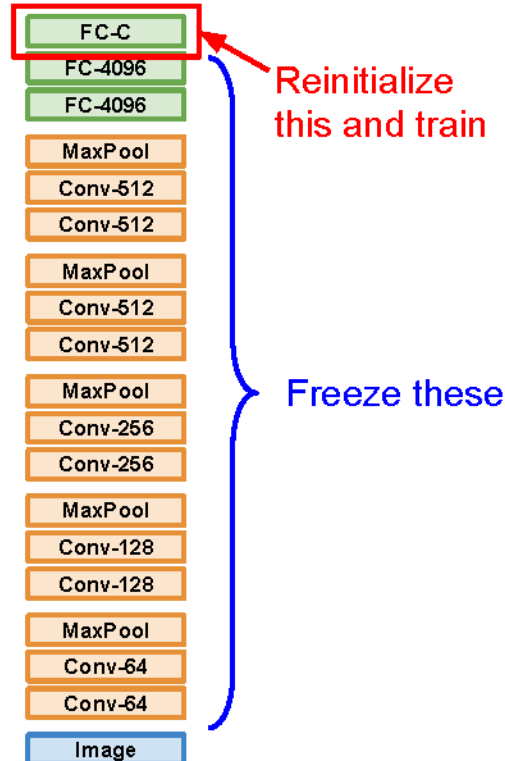
# Transfer Learning with CNNs

Donahue et al, "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014  
Razavian et al, "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014

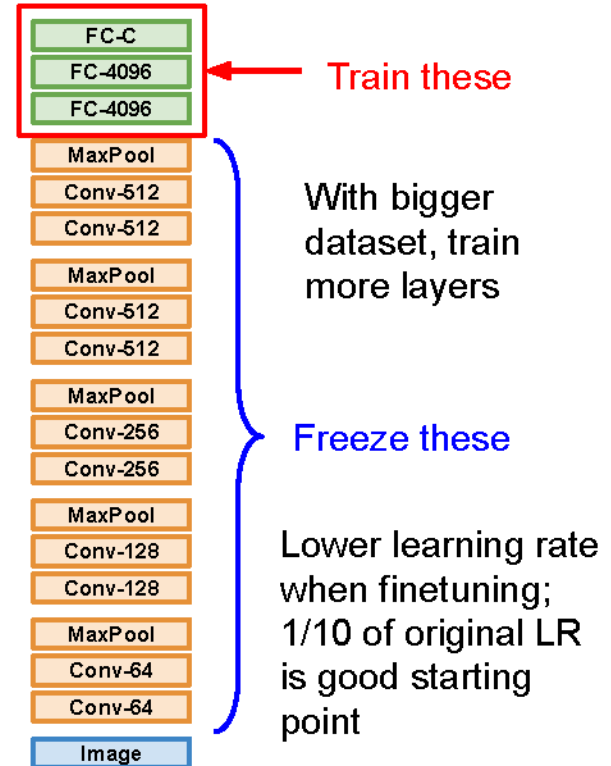
## 1. Train on Imagenet



## 2. Small Dataset (C classes)



## 3. Bigger dataset

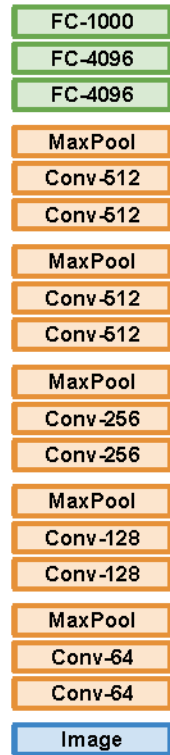




More specific

More generic

	<b>very similar dataset</b>	<b>very different dataset</b>
<b>very little data</b>	?	?
<b>quite a lot of data</b>	?	?



More specific

More generic

	<b>very similar dataset</b>	<b>very different dataset</b>
<b>very little data</b>	Use Linear Classifier on top layer	?
<b>quite a lot of data</b>	Finetune a few layers	?



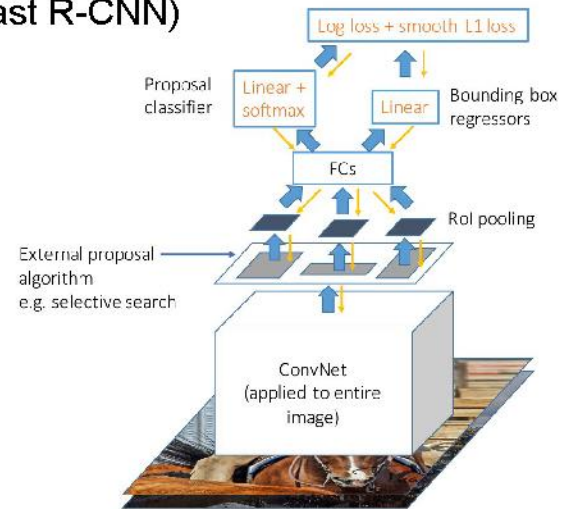
More specific

More generic

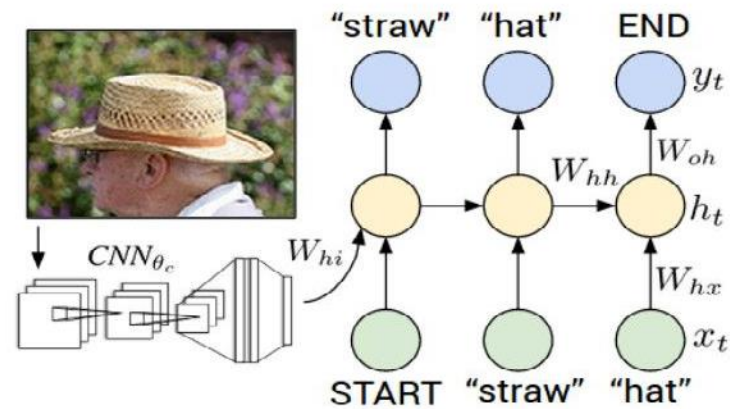
	<b>very similar dataset</b>	<b>very different dataset</b>
<b>very little data</b>	Use Linear Classifier on top layer	You're in trouble... Try linear classifier from different stages
<b>quite a lot of data</b>	Finetune a few layers	Finetune a larger number of layers

# Transfer learning with CNNs is pervasive... (it's the norm, not an exception)

## Object Detection (Fast R-CNN)

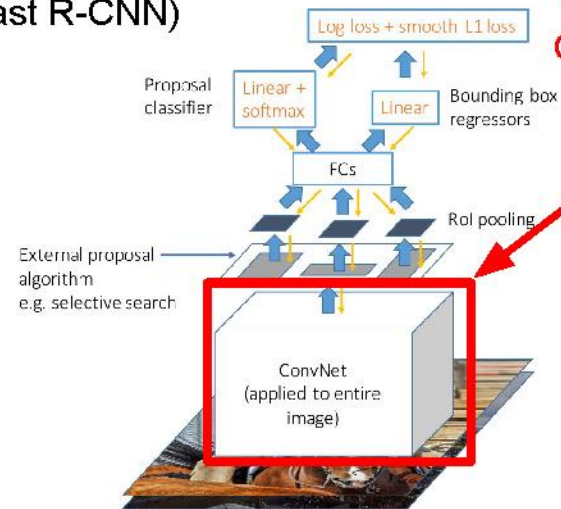


## Image Captioning: CNN + RNN



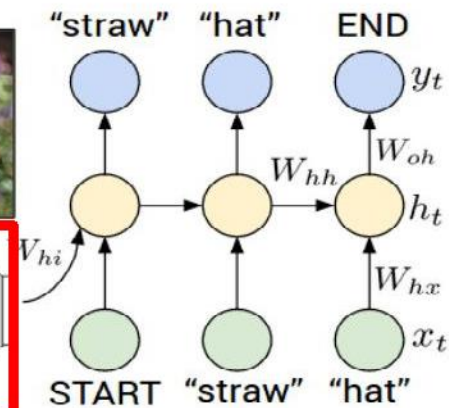
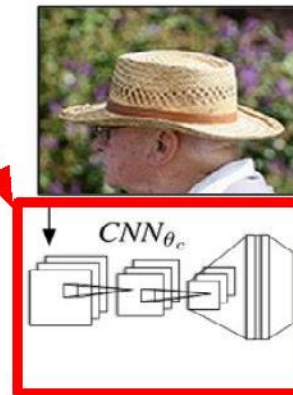
# Transfer learning with CNNs is pervasive... (it's the norm, not an exception)

## Object Detection (Fast R-CNN)



CNN pretrained  
on ImageNet

## Image Captioning: CNN + RNN

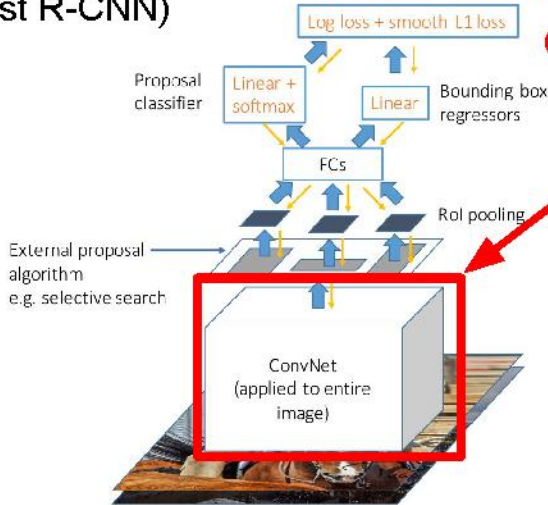




# Transfer learning with CNNs is pervasive...

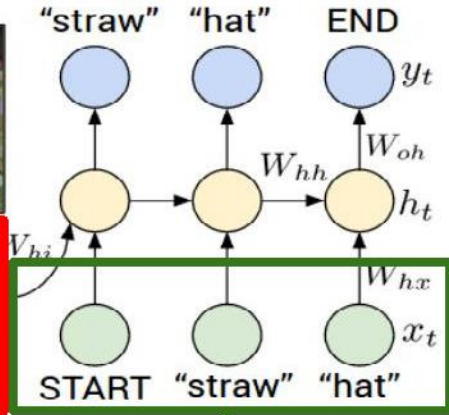
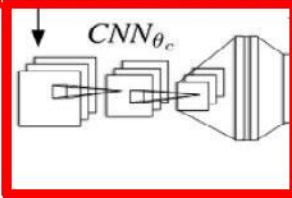
(it's the norm, not an exception)

Object Detection  
(Fast R-CNN)



CNN pretrained on ImageNet

Image Captioning: CNN + RNN



Word vectors pretrained with word2vec

Karpathy and Fei-Fei, "Deep Visual-Semantic Alignments for Generating Image Descriptions", CVPR 2015  
Figure copyright IEEE, 2015. Reproduced for educational purposes.

Girshick, "Fast R-CNN", ICCV 2015  
Figure copyright Ross Girshick, 2015. Reproduced with permission.

# Takeaway for your projects and beyond:

Have some dataset of interest but it has  $< \sim 1\text{M}$  images?

1. Find a very large dataset that has similar data, train a big ConvNet there
2. Transfer learn to your dataset

Deep learning frameworks provide a “Model Zoo” of pretrained models so you don’t need to train your own

TensorFlow: <https://github.com/tensorflow/models>

PyTorch: <https://github.com/pytorch/vision>

Common modern approach:  
start with a ResNet  
architecture pre-trained on  
ImageNet, and fine-tune on  
your (smaller) dataset

**Questions?**