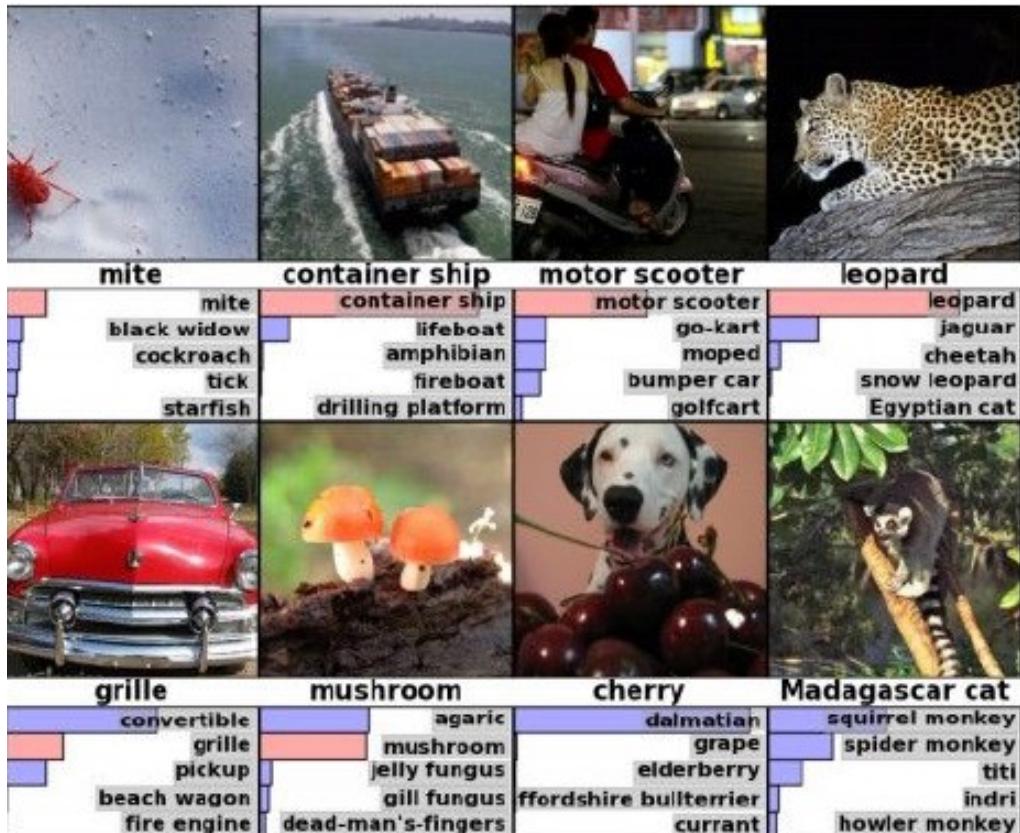


Applications of Deep Learning

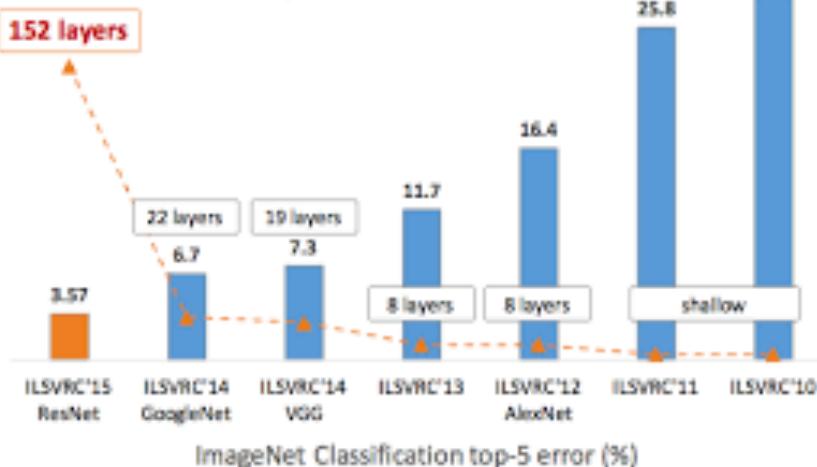
Slides from

- 1) Stanford CS231n Lecture Notes
- 2) U Michigan EECS 498 Lecture Notes

ImageNet Challenge



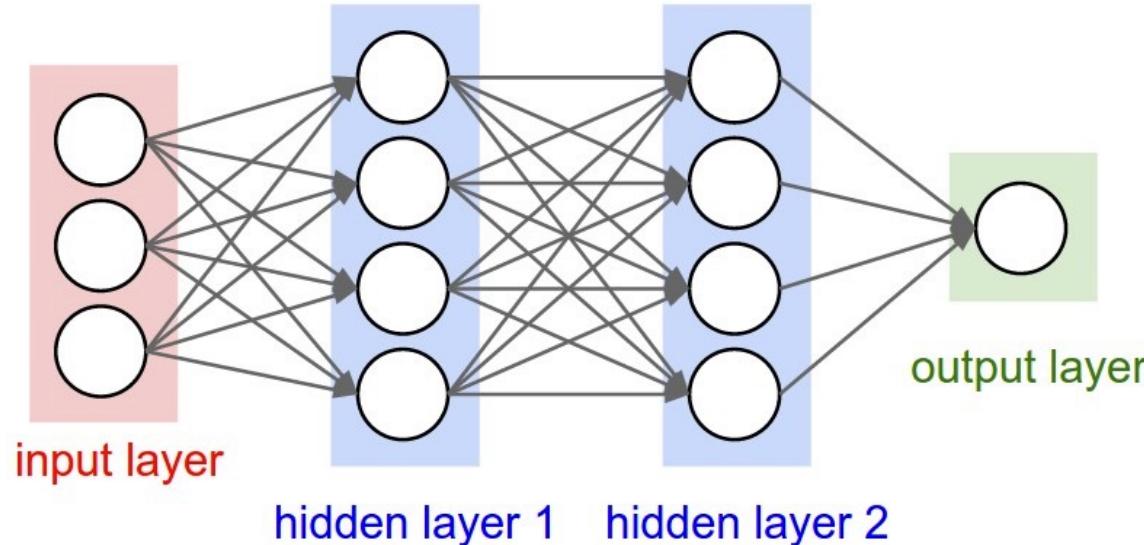
Revolution of Depth



Convolution Neural Network

Natural Language Processing

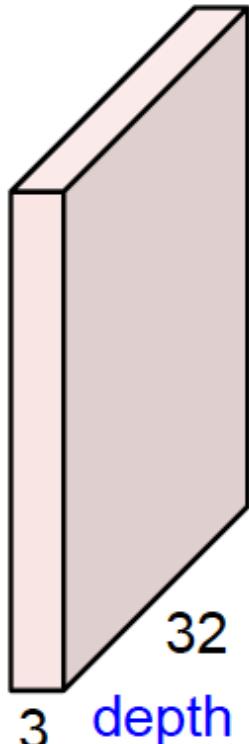
Fully Connected Layers



- ImageNet: 224x224x3 pixels
- One hidden Layer with 5,000 nodes
 - Input -> Hidden Layer: 750 Million weights $224 \times 224 \times 3 \times 5000 = 750M$
 - Hidden Layer -> Output: 5 Million weights

Signal Processing

Filter Parameters



Input Image

$5 \times 5 \times 3$



Filter

Learn parameters for
filter

Input Image

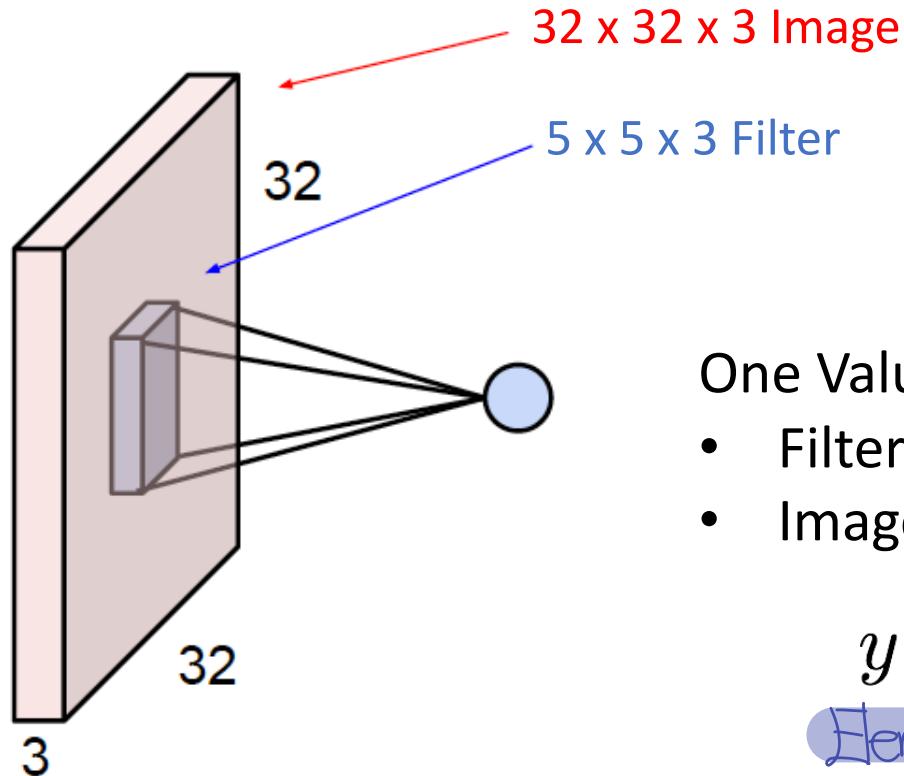
- Height
- Width
- Depth

Filter

- Height
- Width
- Depth

=

Convolution Operation



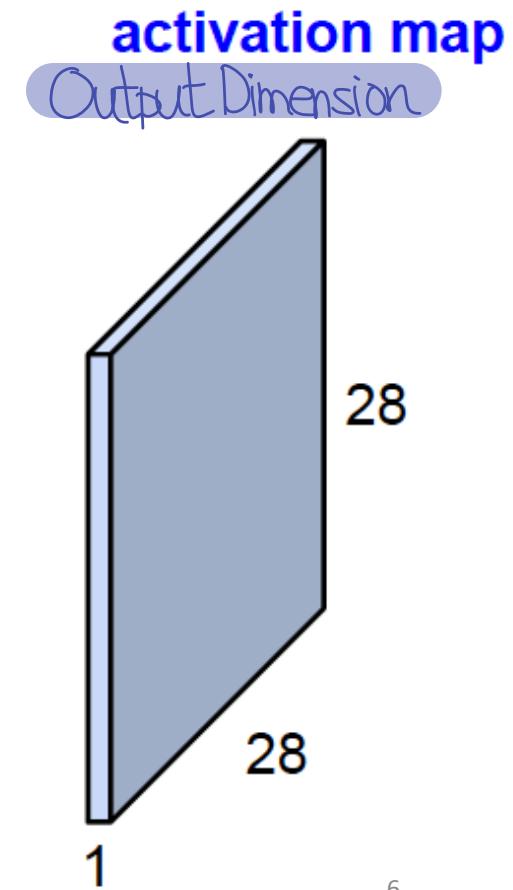
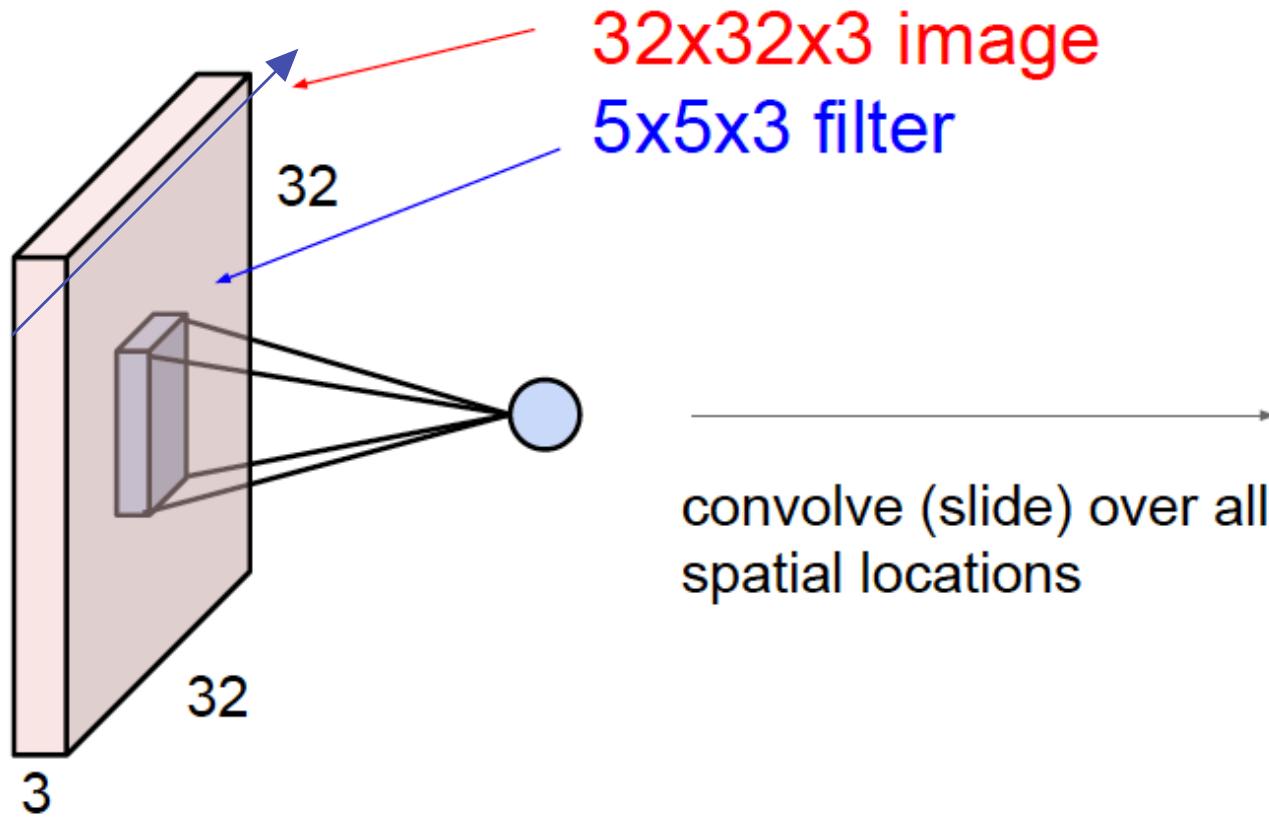
One Value obtained from this operation:

- Filter: \mathbf{w}
- Image Patch: \mathbf{x}

$$y = \mathbf{w}^T \mathbf{x} + b$$

Elementwise Multiplication

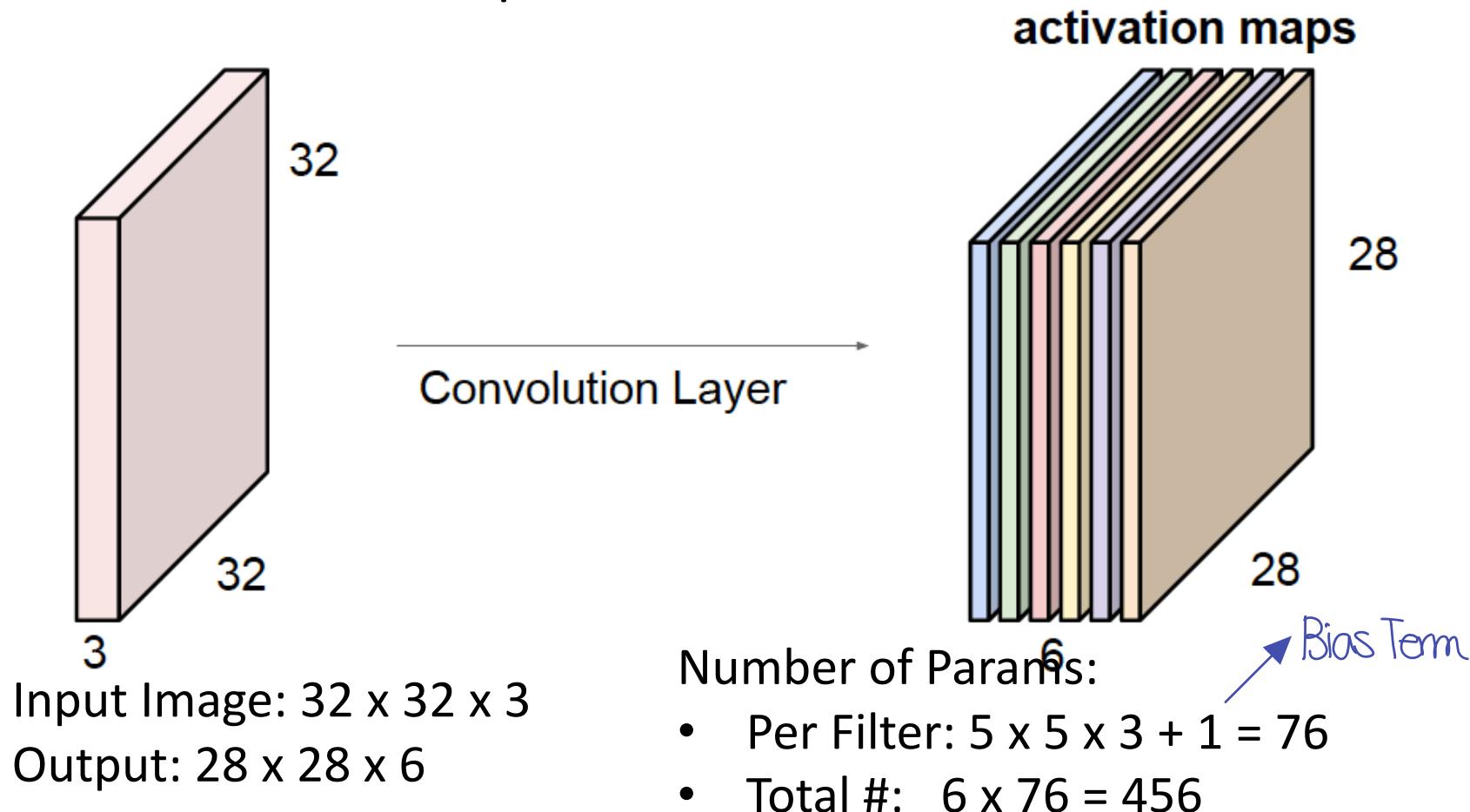
Convolution Operation - II



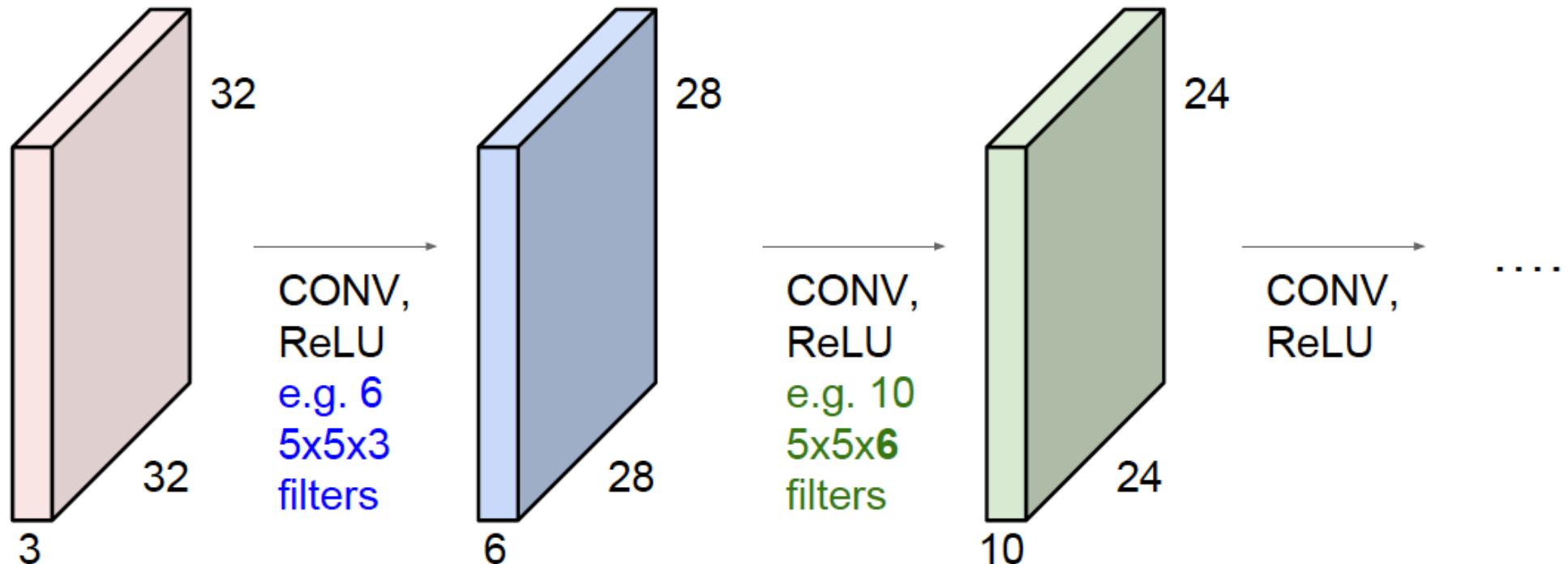
Convolution Operation - III



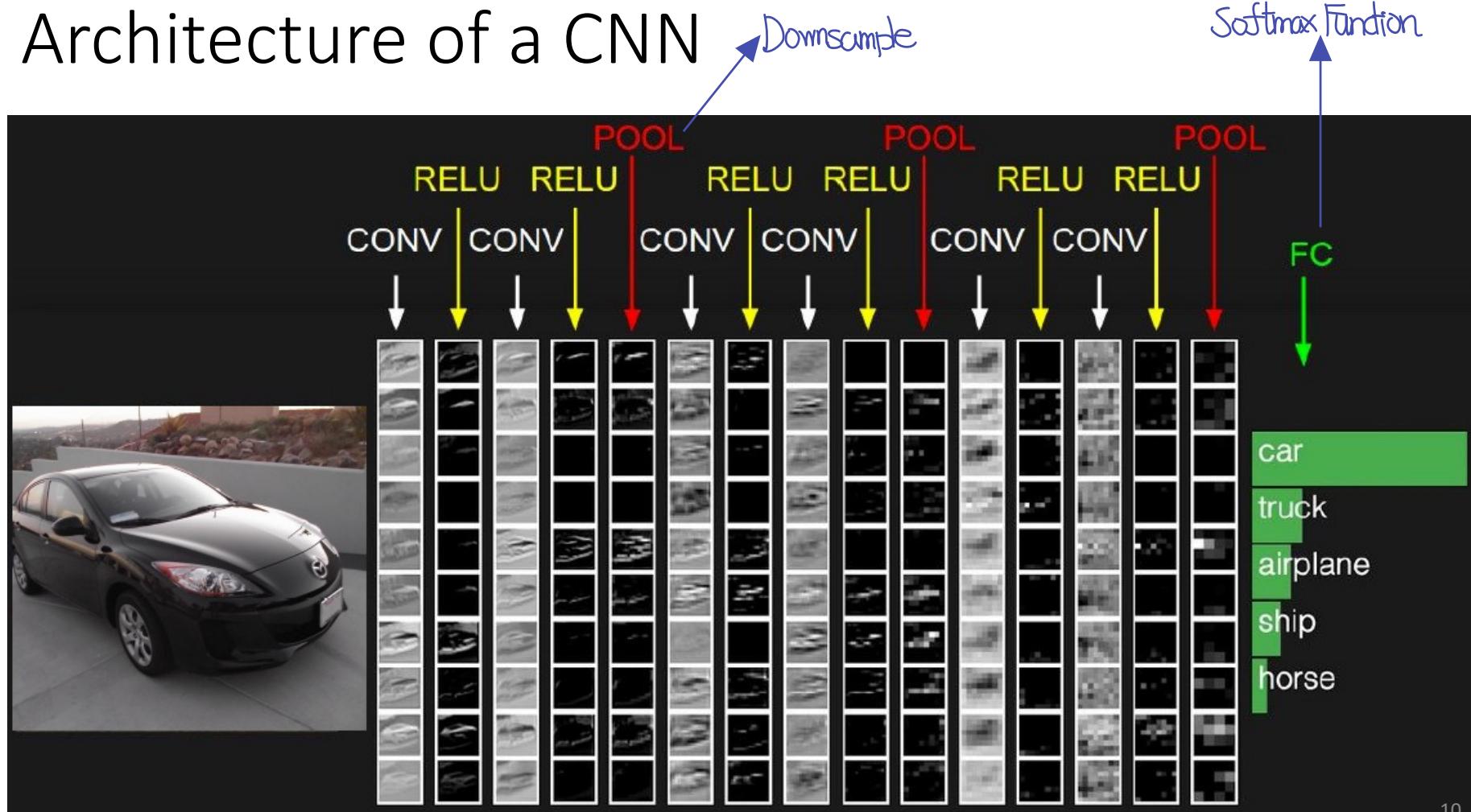
Convolution Operation - IV



Convolutional Neural Network (CNN)

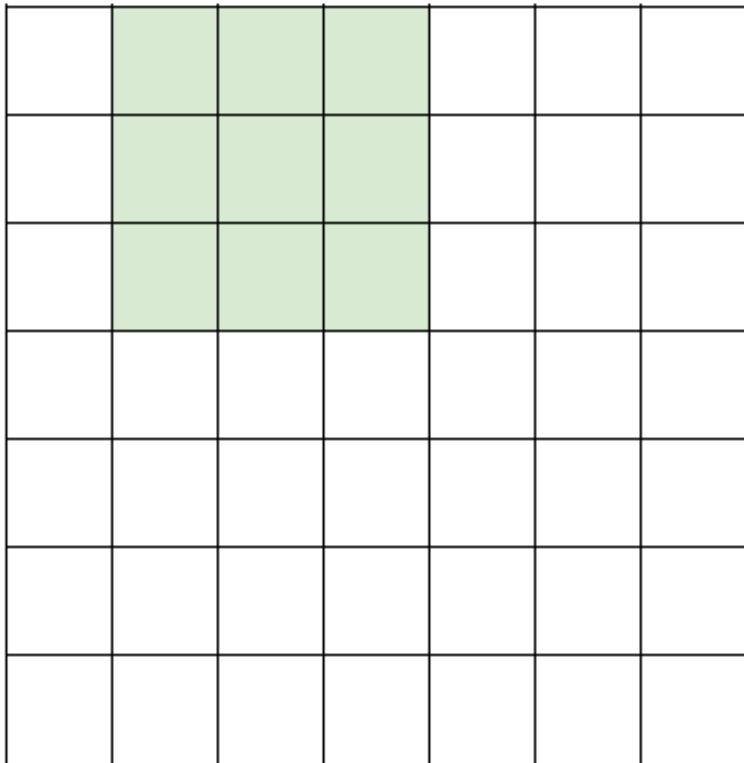


Architecture of a CNN



Closer Look at Convolution Operation

7

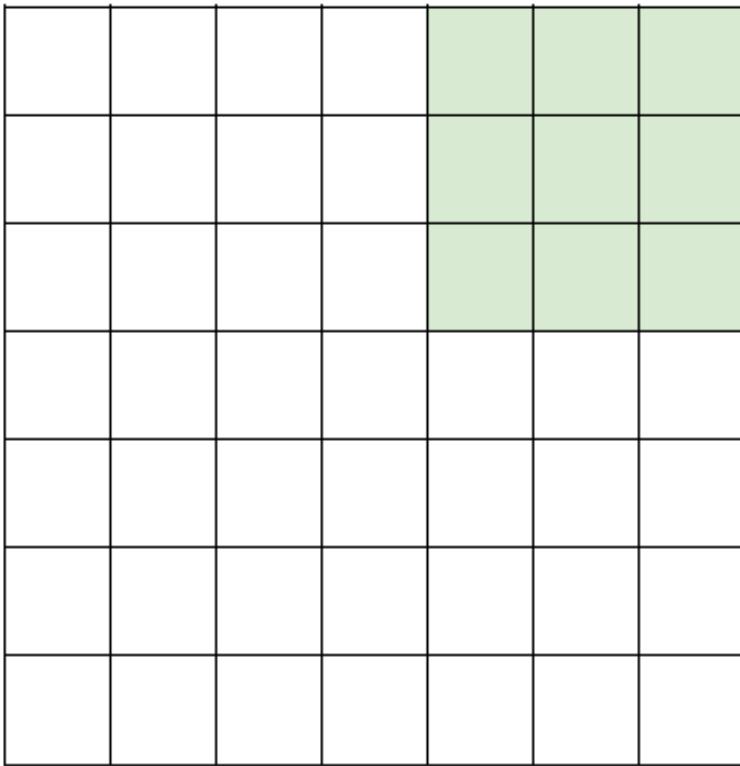


7

- Input: 7×7
- Filter: 3×3
- Stride: 1
- Output: 5×5

Closer Look at Convolution Operation

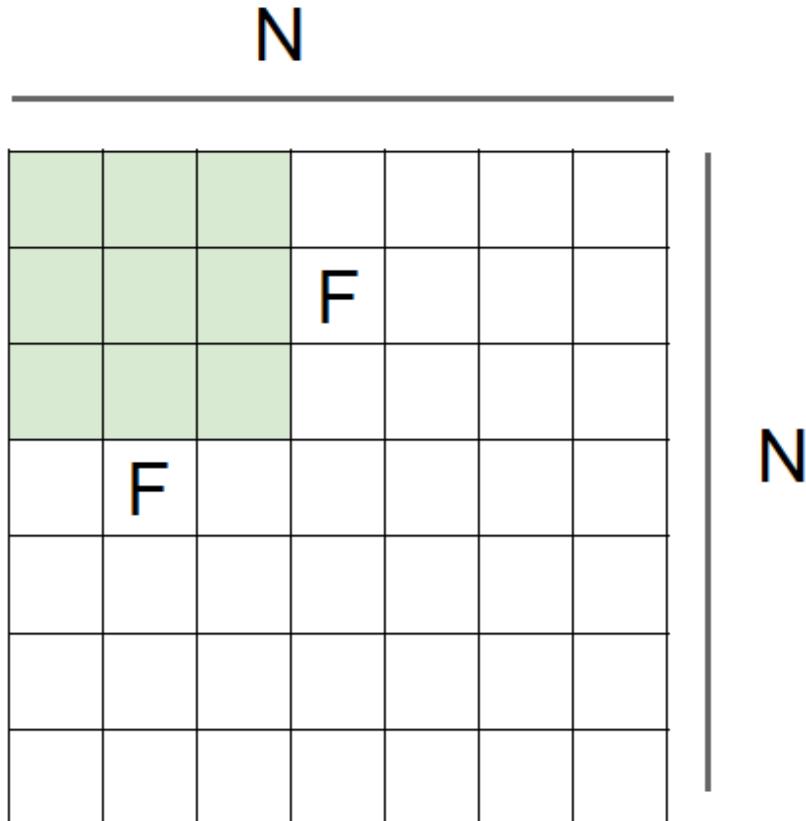
7



7

- Input: 7×7
- Filter: 3×3
- Stride: 2 (*Jump By 2*)
- Output: 3×3

Output Dimensions



Output size:

$$(N - F) / \text{stride} + 1$$

N (Input Size), F (Filter Size)

e.g. $N = 7, F = 3$:

stride 1 => $(7 - 3)/1 + 1 = 5$

stride 2 => $(7 - 3)/2 + 1 = 3$

~~stride 3 => $(7 - 3)/3 + 1 = 2.33$~~

Zero Padding Input Layers

0	0	0	0	0	0		
0							
0							
0							
0							

- Input 7×7
- Filter: 3×3
- Stride: 1
- Zero Pad: 1 pixel border
- Output: 7×7

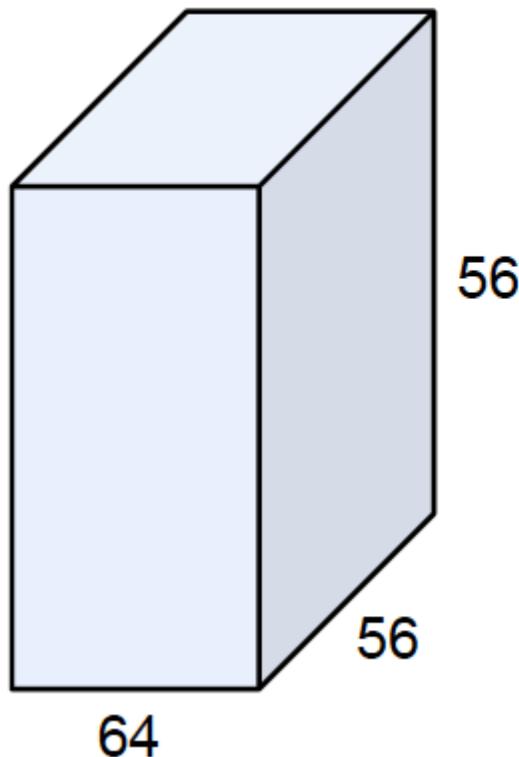
Preserving Image Size

- Filter Size $F \times F$, Stride = 1
- Zero Pad with $(F-1)/2$ pixel borders
 - $3 \times 3 \Rightarrow 1$ pixel border
 - $5 \times 5 \Rightarrow 2$ pixel border
 - $7 \times 7 \Rightarrow 3$ pixel border

Example

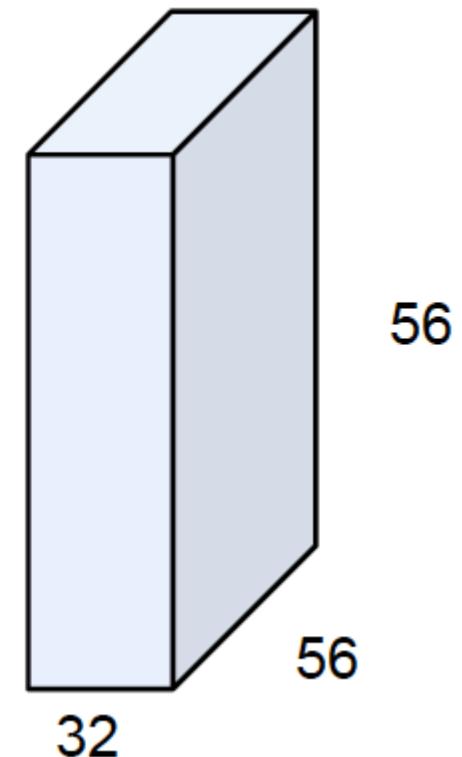
- Input: $32 \times 32 \times 3$
- 10 Filters, each $5 \times 5 \times 3$
- Stride: 1
- Zero Pad: 2 pixels on each border
- Output Dimensions: $32 \times 32 \times 10$ → Ten Filters
- Number of Parameters: $10 * (5 \times 5 \times 3 + 1) = 760$

1×1 Filters (Changing Depth)



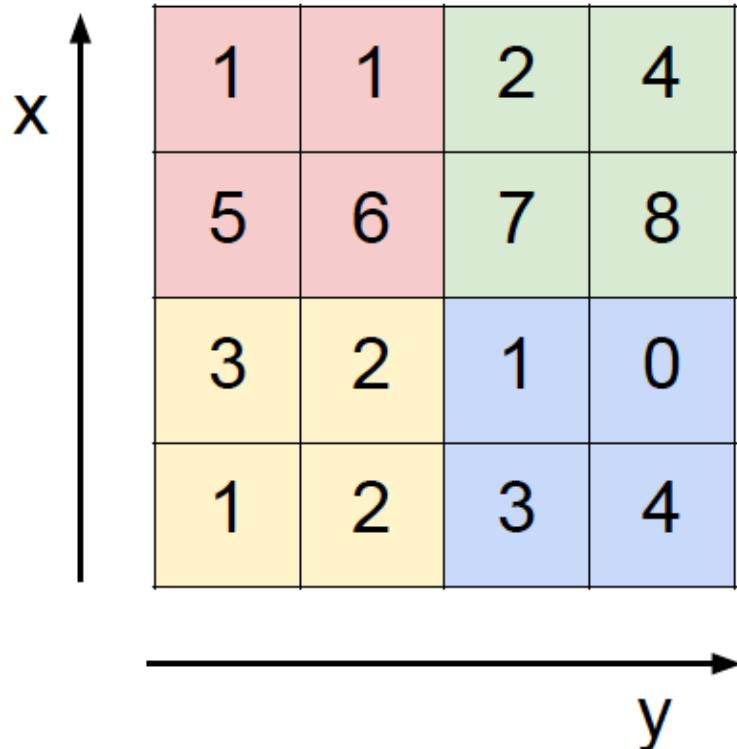
1×1 CONV
with 32 filters

(each filter has size
 $1 \times 1 \times 64$, and performs a
64-dimensional dot product)



Max Pooling

Single depth slice

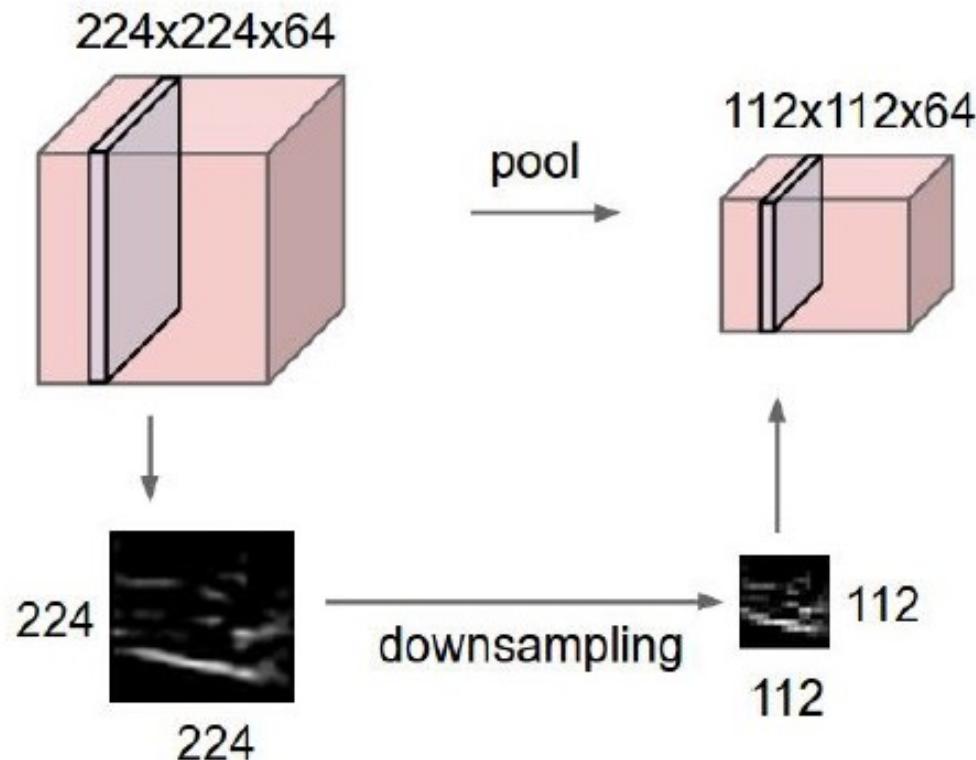


max pool with 2x2 filters
and stride 2

Maximum in each 2x2 square

6	8
3	4

Pooling Layer: Down-Sampling Operation



Alexnet

Alex Krizhevsky, Ilya Sutskever, Geoffrey E. Hinton (2012)

Architecture:

CONV1

MAX POOL1

NORM1

CONV2

MAX POOL2

NORM2

CONV3

CONV4

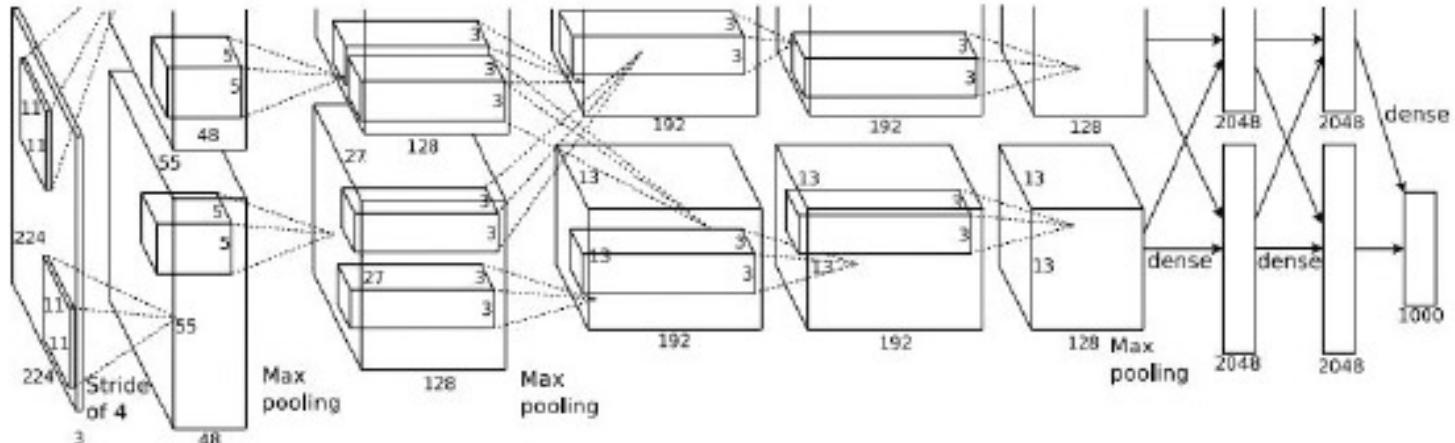
CONV5

Max POOL3

FC6

FC7

FC8



Layer 1 Parameters:

- Input: 227x227x3 images
- Layer 1: 96 11x11 filters applied at stride 4
- Output: 55 x 55 x 96
- Number of Parameters $(11*11*3)*96 = 35K$

Alexnet

Alex Krizhevsky, Ilya Sutskever, Geoffrey E. Hinton (2012)

MAXPOOL Slide Window

Full (simplified) AlexNet architecture:

[227x227x3] INPUT

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

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

[27x27x96] NORM1: Normalization layer

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

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

[13x13x256] NORM2: Normalization layer

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

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

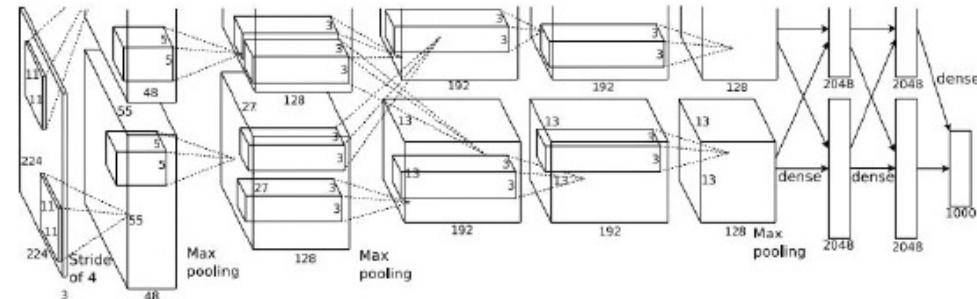
[13x13x256] CONV5: 256 3x3 filters at stride 1,

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

[4096] FC6: 4096 neurons

[4096] FC7: 4096 neurons

[1000] FC8: 1000 neurons (class scores)

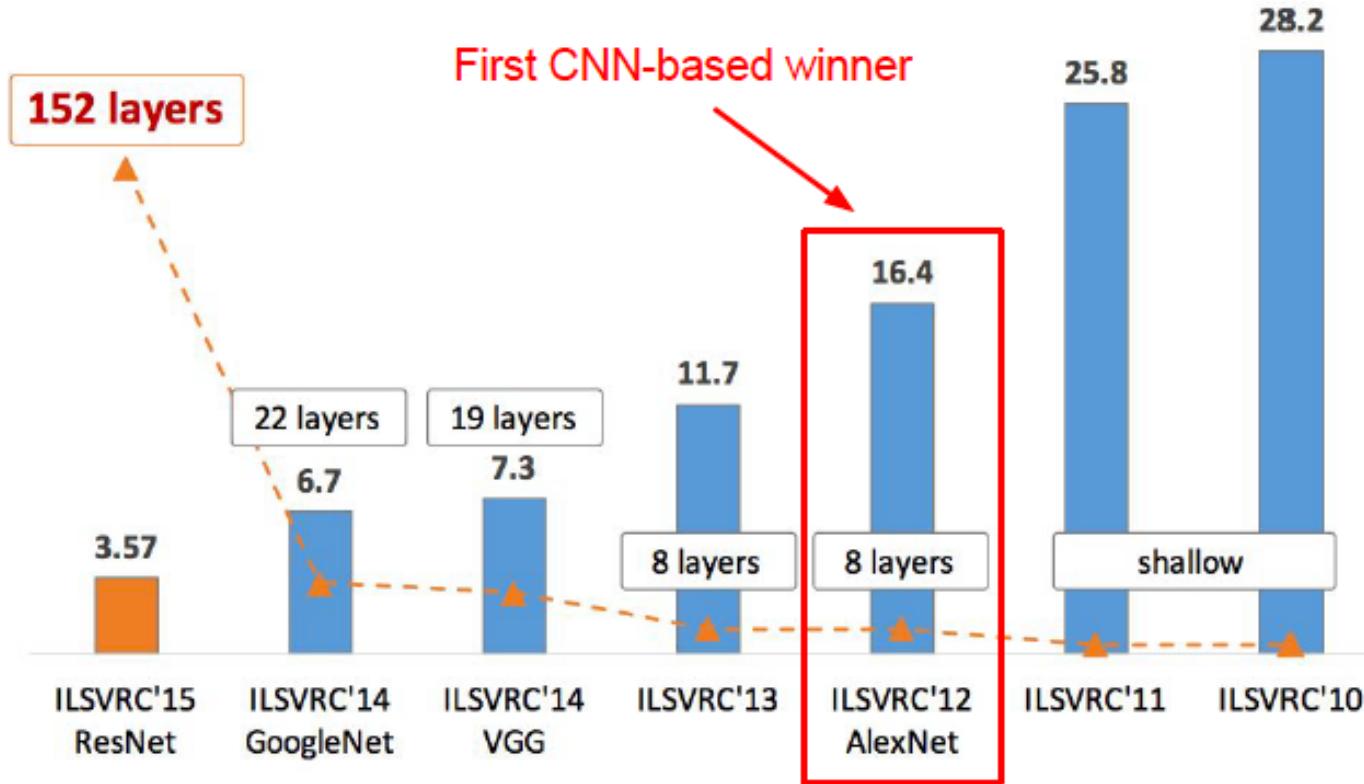


Details/Retrospectives:

- first use of ReLU
- used Norm layers (not common anymore)
- heavy data augmentation
- dropout 0.5
- batch size 128
- SGD Momentum 0.9
- Learning rate 1e-2, reduced by 10 manually when val accuracy plateaus
- L2 weight decay 5e-4
- 7 CNN ensemble: 18.2% -> 15.4%

Ensemble Learning

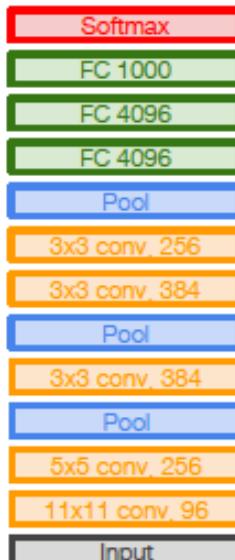
ImageNet Large Scale Visual Recognition Challenge



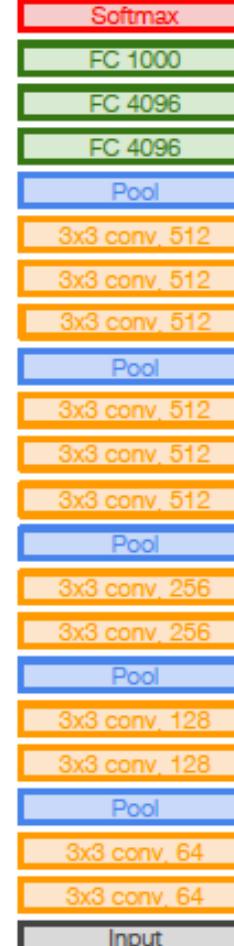
VGG Net

Karen Simonyan & Andrew Zisserman (2014)

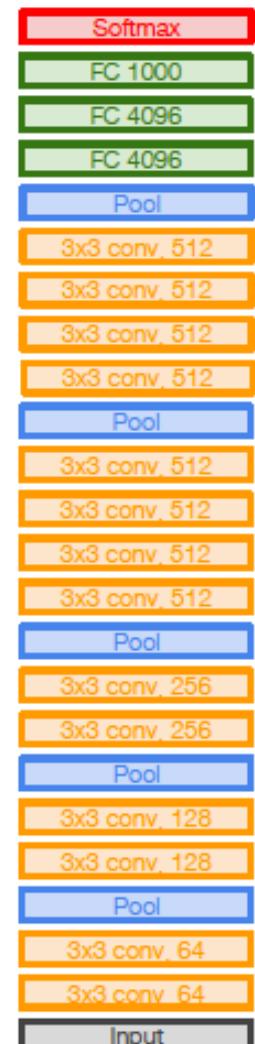
- 8 layers (AlexNet)
- 16 - 19 layers (VGGNet)
- 3x3 CONV stride 1, pad 1
- 2x2 MAX POOL stride 2
- 7.3% top 5 error in ILSVRC'14 vs
- 16.4% error in AlexNet
- Total Number of Parameters:
~138Million in VGGNet
- ~62 Million in Alexnet
Deeper the better



AlexNet

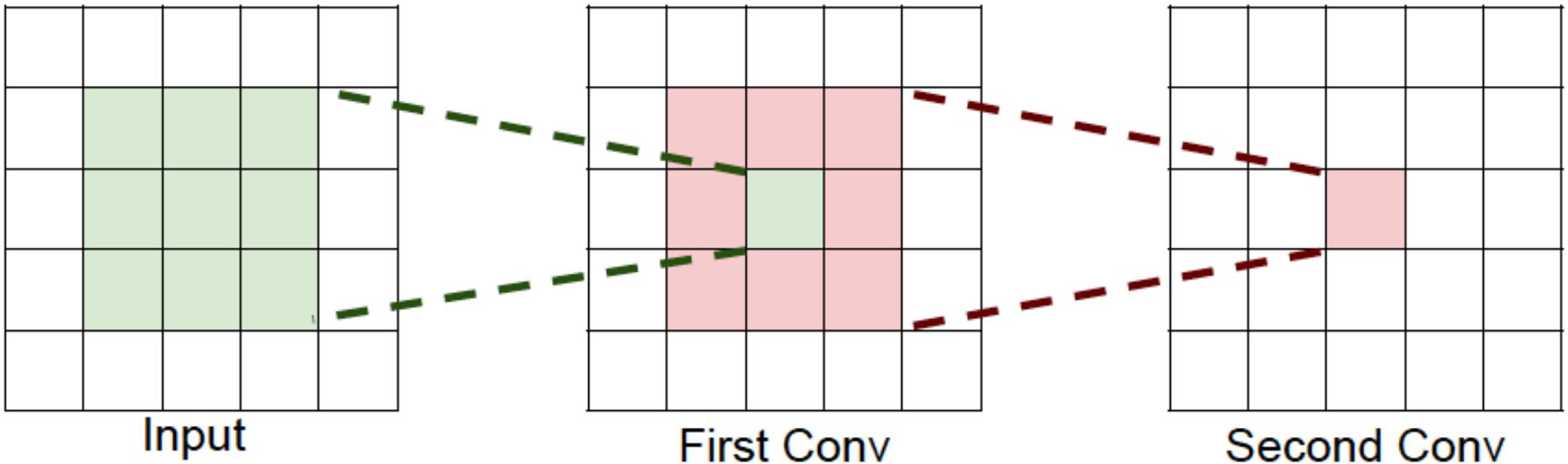


VGG16



VGG19

Why multiple 3x3 layers are preferred?



- Receptor Field of First Layer: 3×3
- Effective Receptor Field of 2nd Layer: 5×5
- Effective Receptor Field of 3rd Layer: 7×7

Why multiple 3x3 layers are preferred?

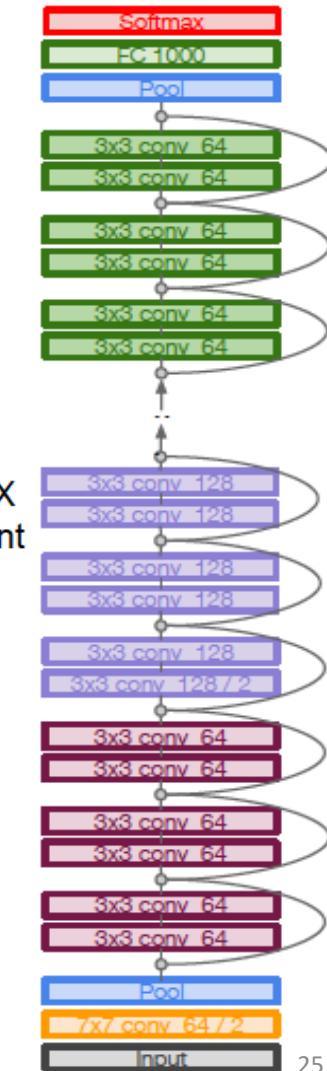
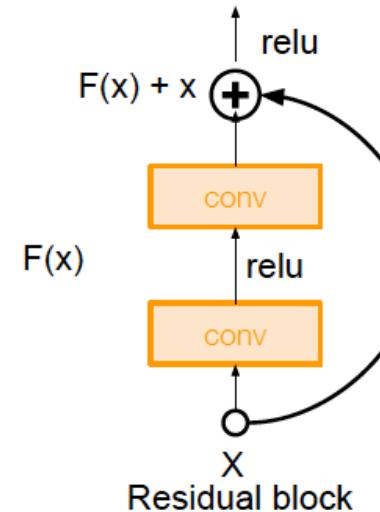
- A stack of three layers with 3×3 filters has the same receptive field as a single layer with 7×7 filters
- Number of Parameters:
 - Three 3×3 filters: 27 (smaller)
 - One 7×7 filter: 49

Same receptive field, less parameter

ResNet

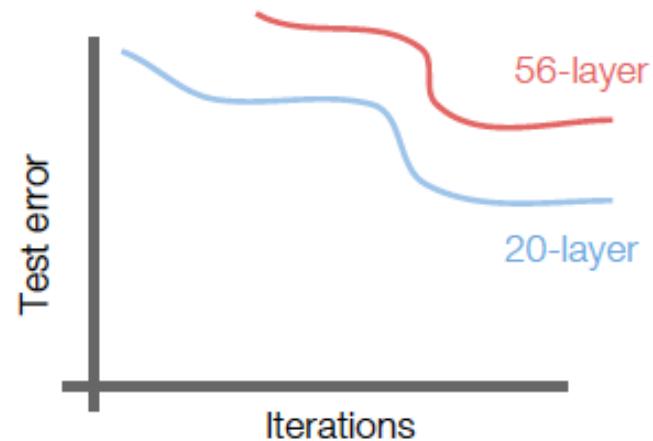
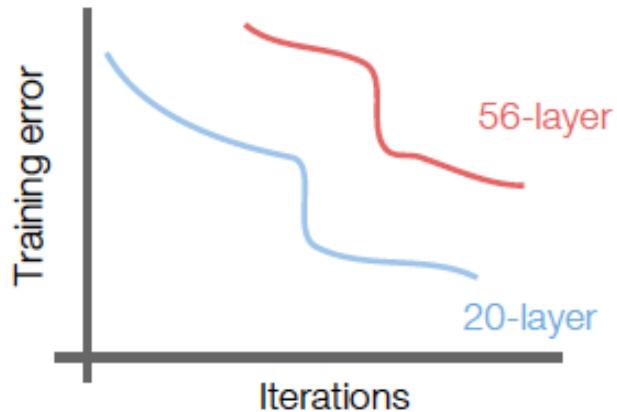
He et. al. (2015)

- 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!



Vanishing Gradient Problem

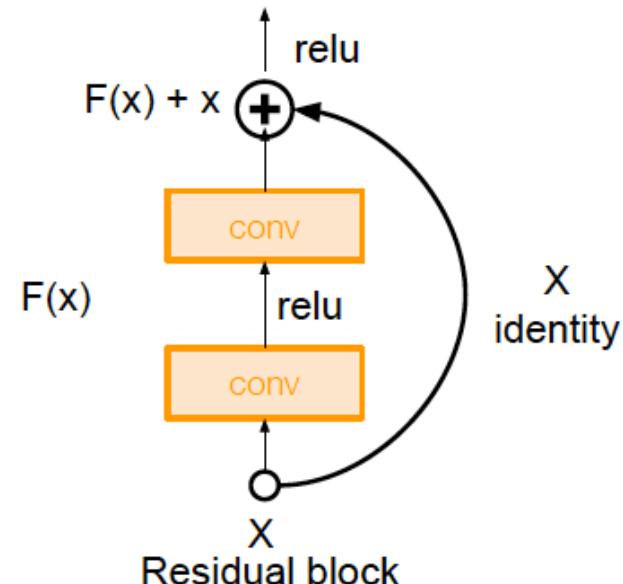
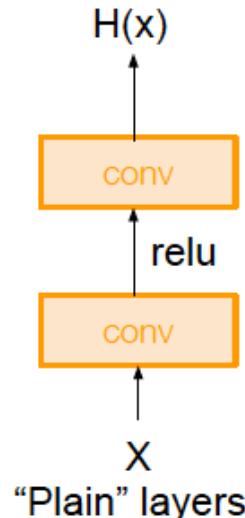
Effect of having more layers



Key Idea

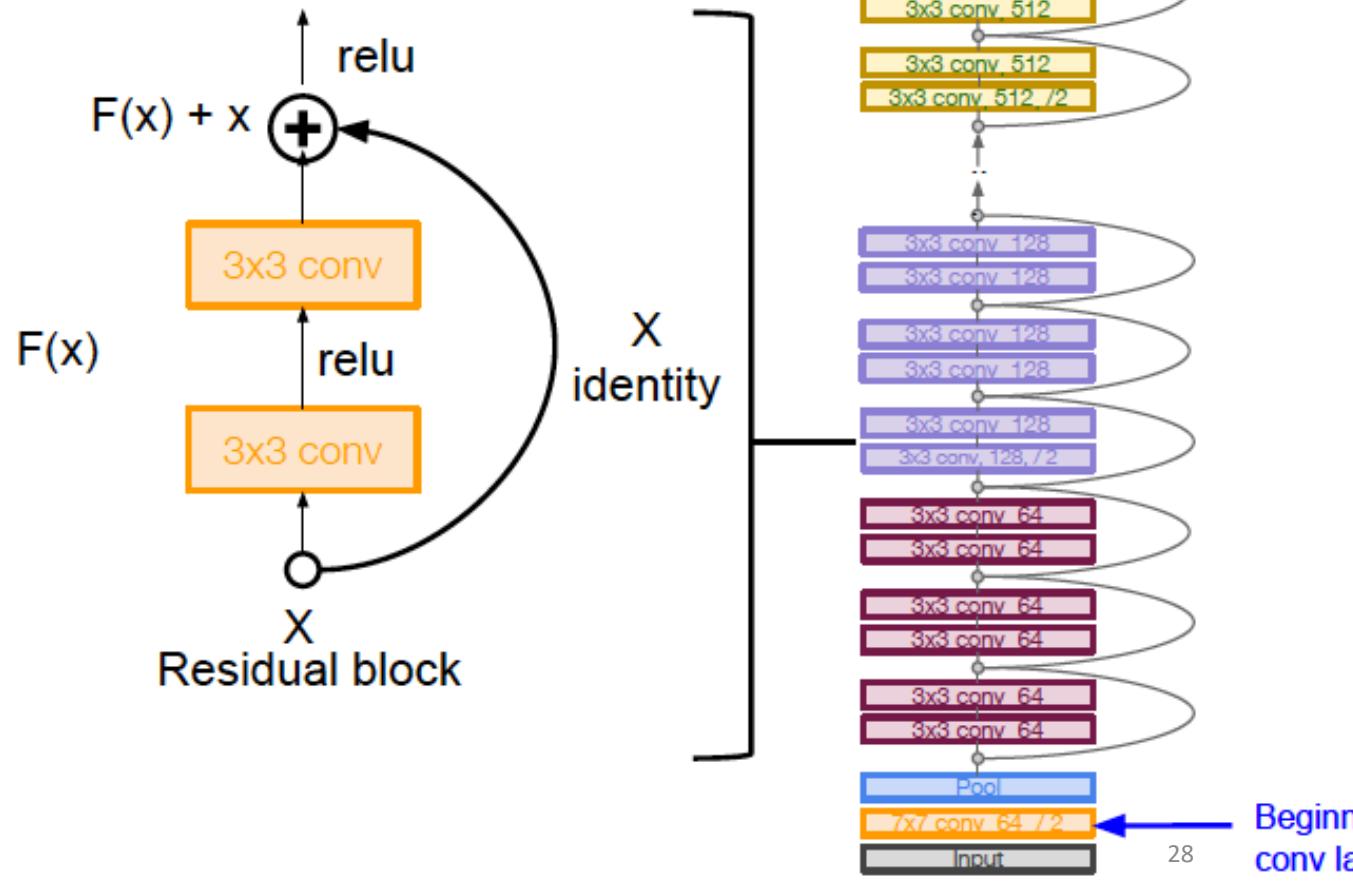
Skip Connection

- Deeper Model should have lower training error than Shallower model if the optimization is done properly.



ResNet Architecture

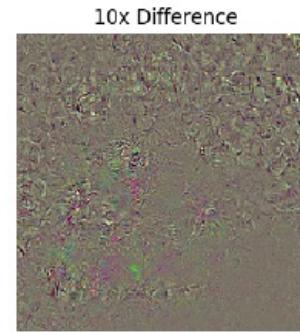
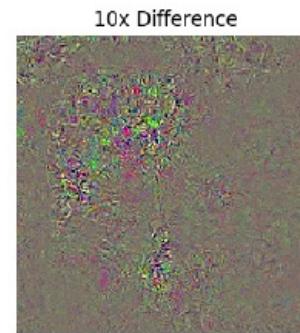
- Full ResNet architecture:
- - Stack residual blocks
- - Every residual block has
- two 3x3 conv layers
- - Periodically, double # of filters and downsample
- spatially using stride 2
- (/2 in each dimension)
- - Additional conv layer at the beginning
- - No FC layers at the end
- (only FC 1000 to output classes)



ResNet Architecture II

- Training ResNet in practice:
 - - Batch Normalization after every CONV layer
 - - Xavier/2 initialization from He et al.
- - SGD + Momentum (0.9)
- - Learning rate: 0.1, divided by 10 when validation error plateaus
- - Mini-batch size 256
- - Weight decay of 1e-5
- - No dropout used

Adversarial Examples



[cat image](#) is CC0 public domain
[elephant image](#) is CC0 public domain

Adversarial Examples

1. Start from an arbitrary image
2. Pick an arbitrary category
3. Modify the image (via gradient ascent)
to maximize the class score
4. Stop when the network is fooled

Feature Inversion

Given a CNN feature vector for an image, find a new image that:

- Matches the given feature vector
- “looks natural” (image prior regularization)

$$\mathbf{x}^* = \underset{\mathbf{x} \in \mathbb{R}^{H \times W \times C}}{\operatorname{argmin}} \ell(\Phi(\mathbf{x}), \Phi_0) + \lambda \mathcal{R}(\mathbf{x})$$

Given feature vector

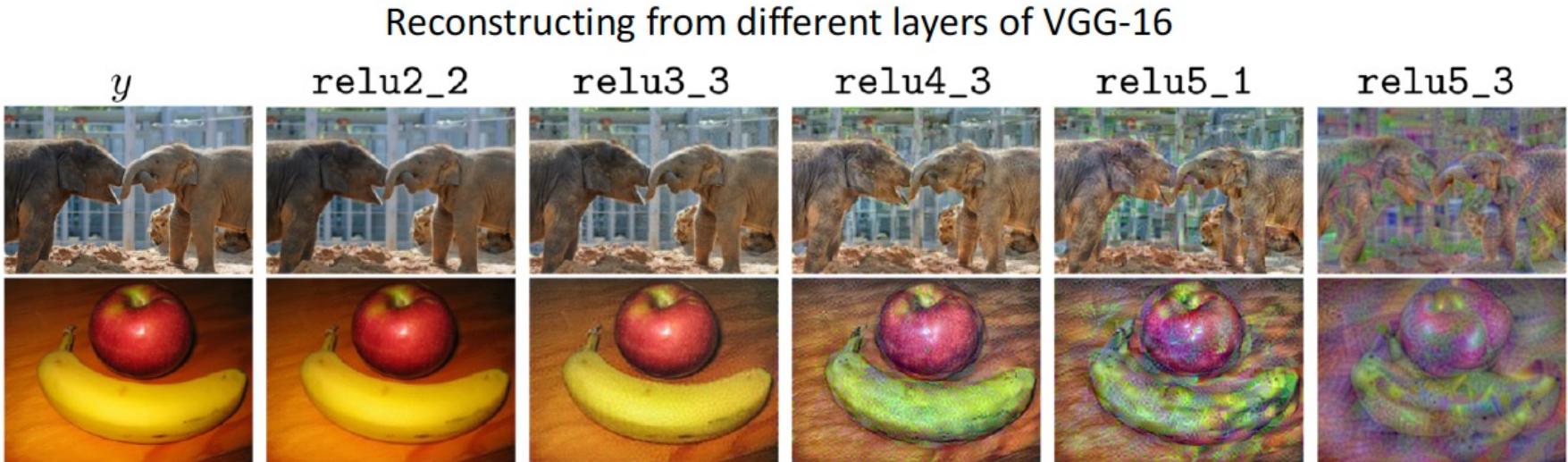
$$\ell(\Phi(\mathbf{x}), \Phi_0) = \|\Phi(\mathbf{x}) - \Phi_0\|^2$$

Features of new image

$$\mathcal{R}_{V^\beta}(\mathbf{x}) = \sum_{i,j} \left((x_{i,j+1} - x_{ij})^2 + (x_{i+1,j} - x_{ij})^2 \right)^{\frac{\beta}{2}}$$

Total Variation regularizer (encourages spatial smoothness)

Feature Inversion



Mahendran and Vedaldi. "Understanding Deep Image Representations by Inverting Them". CVPR 2015

Neural Style Transfer

Content Image



Style Image



+

=

Output Image

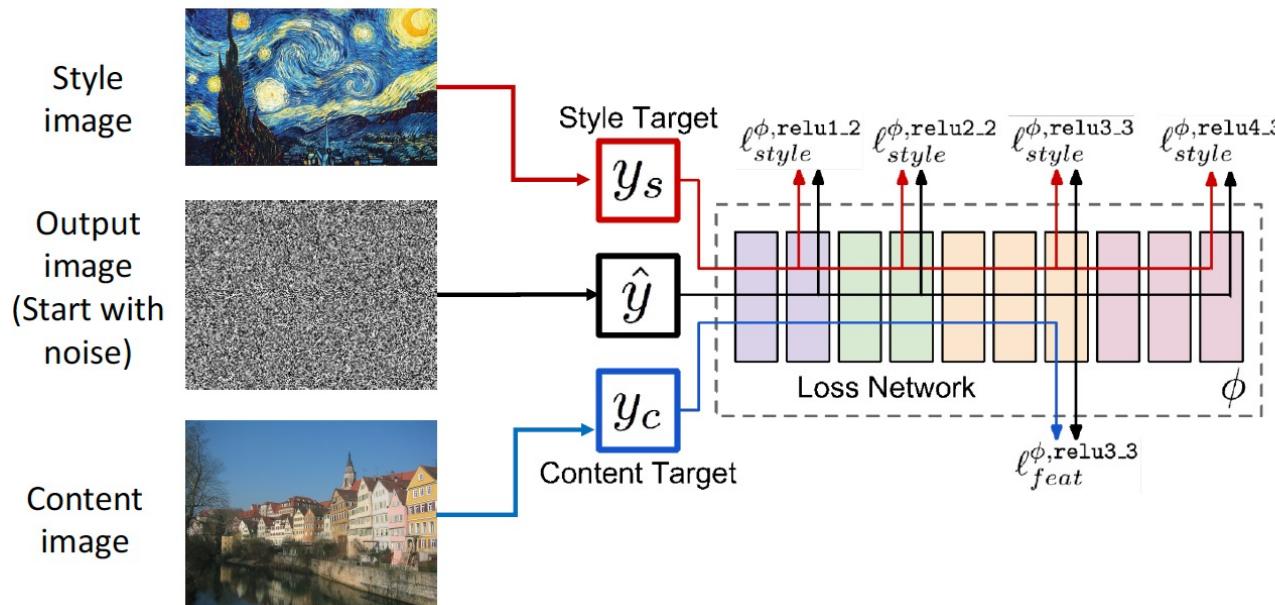


[This image](#) is licensed under CC-BY 3.0

[Starry Night](#) by Van Gogh is in the public domain

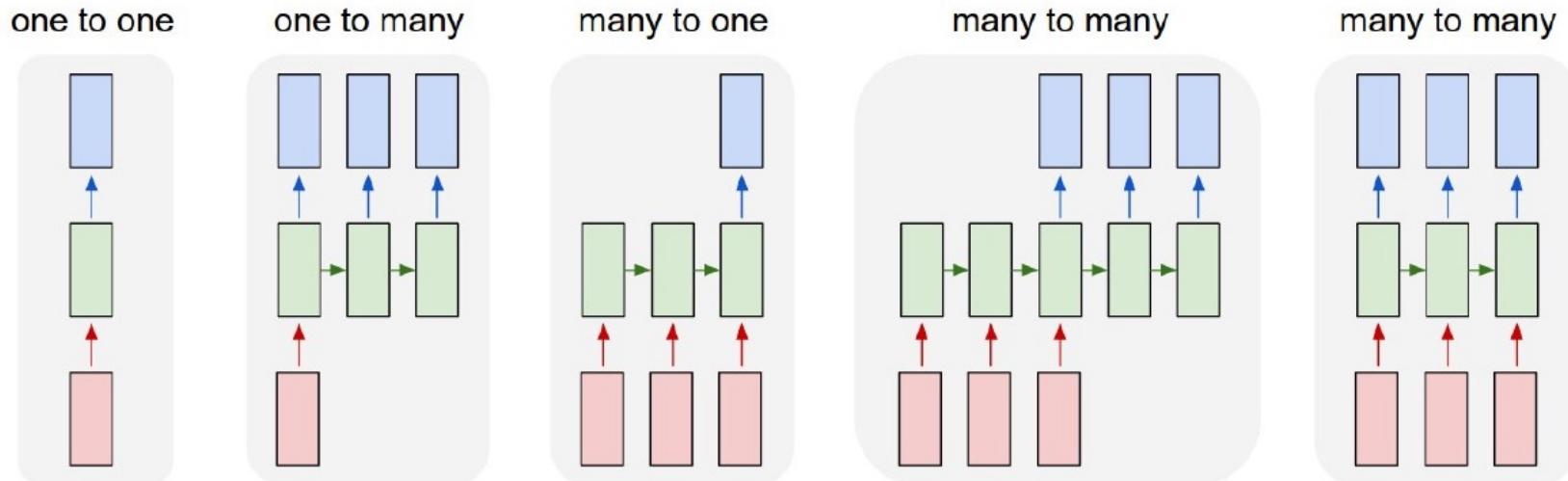
[This image](#) copyright Justin Johnson, 2015. Reproduced with permission.

Neural Style Transfer

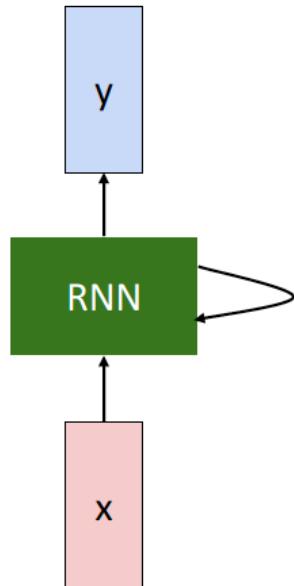


$$\arg \min_y \{ \beta \ell(\Phi(y), \Phi_1) + (1 - \beta) \ell(\Phi(y), \Phi_2) + \lambda \mathcal{R}(y) \}$$

Recurrent Neural Networks



RNN Architecture



The state consists of a single “*hidden*” vector \mathbf{h} :

$$h_t = f_W(h_{t-1}, x_t)$$

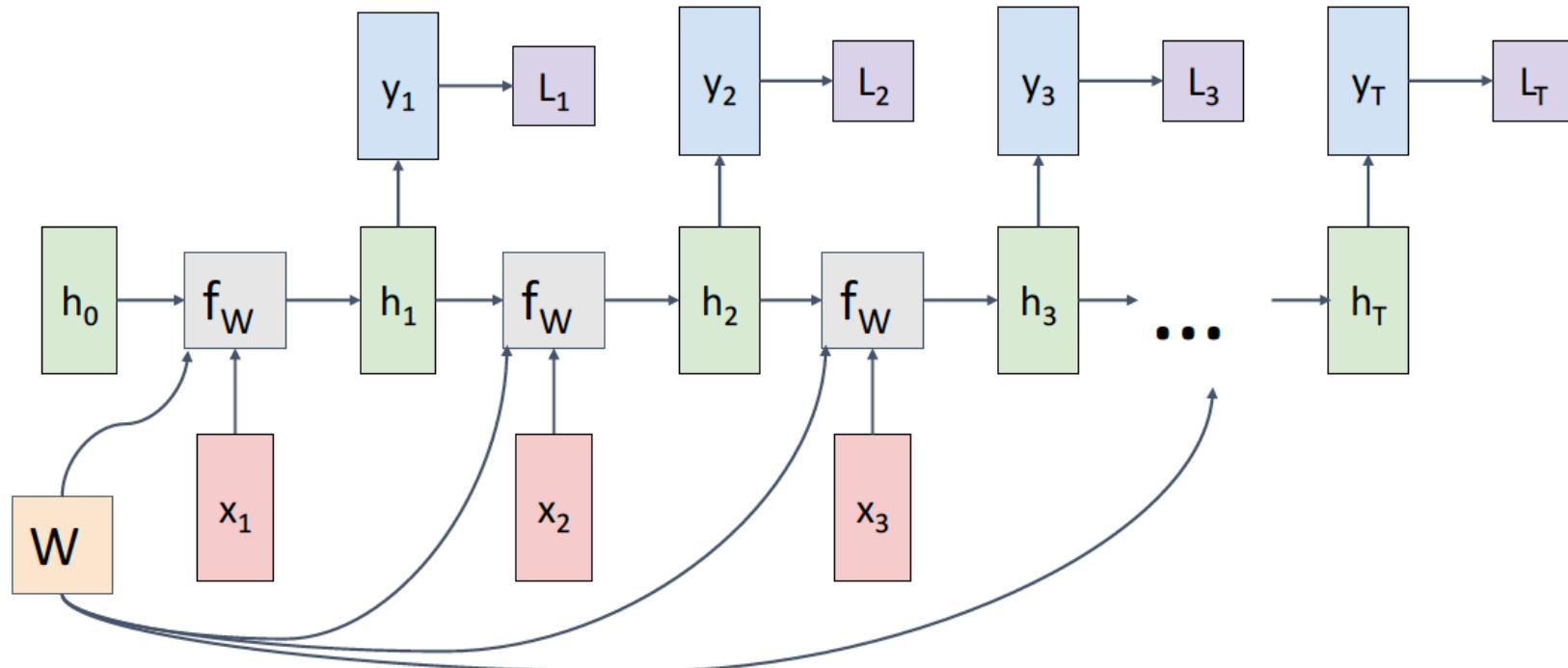
(also bias term)

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$y_t = W_{hy}h_t$$

Sometimes called a “Vanilla RNN” or an
“Elman RNN” after Prof. Jeffrey Elman

RNN Computational Graph



Application

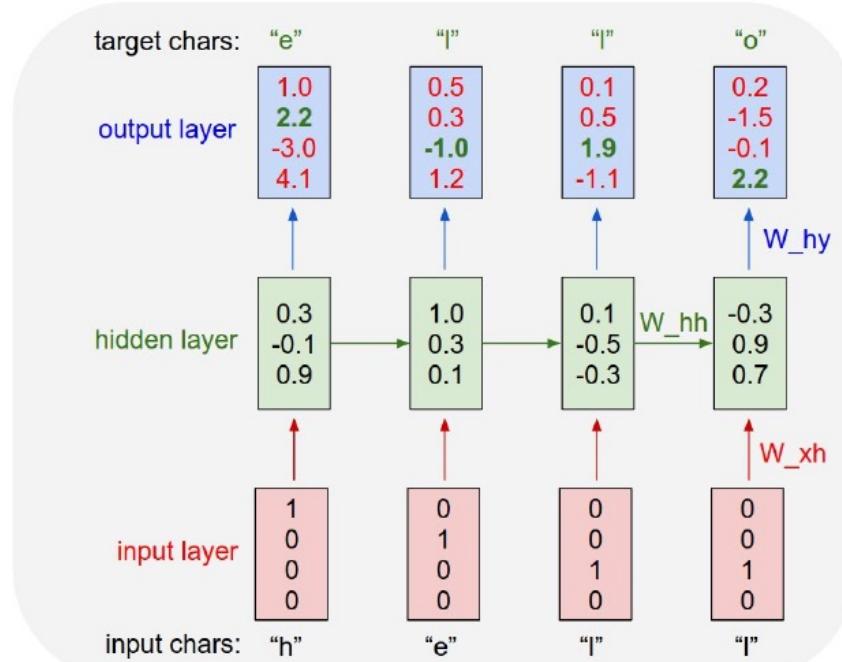
Example: Language Modeling

Given characters 1, 2, ..., t,
model predicts character t

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

Training sequence: "hello"

Vocabulary: [h, e, l, o]



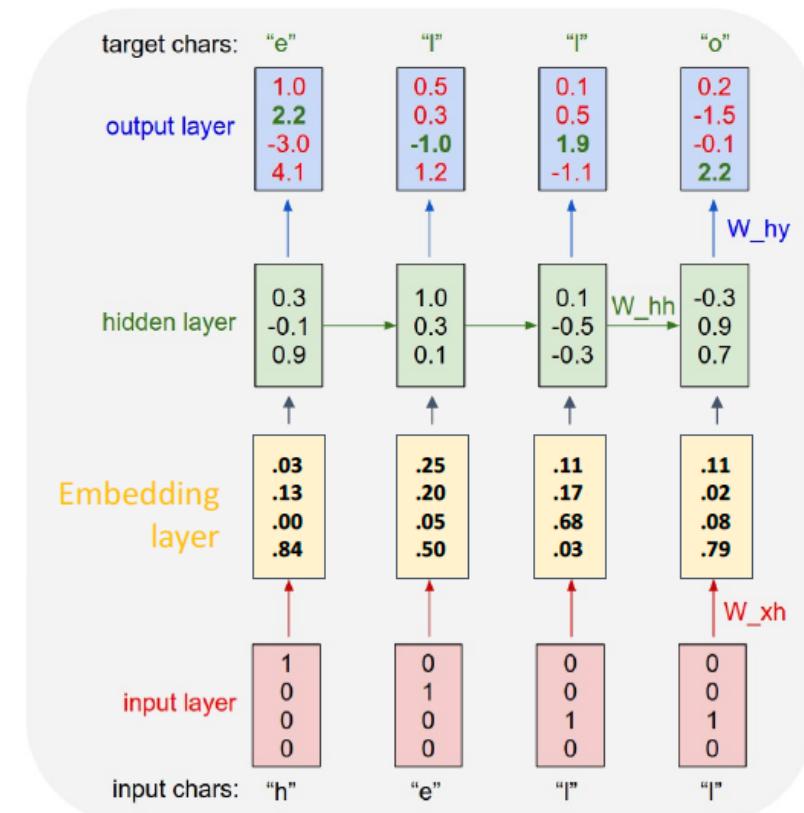
Application

Example: Language Modeling

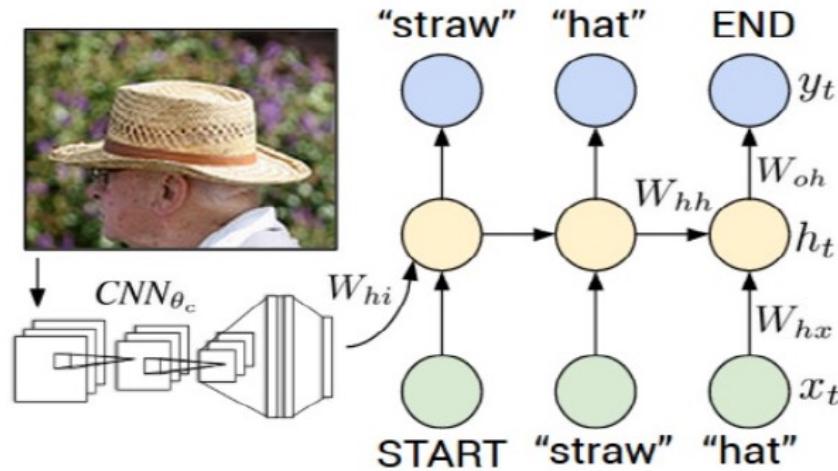
So far: encode inputs as **one-hot-vector**

$$\begin{bmatrix} w_{11} & w_{12} & w_{13} & w_{14} \end{bmatrix} [1] \quad [w_{11}] \\ \begin{bmatrix} w_{21} & w_{22} & w_{23} & w_{14} \end{bmatrix} [0] = [w_{21}] \\ \begin{bmatrix} w_{31} & w_{32} & w_{33} & w_{14} \end{bmatrix} [0] \quad [w_{31}] \\ [0] \end{math>$$

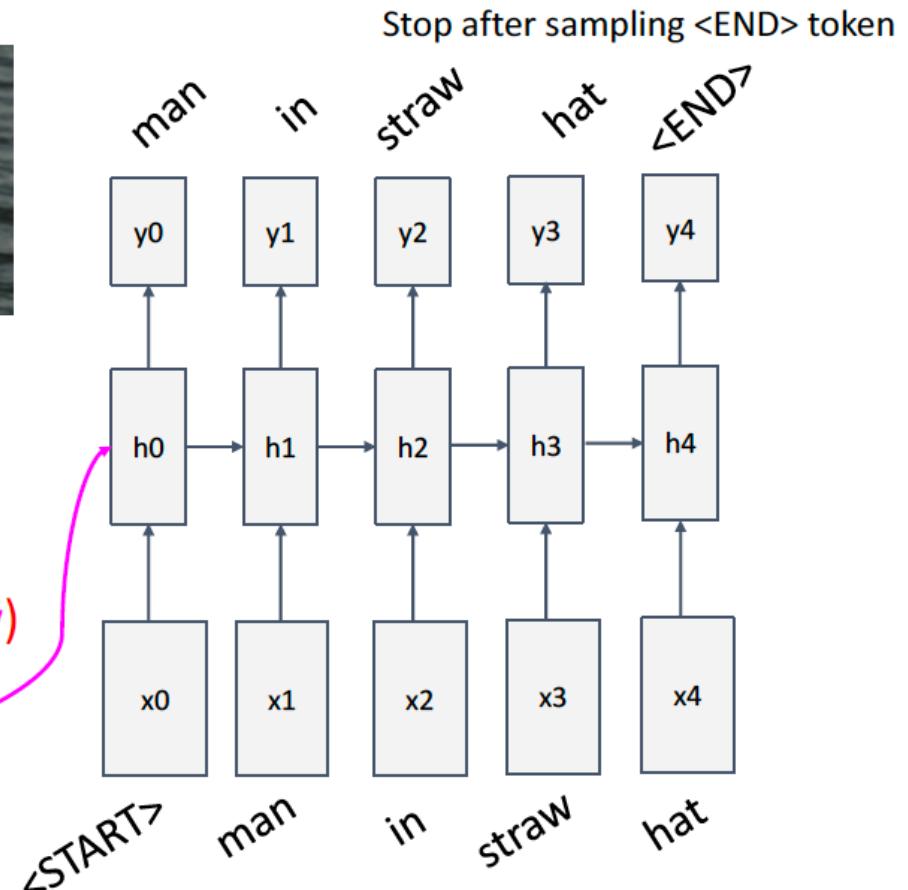
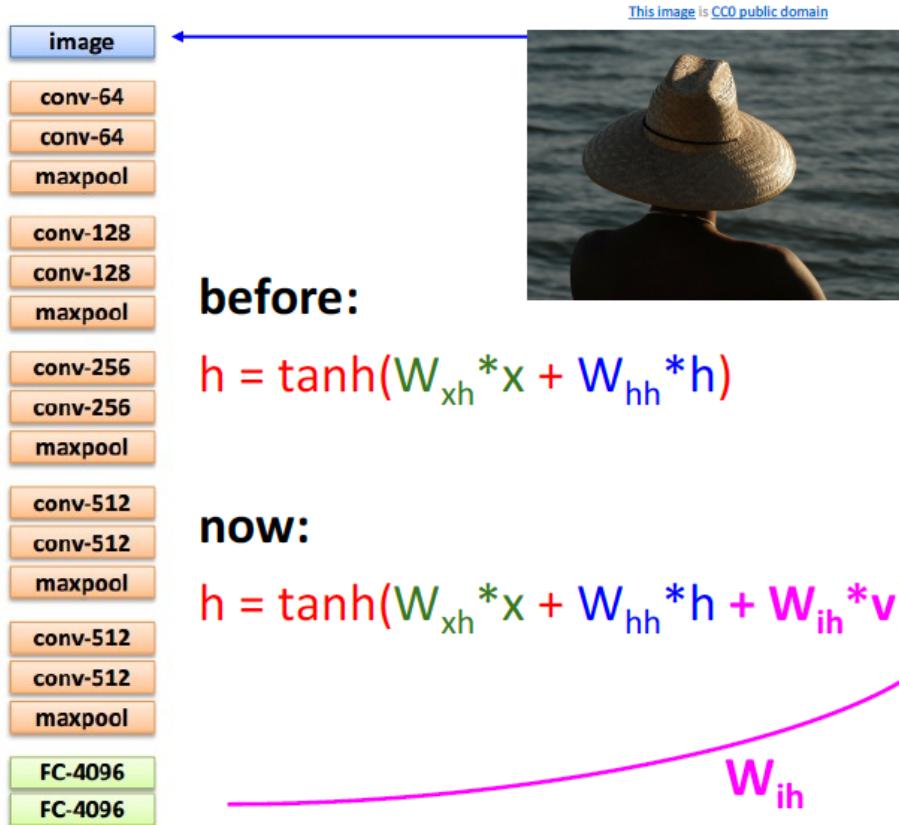
Matrix multiply with a one-hot vector just extracts a column from the weight matrix.
Often extract this into a separate **embedding layer**



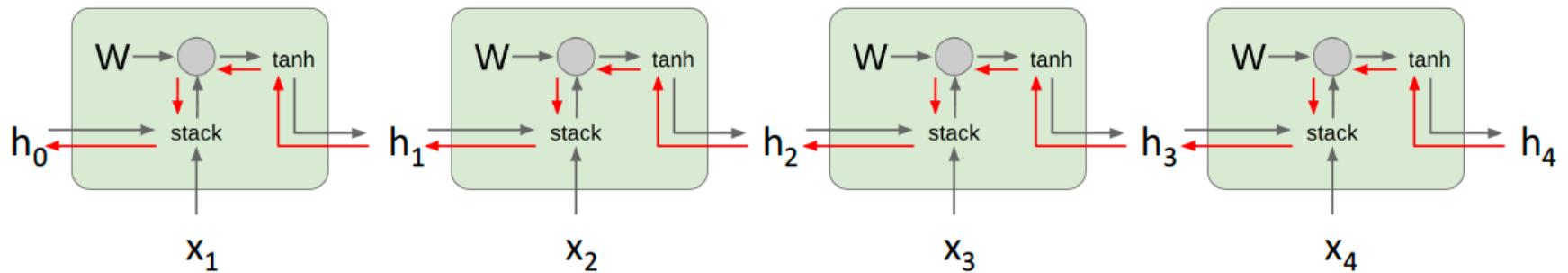
Example: Image Captioning



Application: Image Captioning



Vanishing Gradient Problem



Computing gradient of
 h_0 involves many
factors of W
(and repeated tanh)

Largest singular value > 1 :
Exploding gradients

Largest singular value < 1 :
Vanishing gradients

Long Short Term Memory

Vanilla RNN

$$h_t = \tanh \left(W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix} \right)$$

Compute four **gates**
at each timestep

LSTM

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} W \begin{pmatrix} h_{t-1} \\ x_t \end{pmatrix}$$

$$c_t = f \odot c_{t-1} + i \odot g$$

$$h_t = o \odot \tanh(c_t)$$

Unsupervised Learning (Not On Midterm)

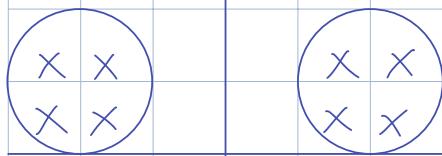
Training Set

$$D = \{x_1, x_2, \dots, x_N\}$$

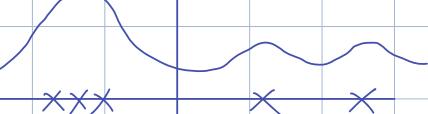
No Labels: y_1, y_2, \dots, y_n

- 1) Clustering
- 2) Density Estimation
- 3) Dimensionality Reduction

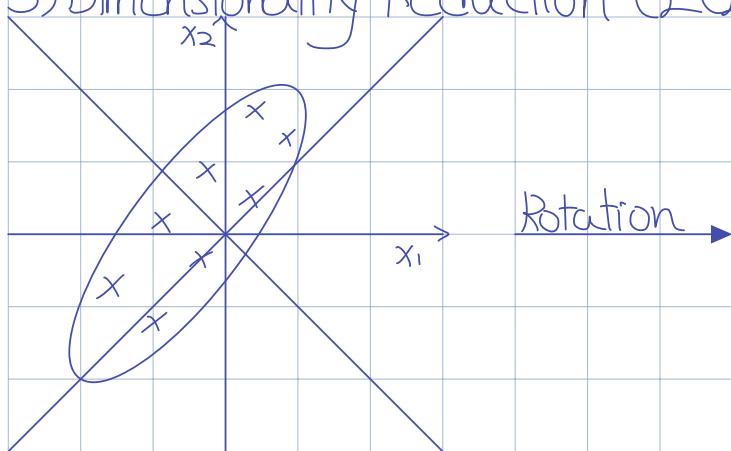
1) Clustering



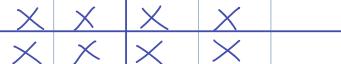
2) Density Estimation
 $f(x)$



3) Dimensionality Reduction (ECE 367)



(Matrix Analysis & Optimization)



Clustering

Partition training dataset D into k disjoint clusters such that elements within each cluster are close to one another.

Input, $D = \{x_1, x_2, \dots, x_n\}$

$k = \# \text{ of clusters } (k \ll N)$

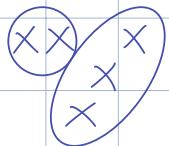
Output, Clusters S_1, S_2, \dots, S_k

$$S_i \subseteq D$$

$$S_i \cap S_j = \emptyset \quad (i \neq j)$$

$$S_1 \cup S_2 \cup \dots \cup S_k = D$$

$$K=2$$



Cluster Center

Let μ_i be the cluster center for S_i

Approximation error for cluster S_j with cluster center μ_j

$$E_j(S_j, \mu_j) = \sum_{x_n \in S_j} \|x_n - \mu_j\|^2$$

Total Approximation Error

$$\begin{aligned} E_n(S_1, S_2, \dots, S_k, \underline{\mu_1, \mu_2, \dots, \mu_k}) &= \sum_{j=1}^K E_j(S_j, \mu_j) \\ &= \sum_{\substack{j=1 \\ x_n \in S_j}}^K \|x_n - \mu_j\|^2 \end{aligned}$$

$$= \sum_{n=1}^N \|x_n - \mu(x_n)\|^2$$



Center of cluster associated with x_n

Optimal Clustering

$$\underset{\substack{S_1, \dots, S_k \\ \mu_1, \dots, \mu_k}}{\text{minimize}} E_n(S_1, S_2, \dots, S_k, \underline{\mu_1, \mu_2, \dots, \mu_k})$$

NP Hard (Bad News)

Excellent Heuristics (Good News)

Find Local Minimum

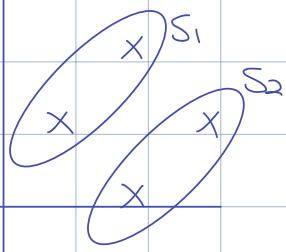
K Means Algorithm

Problem 1, Given S_1, S_2, \dots, S_k ,

Find $\mu_1, \mu_2, \dots, \mu_k$ to minimize approximation error

Select μ_j to minimize

$$E_j(S_j, \mu_j) = \sum_{x_n \in S_j} \|x_n - \mu_j\|^2$$



$$E_1(S_1, \mu_1) + E_2(S_2, \mu_2)$$

$$\mu_1 = \frac{1}{2}(x_1 + x_2)$$

$$\nabla_{\mu_j} E_j(S_j, \mu_j) = 0$$

$$\nabla_{\mu_j} \sum_{x_n \in S_j} \|x_n - \mu_j\|^2 = 0$$

$$\mu_j = \frac{1}{|S_j|} \sum_{x_n \in S_j} x_n$$

number of points in S_j

Given $\mu_1, \mu_2, \dots, \mu_K$ (Problem 2)

Find S_1, S_2, \dots, S_K to minimize

$$\begin{aligned} E_{\text{tot}}(S_1, S_2, \dots, S_K, \mu_1, \mu_2, \dots, \mu_K) &= \sum_{j=1}^K E_j(S_j, \mu_j) \\ &= \sum_{n=1}^N \|x_n - \mu(x_n)\|^2 \end{aligned}$$

Select closest point

$$x_n \in S_j^*$$

$$\text{if } j^* = \operatorname{argmin}_{\substack{j \\ j \in \{1, 2, \dots, k\}}} \|x_n - \mu_j\|^2$$