

# Introduction to Deep Learning

DO VAN NGUYEN

---

Many Contents From:

*Fei-Fei Li, Justin Johnson & Serena Yeung: Stanford Course on “CNN for Visual Recognition”.*

# Parts of presentation

---

## Neural Network

- Linear Classification and Perceptron
- Model optimization
- Multiple Perceptron: Neural network
- Training with gradient descent

## Deep Learning: Architectures and Learning Methods

- Supervised Learning with Convolutional Neural Network
- Sequence Learning with Recurrent Neural Network
- Unsupervised learning - Generative Models
- Reinforcement learning (Just Introduction)

# Neural Network

---



[www.image-net.org](http://www.image-net.org)

**22K** categories and **14M** images

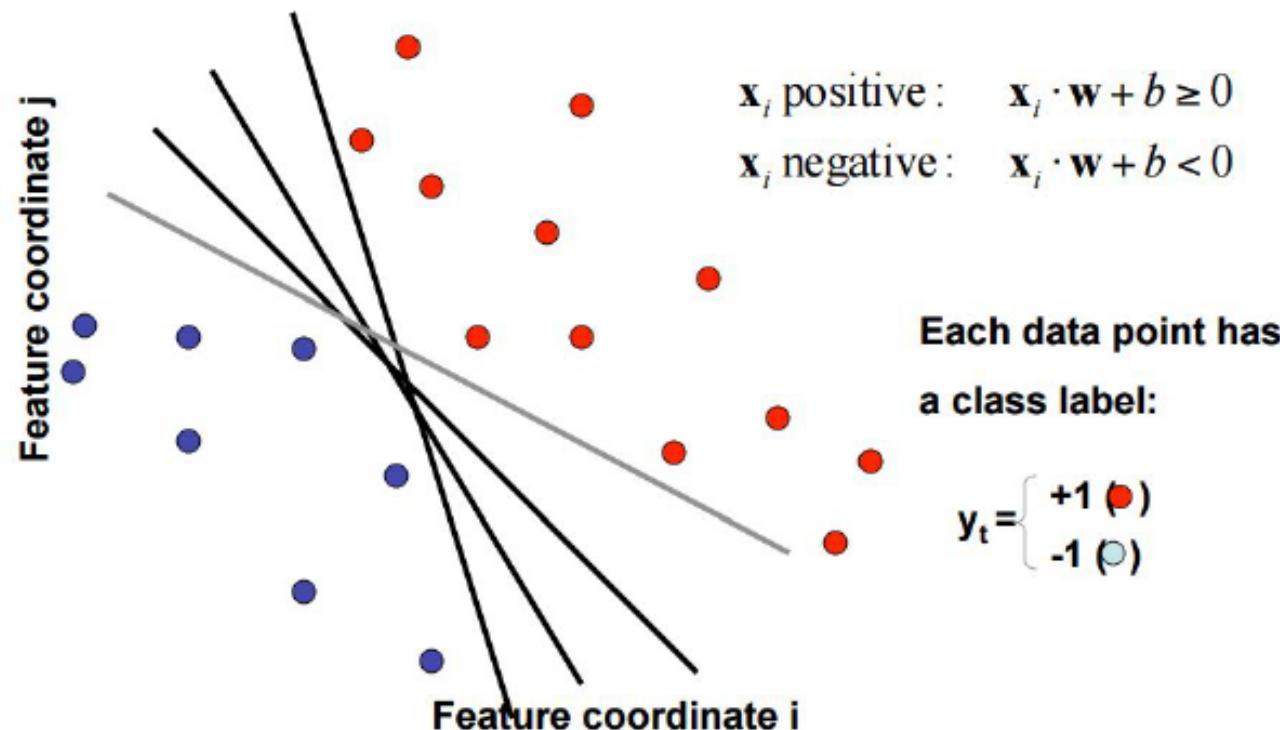
- Animals
  - Bird
  - Fish
  - Mammal
  - Invertebrate
- Plants
  - Tree
  - Flower
  - Food
  - Materials
- Structures
  - Artifact
  - Tools
  - Appliances
  - Structures
- Person
- Scenes
  - Indoor
  - Geological Formations
- Sport Activities



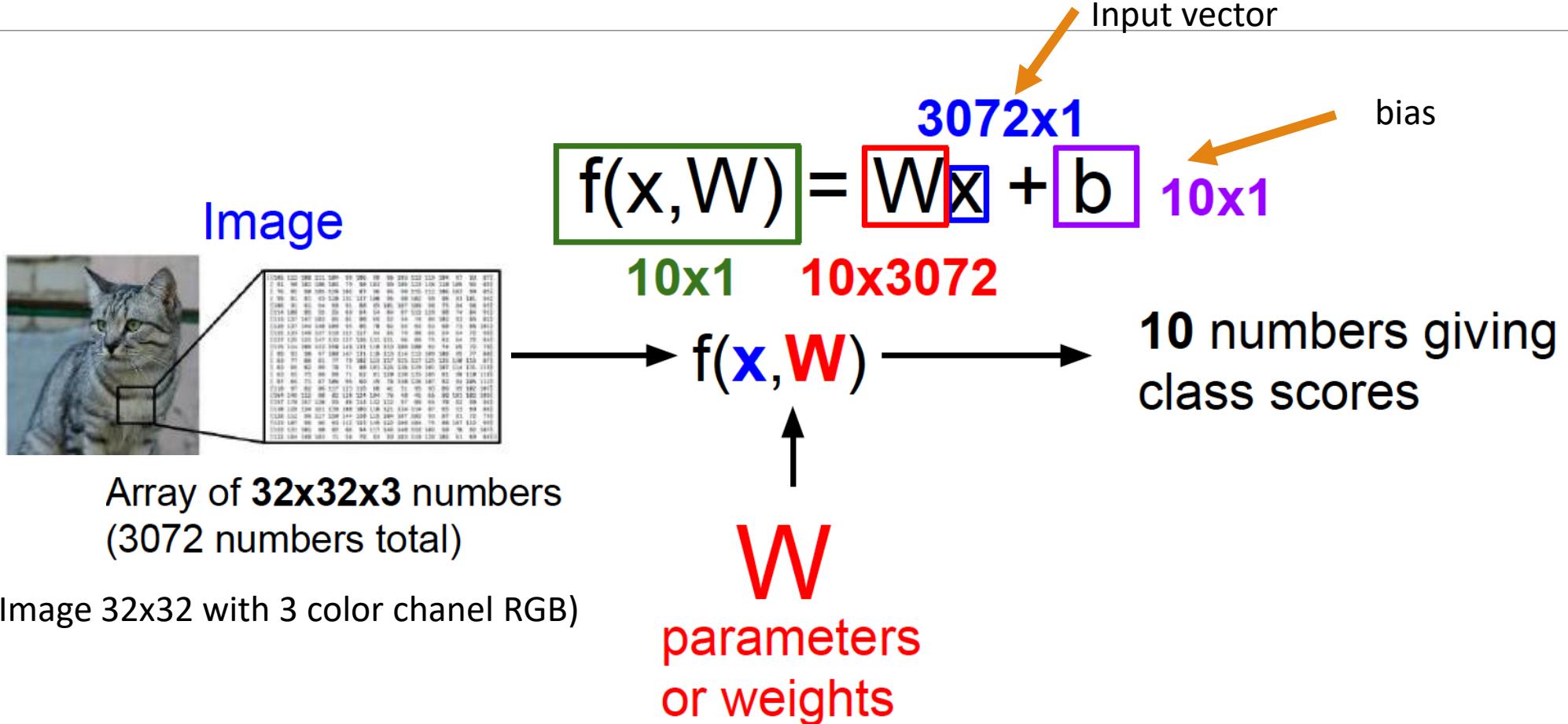
Deng, Dong, Socher, Li, Li, & Fei-Fei, 2009

# Linear Classification

- Find linear expression (*hyperplane*) to separate positive and negative examples

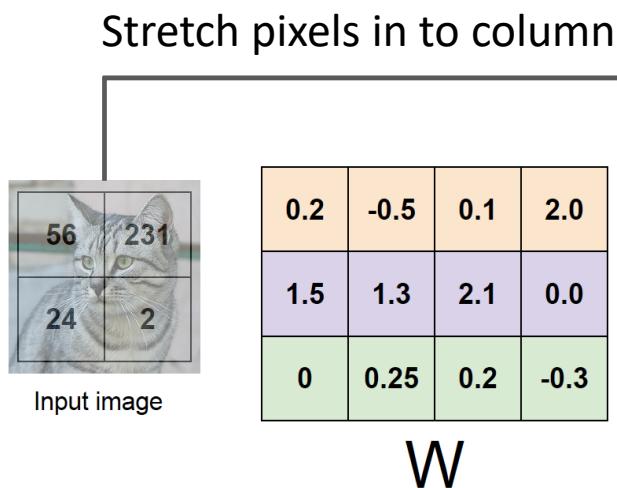


# Linear Classification Recall

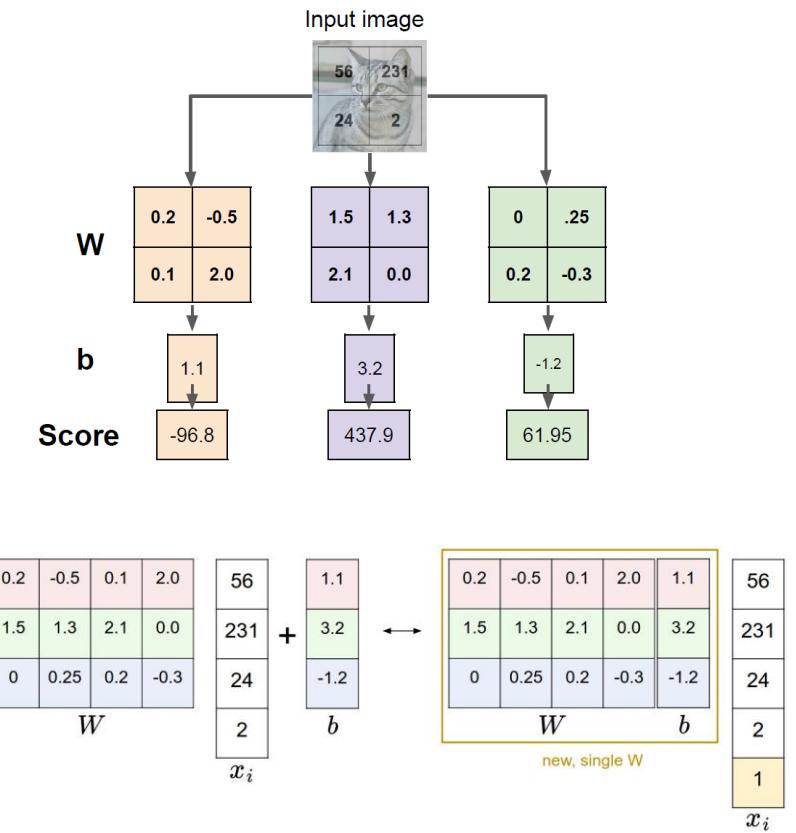


# Linear Classification

Example with an image classification to 3 classes (cat/dog/ship)

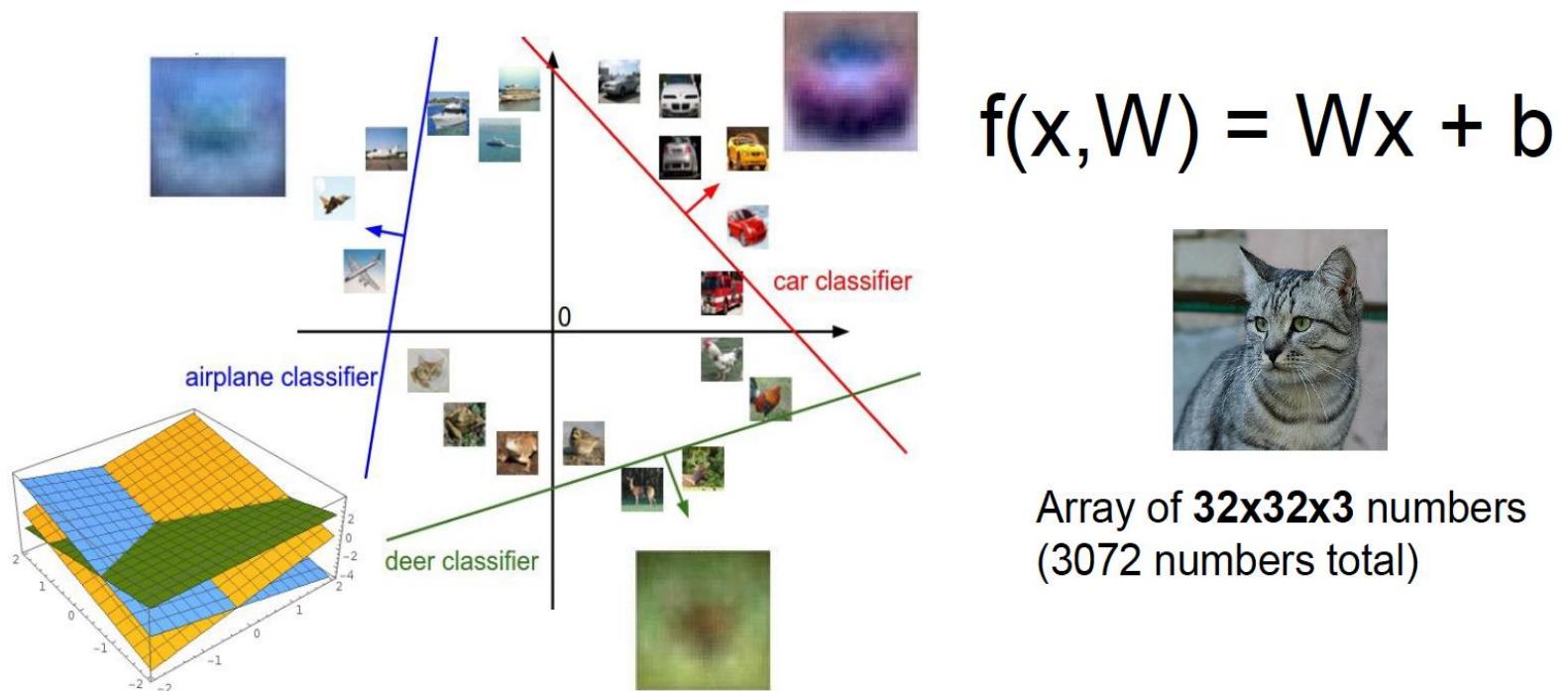


$$f(x, W) = Wx + b$$



# Linear Classification Recall

## Interpreting a Linear Classifier



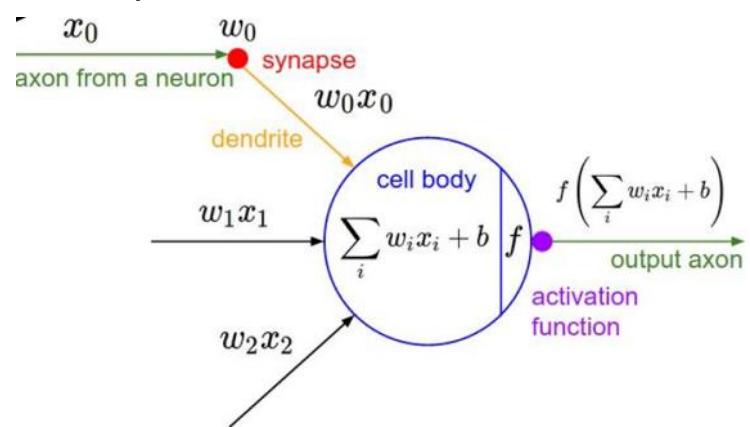
# Neural Network

Linear classification

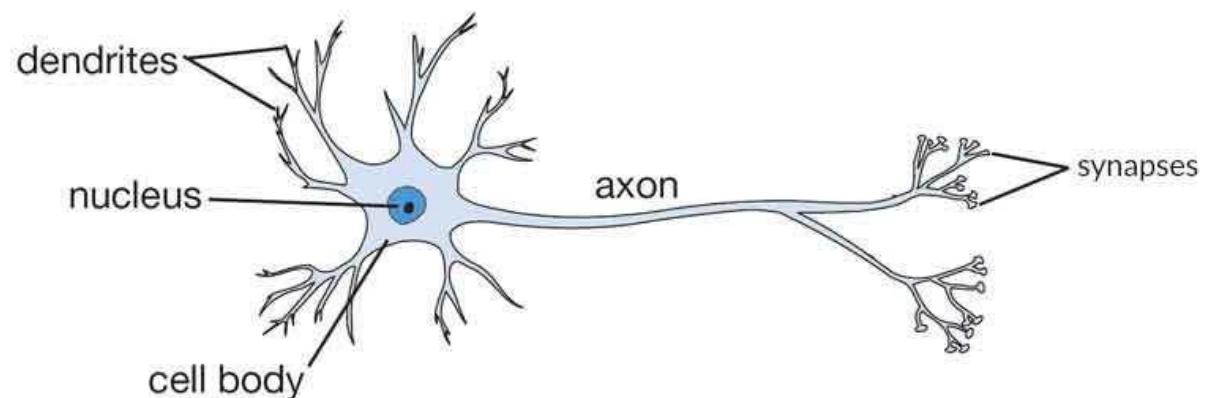
$$\text{output} = \begin{cases} 0 & \text{if } w \cdot x + b \leq 0 \\ 1 & \text{if } w \cdot x + b > 0 \end{cases}$$

$$\text{Output} = \text{sign}(Wx+b)$$

Perceptron

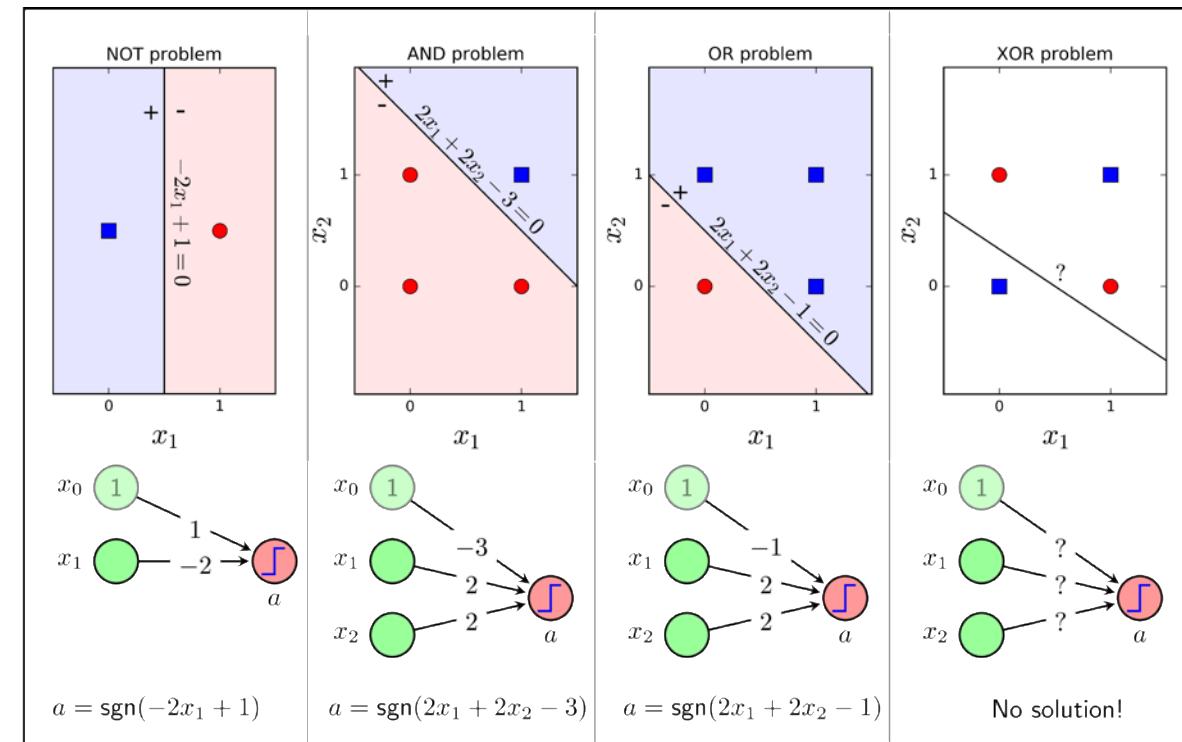


Biological Neuron



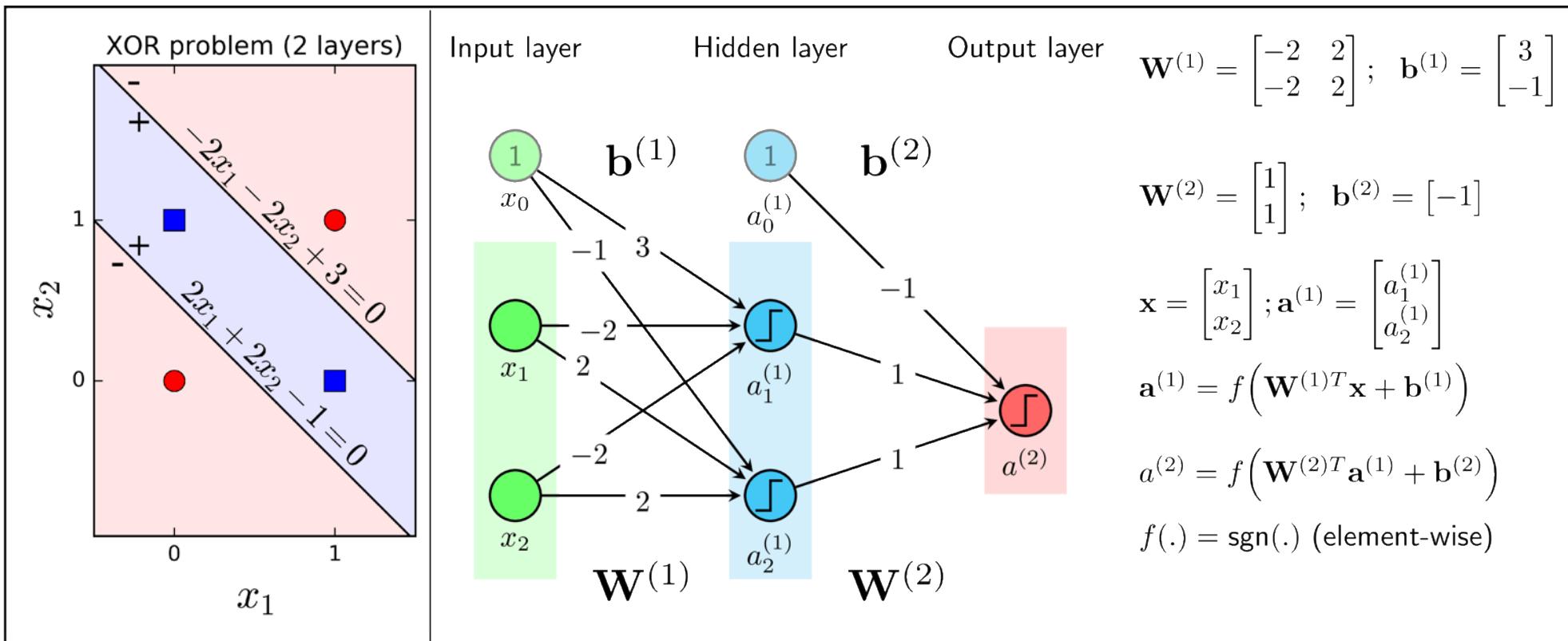
# Neural Network

- ❖ Linear Classification problem
  - ❖ Some data points can be separated by lines
  - ❖ XOR problem: non-linearly separable



# Neural Network

## Multiple perceptron

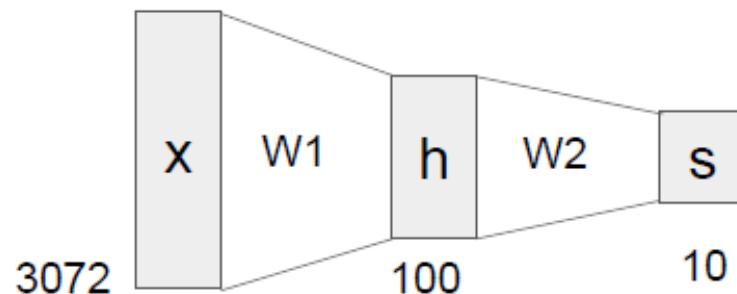


# Multi-perceptron

With 1 hidden layer  
With 2 hidden layer

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

Activate function



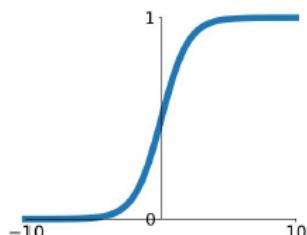
With 3 hidden layer

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

# Activate functions

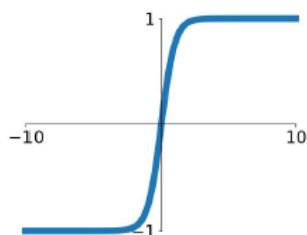
## Sigmoid

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



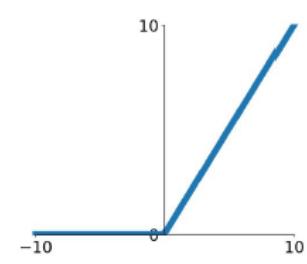
## tanh

$$\tanh(x)$$



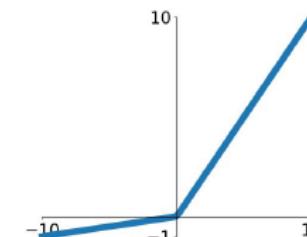
## ReLU

$$\max(0, x)$$



## Leaky ReLU

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

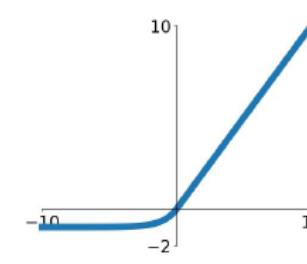


## Maxout

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

## ELU

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

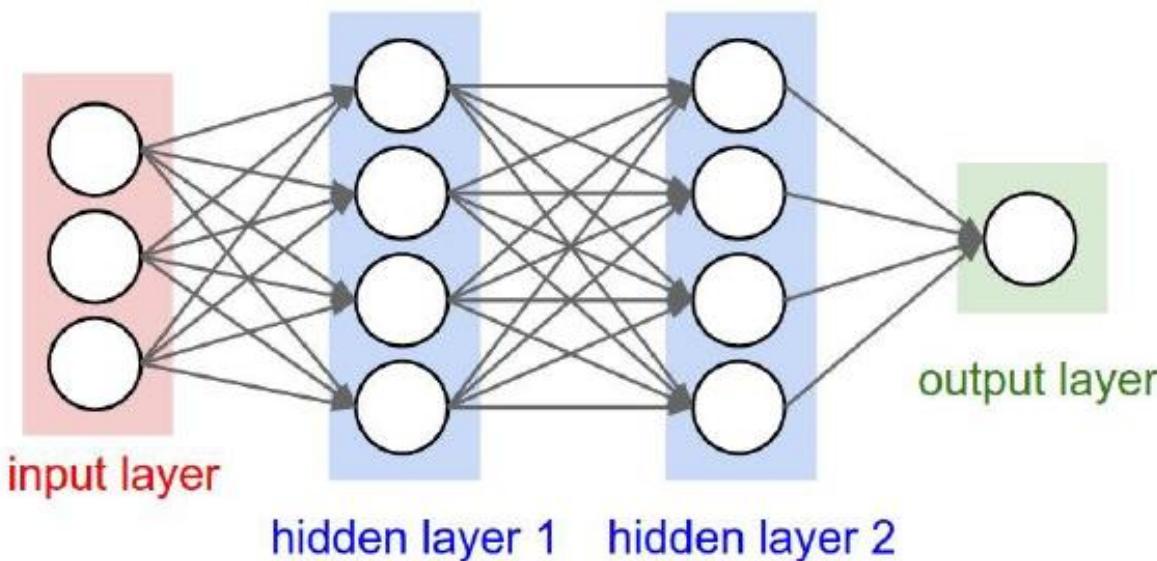


Most practically use

(We will discuss pros and cons in the next lecture)

# Build a Neural Network

---



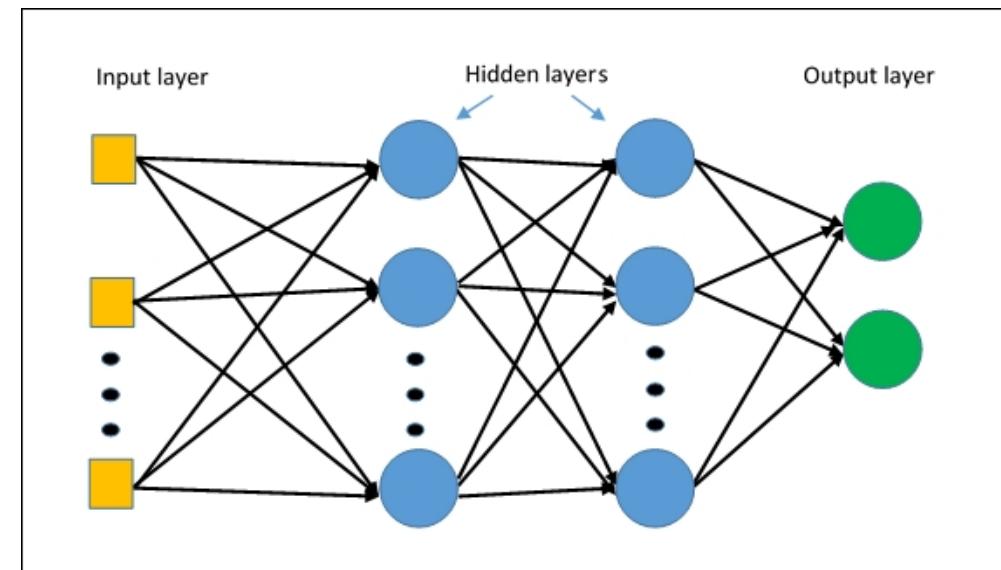
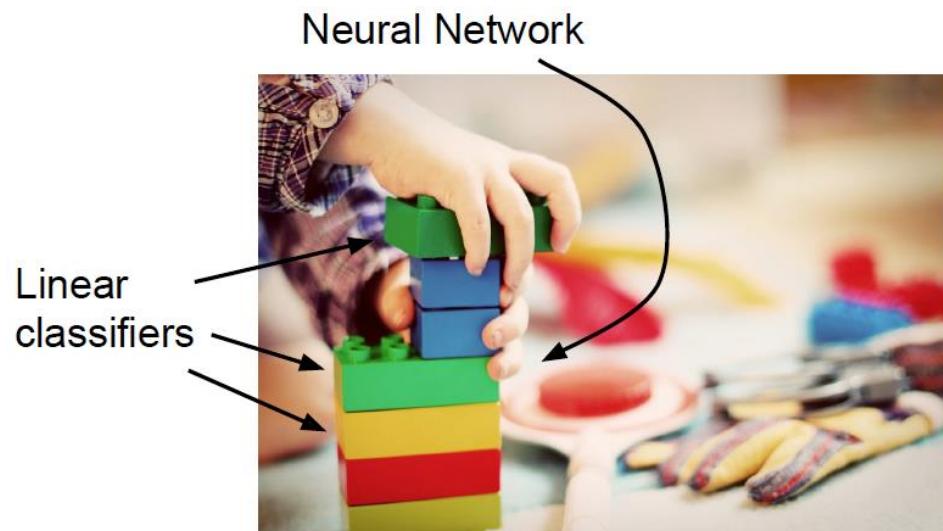
```
# forward-pass of a 3-layer neural network:  
f = lambda x: 1.0/(1.0 + np.exp(-x)) # activation function (use sigmoid)  
x = np.random.randn(3, 1) # random input vector of three numbers (3x1)  
h1 = f(np.dot(W1, x) + b1) # calculate first hidden layer activations (4x1)  
h2 = f(np.dot(W2, h1) + b2)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

# Neural Network

---

## Multi-layers Perceptron

- Layers: Bricks with many types
- Network design: Construction and connection

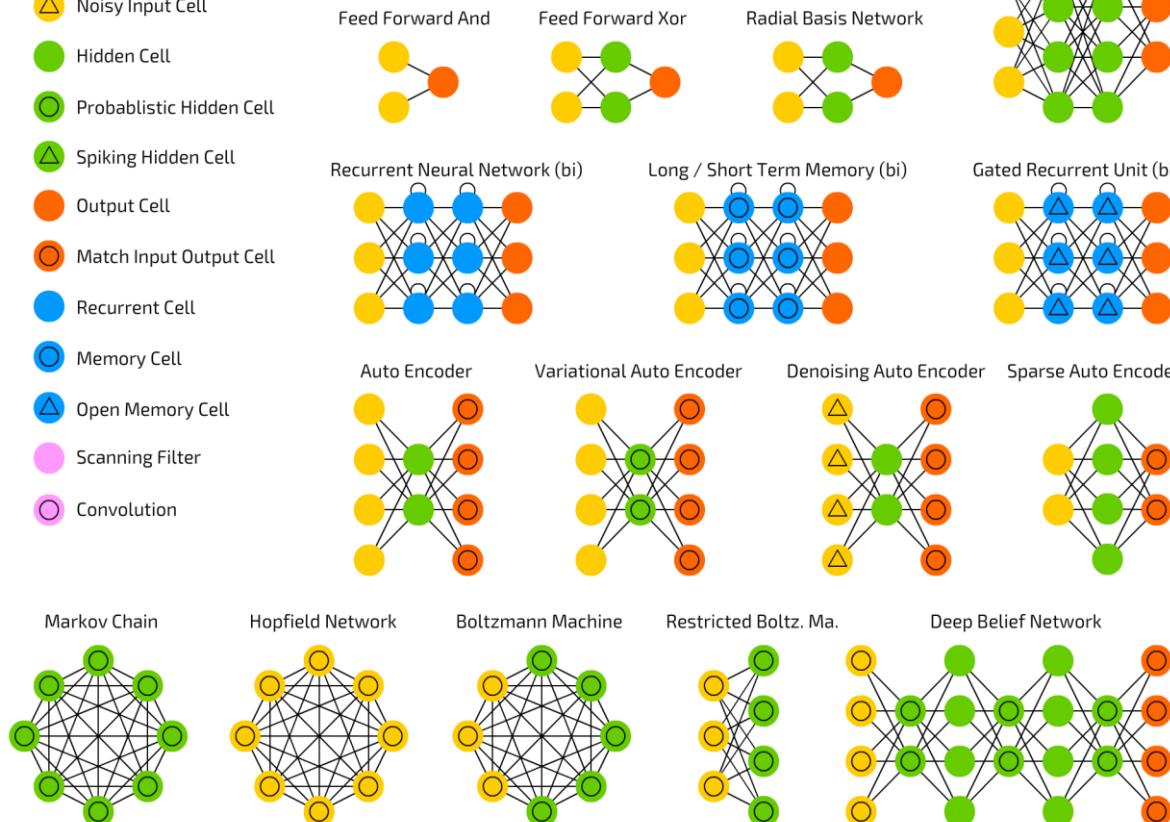


# Neural Network Architects

- Backfed Input Cell
- Input Cell
- △ Noisy Input Cell
- Hidden Cell
- Probabilistic Hidden Cell
- △ Spiking Hidden Cell
- Output Cell
- Match Input Output Cell
- Recurrent Cell
- Memory Cell
- △ Open Memory Cell
- Scanning Filter
- Convolution

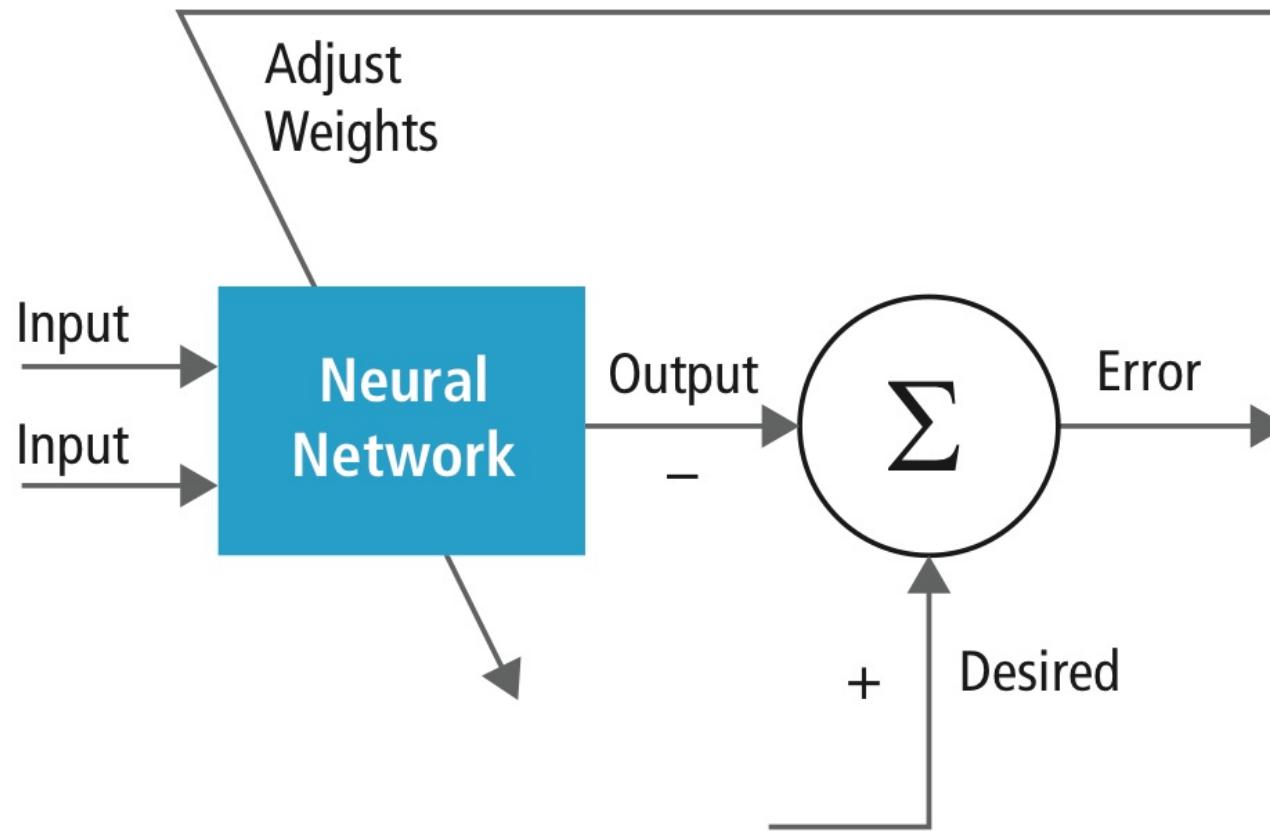
## Neural Networks

A mostly complete chart of architectures  
©2016 Fjodor van Veen



# Training Pipeline

---



# Training Objective

Define a **loss function** that quantifies our unhappiness with the scores across the training data.

Come up with a way of efficiently finding the parameters that minimize the loss function.  
**(optimization)**



airplane	-3.45	-0.51	3.42
automobile	-8.87	<b>6.04</b>	4.64
bird	0.09	5.31	2.65
<b>cat</b>	<b>2.9</b>	-4.22	5.1
deer	4.48	-4.19	2.64
dog	8.02	3.58	5.55
frog	3.78	4.49	<b>-4.34</b>
horse	1.06	-4.37	-1.5
ship	-0.36	-2.09	-4.79
truck	-0.72	-2.93	6.14

Cat image by Nikita is licensed under CC-BY 2.0; Car image is CC0 1.0 public domain; Frog image is in the public domain

# Loss Function

Task: Training set consist of 3 images with 3 classes

Scoring with parameter W in  $f(x, W) = Wx$



cat	<b>3.2</b>	1.3	2.2
car	5.1	<b>4.9</b>	2.5
frog	-1.7	2.0	<b>-3.1</b>
Losses:	<b>2.9</b>	<b>0</b>	<b>12.9</b>

Multiple SVM loss

Training set ( $x_i, y_i$ ):

- $x_i$ : image vectors
- $y_i$ : labels
- Loss on image i

Loss over data set:

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

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

$$\begin{aligned} L &= (2.9 + 0 + 12.9)/3 \\ &= 5.27 \end{aligned}$$

# Softmax classifier

From score to probabilities



$$s = f(x_i; W)$$

$$P(Y = k|X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

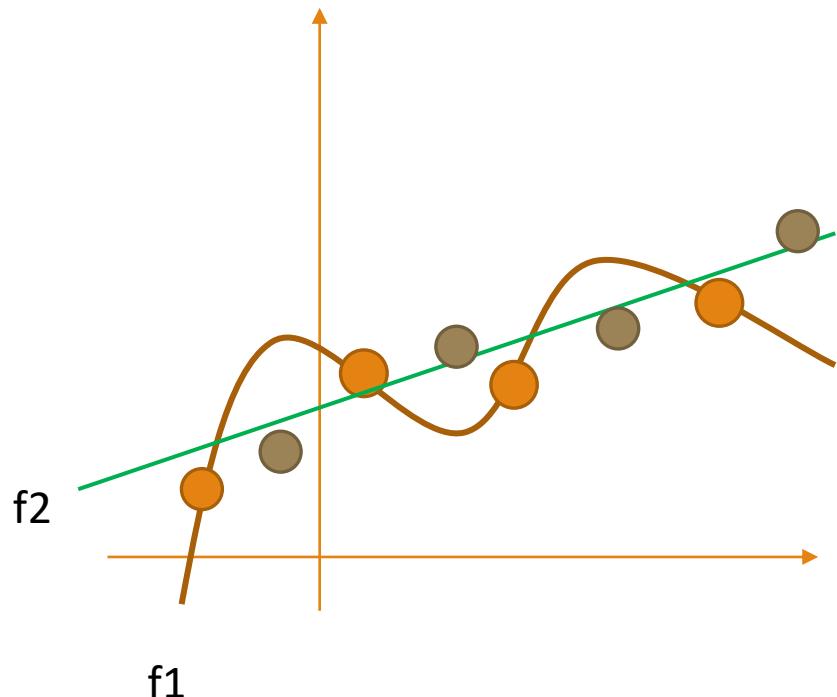
Softmax  
Function

$$L_i = -\log P(Y = y_i|X = x_i)$$

	Score (Unnormalized)			Probs (P)	Compare (By KL Divergence)	Correct Probs (Q)
cat	3.2	$\xrightarrow{\text{exp}}$	24.5	0.13	1.00	0.00
car	5.1	$\xrightarrow{\text{exp}}$	164.0	0.87	0.00	0.00
frog	-1.7		0.18	0.00	0.00	0.00

**Cross Entropy  $H(P,Q)$**

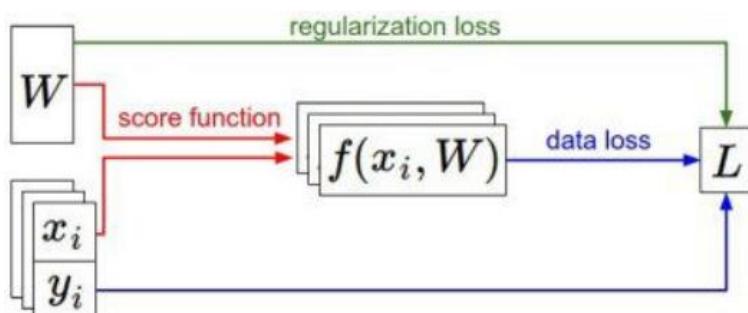
# Regularization



$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(f(x_i, W), y_i) + \lambda R(W)$$

**Data loss:**  
From the model

**Regularization:**  
To reduce the complexity of the model



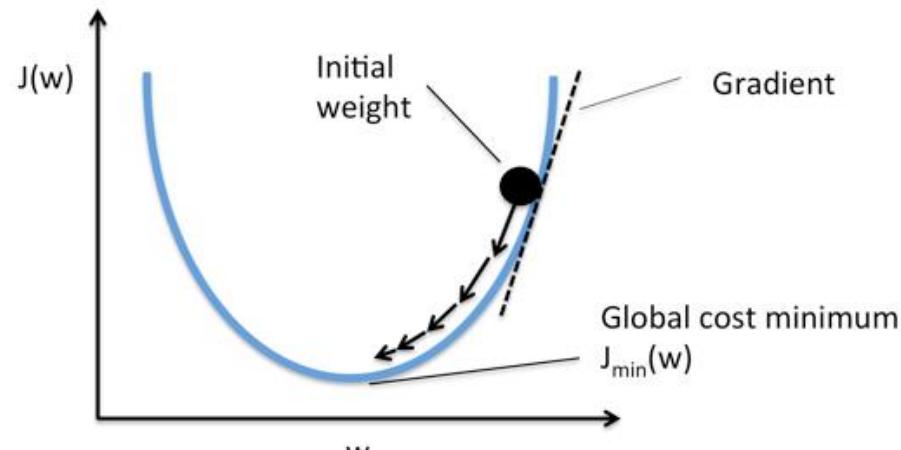
**Reference:**

L2 regularization:  $R(W) = \sum_k \sum_l W_{k,l}^2$

L1 regularization:  $R(W) = \sum_k \sum_l |W_{k,l}|$

Elastic net (L1 + L2):  $R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$

# Gradient Descent (More in supplements)

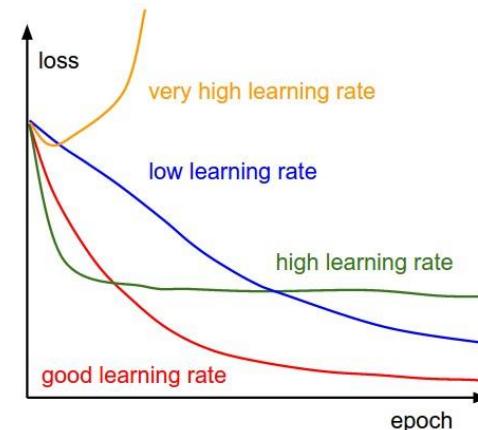
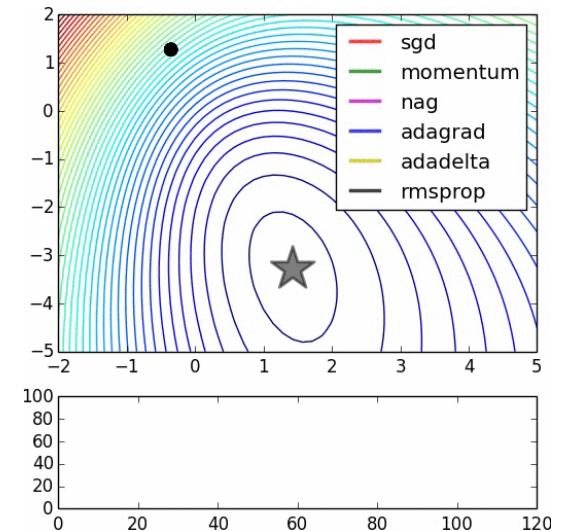
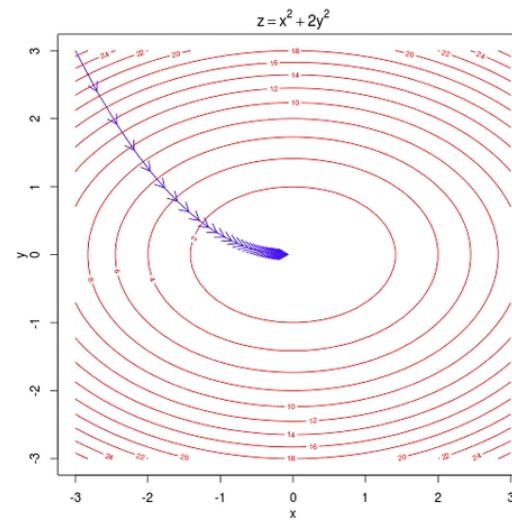


**Algorithm:** Gradient Descent

```

Input:  $Y, \Theta, X, \alpha$ , tolerance, max iterations
Output:  $\Theta$ 
1 for  $i = 0; i < \text{max iterations}; i++$  do
2   current cost =  $\text{Cost}(Y, X, \Theta)$ 
3   if current cost < tolerance then
4     break
5   else
6     gradient =  $\text{Gradient}(Y, X, \Theta)$ 
7      $\theta_j \leftarrow \theta_j - \alpha \cdot \text{gradient}$ 

```

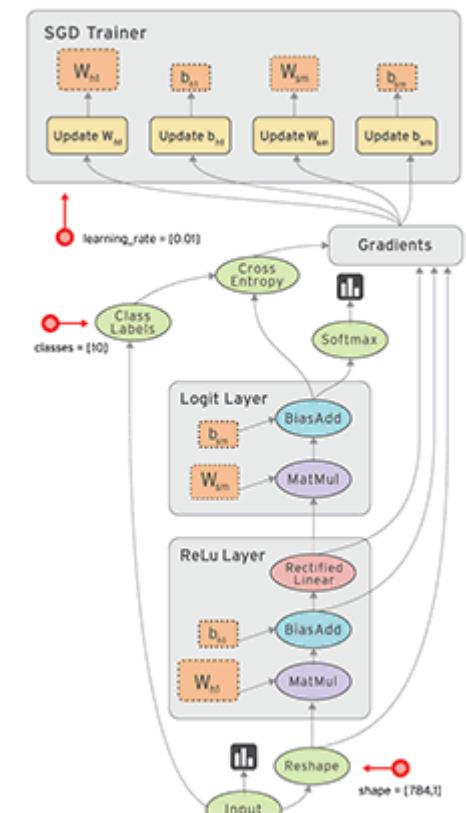
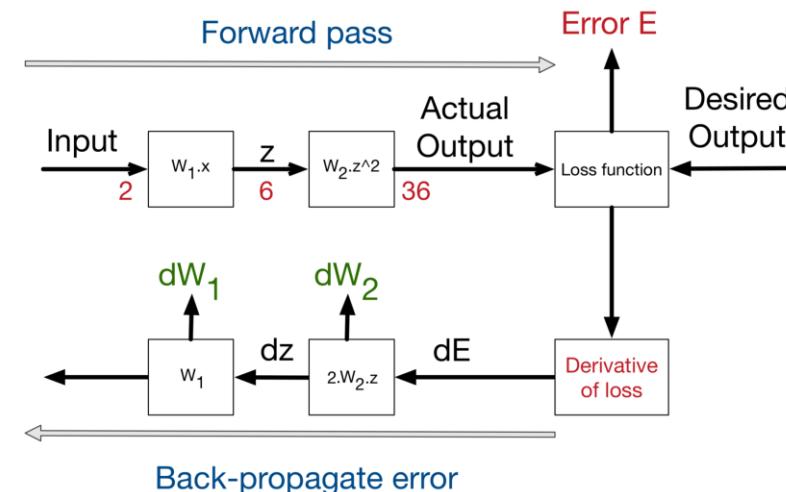
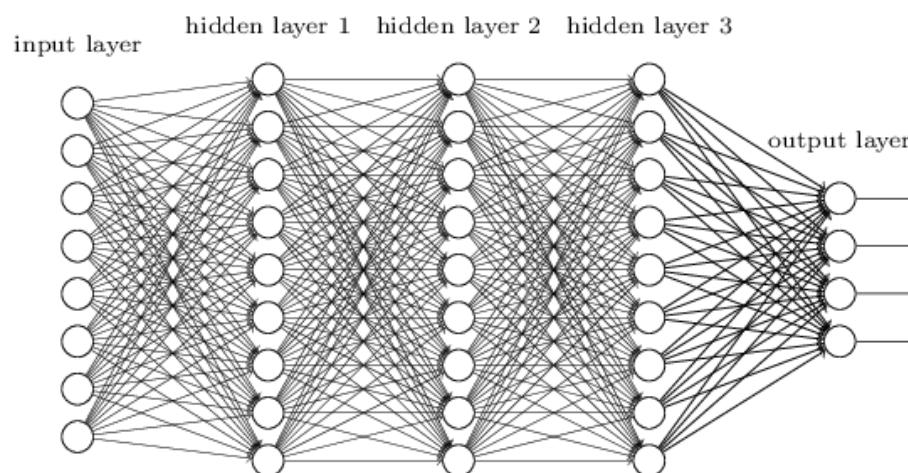


- (Mini) Batch SGD: Using (small) batch instead of each image

# Back propagation (More in Supplements)

Calculating gradient for any  $w_i$  in multiple layer perceptrons

(Considering as black boxes if you do not want to discover )



# Deep Learning

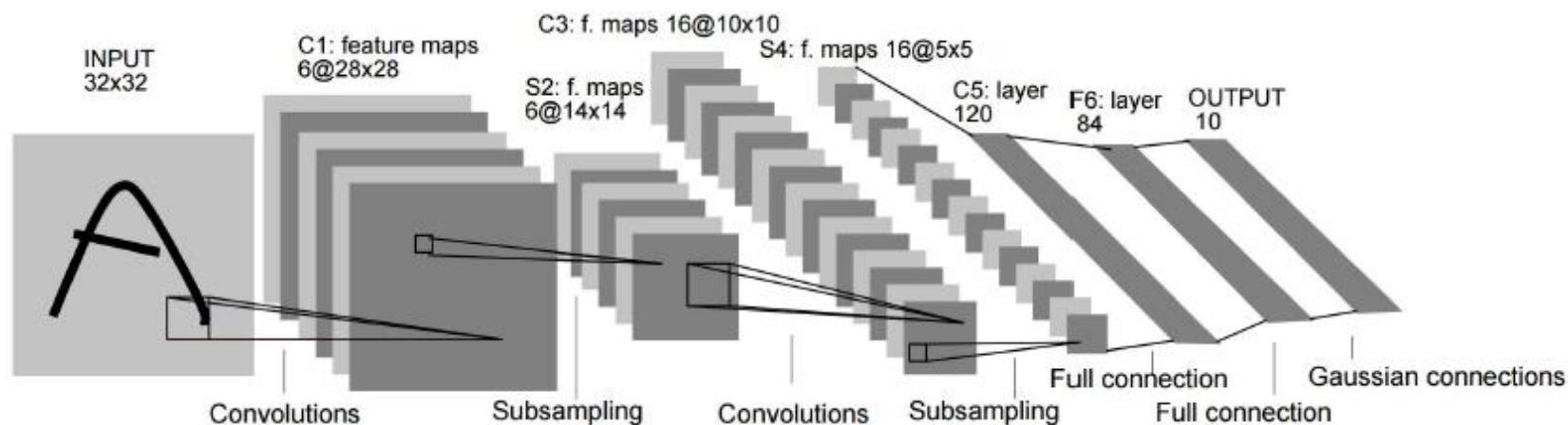
---

NETWORK STRUCTURES AND LEARNING METHODS

# Convolutional Neural Network

---

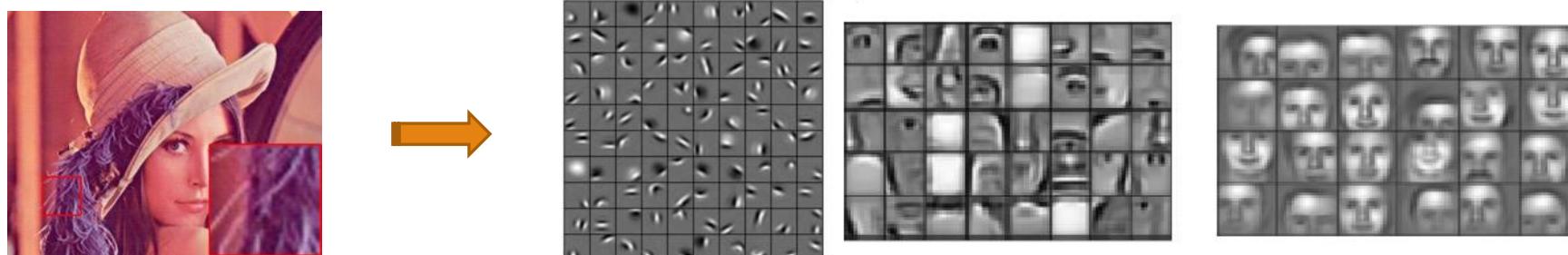
LeCun, Yann, et al. "Gradient-based learning applied to document recognition." Proceedings of the IEEE 86.11 (1998): 2278-2324.



# Image Classification Challenge

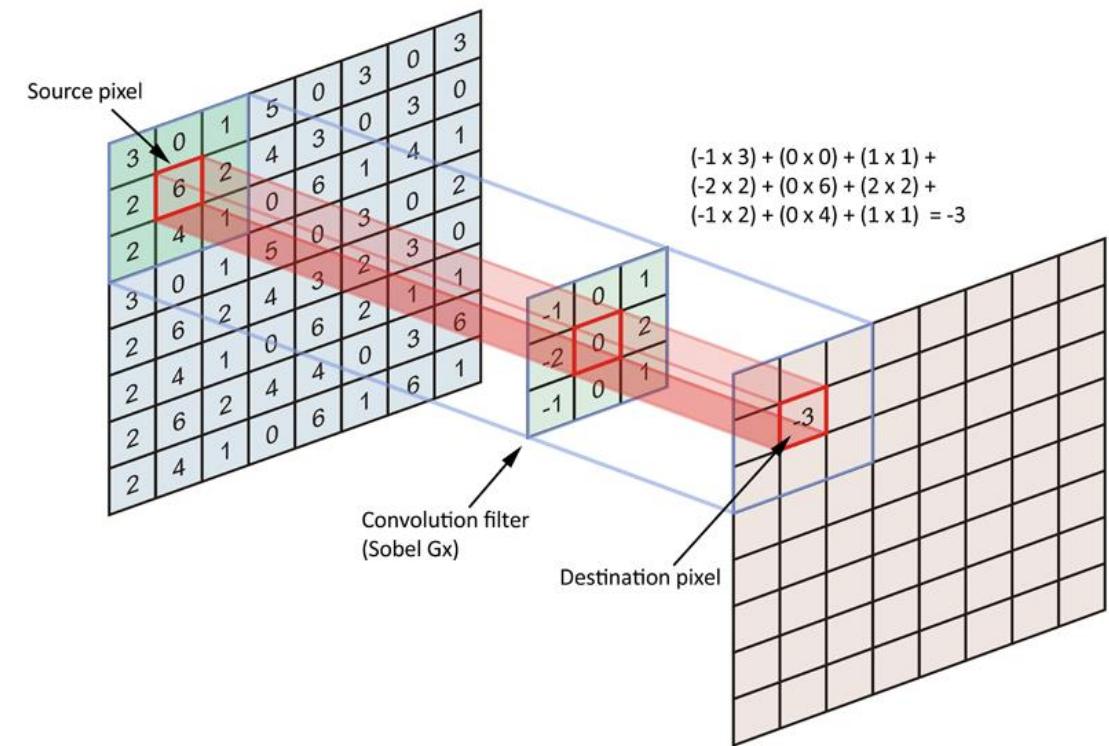
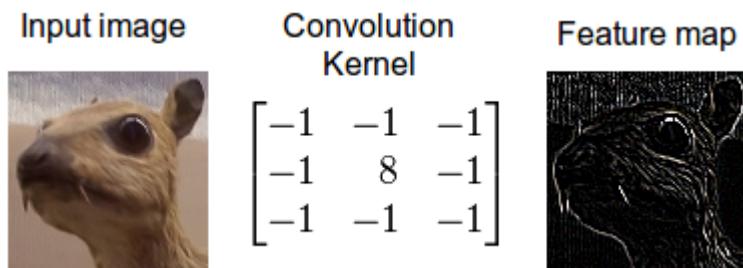
---

- Locality: objects tend to have a local spatial support
- Translation invariance: object appearance is independent of location
  - Can define these properties since an image lies on a grid/lattice
  - ConvNet machinery applicable to other data with such properties, e.g. audio/text
- Recognition: Supporting from hidden feature representation



# Locality

- Make fully-connected layer locally-connected
- Each unit/neuron is connected to a local rectangular area – receptive field
- Different units connected to different locations
- output (“feature map”) lies on a grid itself



# Convolutional Kernels

Identity



0	0	0
0	1	0
0	0	0

Blur

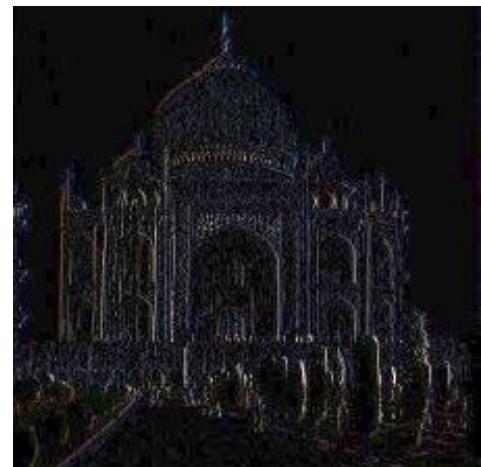
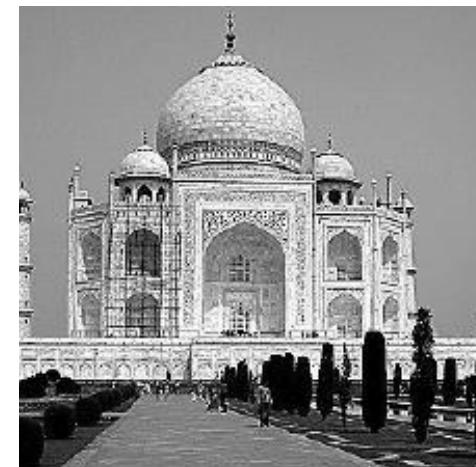
1	1	1
1	1	1
1	1	1

Sharpen

0	-1	0
-1	4	-1
0	-1	0

Edge detector

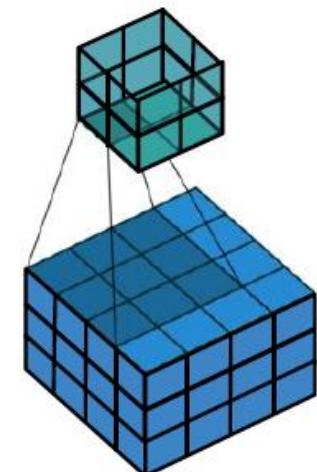
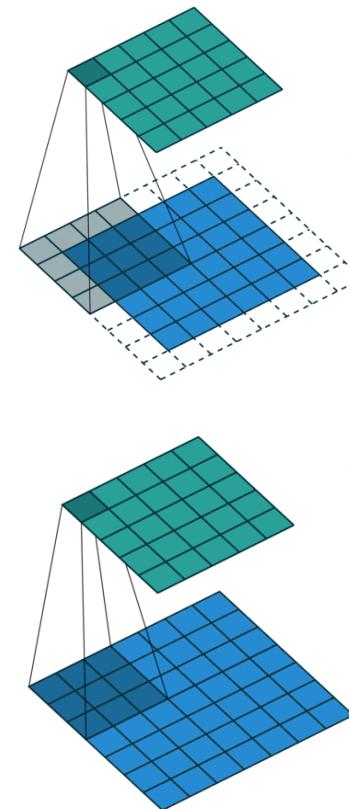
0	0	0
-1	1	0
0	0	0



# Translation Invariance

---

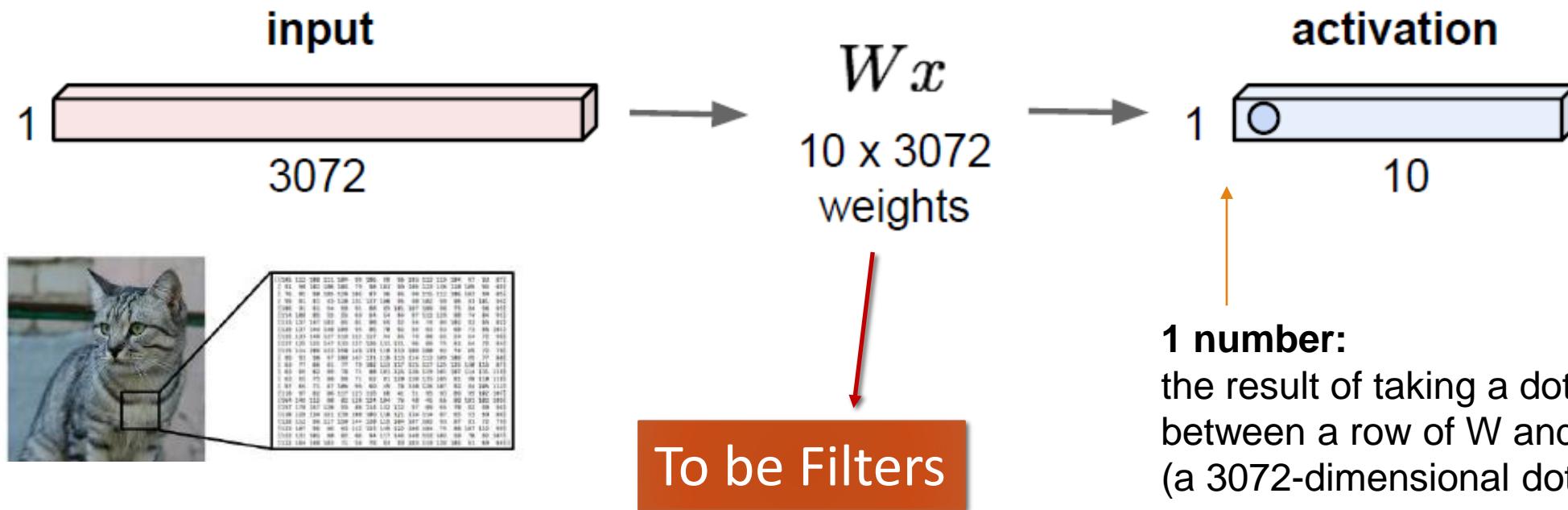
- Weight sharing
  - units connected to different locations have the same weights
  - equivalently, each unit is applied to all locations
- Convolutional layer – locally-connected layer with weight sharing (translation invariance)
- The weights are invariant, the output is equivariant
- Convolutional layer input and output can have multiple channels



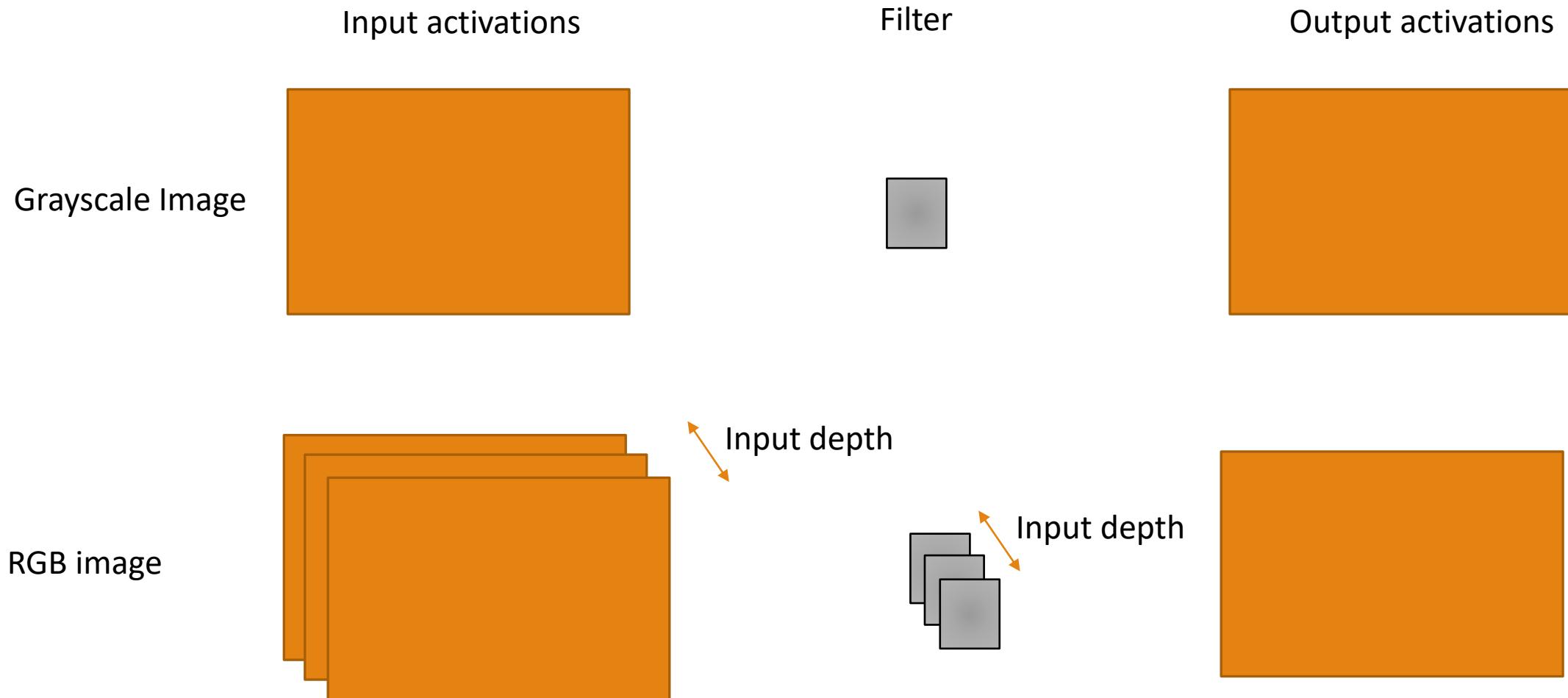
feature maps are 3-D tensors  
height × width × channels

# Fully connected to Convolution Kernels

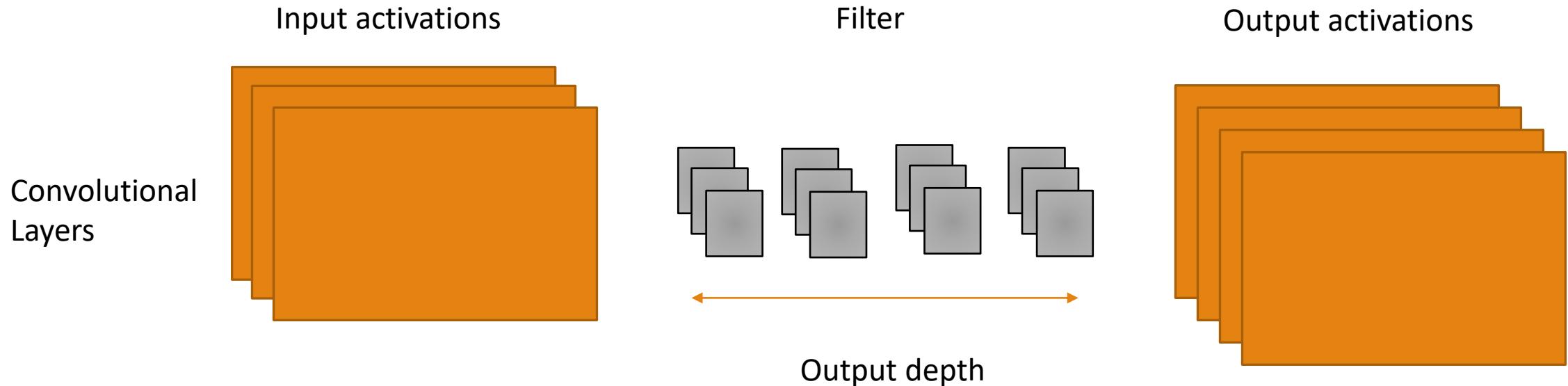
32x32x3 image -> stretch to 3072 x 1



# Generalizing Convolutions

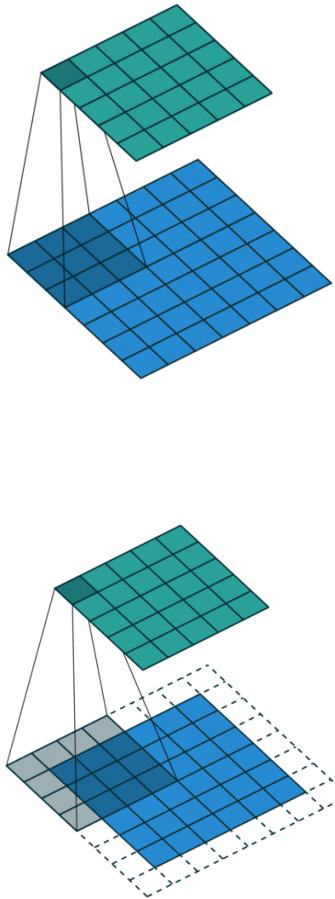


# Generalizing Convolutions



- Input and output depth are arbitrary parameters and not equal.

# Strides and Padding



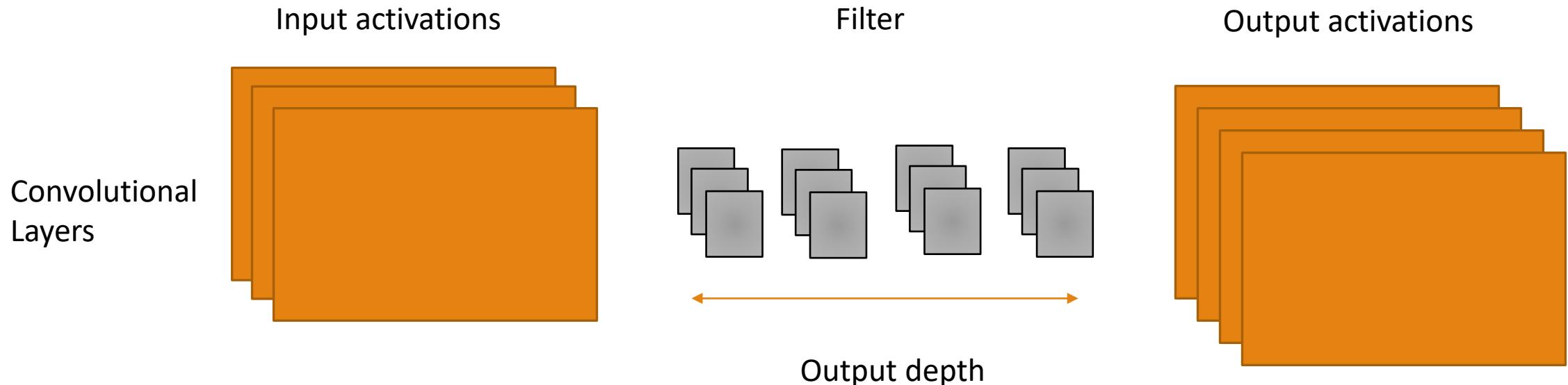
stride: # pixels to shift  
when applying a kernel

padding: what to do at the  
edges: valid or same

0 <sub>2</sub>	0 <sub>0</sub>	0 <sub>1</sub>	0	0	0	0
0 <sub>1</sub>	2 <sub>0</sub>	2 <sub>0</sub>	3	3	3	0
0 <sub>0</sub>	0 <sub>1</sub>	1 <sub>1</sub>	3	0	3	0
0	2	3	0	1	3	0
0	3	3	2	1	2	0
0	3	3	0	2	3	0
0	0	0	0	0	0	0

1 6 5  
7 10 9  
7 10 8

# Generalizing Convolutions



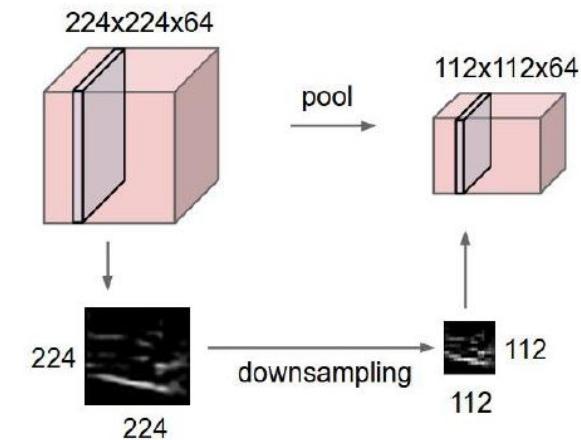
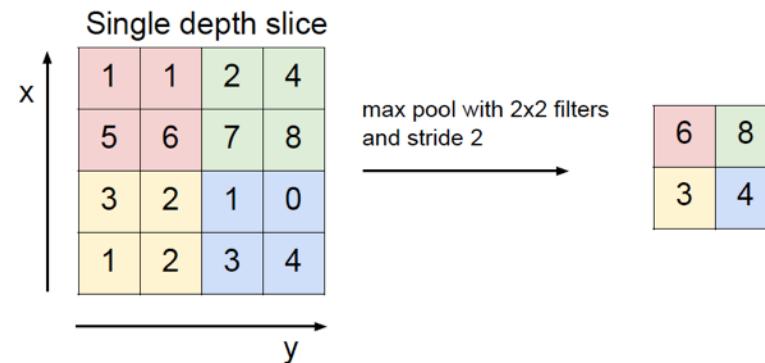
- How many parameters are in a single layer?  
 $(\text{filter width} \times \text{filter height}) \times (\text{input depth}) \times (\text{output depth})$
- How much computational cost in a single layer?  
 $(\text{filter width} \times \text{filter height}) \times (\text{input depth}) \times (\text{output depth}) \times (\text{input width / stride}) \times (\text{input height / stride})$

# Recognition from hidden features

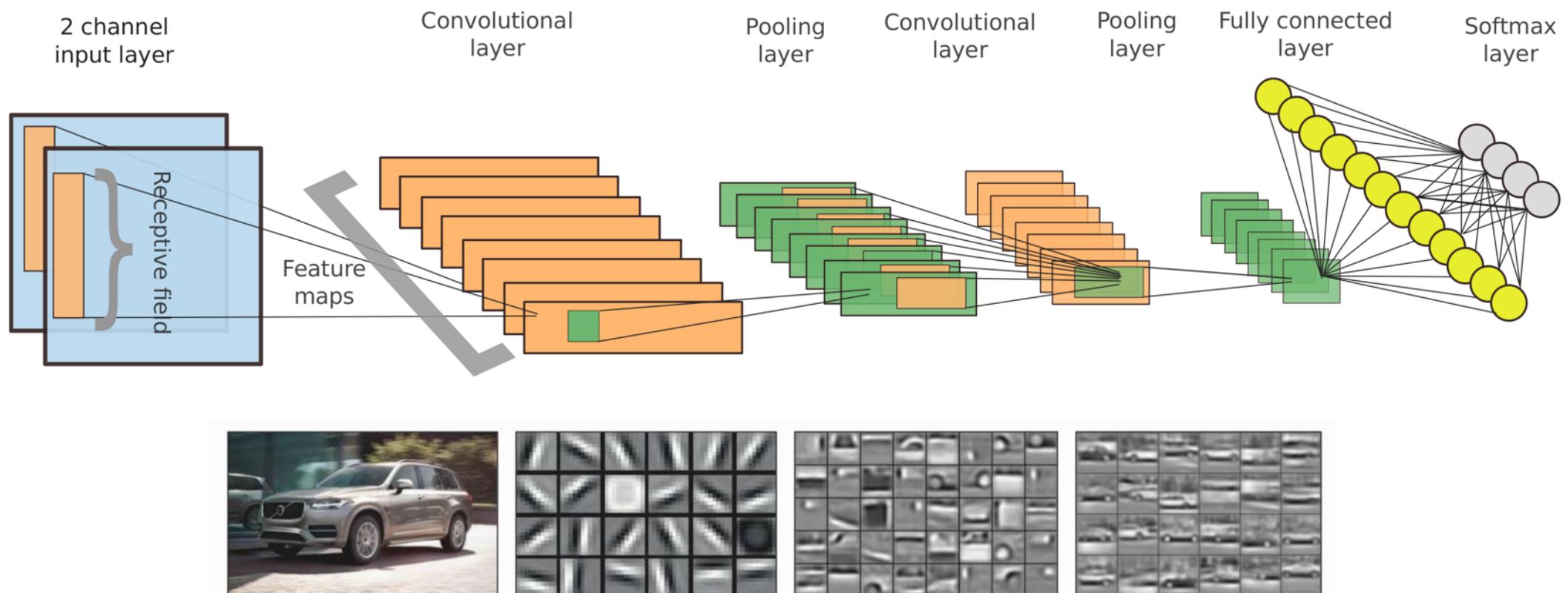
## Pooling layer:

- makes the representations smaller and more manageable
- operates over each activation map independently

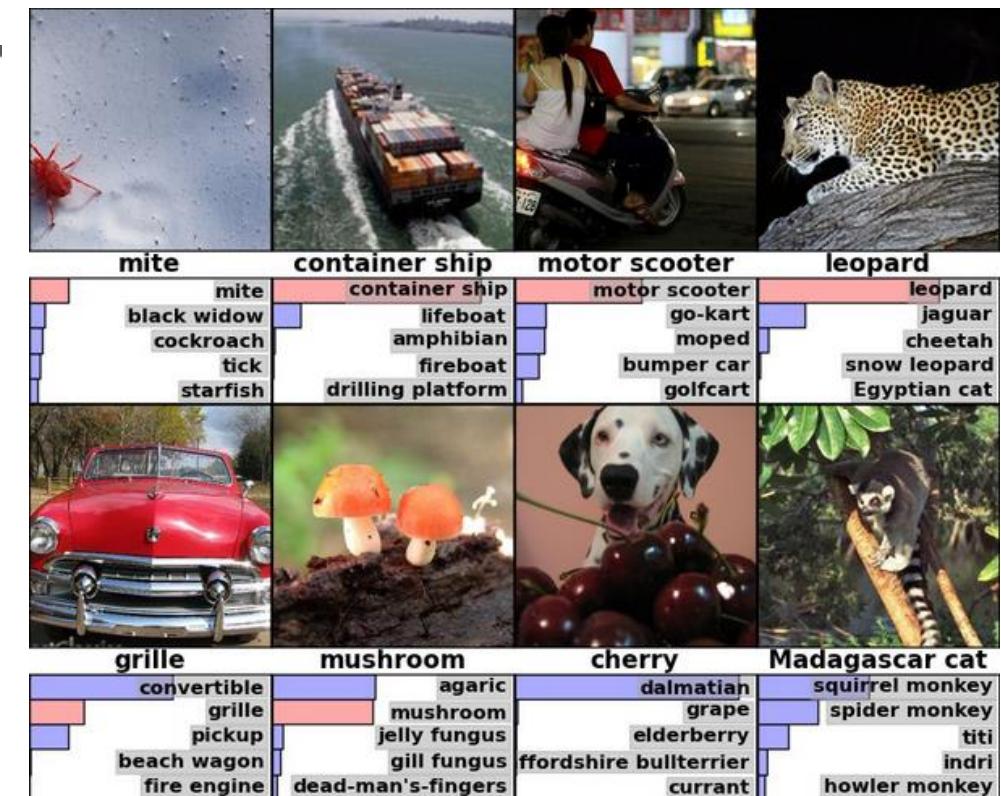
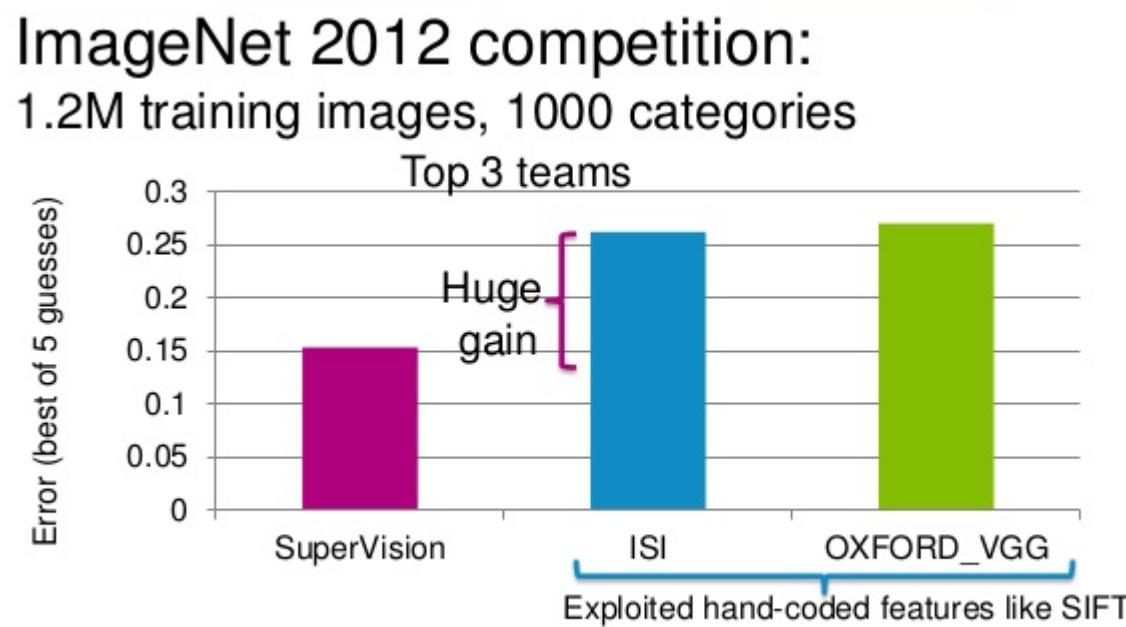
## Max pooling:



# Typical structure for image classification



# Object Recognition Show



# Case Study: AlexNet

Full (simplified) AlexNet architecture:

[227x227x3] INPUT

[55x55x96] CONV1: 96 11x11 filters at stride 4, pad 0

[27x27x96] MAX POOL1: 3x3 filters at stride 2

[27x27x96] NORM1: Normalization layer

[27x27x256] CONV2: 256 5x5 filters at stride 1, pad 2

[13x13x256] MAX POOL2: 3x3 filters at stride 2

[13x13x256] NORM2: Normalization layer

[13x13x384] CONV3: 384 3x3 filters at stride 1, pad 1

[13x13x384] CONV4: 384 3x3 filters at stride 1, pad 1

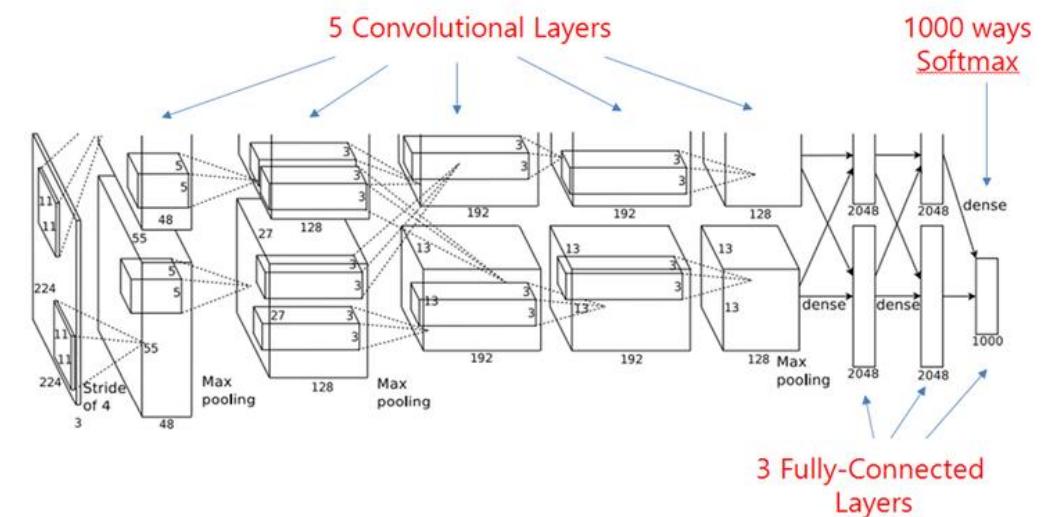
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1

[6x6x256] MAX POOL3: 3x3 filters at stride 2

[4096] FC6: 4096 neurons

[4096] FC7: 4096 neurons

[1000] FC8: 1000 neurons (class scores)



[Krizhevsky et al. 2012]

# Case Study: VGGNet

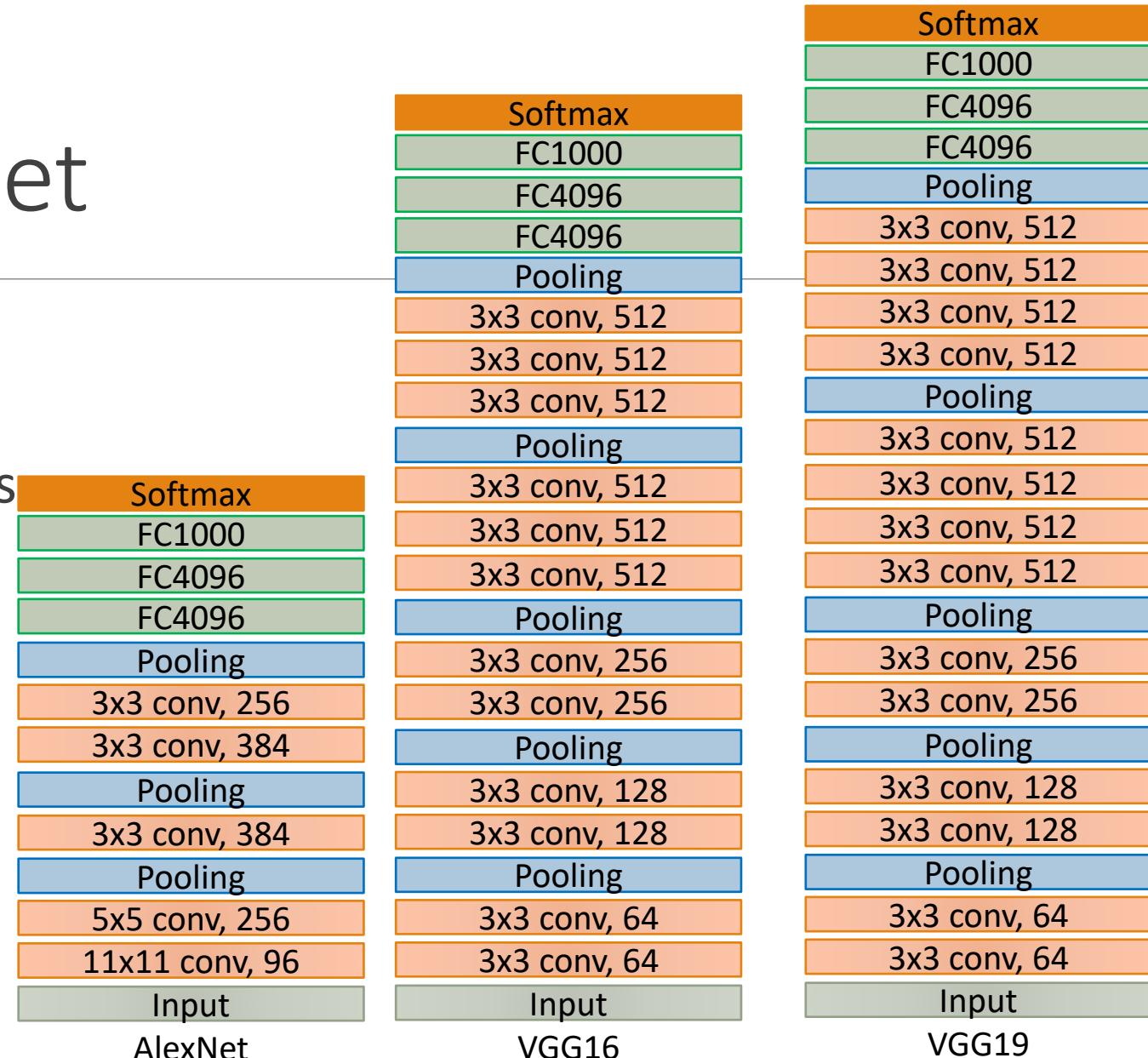
Small filters, Deeper networks  
[Simonyan and Zisserman, 2014]

8 layers (AlexNet) -> 16 - 19 layers  
(VGG16Net)

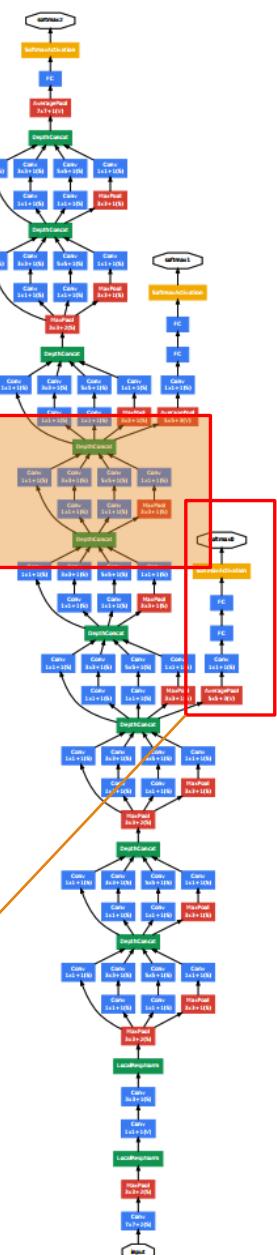
Stack of three 3x3 conv (stride 1)  
layers has same effective  
receptive field as one 7x7 conv  
layer

Deeper, more non-linearities

And fewer parameters

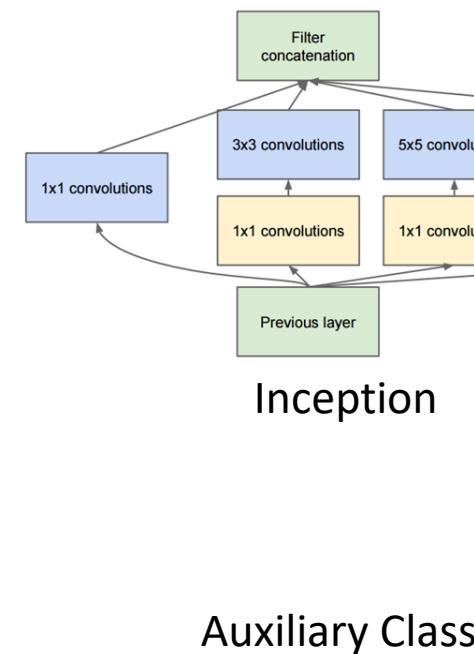


# Case Study: GoogleNet [Szegedy et al., 2014]



## Inception Module: Network in Network

- Apply parallel filter operations on the input from previous layer:
  - Multiple receptive field sizes for convolution (1x1, 3x3, 5x5)
  - Pooling operation (3x3)
- Concatenate all filter outputs together depth-wise
- Auxiliary Classifier: Auxiliary classification outputs to inject additional gradient at lower layers



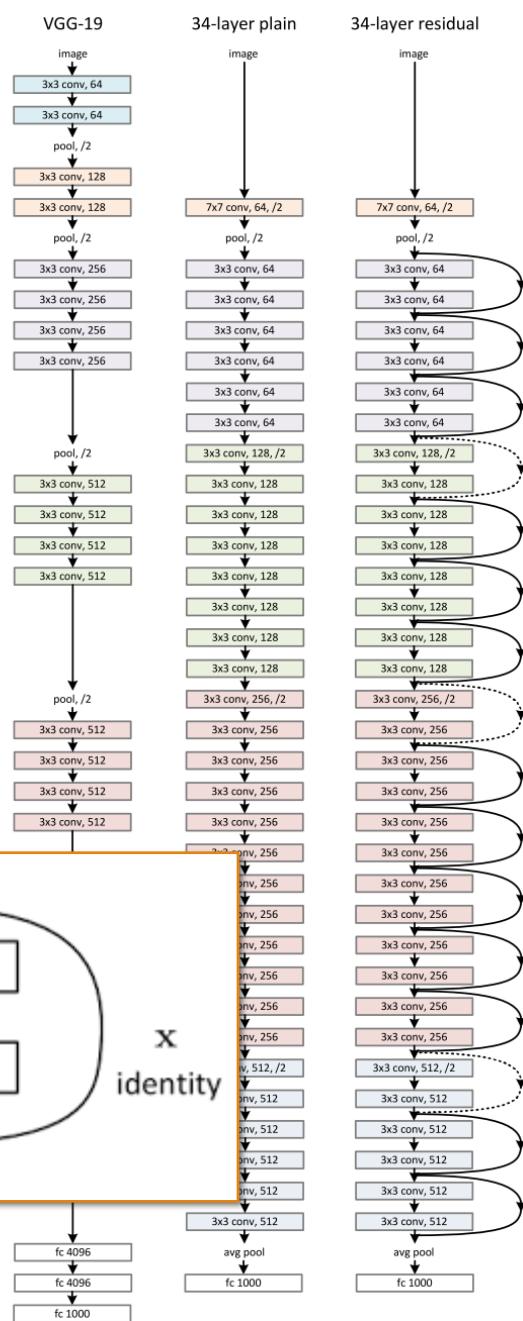
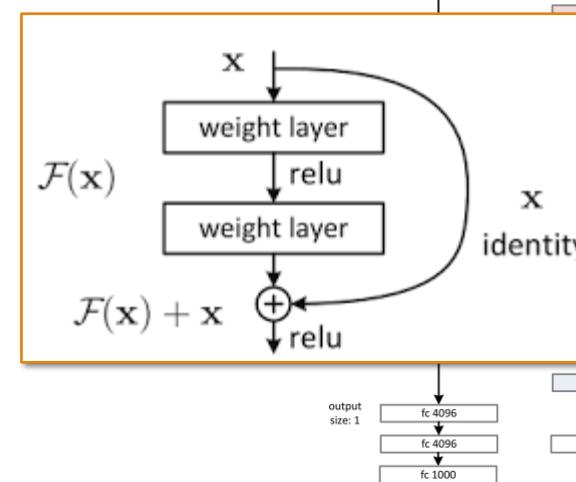
Inception

Auxiliary Classifier

# Case Study: ResNet

Very deep networks using residual connections

- 152-layer model for ImageNet
- ILSVRC'15 classification winner (3.57% top 5 error)
- Swept all classification and detection competitions in ILSVRC'15 and COCO'15!
- Optimization problem: The deeper model should be able to perform at least as well as the shallower model.
- Many approaches to deal



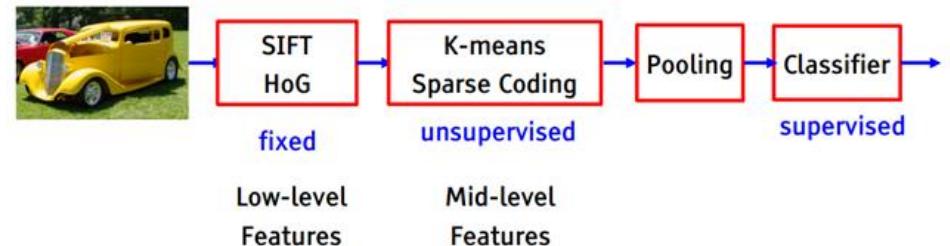
# Why Deep Learning

2010-11: hand-crafted computer vision pipelines

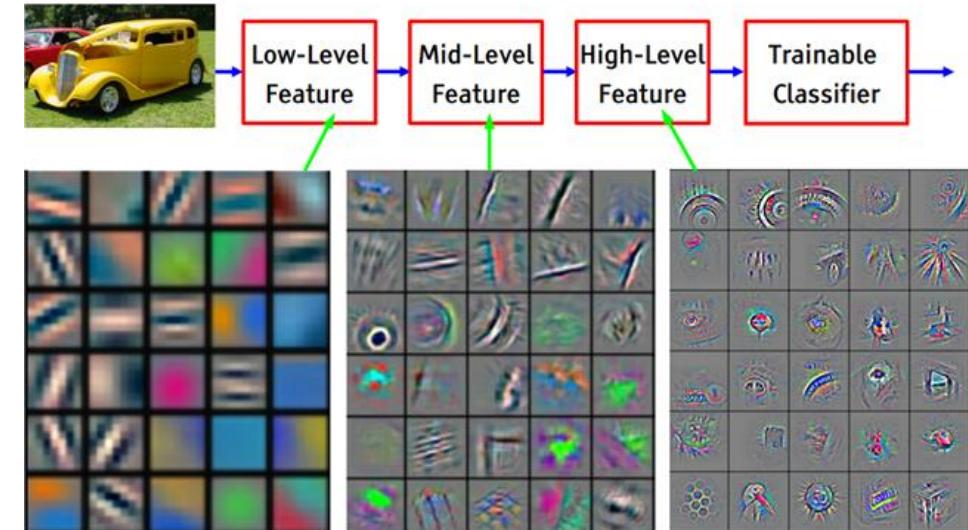
2012-2016: ConvNets, now tend to feature learning

- 2012: AlexNet : major deep learning success
- 2013: ZFNet: improvements over AlexNet
- 2014:
  - VGGNet: deeper, simpler
  - InceptionNet: deeper, faster
- 2015: ResNet: even deeper
- 2016: ensembled networks
- 2017: Squeeze and Excitation Network

Object recognition 2006-2012

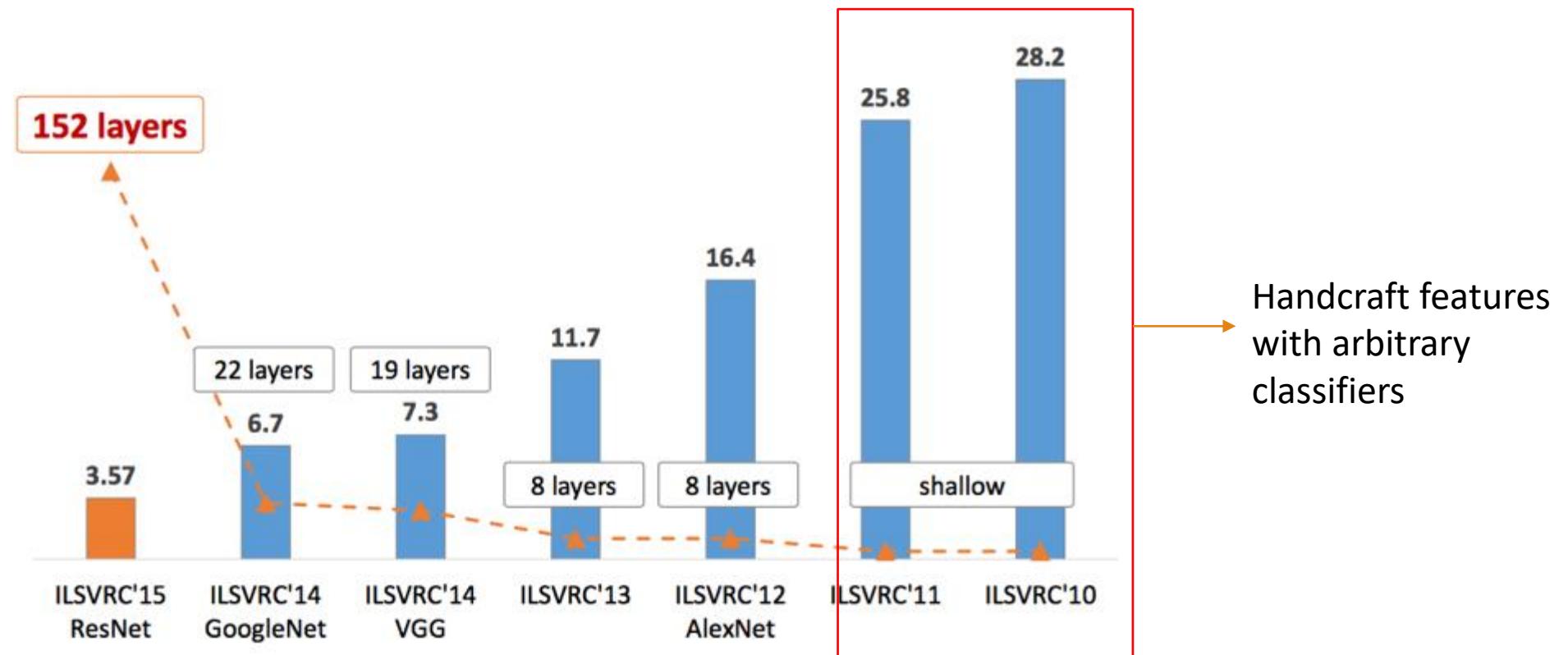


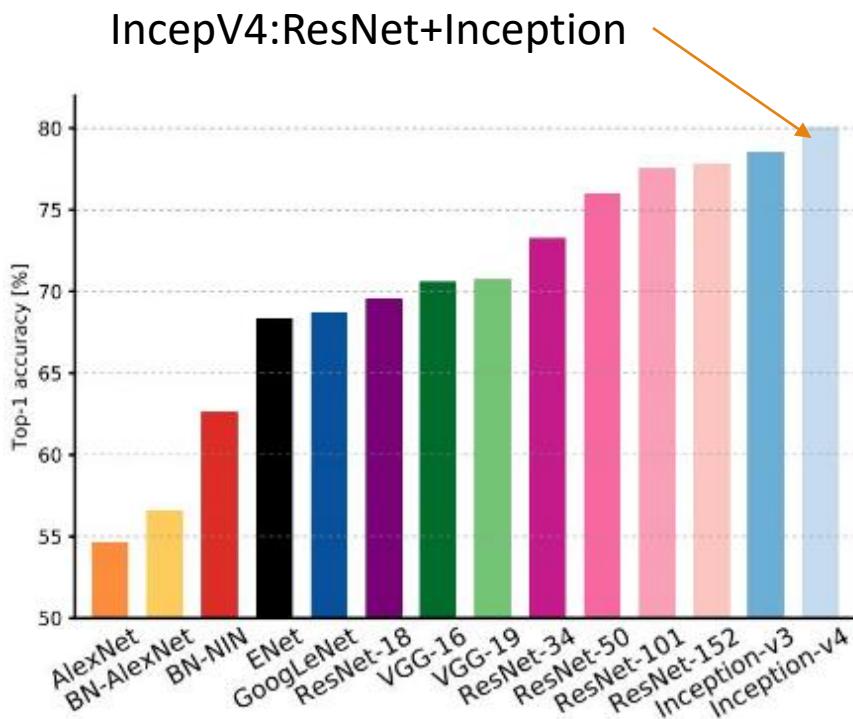
State of the art object recognition using CNNs



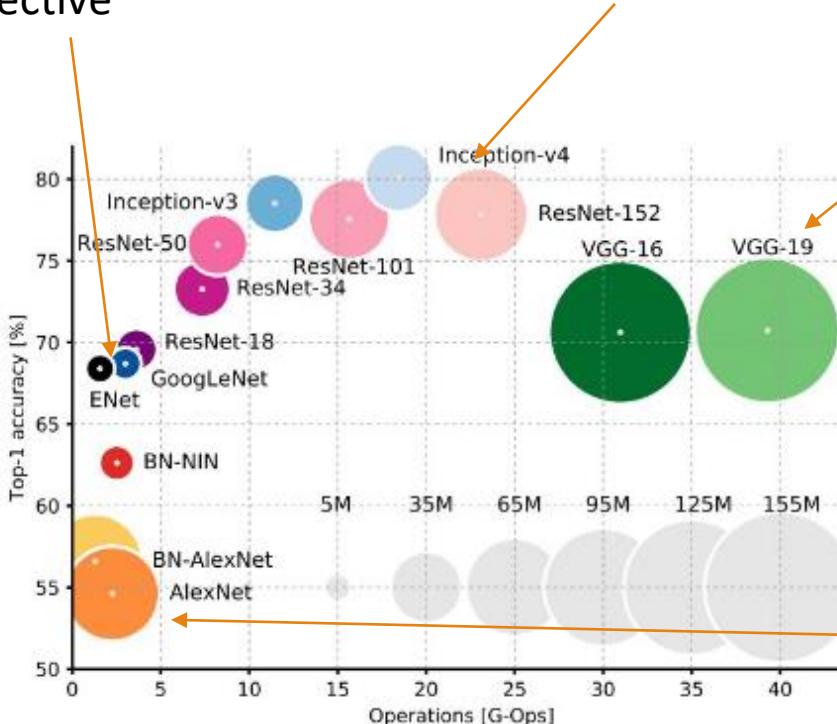
# Convolutional Network in Competitions

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners





GoogleNet:  
effective



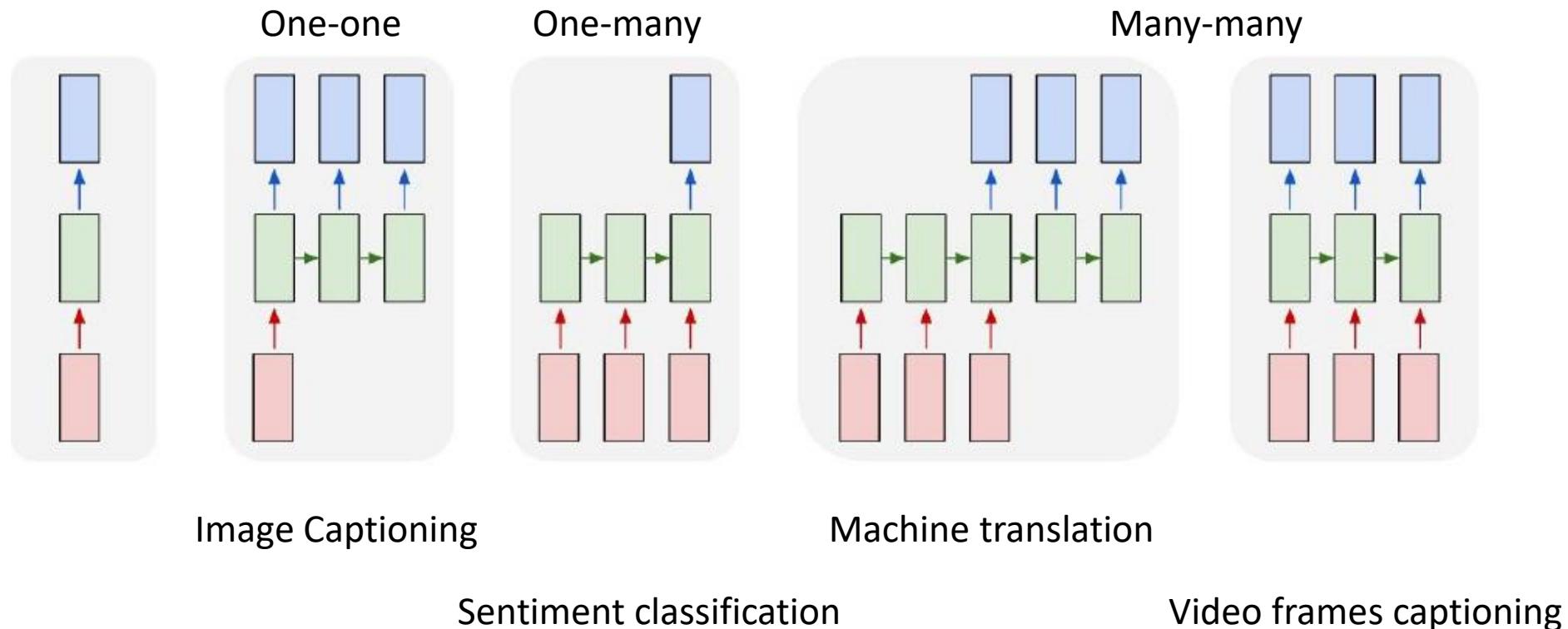
ResNet: Moderate  
efficiency depending  
on model, highest accuracy

VGG: Highest  
memory, most  
operations

AlexNet: Smaller  
compute, still  
memory heavy,  
lower accuracy

An Analysis of Deep Neural Network Models for Practical Applications, 2017.

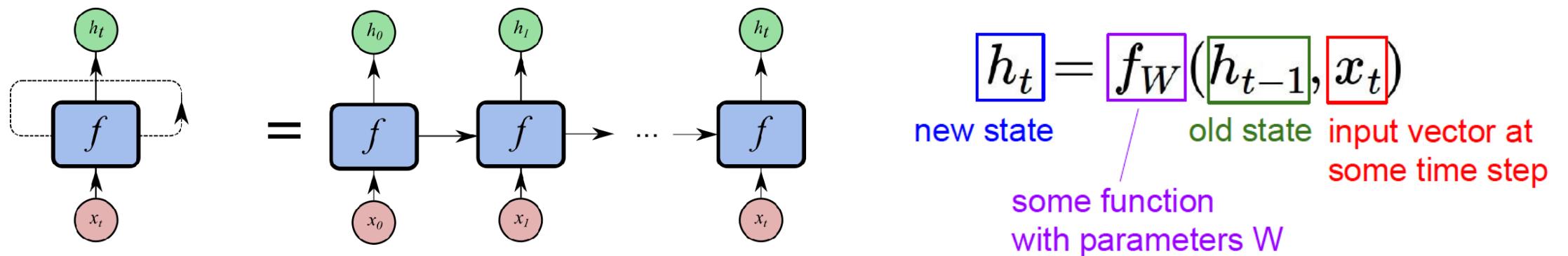
# Sequence Learning



# Recurrent Neural Network

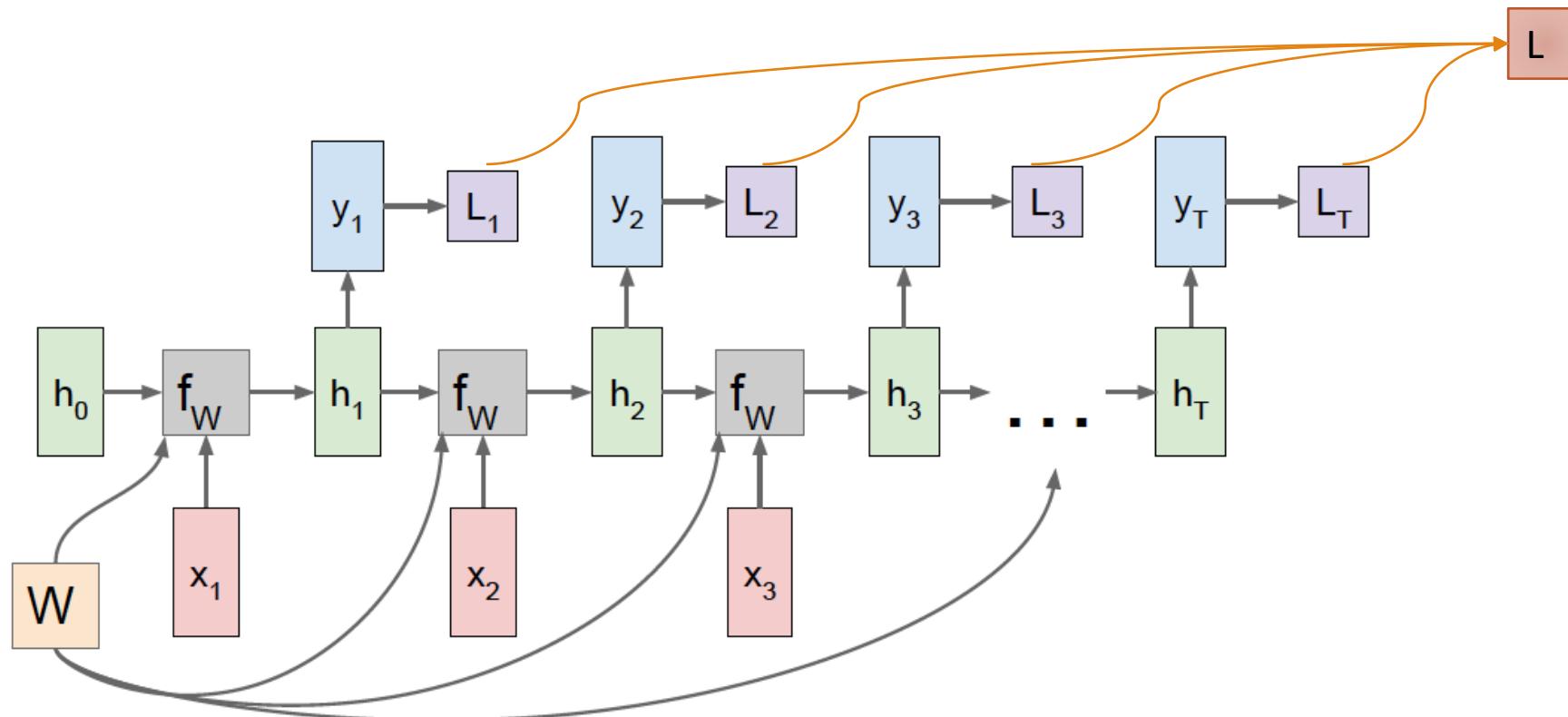
---

- Usually want to predict a vector at some time steps
- Process a sequence of vector  $x$  by applying a recurrence formula at every step



# RNN: Computational Graph

---

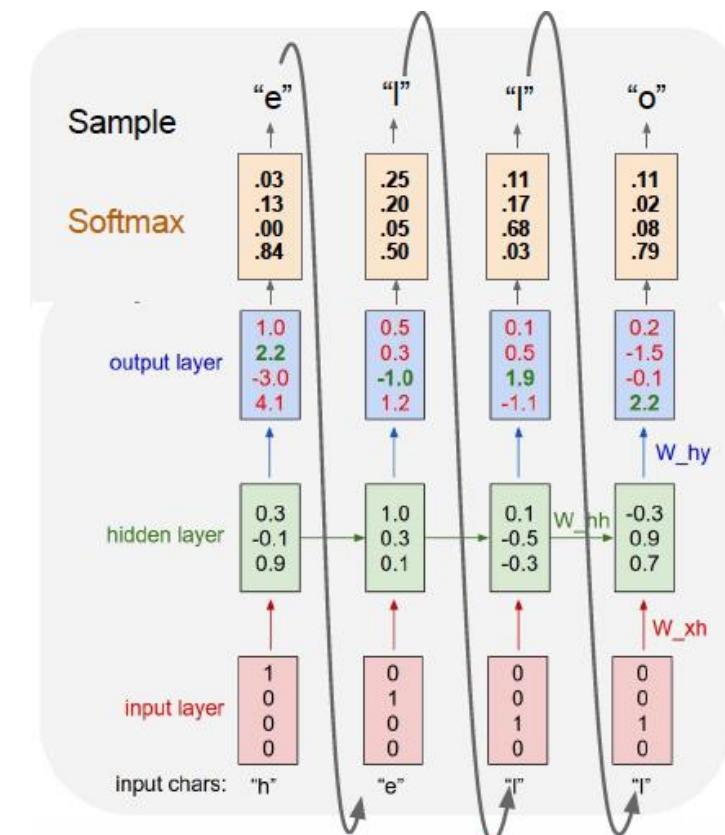


# RNN: Character-level

Character level Language Sampling

Vocabulary: [h,e,l,o]

At least-time sample characters  
one at a time feed back to the  
model

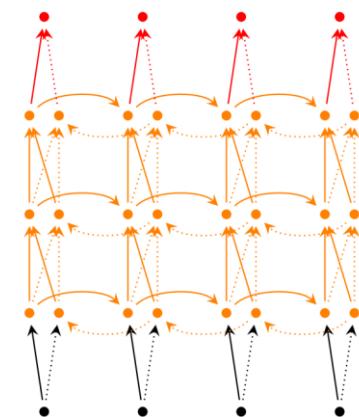
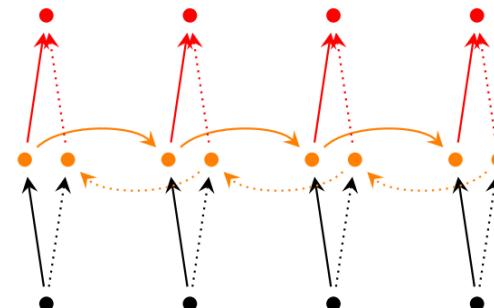


# More RNN Structures

---

Bidirectional RNNs:

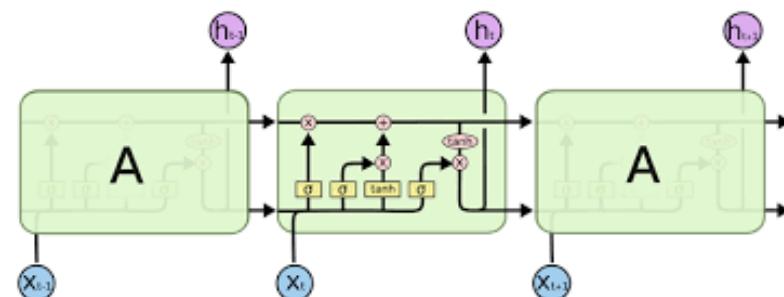
- Output at time  $t$  depend on:
  - the previous elements in the sequence.
  - future elements.



Deep (Bidirectional) RNNs:

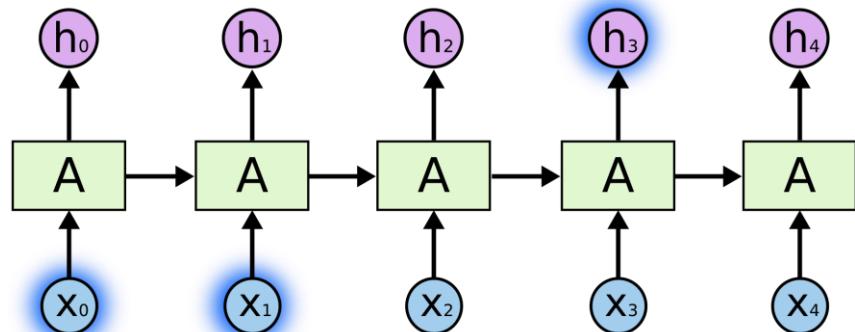
- Similar to Bidirectional RNNs.
- Have multiple layers per time step.

Long short-term memory (next slides)

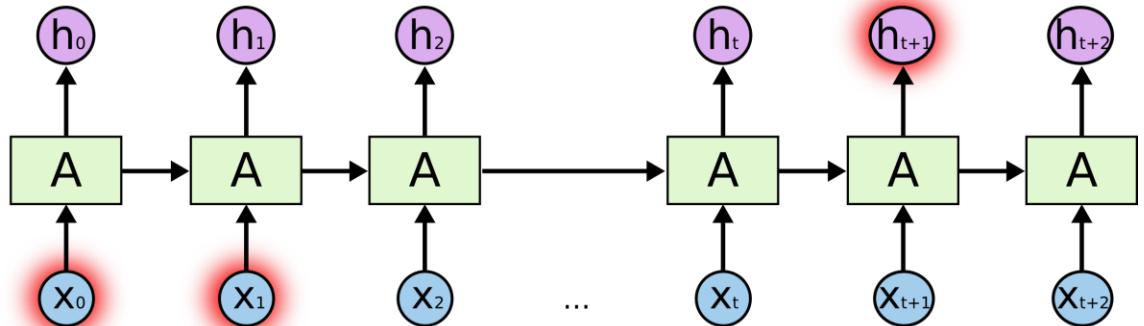


# Long short-term Memory

## Long-Term Dependencies Problem

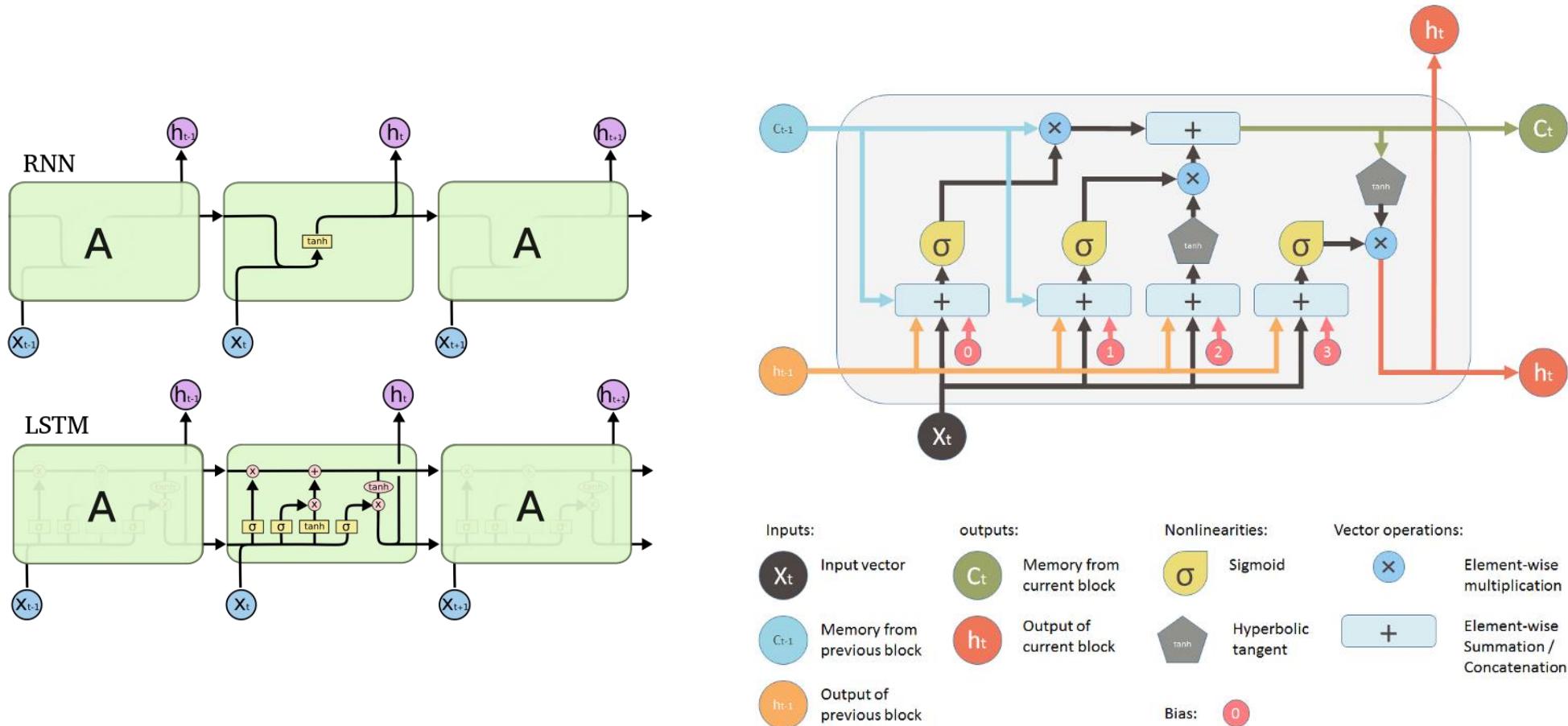


RNN looks at recent information



RNN unable to learn to connect long term information

# Long short-term Memory

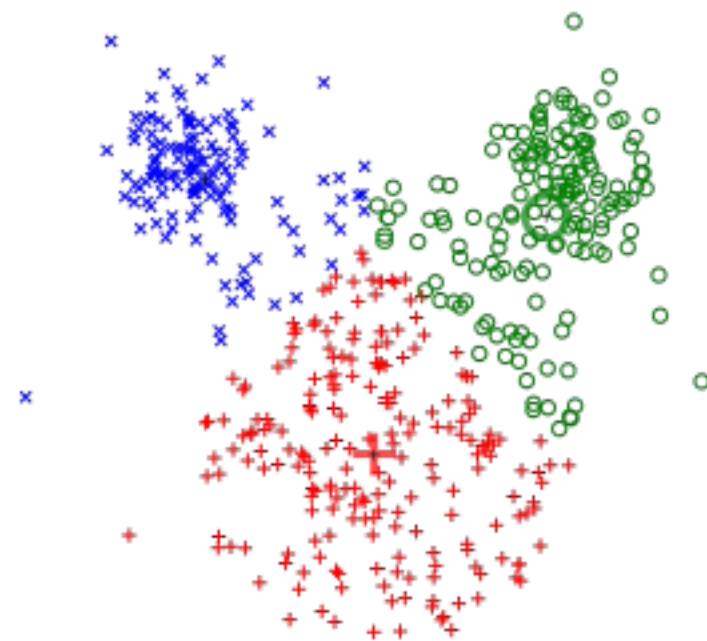


# Generative Model

---

## Unsupervised learning

- Data:  $x$
- Just data, no labels!
- Goal: Learn some underlying hidden structure of the data
- Examples: Clustering, dimensionality reduction, feature learning, density estimation, etc.



K-means clustering

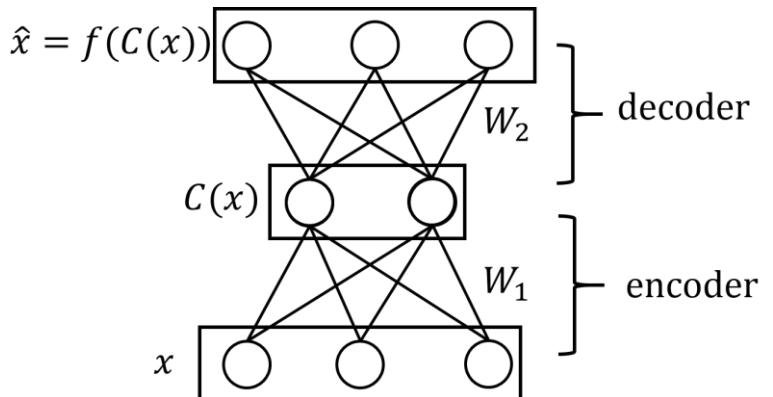
# Auto Encoder

---

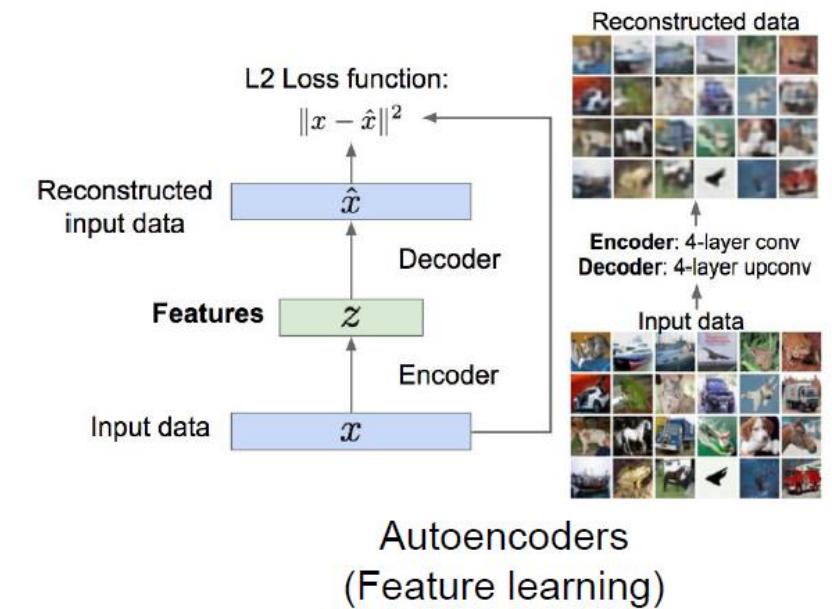
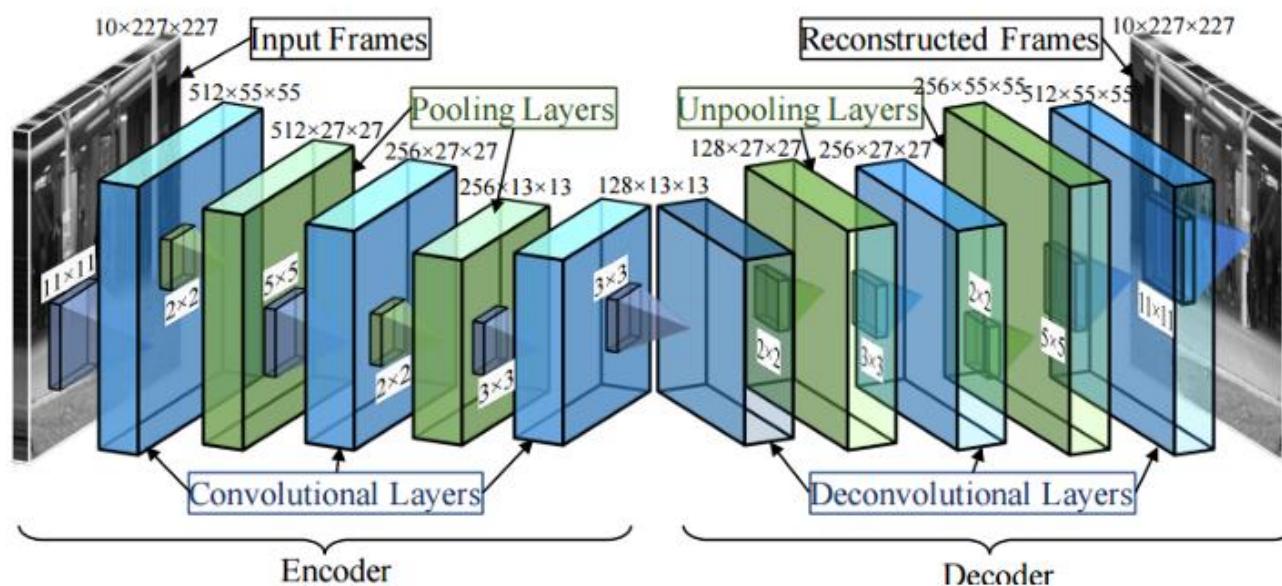
Unsupervised learning: no labels for training data

## Unsupervised Encode Input Vector

- Reconstruct input: Numbers neural in input layer and output layer are equal
- The middle layer may have:
  - less neurons: undercomplete autoencoder
  - more neurons: overcomplete autoencoder
  - $W_1$  and  $W_2$  are usually (but not always) tied, i.e.  $W_2 = W_1^\top$



# Deep Convolutional Autoencoder



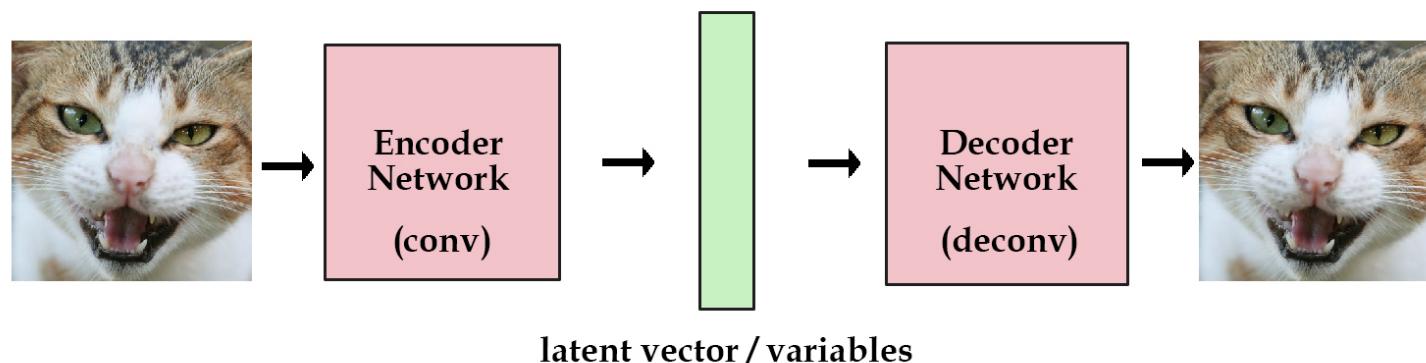
# Variational Autoencoders

---

Encoder to represent latent feature of image

Decoder:

- Use as generator to generate images from a vector
- Problem: Just remember learned images, not gaussian distribution

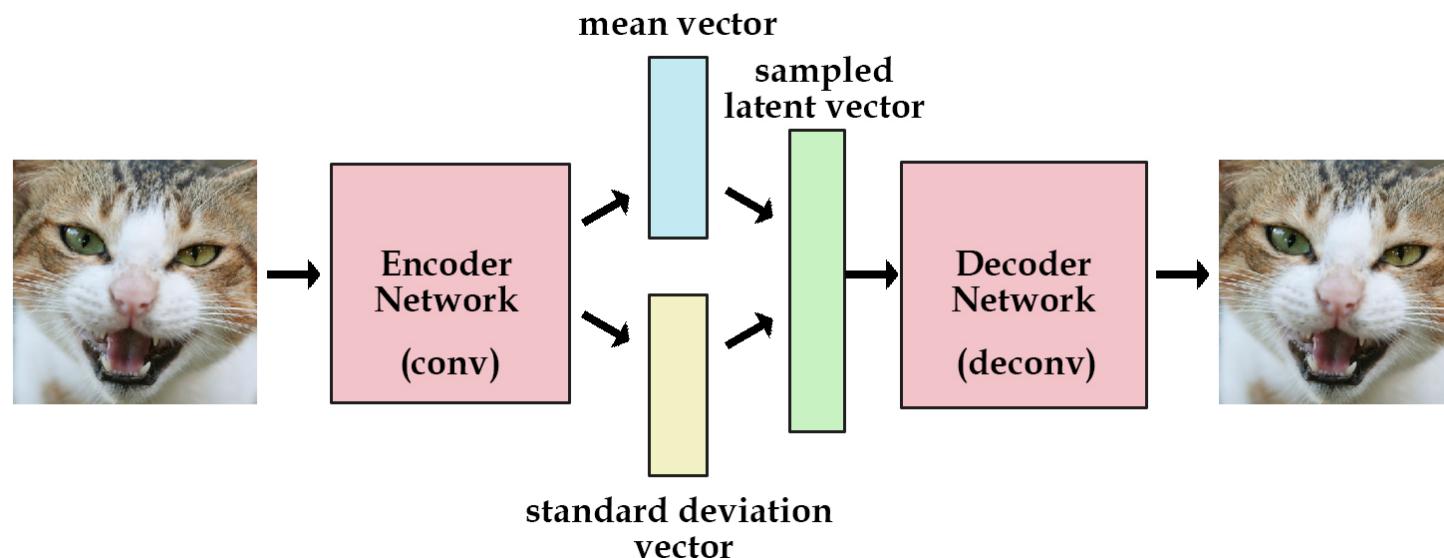


# Variational Autoencoder

---

Solution: separate two losses

- Generative loss: mean squared error of accurately reconstructed images
- Latent loss: measure how closely latent variable match a unit gaussian

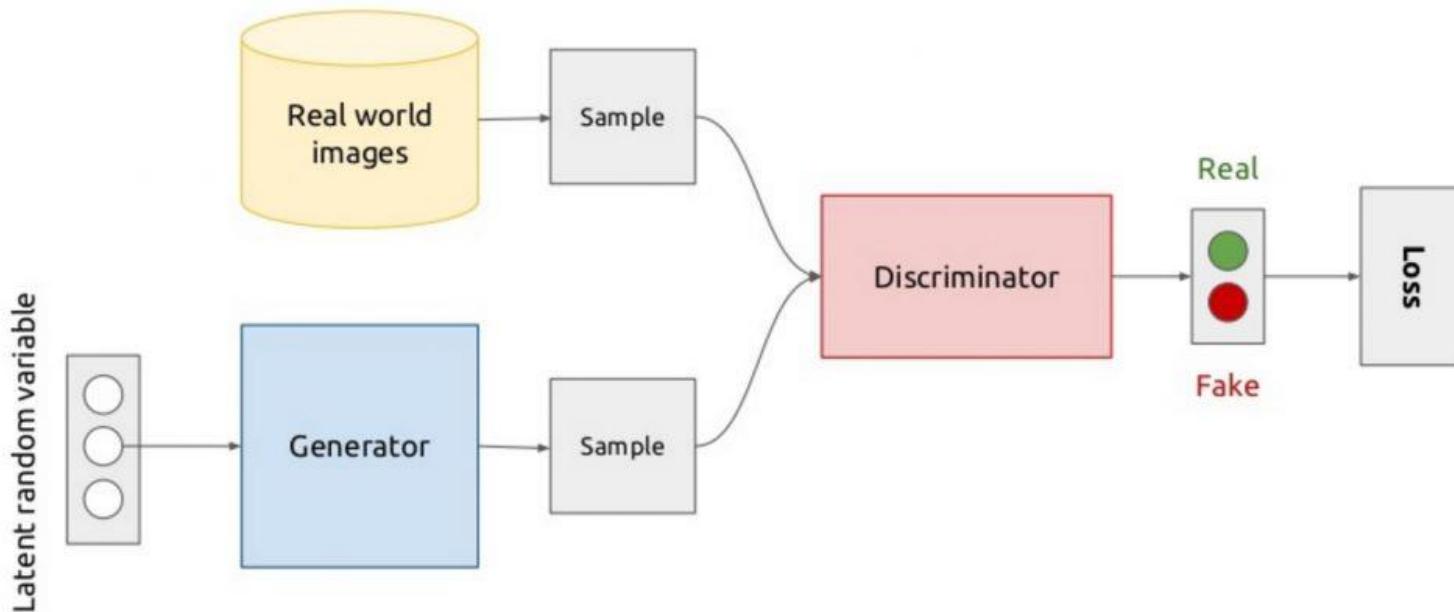


# Generative Adversarial Network

---

Generator to create fake image

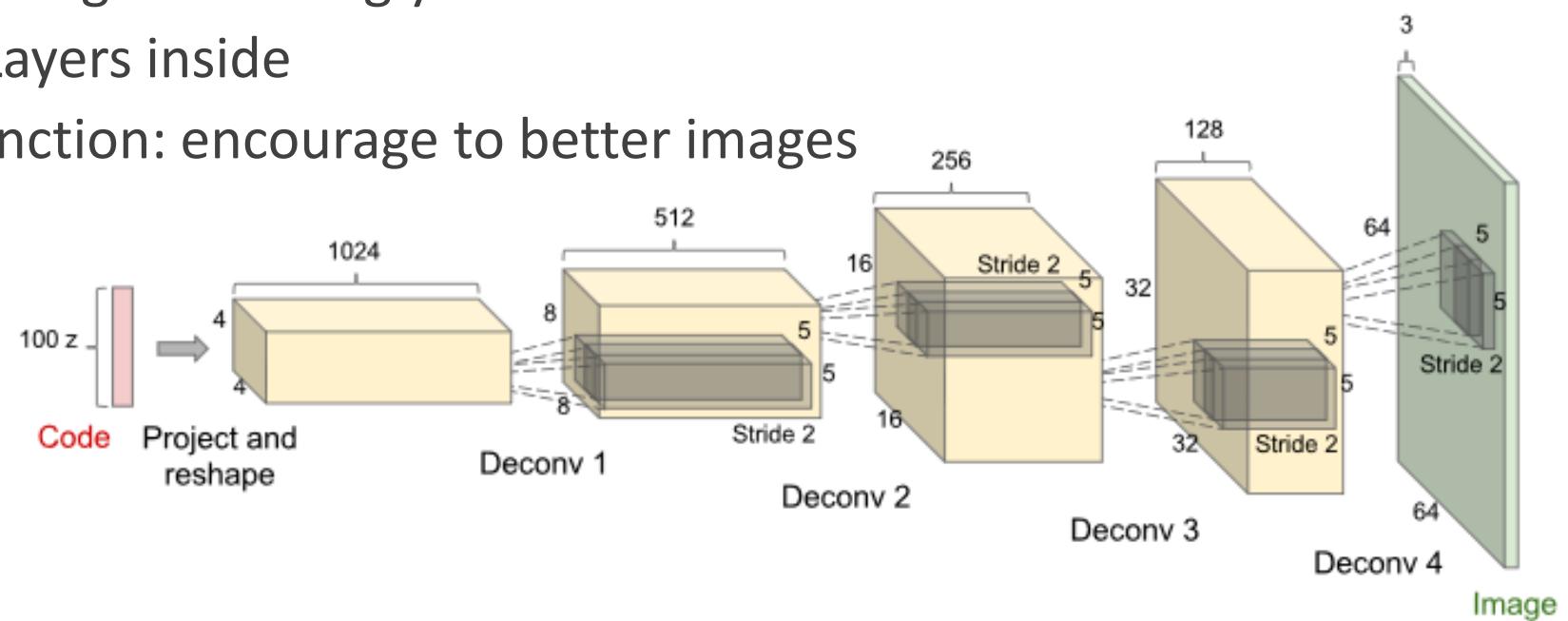
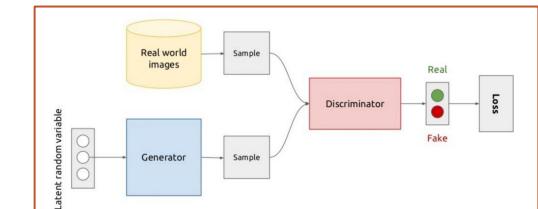
Discriminator to classify real or fake image



# Deconvolutional Generator

## Generator

- Takes code: parameter of noise
- Output synthetic image accordingly
- Deconvolutional Layers inside
- Generator loss Function: encourage to better images



# Reinforcement learning

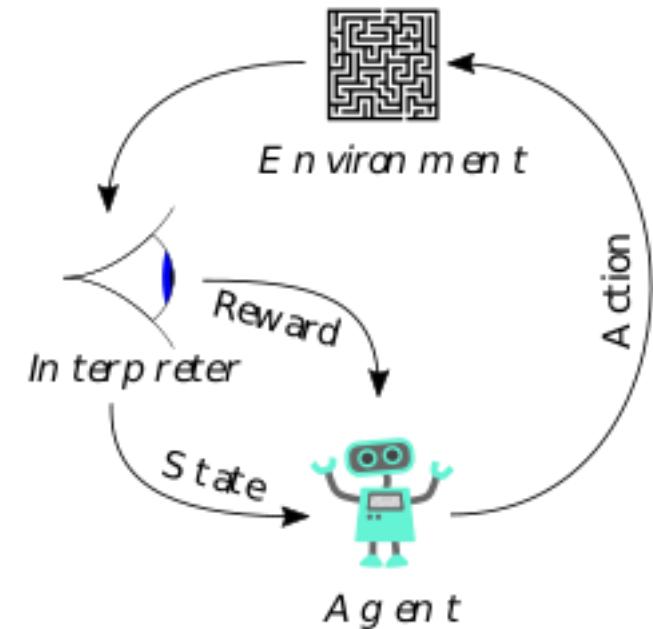
---

Making good decision to do new task: fundamental challenge in Artificial Intelligence and Machine learning

Learn to make good sequence of decisions

Intelligent agents learning and acting

- Learning by trial-and-error, in real time
- Improve with experience
- Inspired by psychology:
  - Agents + environment
  - Agents select action to maximize *cumulative* rewards



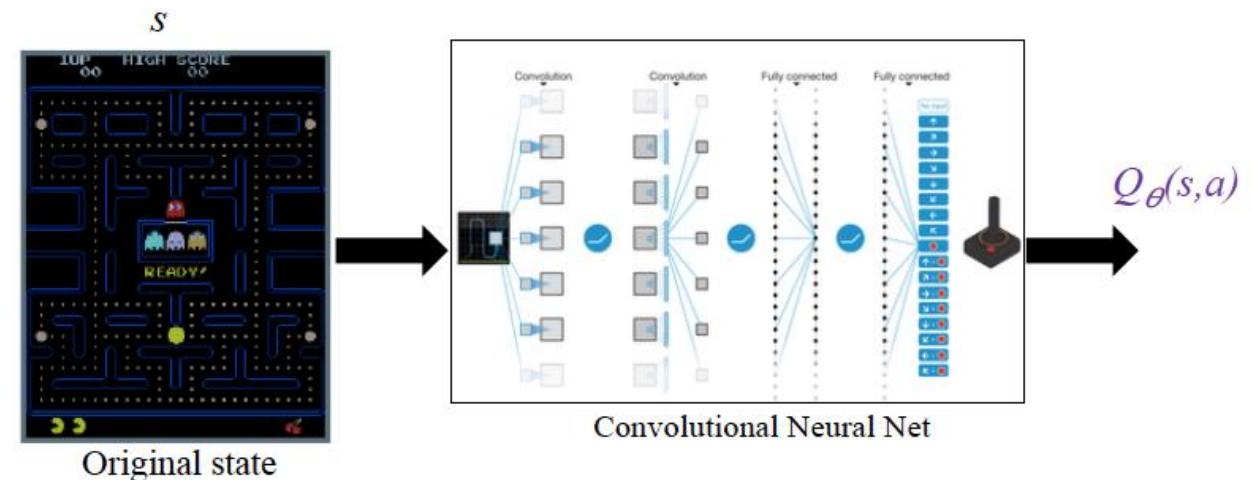
# Deep Reinforcement Learning: DeepNN + RL

Real applications: very large environment, many states, noise

- Use deep learning to learn policies, values or models to use in a reinforcement learning domain
- Function estimation:

$$\hat{v}(s, \mathbf{w}) \approx v_\pi(s)$$

or  $\hat{q}(s, a, \mathbf{w}) \approx q_\pi(s, a)$



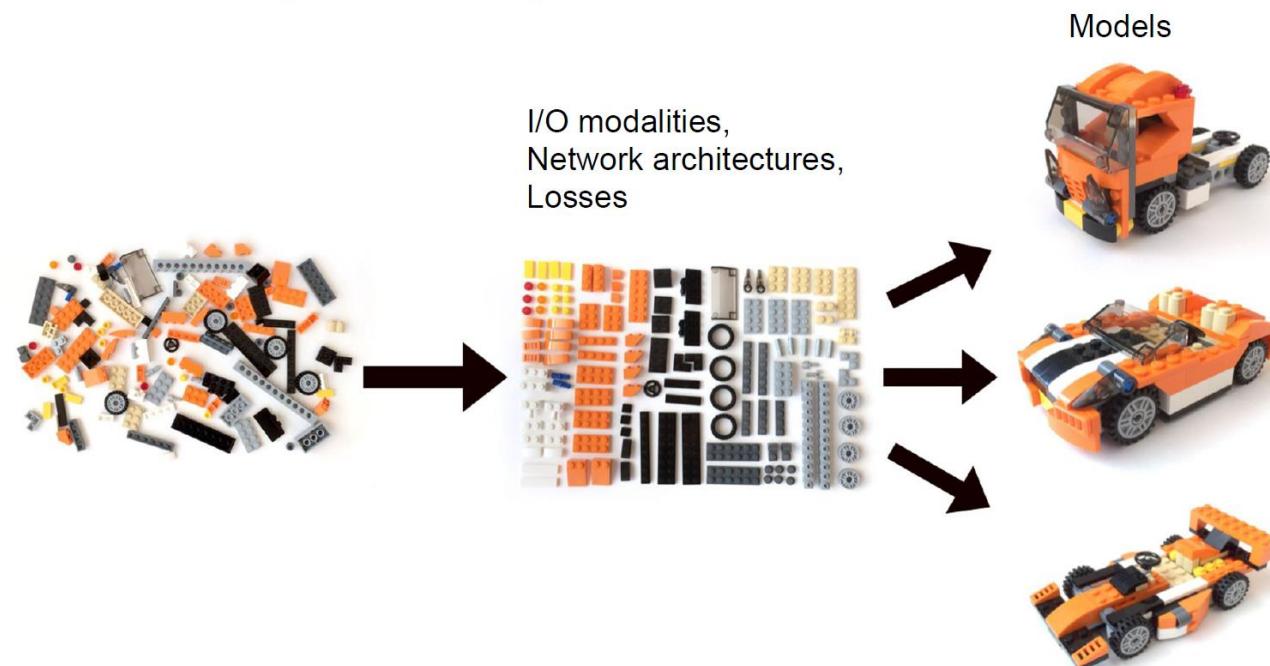
Deep Q-Network trained with stochastic gradient descent.

[DeepMind: Mnih et al., 2015].

# Deep Learning Recap

---

## Deep Learning Building Blocks



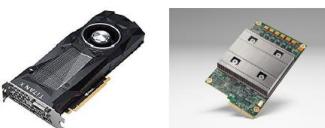
*Image: Nagel, Wolfram. Multiscreen UX Design: Developing for a Multitude of Devices. Morgan Kaufmann, 2015.*

# Deep Learning Recap

## Deep Learning: Zooming Out



Platforms



Frameworks



Datasets



Caltech 101 MGENET



WMT Workshop 2014

6

# Deep Learning Recap

## Deep Learning: Zooming In

### Non-Linearities

Relu  
Sigmoid  
Tanh  
GRU  
LSTM  
Linear  
...

### Optimizer

SGD  
Momentum  
RMSProp  
Adagrad  
Adam  
Second Order (KFac)  
...

### Connectivity Pattern

Fully connected  
Convolutional  
Dilated  
Recurrent  
Recursive  
Skip / Residual  
Random  
...

### Loss

Cross Entropy  
Adversarial  
Variational  
Max. Likelihood  
Sparse  
L2 Reg  
REINFORCE  
...



### Hyper Parameters

Learning Rate  
Decay  
Layer Size  
Batch Size  
Dropout Rate  
Weight init  
Data augmentation  
Gradient Clipping  
Beta  
Momentum



7

# Q&A

---

THANK YOU!