



Carnegie Mellon University

Introduction to Deep Learning for Engineers

Spring 2025, Deep Learning for Engineers
Feb 8, 2025, Eighth Session

Amir Barati Farimani

*Associate Professor of Mechanical Engineering and Bio-Engineering
Carnegie Mellon University*

Story so far on CNN:

1. Concept of Representation + Learning

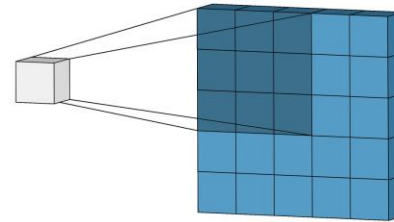
We need a **UNIVERSAL REPRESENTATION LEARNER** to ease the learning

2. How to learn robust representation?

We need a representation learner that can automatically learn the features needed for the task

3. How to learn spatial, high level features (Swan example)

Since high level features are spatial (not pixels) we need a scanner to scan patches of data instead of single pixels



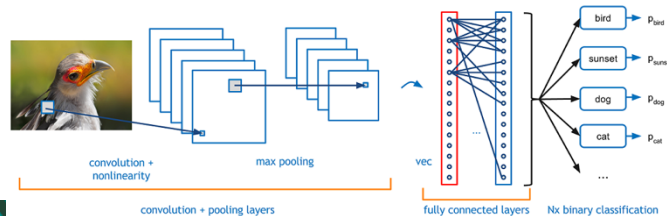
Story so far on CNN:

4. How to build a scanner for feature learning? what should be the properties of this scanner?

1. It should be numbers (a matrix) because it should be machine readable
2. It should be learnable
3. It should be flexible in size and dimension
4. Should be pluggable to Neural Networks

5. Can we design the scanners based on the learning tasks?

Yes, and we should. Because the mode of data might be different (sound, image, video) and features are needed based on the task to make a good model, the scanner should **ONLY** learn the relevant features connecting them to the output



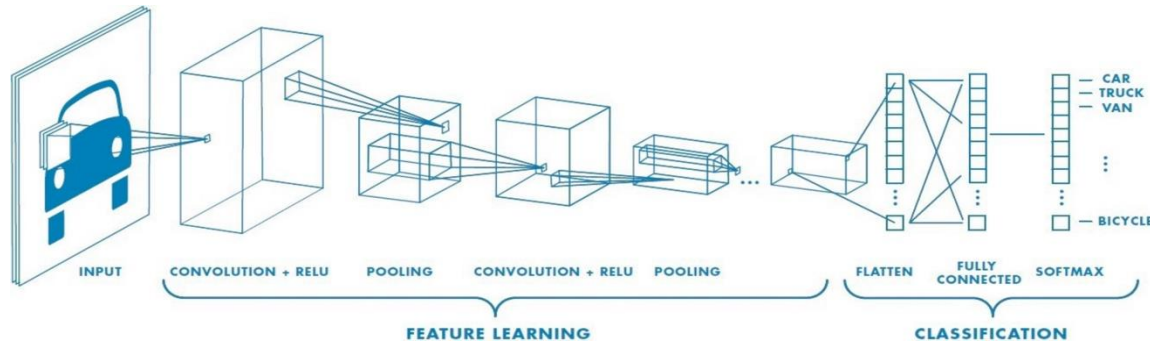
Story so far on CNN:

6. How can the scanners learn?

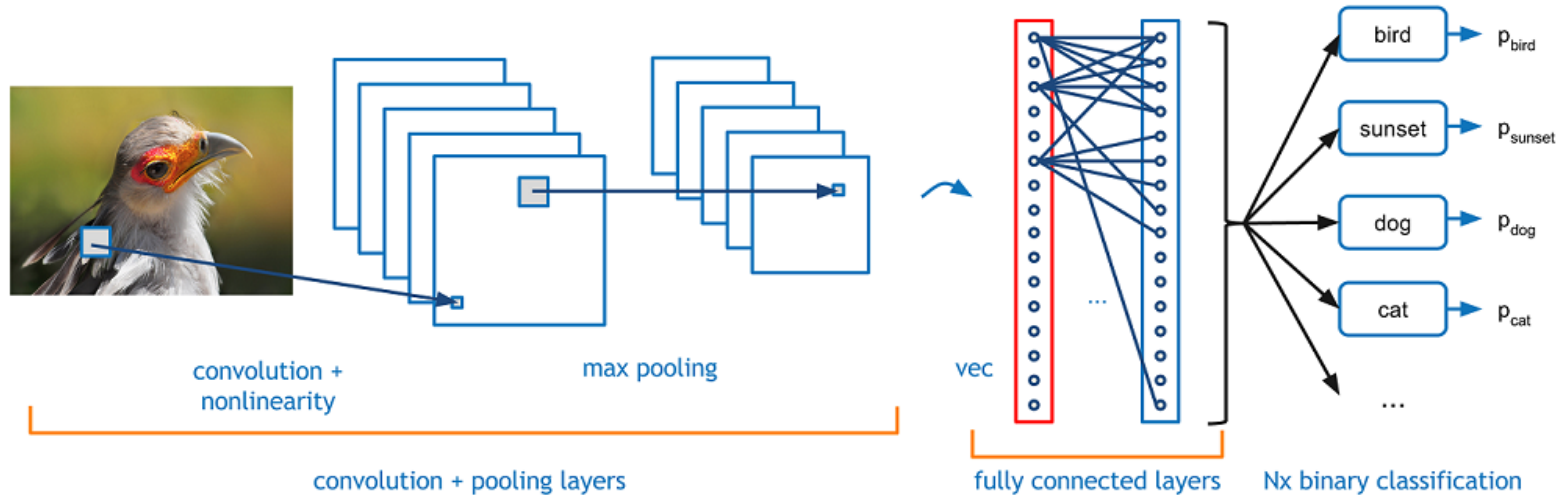
Inspired by iterative optimization and backpropagation in neural networks, we can iteratively learn the initialized weights of scanners (remember these are numbers)

7. How can we plug in the scanners into Neural networks?

By flattening the output of the last convolved map and passing it to the FC layer, we can forward propagate, and we can backpropagate to learn the parameters of a filter (scanner)



Recap: CNN Overall Architecture



CNNs are Automatic Feature Detectors

Sharing
Weights

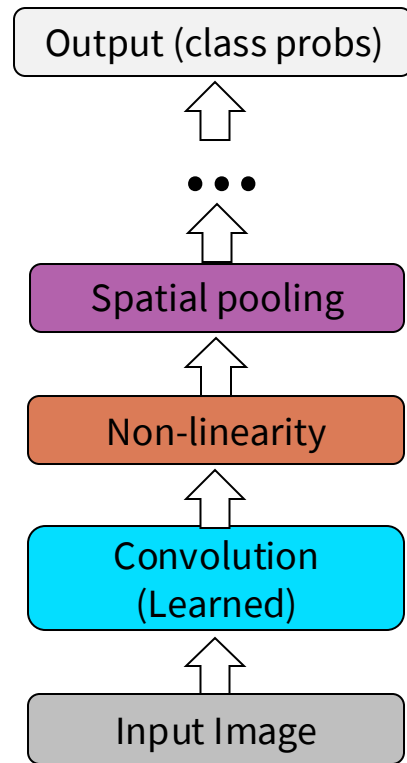
Feature
Detection

Spatial Local
Features

Translation
In-variant

Convolutional Neural Networks (CNN)

- FEED-FORWARD FEATURE EXTRACTION:
 - 1. Convolve input with learned filters**
 - 2. Apply non-linearity**
 - 3. Spatial pooling (downsample)**
- SUPERVISED TRAINING OF CONVOLUTIONAL FILTERS BY BACK-PROPAGATING CLASSIFICATION ERROR



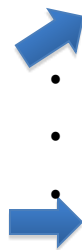
Convolution

- APPLY LEARNED FILTER WEIGHTS
- ONE FEATURE MAP PER FILTER
- STRIDE CAN BE GREATER THAN 1
(FASTER, LESS MEMORY)



Input

Adapted from Rob Fergus



Feature Map

The Convolution operation

ASSUME A 6X6 MATRIX **M** AS INPUT. THE 2D CONVOLUTION OF **M** WITH *FILTER (OR KERNEL) F* AND *STRIDE 1* IS A 4X4 MATRIX **CM** (SOMETIMES CALLED *FEATURE MAP*) COMPUTED AS FOLLOWS:

3	0	1	2	4	7
1	5	8	9	3	1
2	7	2	5	1	3
0	1	3	1	7	8
4	2	1	6	2	8
2	4	5	2	3	9

M

*

1	0	-1
1	0	-1
1	0	-1

F

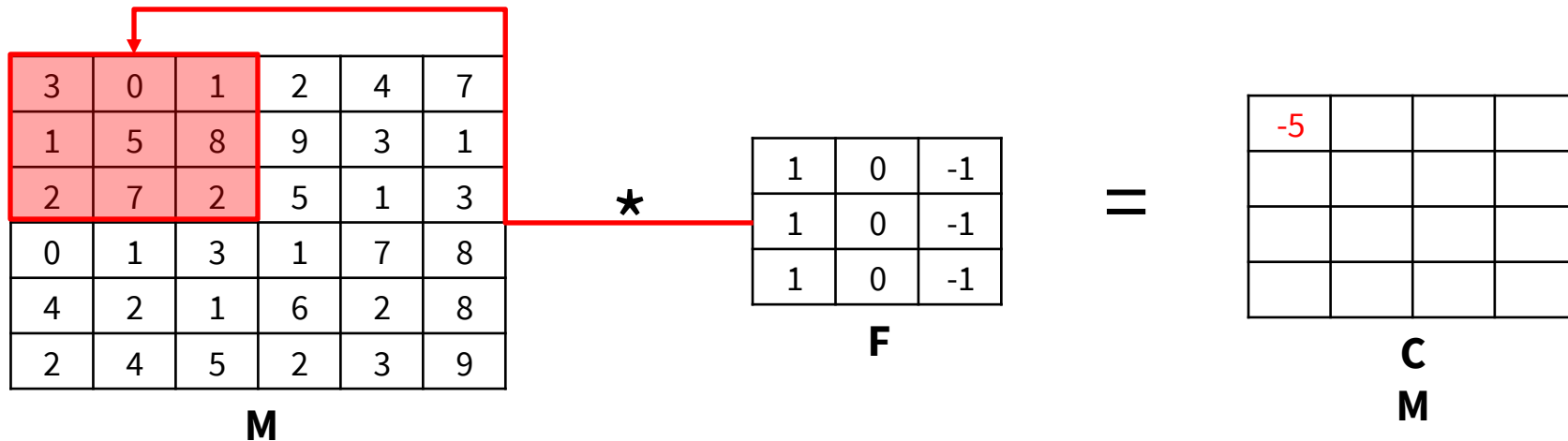
=

CM

convoluti
on
operator

The Convolution operation

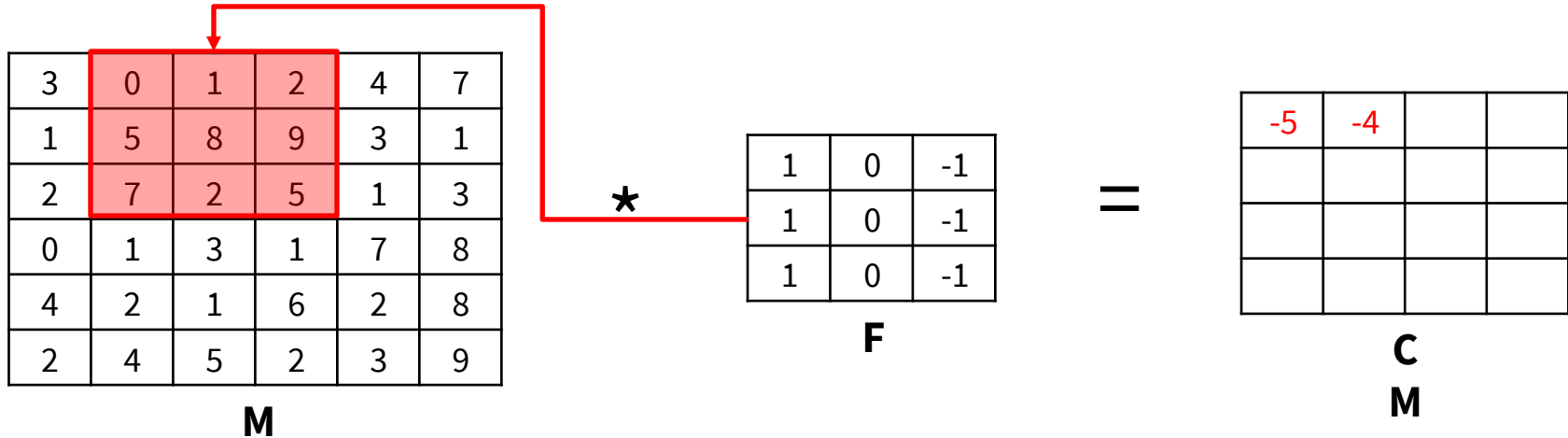
ASSUME A 6X6 MATRIX **M** AS INPUT. THE 2D CONVOLUTION OF **M** WITH *FILTER (OR KERNEL) F* AND *STRIDE 1* IS A 4X4 MATRIX **CM** (SOMETIMES CALLED *FEATURE MAP*) COMPUTED AS FOLLOWS:



$$3*1 + 1*1 + 2*1 + 0*0 + 5*0 + 7*0 + 1*(-1) + 8*(-1) + 2*(-1)$$

The Convolution operation

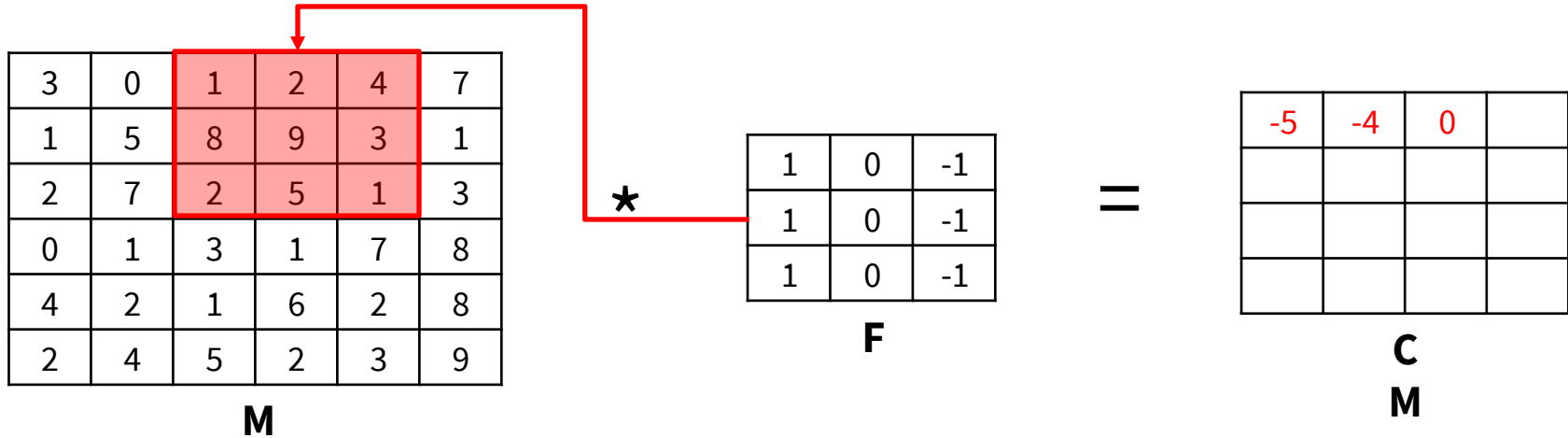
ASSUME A 6X6 MATRIX **M** AS INPUT. THE 2D CONVOLUTION OF **M** WITH *FILTER (OR KERNEL) F* AND *STRIDE 1* IS A 4X4 MATRIX **CM** (SOMETIMES CALLED *FEATURE MAP*) COMPUTED AS FOLLOWS:



$$0*1 + 5*1 + 7*1 + 1*0 + 8*0 + 2*0 + 2*(-1) + 9*(-1) + 5*(-1)$$

The Convolution operation

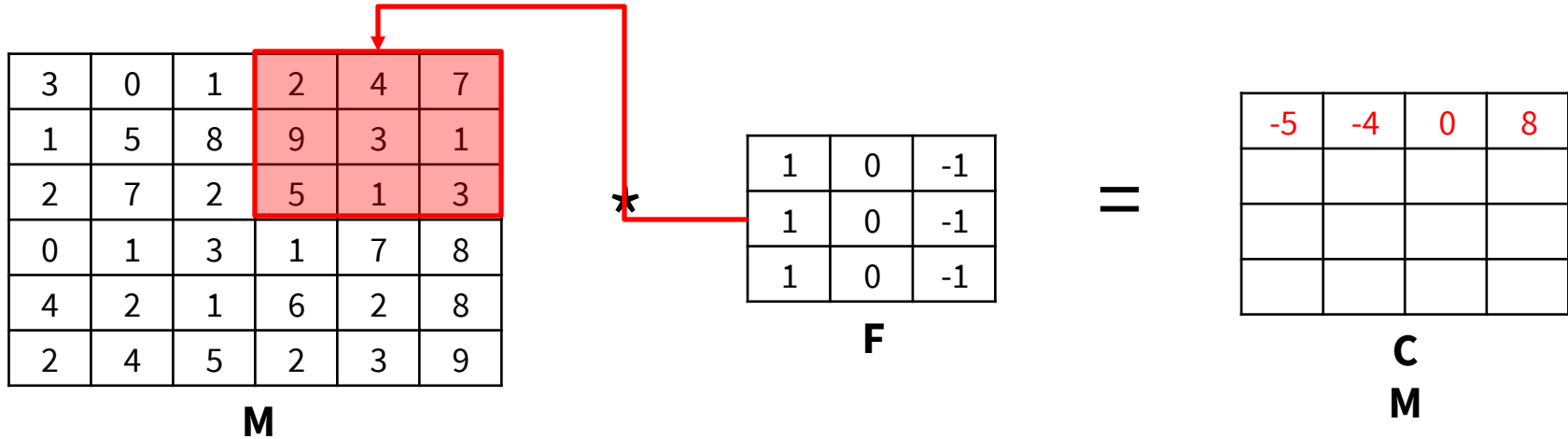
ASSUME A 6X6 MATRIX **M** AS INPUT. THE 2D CONVOLUTION OF **M** WITH *FILTER (OR KERNEL) F* AND *STRIDE 1* IS A 4X4 MATRIX **CM** (SOMETIMES CALLED *FEATURE MAP*) COMPUTED AS FOLLOWS:



$$1*1 + 8*1 + 2*1 + 2*0 + 9*0 + 5*0 + 4*(-1) + 3*(-1) + 1*(-1)$$

The Convolution operation

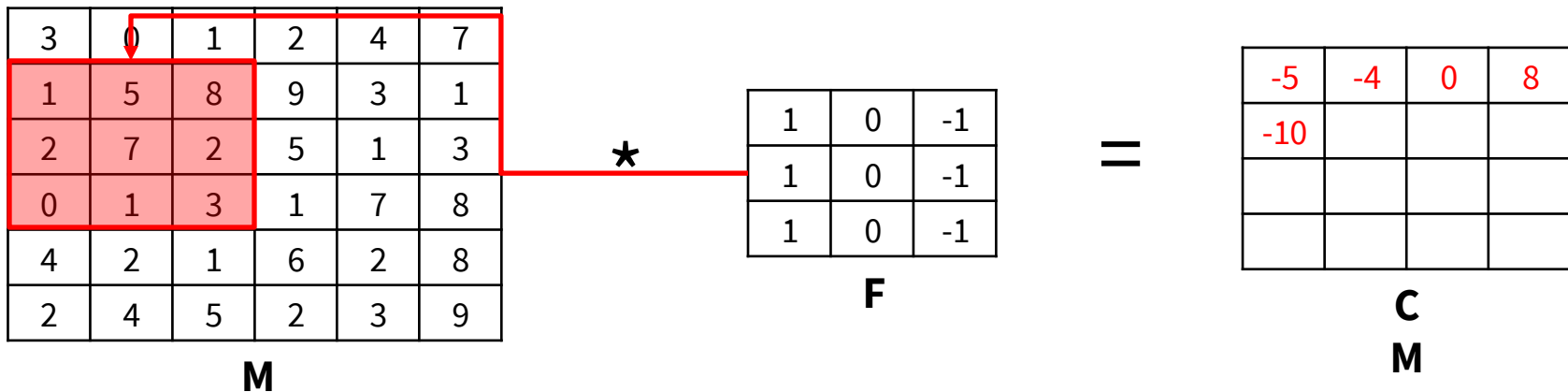
ASSUME A 6X6 MATRIX **M** AS INPUT. THE 2D CONVOLUTION OF **M** WITH *FILTER (OR KERNEL) F* AND *STRIDE 1* IS A 4X4 MATRIX **CM** (SOMETIMES CALLED *FEATURE MAP*) COMPUTED AS FOLLOWS:



$$2*1 + 9*1 + 5*1 + 4*0 + 3*0 + 1*0 + 7*(-1) + 1*(-1) + 3*(-1)$$

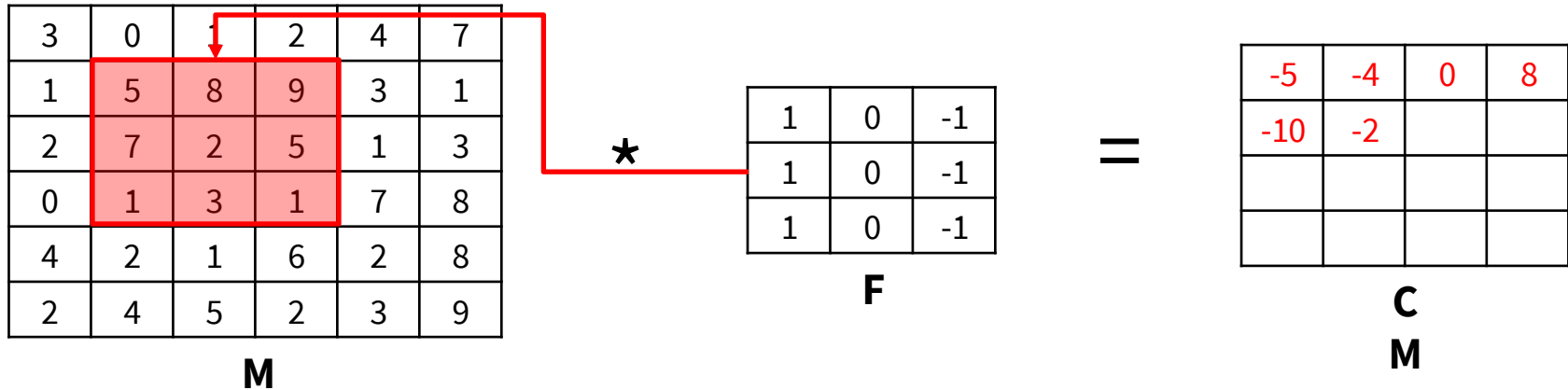
The Convolution operation

ASSUME A 6X6 MATRIX **M** AS INPUT. THE 2D CONVOLUTION OF **M** WITH *FILTER (OR KERNEL) F* AND *STRIDE 1* IS A 4X4 MATRIX **CM** (SOMETIMES CALLED *FEATURE MAP*) COMPUTED AS FOLLOWS:



The Convolution operation

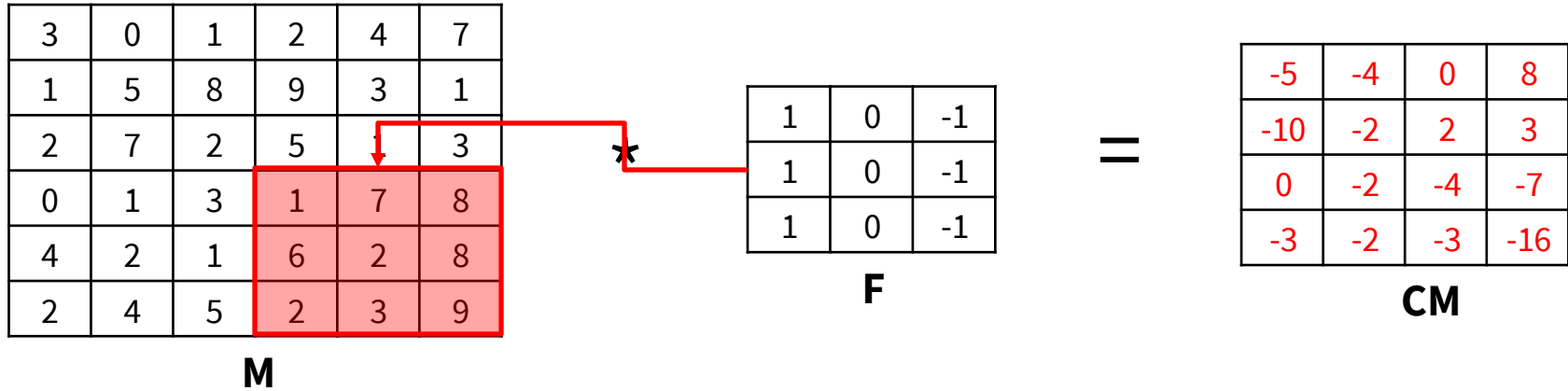
ASSUME A 6X6 MATRIX **M** AS INPUT. THE 2D CONVOLUTION OF **M** WITH *FILTER (OR KERNEL)* **F** AND *STRIDE 1* IS A 4X4 MATRIX **CM** (SOMETIMES CALLED *FEATURE MAP*) COMPUTED AS FOLLOWS:



$$5*1 + 7*1 + 1*1 + 8*0 + 2*0 + 3*0 + 9*(-1) + 5*(-1) + 1*(-1)$$

The Convolution operation¹

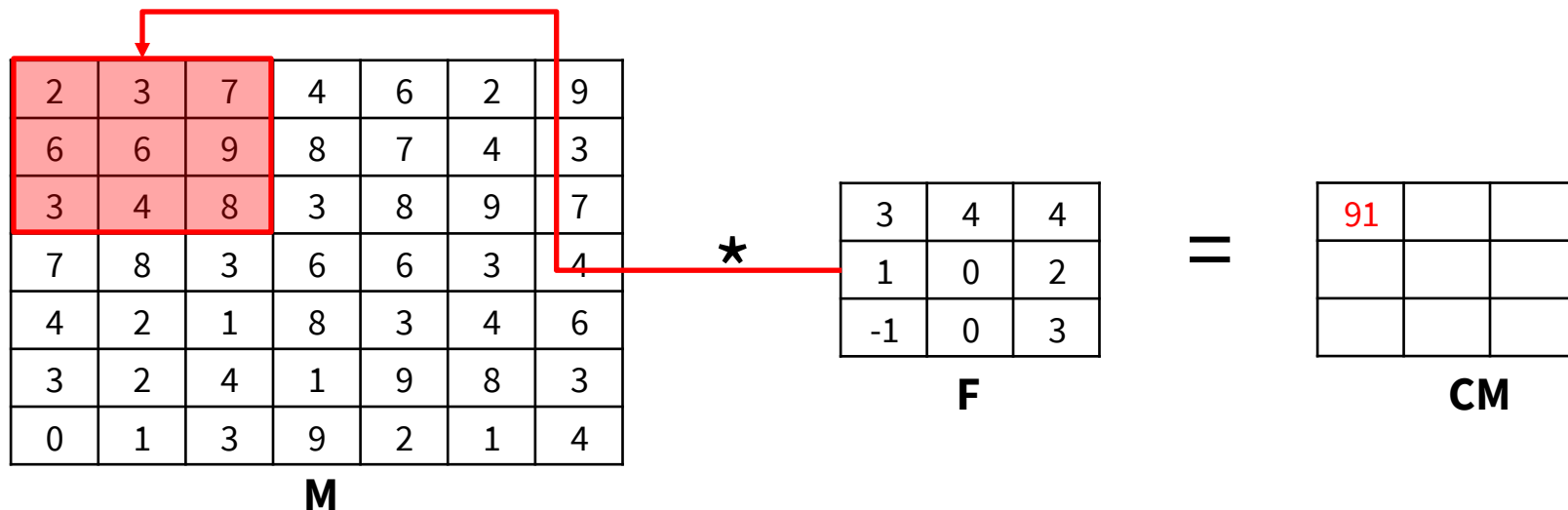
ASSUME A 6X6 MATRIX **M** AS INPUT. THE 2D CONVOLUTION OF **M** WITH *FILTER (OR KERNEL)* **F** AND *STRIDE 1* IS A 4X4 MATRIX **CM** (SOMETIMES CALLED *FEATURE MAP*) COMPUTED AS FOLLOWS:



$$1*1 + 6*1 + 2*1 + 7*0 + 2*0 + 3*0 + 8*(-1) + 8*(-1) + 9*(-1)$$

Larger Strides

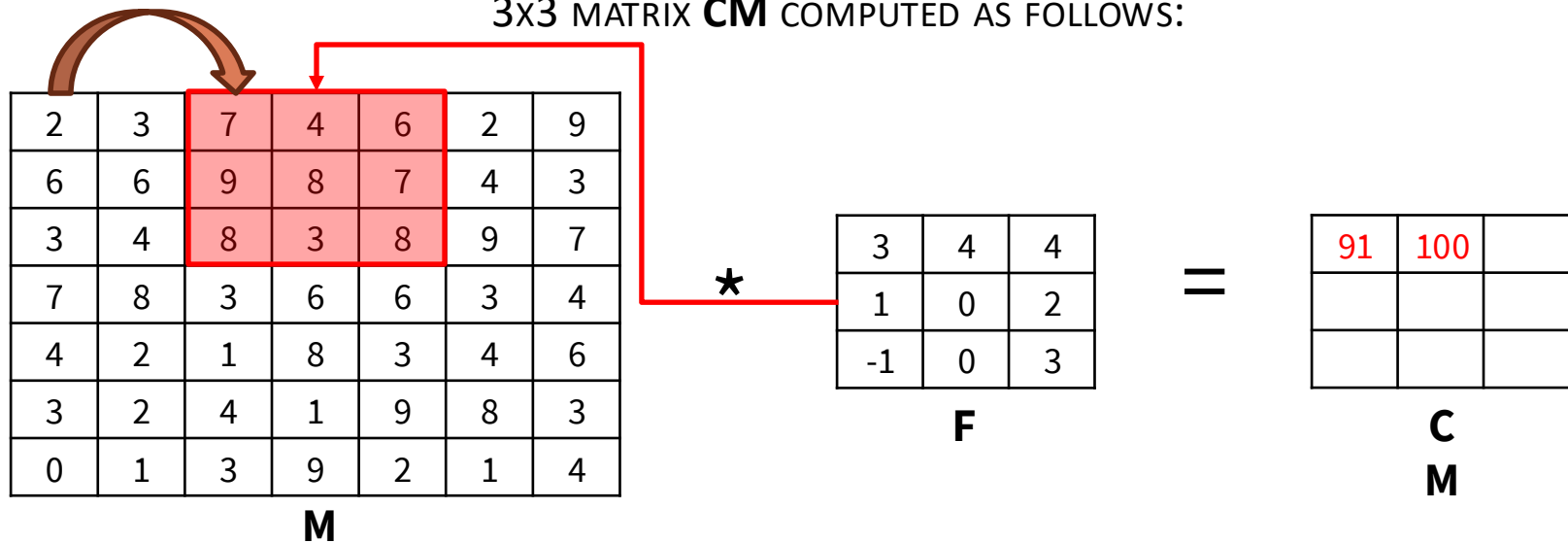
ASSUME A M , A 7×7 MATRIX. THE 2D CONVOLUTION OF M WITH *FILTER* F AND **STRIDE 2** IS A 3×3 MATRIX **CM** COMPUTED AS FOLLOWS:



$$2*3 + 6*1 + 3*(-1) + 3*4 + 6*0 + 4*0 + 7*4 + 9*2 + 8*3$$

Larger Strides

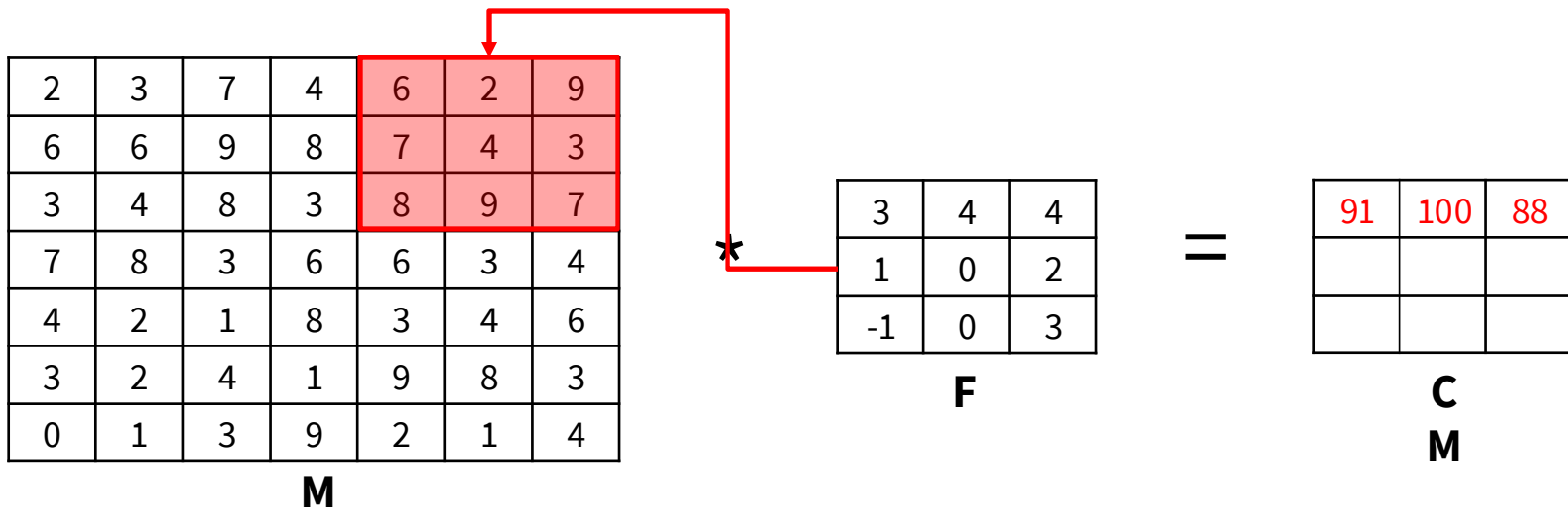
ASSUME A M , A 7×7 MATRIX. THE 2D CONVOLUTION OF M WITH *FILTER* F AND **STRIDE 2** IS A 3×3 MATRIX **CM** COMPUTED AS FOLLOWS:



$$7*3 + 9*1 + 8*(-1) + 4*4 + 8*0 + 3*0 + 6*4 + 7*2 + 8*3$$

Larger Strides

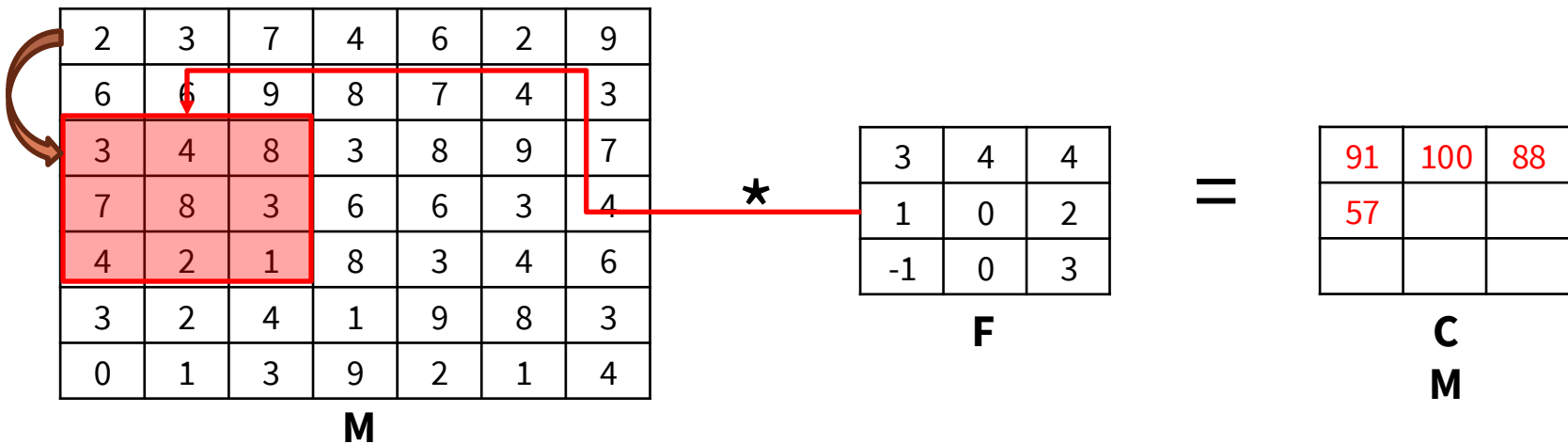
ASSUME A **M**, A 7x7 MATRIX. THE 2D CONVOLUTION OF **M** WITH *FILTER F* AND **STRIDE 2** IS A 3x3 MATRIX **CM** COMPUTED AS FOLLOWS:



$$6*3 + 7*1 + 8*(-1) + 2*4 + 4*0 + 9*0 + 9*4 + 3*2 + 7*3$$

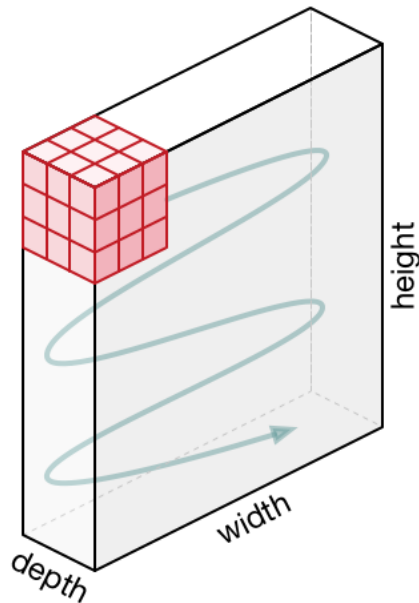
Larger Strides

ASSUME A M , A 7×7 MATRIX. THE 2D CONVOLUTION OF M WITH *FILTER* F AND **STRIDE 2** IS A 3×3 MATRIX **CM** COMPUTED AS FOLLOWS:



$$3 \cdot 3 + 7 \cdot 1 + 4 \cdot (-1) + 4 \cdot 4 + 8 \cdot 0 + 2 \cdot 0 + 8 \cdot 4 + 3 \cdot 2 + 1 \cdot 3$$

Recap: How it works for Images?



0	0	0	0	0	0	...
0	156	155	156	158	158	...
0	153	154	157	159	159	...
0	149	151	155	158	159	...
0	146	146	149	153	158	...
0	145	143	143	148	158	...
...

Input Channel #1 (Red)

-1	-1	1
0	1	-1
0	1	1

Kernel Channel #1



308

+

0	0	0	0	0	0	...
0	167	166	167	169	169	...
0	164	165	168	170	170	...
0	160	162	166	169	170	...
0	156	156	159	163	168	...
0	155	153	153	158	168	...
...

Input Channel #2 (Green)

1	0	0
1	-1	-1
1	0	-1

Kernel Channel #2



-498

+

0	0	0	0	0	0	...
0	163	162	163	165	165	...
0	160	161	164	166	166	...
0	156	158	162	165	166	...
0	155	155	158	162	167	...
0	154	152	152	157	167	...
...

Input Channel #3 (Blue)

0	1	1
0	1	0
1	-1	1

Kernel Channel #3



164

+

+ 1 = -25

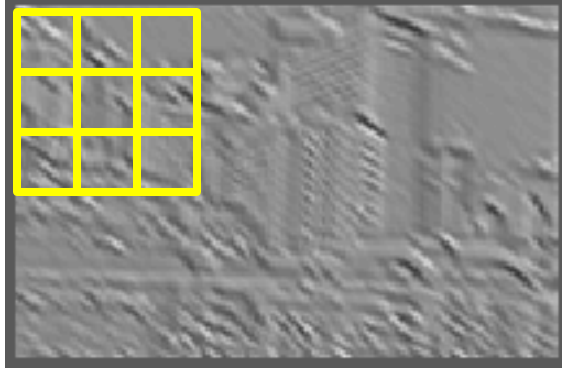
↑
Bias = 1

Output

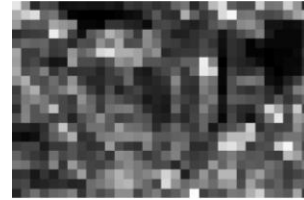
-25			...
			...
			...
			...
...

Spatial Pooling

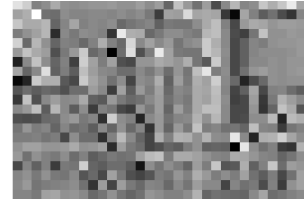
- SUM OR MAX OVER NON-OVERLAPPING / OVERLAPPING REGIONS
- ROLE OF POOLING:
 - Invariance to small transformations
 - Larger receptive fields (neurons see more of input)



Max



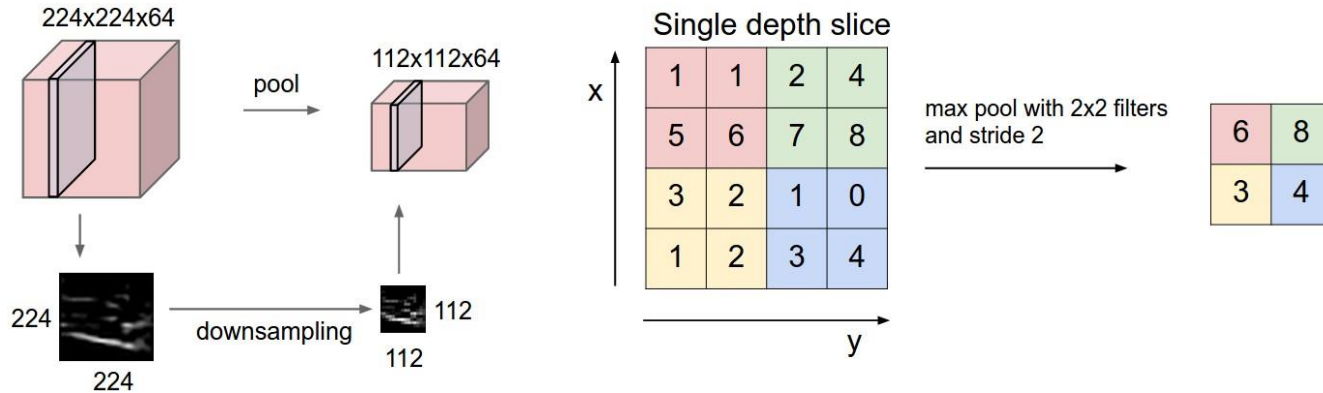
Sum



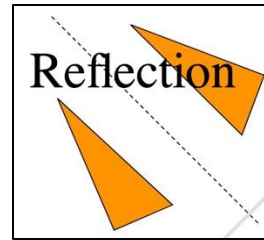
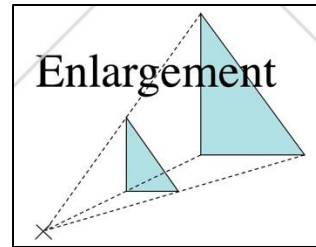
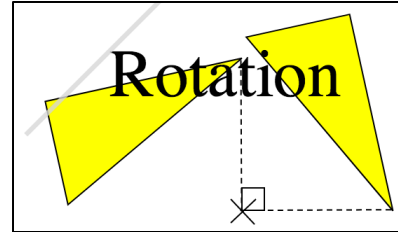
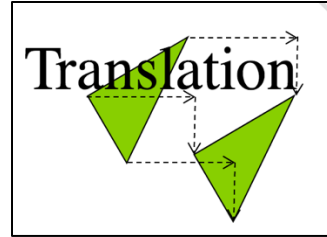
Adapted from Rob Fergus

Spatial Pooling

- SUM OR MAX OVER NON-OVERLAPPING / OVERLAPPING REGIONS
- ROLE OF POOLING:
 - Invariance to small transformations
 - Larger receptive fields (neurons see more of input)



Some Types of Linear Transformations



Types of Pooling functions

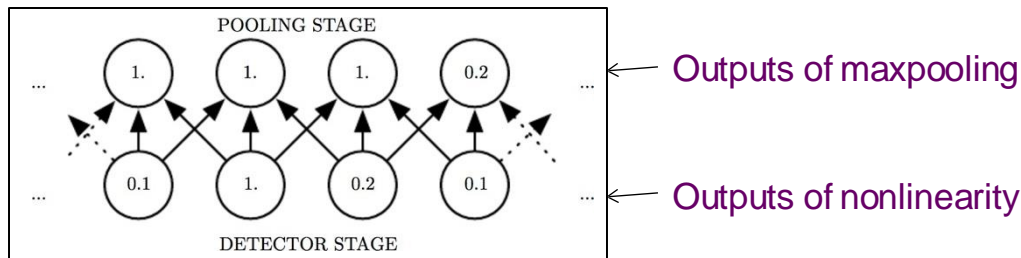
- A pooling function replaces the output of the net at a certain location with a summary statistic of the nearby inputs
- Popular pooling functions:
 1. *max pooling* operation reports the maximum output within a rectangular neighborhood
 2. Average of a rectangular neighborhood
 3. L^2 norm of a rectangular neighborhood
 4. Weighted average based on the distance from the central pixel

Pooling Causes Translation Invariance

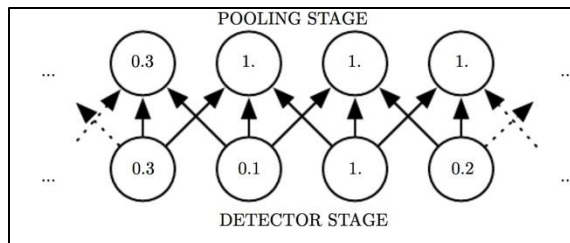
- In all cases, pooling helps make the representation become approximately *invariant* to small translations of the input
 - If we translate the input by a small amount values of most of the outputs does not change (example next slide)
 - Pooling can be viewed as adding a strong prior that the function the layer learns must be invariant to small translations

Max Pooling Introduces Invariance to Translation

- View of middle of output of a convolutional layer



- Same network after the input has been shifted by one pixel



- Every input value has changed, but only half the values of output have changed because maxpooling units are only sensitive to maximum value in neighborhood not exact value

Importance of Translation Invariance

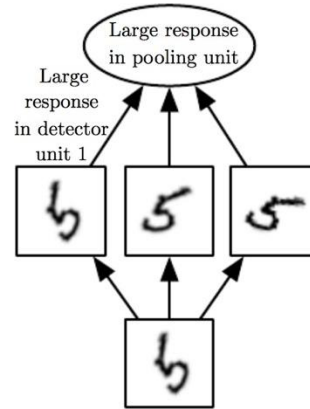
- Invariance to translation is important if we care about whether a feature is present rather than exactly where it is
 - Ex: for detecting a face just need to know that an eye is present in a region, not its exact location
- In other context it is more important to preserve location of a feature
 - Ex: to determine a corner we need to know whether two edges are present and test whether they meet

Learning other Invariances

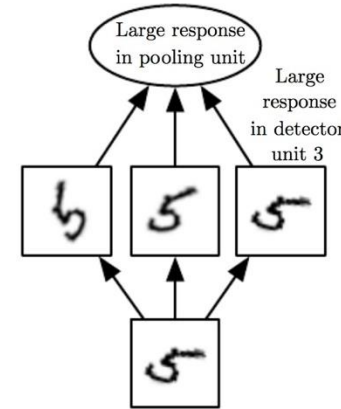
- Pooling over spatial regions produces invariance to translation
- But if we pool over the results of separately parameterized convolutions, the features can learn which transformations to become invariant to

Learning Invariance to Rotation

- A pooling unit that pools over multiple features that are learned with separate parameters can learn to be invariant to transformations of the input



Input tilted left
gets large response
from unit tuned to
left-tilted images



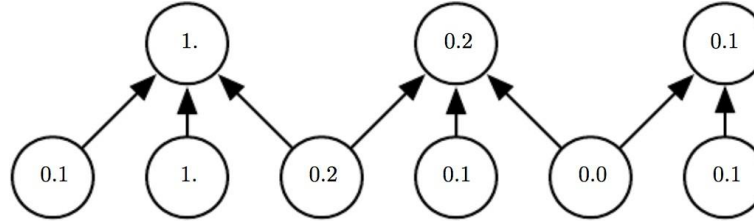
Tilted Right

Using Fewer Pooling Units than Detector Units

- Because pooling summarizes the responses over a whole neighborhood, it is possible to use fewer pooling units than detector units
 - This is due to the effect of reporting summary statistics for pooling regions spaced k pixels apart rather than one pixel apart
 - This improves computational efficiency because the next layer has k times fewer inputs to process
- An example is given next

Pooling with Downsampling

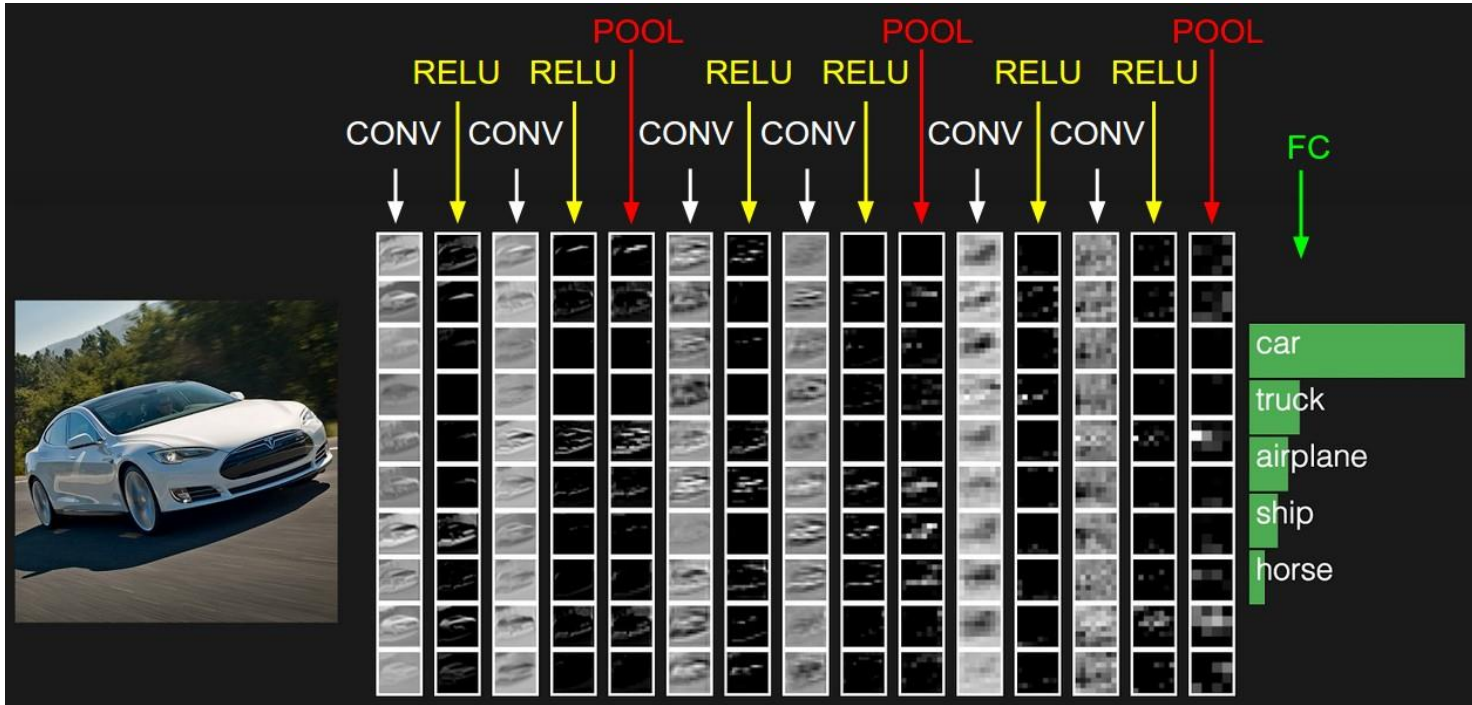
- Max-pooling with a pool width of three and a stride between pools of two



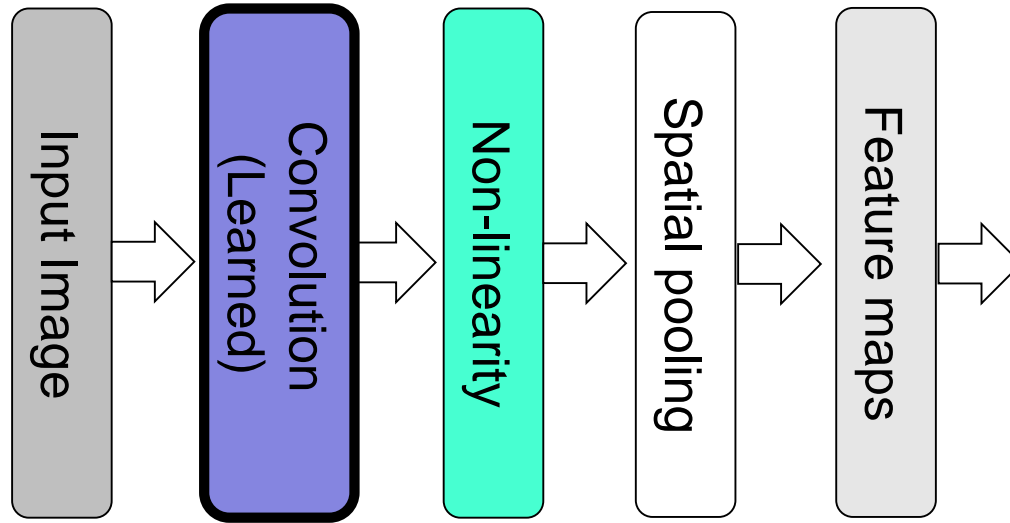
- This reduces representation size by a factor of two
 - Which reduces computational burden of next layer
 - Rightmost pooling region has a smaller size but must be included if we don't want to ignore some of the detector units

Downsampling makes a digital audio signal smaller by lowering its sampling rate or sample size (bits per sample). It decreases the bit rate when transmitting over a limited bandwidth or to convert to a more limited audio format

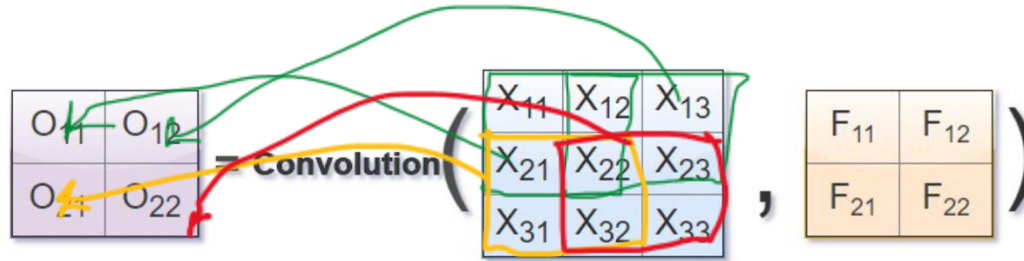
Feature Map (Convolution Layer)



CNN Components



Backprop in Convolution



$$O_{11} = F_{11}X_{11} + F_{12}X_{12} + F_{21}X_{21} + F_{22}X_{22}$$

$$O_{12} = F_{11}X_{12} + F_{12}X_{13} + F_{21}X_{22} + F_{22}X_{23}$$

$$O_{21} = F_{11}X_{21} + F_{12}X_{22} + F_{21}X_{31} + F_{22}X_{32}$$

$$O_{22} = F_{11}X_{22} + F_{12}X_{23} + F_{21}X_{32} + F_{22}X_{33}$$

Backprop in Convolution

$$\begin{bmatrix} O_{11} & O_{12} \\ O_{21} & O_{22} \end{bmatrix} = \text{Convolution} \left(\begin{bmatrix} X_{11} & X_{12} & X_{13} \\ X_{21} & X_{22} & X_{23} \\ X_{31} & X_{32} & X_{33} \end{bmatrix}, \begin{bmatrix} F_{11} & F_{12} \\ F_{21} & F_{22} \end{bmatrix} \right)$$

$$O_{11} = F_{11}X_{11} + F_{12}X_{12} + F_{21}X_{21} + F_{22}X_{22}$$

$$O_{12} = F_{11}X_{12} + F_{12}X_{13} + F_{21}X_{22} + F_{22}X_{23}$$

$$O_{21} = F_{11}X_{21} + F_{12}X_{22} + F_{21}X_{31} + F_{22}X_{32}$$

$$O_{22} = F_{11}X_{22} + F_{12}X_{23} + F_{21}X_{32} + F_{22}X_{33}$$

$$\frac{\partial E}{\partial F_{11}} = \frac{\partial E}{\partial O_{11}} \frac{\partial O_{11}}{\partial F_{11}} + \frac{\partial E}{\partial O_{12}} \frac{\partial O_{12}}{\partial F_{11}} + \frac{\partial E}{\partial O_{21}} \frac{\partial O_{21}}{\partial F_{11}} + \frac{\partial E}{\partial O_{22}} \frac{\partial O_{22}}{\partial F_{11}}$$

$$\frac{\partial E}{\partial F_{12}} = \frac{\partial E}{\partial O_{11}} \frac{\partial O_{11}}{\partial F_{12}} + \frac{\partial E}{\partial O_{12}} \frac{\partial O_{12}}{\partial F_{12}} + \frac{\partial E}{\partial O_{21}} \frac{\partial O_{21}}{\partial F_{12}} + \frac{\partial E}{\partial O_{22}} \frac{\partial O_{22}}{\partial F_{12}}$$

$$\frac{\partial E}{\partial F_{21}} = \frac{\partial E}{\partial O_{11}} \frac{\partial O_{11}}{\partial F_{21}} + \frac{\partial E}{\partial O_{12}} \frac{\partial O_{12}}{\partial F_{21}} + \frac{\partial E}{\partial O_{21}} \frac{\partial O_{21}}{\partial F_{21}} + \frac{\partial E}{\partial O_{22}} \frac{\partial O_{22}}{\partial F_{21}}$$

$$\frac{\partial E}{\partial F_{22}} = \frac{\partial E}{\partial O_{11}} \frac{\partial O_{11}}{\partial F_{22}} + \frac{\partial E}{\partial O_{12}} \frac{\partial O_{12}}{\partial F_{22}} + \frac{\partial E}{\partial O_{21}} \frac{\partial O_{21}}{\partial F_{22}} + \frac{\partial E}{\partial O_{22}} \frac{\partial O_{22}}{\partial F_{22}}$$

Backprop in Convolution

Which evaluates to-

$$O_{11} = F_{11}X_{11} + F_{12}X_{12} + F_{21}X_{21} + F_{22}X_{22}$$

$$O_{12} = F_{11}X_{12} + F_{12}X_{13} + F_{21}X_{22} + F_{22}X_{23}$$

$$O_{21} = F_{11}X_{21} + F_{12}X_{22} + F_{21}X_{31} + F_{22}X_{32}$$

$$O_{22} = F_{11}X_{22} + F_{12}X_{23} + F_{21}X_{32} + F_{22}X_{33}$$

$$\begin{bmatrix} O_{11} & O_{12} \\ O_{21} & O_{22} \end{bmatrix} = \text{Convolution} \left(\begin{bmatrix} X_{11} & X_{12} & X_{13} \\ X_{21} & X_{22} & X_{23} \\ X_{31} & X_{32} & X_{33} \end{bmatrix}, \begin{bmatrix} F_{11} & F_{12} \\ F_{21} & F_{22} \end{bmatrix} \right)$$

$$\frac{\partial E}{\partial F_{11}} = \frac{\partial E}{\partial O_{11}} X_{11} + \frac{\partial E}{\partial O_{12}} X_{12} + \frac{\partial E}{\partial O_{21}} X_{21} + \frac{\partial E}{\partial O_{22}} X_{22}$$

$$\frac{\partial E}{\partial F_{12}} = \frac{\partial E}{\partial O_{11}} X_{12} + \frac{\partial E}{\partial O_{12}} X_{13} + \frac{\partial E}{\partial O_{21}} X_{22} + \frac{\partial E}{\partial O_{22}} X_{23}$$

$$\frac{\partial E}{\partial F_{21}} = \frac{\partial E}{\partial O_{11}} X_{21} + \frac{\partial E}{\partial O_{12}} X_{22} + \frac{\partial E}{\partial O_{21}} X_{31} + \frac{\partial E}{\partial O_{22}} X_{32}$$

$$\frac{\partial E}{\partial F_{22}} = \frac{\partial E}{\partial O_{11}} X_{22} + \frac{\partial E}{\partial O_{12}} X_{23} + \frac{\partial E}{\partial O_{21}} X_{32} + \frac{\partial E}{\partial O_{22}} X_{33}$$

Backprop in Convolution

If we look closely the previous equation can be written in form of our convolution operation.

$$\begin{array}{|c|c|} \hline \partial E / \partial F_{11} & \partial E / \partial F_{12} \\ \hline \partial E / \partial F_{21} & \partial E / \partial F_{22} \\ \hline \end{array} = \text{Convolution} \left(\begin{array}{|c|c|c|} \hline X_{11} & X_{12} & X_{13} \\ \hline X_{21} & X_{22} & X_{23} \\ \hline X_{31} & X_{32} & X_{33} \\ \hline \end{array} , \begin{array}{|c|c|} \hline \partial E / \partial O_{11} & \partial E / \partial O_{12} \\ \hline \partial E / \partial O_{21} & \partial E / \partial O_{22} \\ \hline \end{array} \right)$$

$$\frac{\partial E}{\partial F_{11}} = \frac{\partial E}{\partial O_{11}} X_{11} + \frac{\partial E}{\partial O_{12}} X_{12} + \frac{\partial E}{\partial O_{21}} X_{21} + \frac{\partial E}{\partial O_{22}} X_{22}$$

$$\frac{\partial E}{\partial F_{12}} = \frac{\partial E}{\partial O_{11}} X_{12} + \frac{\partial E}{\partial O_{12}} X_{13} + \frac{\partial E}{\partial O_{21}} X_{22} + \frac{\partial E}{\partial O_{22}} X_{23}$$

$$\frac{\partial E}{\partial F_{21}} = \frac{\partial E}{\partial O_{11}} X_{21} + \frac{\partial E}{\partial O_{12}} X_{22} + \frac{\partial E}{\partial O_{21}} X_{31} + \frac{\partial E}{\partial O_{22}} X_{32}$$

$$\frac{\partial E}{\partial F_{22}} = \frac{\partial E}{\partial O_{11}} X_{22} + \frac{\partial E}{\partial O_{12}} X_{23} + \frac{\partial E}{\partial O_{21}} X_{32} + \frac{\partial E}{\partial O_{22}} X_{33}$$

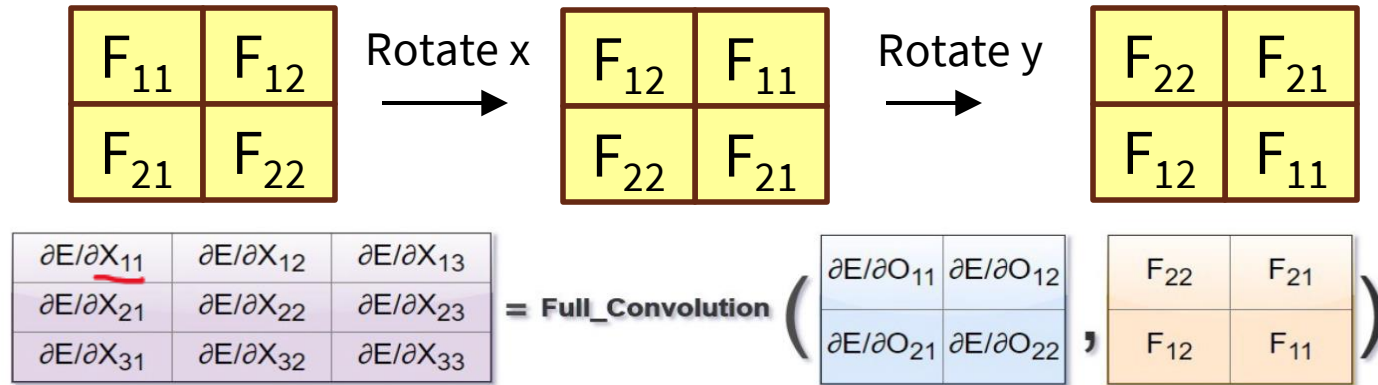
Backprop in Convolution

Similarly we can find the gradients of the error 'E' with respect to the input matrix 'X'.

$$\begin{aligned}\frac{\partial E}{\partial X_{11}} &= \frac{\partial E}{\partial O_{11}} F_{11} + \frac{\partial E}{\partial O_{12}} 0 + \frac{\partial E}{\partial O_{21}} 0 + \frac{\partial E}{\partial O_{22}} 0 \\ \frac{\partial E}{\partial X_{12}} &= \frac{\partial E}{\partial O_{11}} F_{12} + \frac{\partial E}{\partial O_{12}} F_{11} + \frac{\partial E}{\partial O_{21}} 0 + \frac{\partial E}{\partial O_{22}} 0 \\ \frac{\partial E}{\partial X_{13}} &= \frac{\partial E}{\partial O_{11}} 0 + \frac{\partial E}{\partial O_{12}} F_{12} + \frac{\partial E}{\partial O_{21}} 0 + \frac{\partial E}{\partial O_{22}} 0 \\ \frac{\partial E}{\partial X_{21}} &= \frac{\partial E}{\partial O_{11}} F_{21} + \frac{\partial E}{\partial O_{12}} 0 + \frac{\partial E}{\partial O_{21}} F_{11} + \frac{\partial E}{\partial O_{22}} 0 \\ \frac{\partial E}{\partial X_{22}} &= \frac{\partial E}{\partial O_{11}} F_{22} + \frac{\partial E}{\partial O_{12}} F_{21} + \frac{\partial E}{\partial O_{21}} f_{12} + \frac{\partial E}{\partial O_{22}} F_{11} \\ \frac{\partial E}{\partial X_{23}} &= \frac{\partial E}{\partial O_{11}} 0 + \frac{\partial E}{\partial O_{12}} F_{22} + \frac{\partial E}{\partial O_{21}} 0 + \frac{\partial E}{\partial O_{22}} F_{11} \\ \frac{\partial E}{\partial X_{31}} &= \frac{\partial E}{\partial O_{11}} 0 + \frac{\partial E}{\partial O_{12}} 0 + \frac{\partial E}{\partial O_{21}} F_{21} + \frac{\partial E}{\partial O_{22}} 0 \\ \frac{\partial E}{\partial X_{32}} &= \frac{\partial E}{\partial O_{11}} 0 + \frac{\partial E}{\partial O_{12}} 0 + \frac{\partial E}{\partial O_{21}} F_{22} + \frac{\partial E}{\partial O_{22}} F_{21} \\ \frac{\partial E}{\partial X_{33}} &= \frac{\partial E}{\partial O_{11}} 0 + \frac{\partial E}{\partial O_{12}} 0 + \frac{\partial E}{\partial O_{21}} 0 + \frac{\partial E}{\partial O_{22}} F_{22}\end{aligned}$$

The previous computation can be obtained by a different type of convolution operation known as full convolution.

In order to obtain the gradients of the input matrix we need to rotate the filter by 180 degree and calculate the full convolution of the rotated filter by the gradients of the output with respect to error.



F_{22}	F_{21}	
F_{12}	$F_{11}\delta O_{11}$	δO_{12}
	δO_{21}	δO_{22}

$$\delta X_{11} = F_{11}\delta O_{11}$$

F_{22}	F_{21}	
$F_{12}\delta O_{11}$	$F_{11}\delta O_{12}$	
δO_{21}	δO_{22}	

$$\delta X_{12} = F_{12}\delta O_{11} + F_{11}\delta O_{12}$$

	F_{22}	F_{21}
δO_{11}	$F_{12}\delta O_{12}$	F_{11}
δO_{21}	δO_{22}	

$$\delta X_{13} = F_{12}\delta O_{12}$$

F_{22}	$F_{21}\delta O_{11}$	δO_{12}
F_{12}	$F_{11}\delta O_{21}$	δO_{22}

$$\delta X_{21} = F_{21}\delta O_{11} + F_{11}\delta O_{21}$$

$F_{22}\delta O_{11}$	$F_{21}\delta O_{12}$
$F_{12}\delta O_{21}$	$F_{11}\delta O_{22}$

$$\delta X_{22} = F_{22}\delta O_{11} + F_{21}\delta O_{12} + F_{12}\delta O_{21} + F_{11}\delta O_{22}$$

δO_{11}	$F_{22}\delta O_{12}$	F_{21}
δO_{21}	$F_{12}\delta O_{22}$	F_{11}

$$\delta X_{23} = F_{22}\delta O_{12} + F_{12}\delta O_{22}$$

		δO_{11}	δO_{12}
F_{22}	$F_{21}\delta O_{21}$	δO_{22}	
F_{12}	F_{11}		

$$\delta X_{31} = F_{21}\delta O_{21}$$

δO_{11}	δO_{12}
$F_{22}\delta O_{21}$	$F_{21}\delta O_{22}$
F_{12}	F_{11}

$$\delta X_{32} = F_{22}\delta O_{21} + F_{21}\delta O_{22}$$

δO_{11}	δO_{12}	
δO_{21}	$F_{22}\delta O_{22}$	F_{21}
	F_{12}	F_{11}

$$\delta X_{33} = F_{22}\delta O_{22}$$

Backpropagation of max pooling

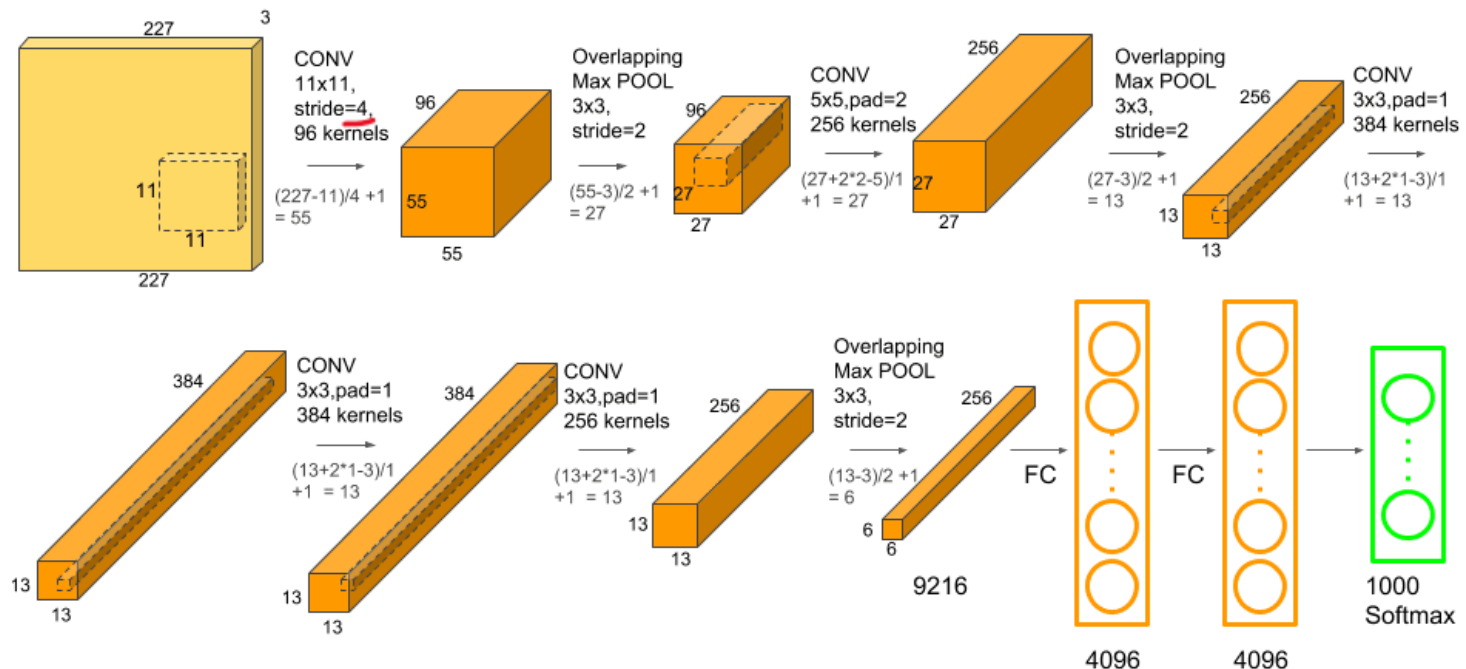
Suppose you have a matrix M of four elements:

a	b
c	d

and $\text{maxpool}(M)$ returns d .

Then, the maxpool function really only depends on d . So, the derivative of maxpool relative to d is 1, and its derivative relative to a, b, c is zero. So you backpropagate 1 to the unit corresponding to d , and you backpropagate zero for the other units.

AlexNet



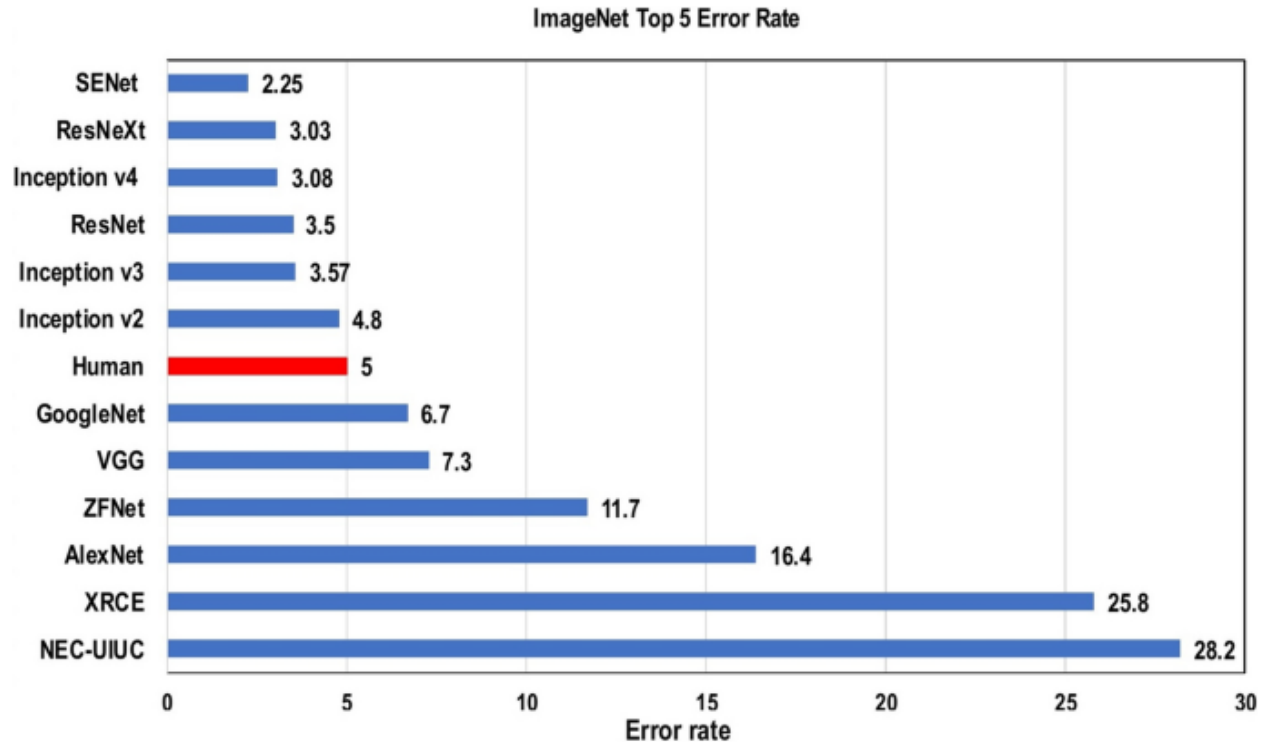
Alex Krizhevsky, ImageNet Classification with Deep Convolutional Neural Networks, 2012



Carnegie Mellon University

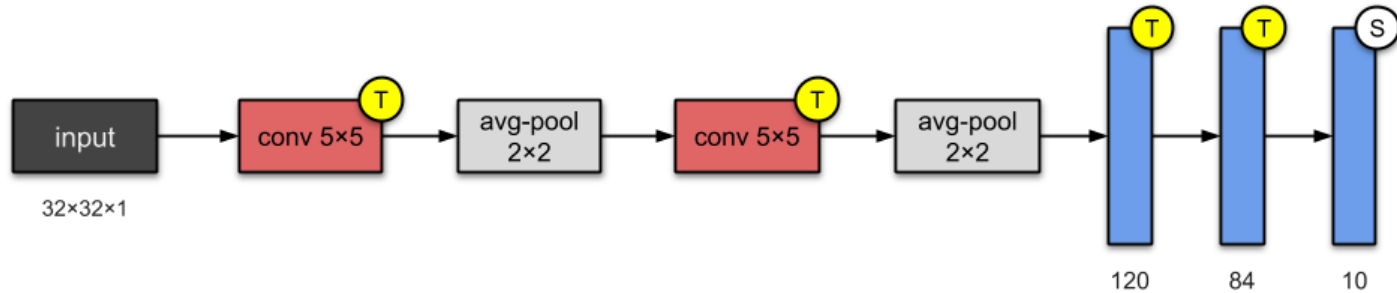
Different CNN Architectures

Notable CNN models¹⁹



Notable CNN models

LENET (1998)⁸



(2) ENDING THE NETWORK WITH ONE OR MORE FULLY-CONNECTED LAYERS

LeNet-5

Gradient Based Learning Applied To Document Recognition

- Y. Lecun, L. Bottou, Y. Bengio, P. Haffner; 1998

Helped establish how we use CNNs today

Replaced manual feature

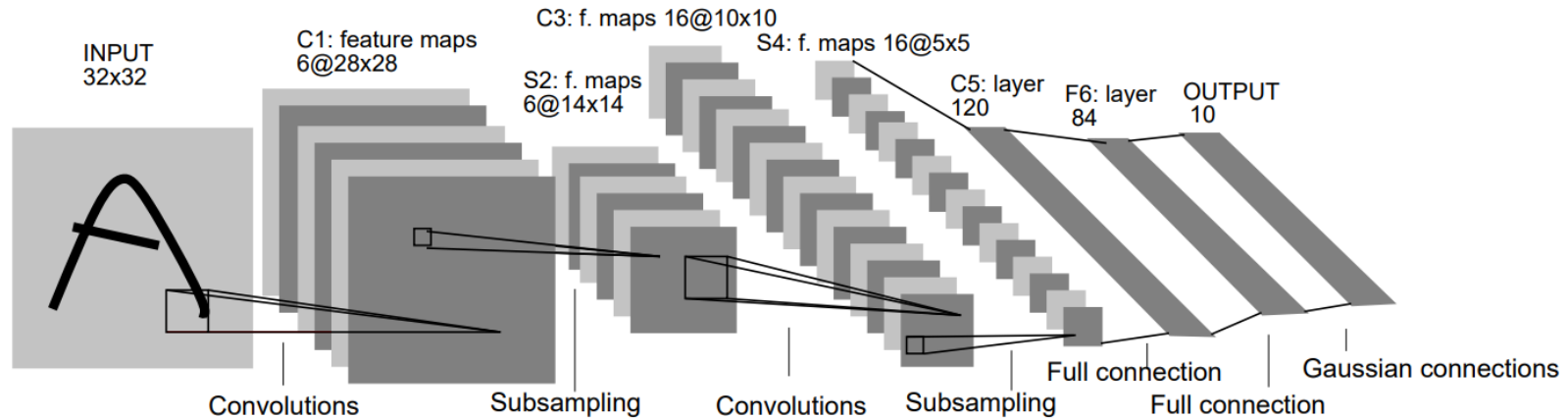
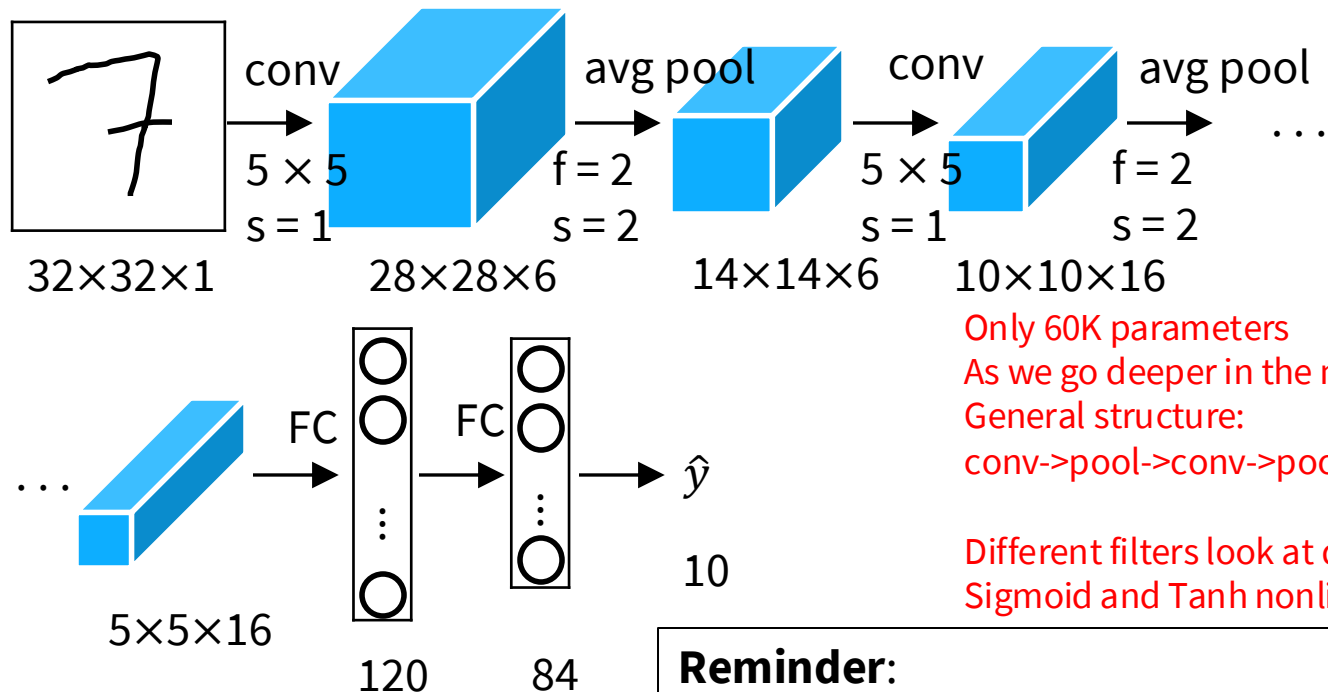


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

LeNet-5



Only 60K parameters

As we go deeper in the network: $N_H \downarrow$, $N_W \downarrow$, $N_C \uparrow$

General structure:

conv->pool->conv->pool->FC->FC->output

Different filters look at different channels

Sigmoid and Tanh nonlinearity

Reminder:

Output size = $(N+2P-F)/\text{stride} + 1$

AlexNet

ImageNet Classification with Deep Convolutional Neural Networks

- Alex Krizhevsky, Ilya Sutskever, Geoffrey E. Hinton; 2012

Facilitated by GPUs, highly optimized convolution implementation and large datasets (ImageNet)

One of the largest CNNs to date

Has 60 Million parameter compared to 60k parameter of LeNet-5

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners

The annual “Olympics” of computer vision.

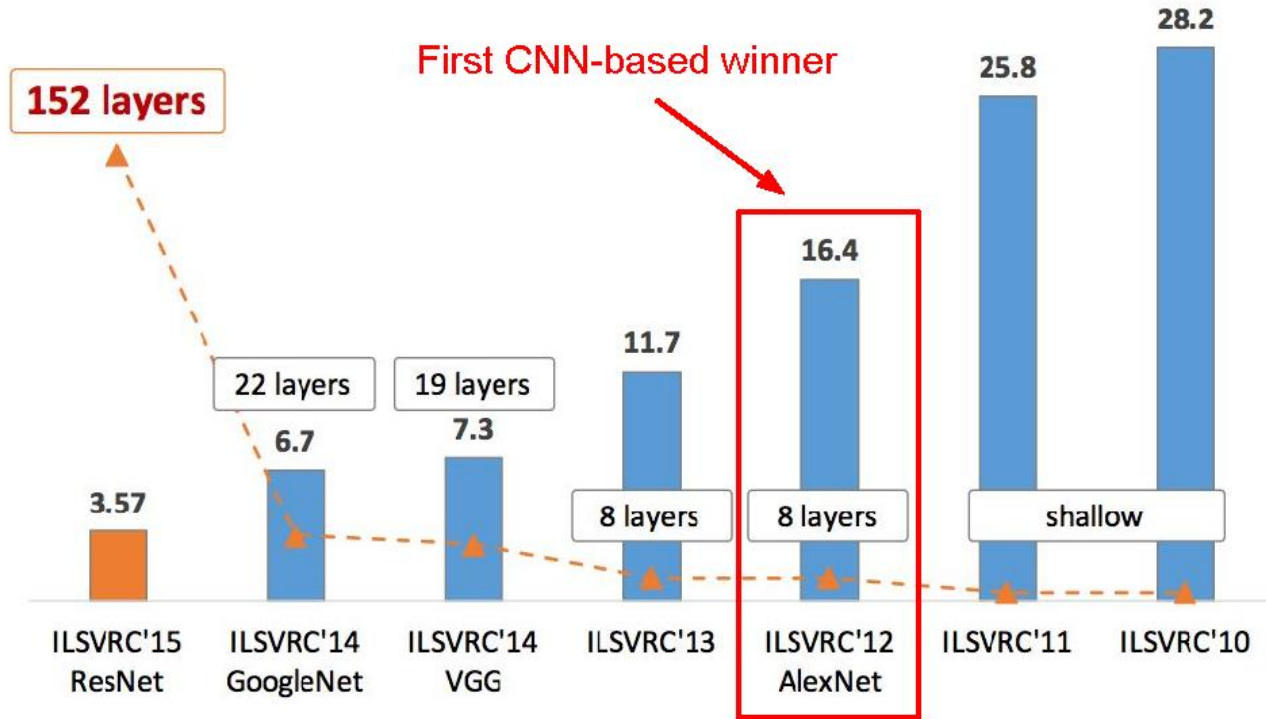
Teams from across the world compete to see who has the best computer vision model for tasks such as classification, localization, detection, and more.

2012 marked **the first year where a CNN was used** to achieve a top 5 test error rate of 15.3%.

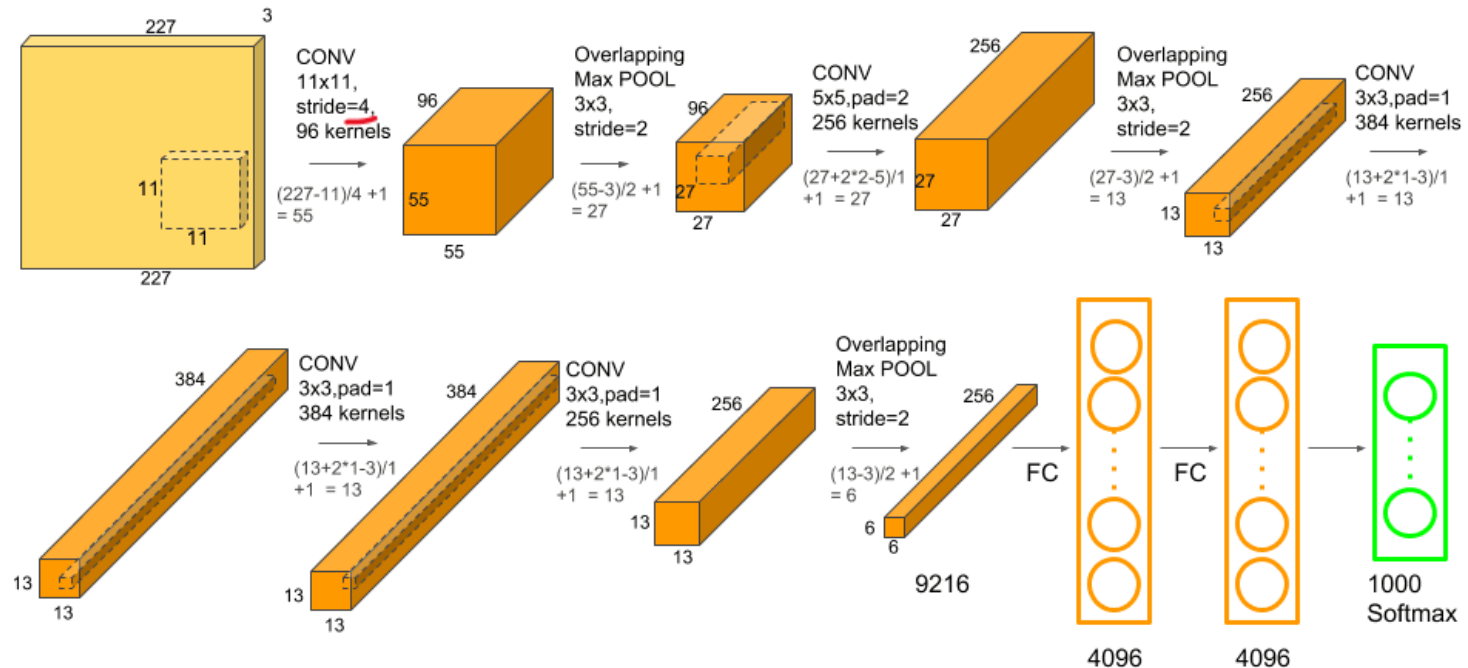
The next best entry achieved an error of 26.2%.



ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



AlexNet



Alex Krizhevsky, ImageNet Classification with Deep Convolutional Neural Networks, 2012

Architecture

CONV1

MAX POOL1

NORM1

CONV2

MAX POOL2

NORM2

CONV3

CONV4

CONV5

Max POOL3

FC6

FC7

FC8

INPUT: 227x227x3 IMAGES (224x224
BEFORE PADDING)

FIRST LAYER: 96 11x11 FILTERS APPLIED AT
STRIDE 4

OUTPUT VOLUME SIZE?

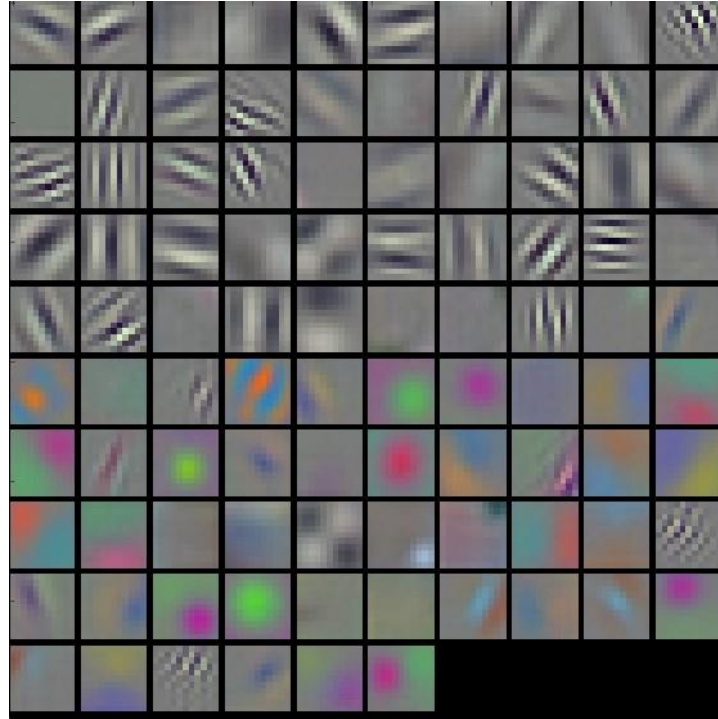
$$(N-F)/s+1 = (227-11)/4+1 = 55 \rightarrow$$

[55x55x96]

NUMBER OF PARAMETERS IN THIS LAYER?

$$(11*11*3)*96 = 35K$$

AlexNet filters



AlexNet

Architecture

CONV1

MAX POOL1

NORM1

CONV2

MAX POOL2

NORM2

CONV3

CONV4

CONV5

Max POOL3

FC6

FC7

FC8

INPUT: 227x227x3 IMAGES (224x224 BEFORE
PADDING)

AFTER CONV1: 55x55x96

SECOND LAYER: 3x3 FILTERS APPLIED AT
STRIDE 2

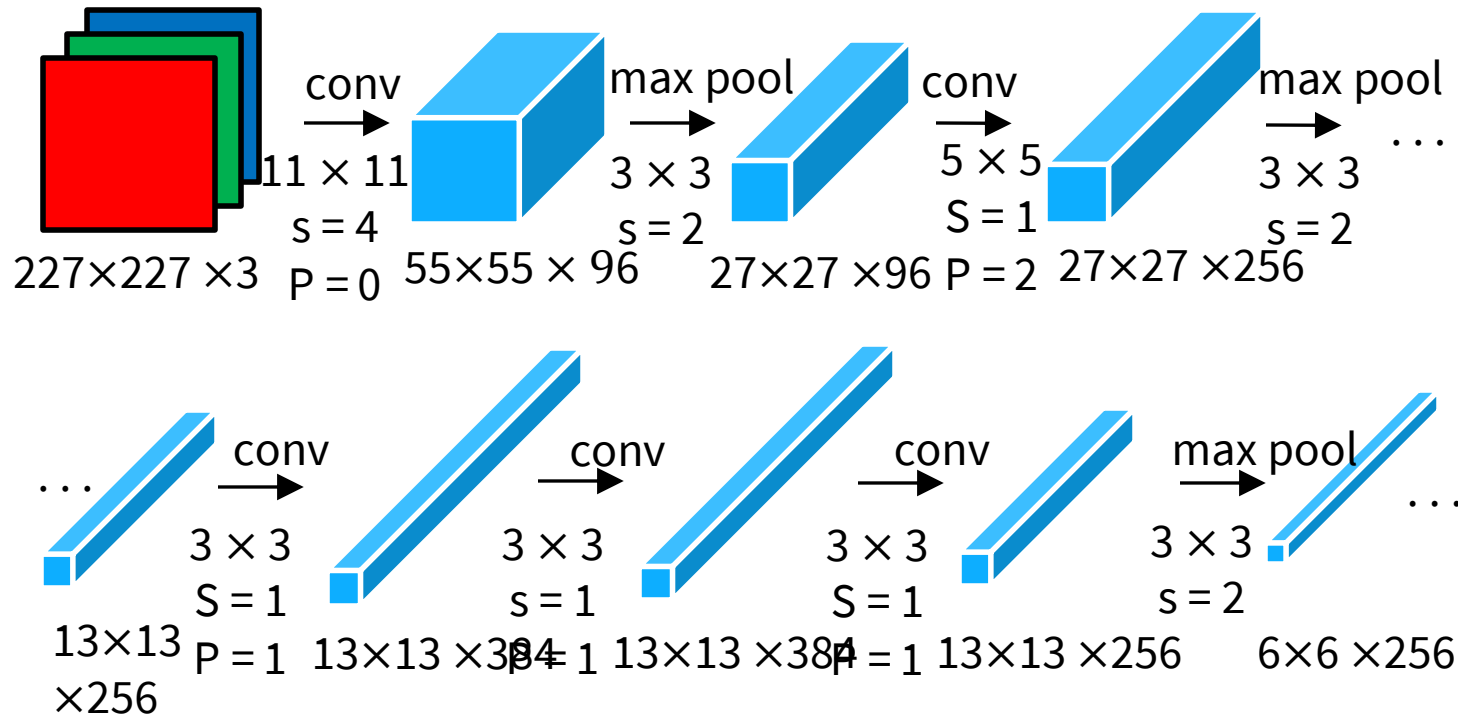
OUTPUT VOLUME SIZE?

$(N-F)/s+1 = (55-3)/2+1 = 27 \rightarrow$
[27x27x96]

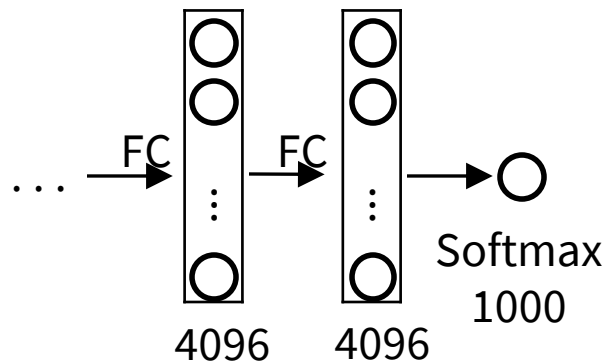
NUMBER OF PARAMETERS IN THIS LAYER?

0!

AlexNet



AlexNet



DETAILS/RETROSPECTIVES:

FIRST USE OF RELU

USED NORM LAYERS (NOT COMMON ANYMORE)

HEAVY DATA AUGMENTATION

DROPOUT 0.5

BATCH SIZE 128

7 CNN ENSEMBLE

TRAINED ON GTX 580 GPU WITH ONLY 3 GB OF MEMORY.

NETWORK SPREAD ACROSS 2 GPUS, HALF THE NEURONS
(FEATURE MAPS) ON EACH GPU.

CONV1, CONV2, CONV4, CONV5:

CONNECTIONS ONLY WITH FEATURE MAPS ON SAME GPU.

CONV3, FC6, FC7, FC8:

CONNECTIONS WITH ALL FEATURE MAPS IN PRECEDING
LAYER, COMMUNICATION ACROSS GPUS.

AlexNet

ALEXNET WAS THE COMING OUT PARTY FOR CNNs IN THE COMPUTER VISION COMMUNITY. THIS WAS **THE FIRST TIME A MODEL PERFORMED SO WELL ON A HISTORICALLY DIFFICULT IMAGENET DATASET**. THIS PAPER ILLUSTRATED THE BENEFITS OF CNNs AND BACKED THEM UP WITH RECORD BREAKING PERFORMANCE IN THE COMPETITION.

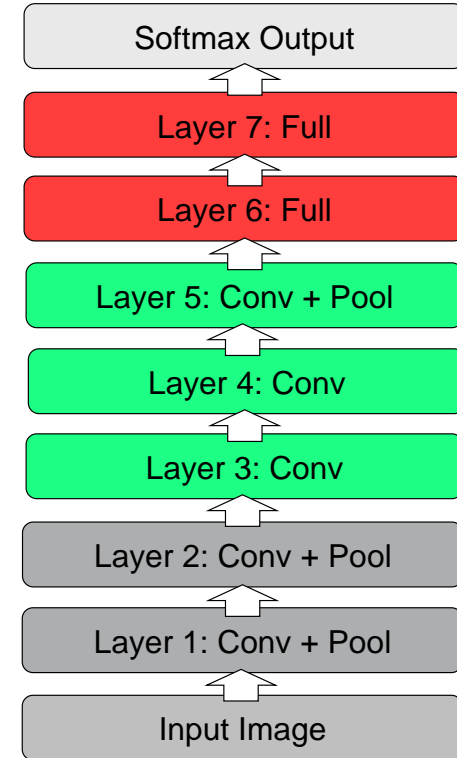
AlexNet

Model	Top-1	Top-5
<i><u>Sparse coding</u> [2]</i>	47.1%	28.2%
<i><u>SIFT</u> + <u>FVs</u> [24]</i>	45.7%	25.7%
CNN	<u>37.5%</u>	<u>17.0%</u>

Alex Krizhevsky, ImageNet Classification with Deep Convolutional Neural Networks, 2012

AlexNet

- „ 8 layers total!
- „ Trained on Imagenet dataset [Deng et al. CVPR'09]!
- „ 18.2% top-5 error !

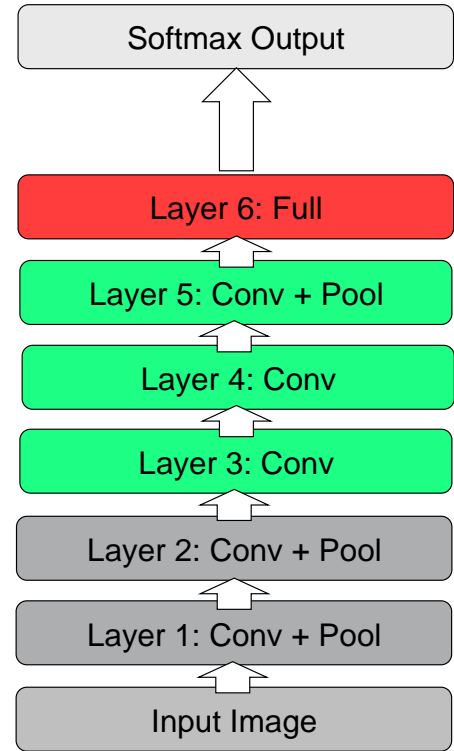


[From Rob Fergus' CIFAR 2016 tutorial] Krizhevsky et al., NIPS 2012

AlexNet

AlexNet

- „ Remove top fully connected layer 7 !
- „ Drop ~16 million parameters!
- „ Only 1.1% drop in performance!!

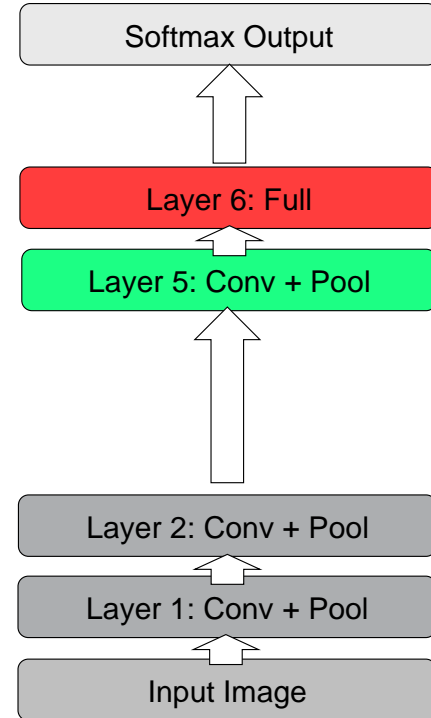


[From Rob Fergus' CIFAR 2016 tutorial] Krizhevsky et al., NIPS 2012

AlexNet

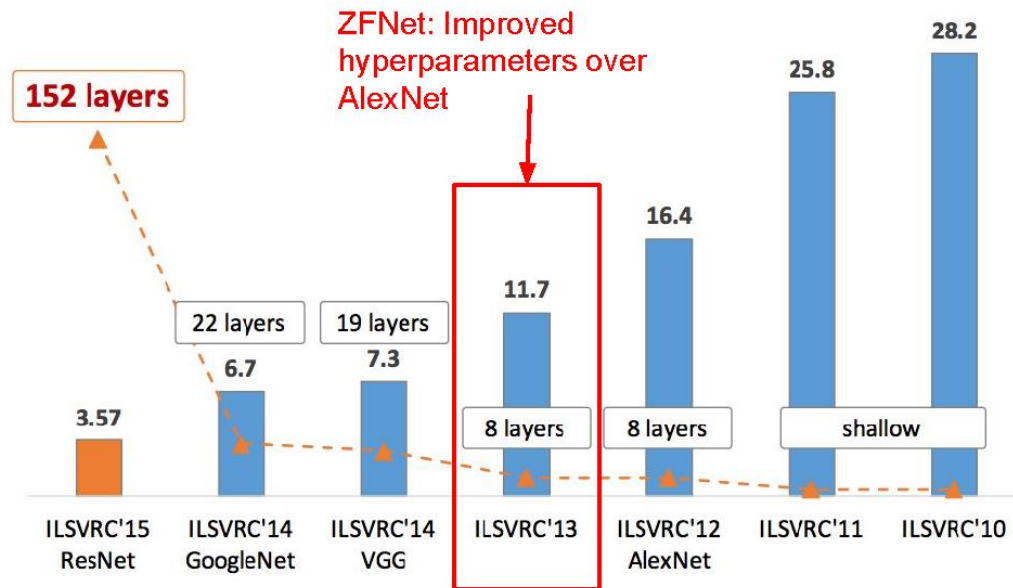
AlexNet

- „ Remove layers 3 4,6 and 7 !
- „ Drop ~50 million parameters!
- „ **33.5%** drop in performance!!
- „ Depth of the network is the key!

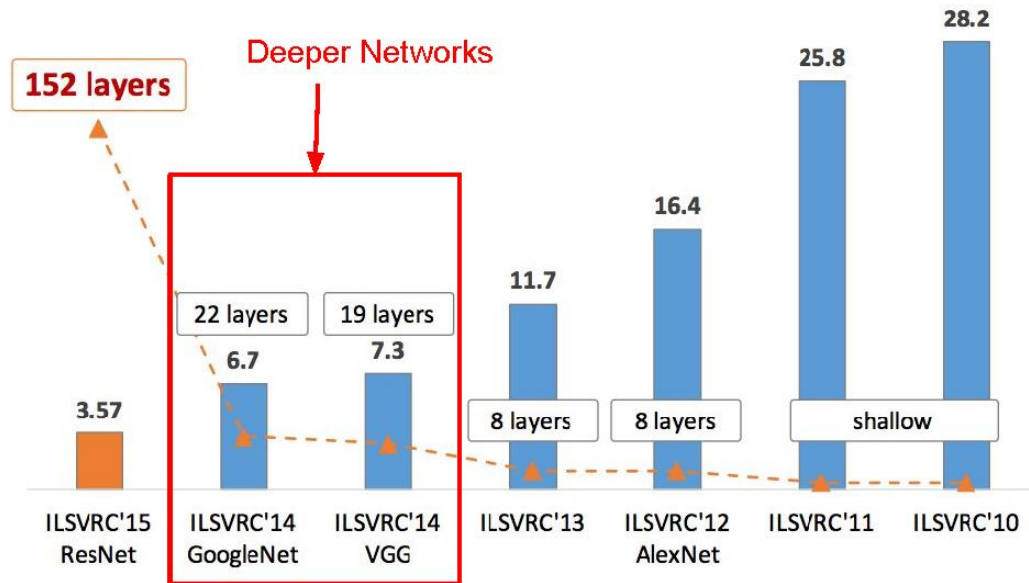


[From Rob Fergus' CIFAR 2016 tutorial] Krizhevsky et al., NIPS 2012

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



VGGNet

*VERY DEEP CONVOLUTIONAL NETWORKS FOR LARGE SCALE
IMAGE RECOGNITION - KAREN SIMONYAN AND ANDREW
ZISSEMAN; 2015*

THE RUNNER-UP AT THE ILSVRC 2014 COMPETITION
SIGNIFICANTLY DEEPER THAN ALEXNET
140 MILLION PARAMETERS

Input

3x3 conv, 64

3x3 conv, 64

Pool 1/2

3x3 conv, 128

3x3 conv, 128

Pool 1/2

3x3 conv, 256

3x3 conv, 256

Pool 1/2

3x3 conv, 512

3x3 conv, 512

3x3 conv, 512

Pool 1/2

3x3 conv, 512

3x3 conv, 512

3x3 conv, 512

Pool 1/2

FC 4096

FC 4096

FC 1000

Softmax

VGGNet

SMALLER FILTERS

ONLY 3X3 CONV FILTERS, STRIDE 1, PAD 1
AND 2X2 MAX POOL, STRIDE 2

DEEPER NETWORK

ALEXNET: 8 LAYERS

VGGNET: 16 - 19 LAYERS

ZFNET: 11.7% TOP 5 ERROR IN ILSVRC'13

VGGNET: 7.3% TOP 5 ERROR IN ILSVRC'14

WHY USE SMALLER FILTERS? (3X3 CONV)

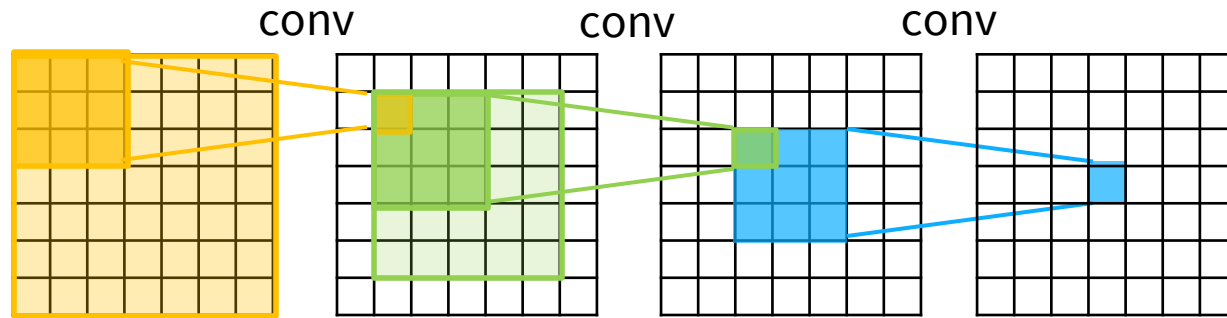
STACK OF THREE 3X3 CONV (STRIDE 1) LAYERS HAS THE SAME EFFECTIVE RECEPTIVE FIELD AS ONE 7X7 CONV LAYER.

WHAT IS THE EFFECTIVE RECEPTIVE FIELD OF THREE 3X3 CONV (STRIDE 1) LAYERS? 7x7

BUT DEEPER, MORE NON-LINEARITIES

AND FEWER PARAMETERS: $3 * (3^2 C^2)$ VS. $7^2 C^2$ FOR C CHANNELS PER LAYER

Reminder: Receptive Field



Input	memory: $224*224*3=150\text{K}$	params: 0
3x3 conv, 64	memory: $224*224*64=3.2\text{M}$	params: $(3*3*3)*64 = 1,728$
3x3 conv, 64	memory: $224*224*64=3.2\text{M}$	params: $(3*3*64)*64 = \underline{36,864}$
Pool	memory: $112*112*64=800\text{K}$	params: 0
3x3 conv, 128	memory: $112*112*128=1.6\text{M}$	params: $(3*3*64)*128 = \underline{73,728}$
3x3 conv, 128	memory: $112*112*128=1.6\text{M}$	params: $(3*3*128)*128 = 147,456$
Pool	memory: $56*56*128=400\text{K}$	params: 0
3x3 conv, 256	memory: $56*56*256=800\text{K}$	params: $(3*3*128)*256 = 294,912$
3x3 conv, 256	memory: $56*56*256=800\text{K}$	params: $(3*3*256)*256 = 589,824$
3x3 conv, 256	memory: $56*56*256=800\text{K}$	params: $(3*3*256)*256 = 589,824$
Pool	memory: $28*28*256=200\text{K}$	params: 0
3x3 conv, 512	memory: $28*28*512=400\text{K}$	params: $(3*3*256)*512 = 1,179,648$
3x3 conv, 512	memory: $28*28*512=400\text{K}$	params: $(3*3*512)*512 = 2,359,296$
3x3 conv, 512	memory: $28*28*512=400\text{K}$	params: $(3*3*512)*512 = 2,359,296$
Pool	memory: $14*14*512=100\text{K}$	params: 0
3x3 conv, 512	memory: $14*14*512=100\text{K}$	params: $(3*3*512)*512 = 2,359,296$
3x3 conv, 512	memory: $14*14*512=100\text{K}$	params: $(3*3*512)*512 = 2,359,296$
3x3 conv, 512	memory: $14*14*512=100\text{K}$	params: $(3*3*512)*512 = 2,359,296$
Pool	memory: $7*7*512=25\text{K}$	params: 0
FC 4096	memory: 4096	params: $7*7*512*4096 = \underline{102,760,448}$
FC 4096	memory: 4096	params: $4096*4096 = \underline{16,777,216}$
FC 1000	memory: 1000	params: $4096*1000 = 4,096,000$

VGGNet

Input

3x3 conv, 64

3x3 conv, 64

Pool

3x3 conv, 128

3x3 conv, 128

Pool

3x3 conv, 256

3x3 conv, 256

3x3 conv, 256

Pool

3x3 conv, 512

3x3 conv, 512

3x3 conv, 512

Pool

3x3 conv, 512

3x3 conv, 512

3x3 conv, 512

Pool

FC 4096

FC 4096

FC 1000

Softmax

VGG16:

TOTAL MEMORY: $24M * 4 \text{ BYTES} \approx 96MB$ / IMAGE

TOTAL PARAMS: 138M PARAMETERS

Input	memory: $224*224*3=150\text{K}$	params: 0
3x3 conv, 64	memory: $224*224*64=3.2\text{M}$	params: $(3*3*3)*64 = 1,728$
3x3 conv, 64	memory: $224*224*64=3.2\text{M}$	params: $(3*3*64)*64 = 36,864$
Pool	memory: $112*112*64=800\text{K}$	params: 0
3x3 conv, 128	memory: $112*112*128=1.6\text{M}$	params: $(3*3*64)*128 = 73,728$
3x3 conv, 128	memory: $112*112*128=1.6\text{M}$	params: $(3*3*128)*128 = 147,456$
Pool	memory: $56*56*128=400\text{K}$	params: 0
3x3 conv, 256	memory: $56*56*256=800\text{K}$	params: $(3*3*128)*256 = 294,912$
3x3 conv, 256	memory: $56*56*256=800\text{K}$	params: $(3*3*256)*256 = 589,824$
3x3 conv, 256	memory: $56*56*256=800\text{K}$	params: $(3*3*256)*256 = 589,824$
Pool	memory: $28*28*256=200\text{K}$	params: 0
3x3 conv, 512	memory: $28*28*512=400\text{K}$	params: $(3*3*256)*512 = 1,179,648$
3x3 conv, 512	memory: $28*28*512=400\text{K}$	params: $(3*3*512)*512 = 2,359,296$
3x3 conv, 512	memory: $28*28*512=400\text{K}$	params: $(3*3*512)*512 = 2,359,296$
Pool	memory: $14*14*512=100\text{K}$	params: 0
3x3 conv, 512	memory: $14*14*512=100\text{K}$	params: $(3*3*512)*512 = 2,359,296$
3x3 conv, 512	memory: $14*14*512=100\text{K}$	params: $(3*3*512)*512 = 2,359,296$
3x3 conv, 512	memory: $14*14*512=100\text{K}$	params: $(3*3*512)*512 = 2,359,296$
Pool	memory: $7*7*512=25\text{K}$	params: 0
FC 4096	memory: 4096	params: $7*7*512*4096 = 102,760,448$
FC 4096	memory: 4096	params: $4096*4096 = 16,777,216$
FC 1000	memory: 1000	params: $4096*1000 = 4,096,000$

DETAILS/RETROSPECTIVES :

ILSVRC'14 2ND IN CLASSIFICATION, 1ST IN LOCALIZATION

SIMILAR TRAINING PROCEDURE AS ALEXNET

NO LOCAL RESPONSE NORMALISATION (LRN)

USE VGG16 OR VGG19 (VGG19 ONLY SLIGHTLY BETTER,
MORE MEMORY)

USE ENSEMBLES FOR BEST RESULTS

FC7 FEATURES GENERALIZE WELL TO OTHER TASKS

TRAINED ON 4 NVIDIA TITAN BLACK GPUS FOR **TWO TO THREE
WEEKS.**

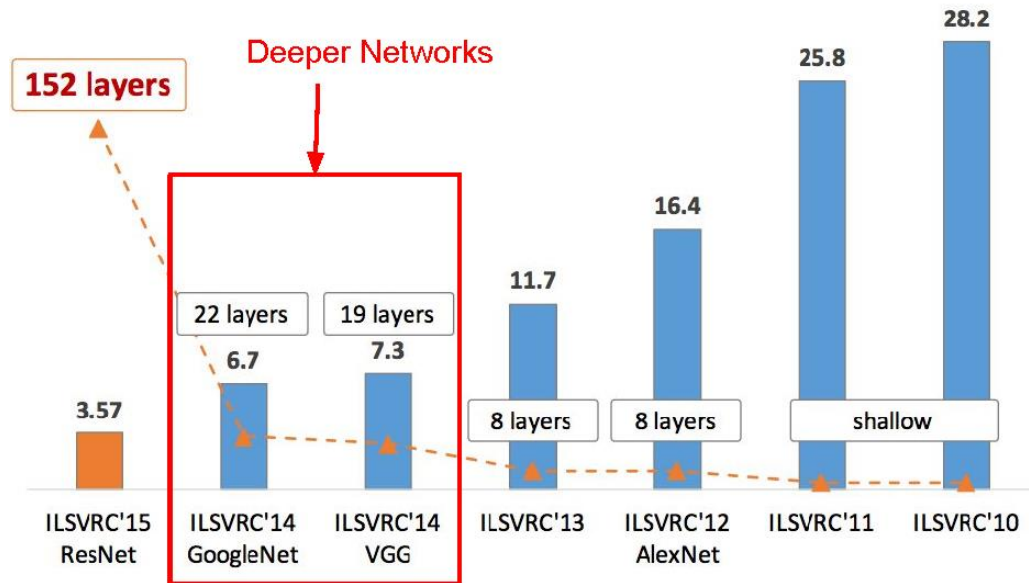


VGG NET REINFORCED THE NOTION
THAT **CONVOLUTIONAL NEURAL NETWORKS HAVE TO HAVE
A DEEP NETWORK OF LAYERS IN ORDER FOR THIS
HIERARCHICAL REPRESENTATION OF VISUAL DATA TO
WORK.**

KEEP IT DEEP.

KEEP IT SIMPLE.

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



GoogleNet

GOING DEEPER WITH CONVOLUTIONS - CHRISTIAN SZEGEDY ET AL.; 2015

ILSVRC 2014 COMPETITION WINNER

ALSO SIGNIFICANTLY DEEPER THAN ALEXNET

X12 LESS PARAMETERS THAN ALEXNET

FOCUSED ON COMPUTATIONAL EFFICIENCY

GoogleNet

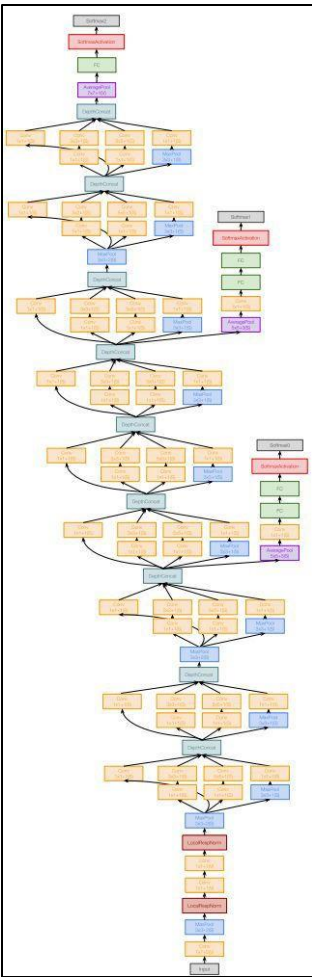
22 LAYERS

EFFICIENT “**INCEPTION**” MODULE - STRAYED FROM
THE GENERAL APPROACH OF SIMPLY STACKING
CONV AND POOLING LAYERS ON TOP OF EACH
OTHER IN A SEQUENTIAL STRUCTURE

NO FC LAYERS

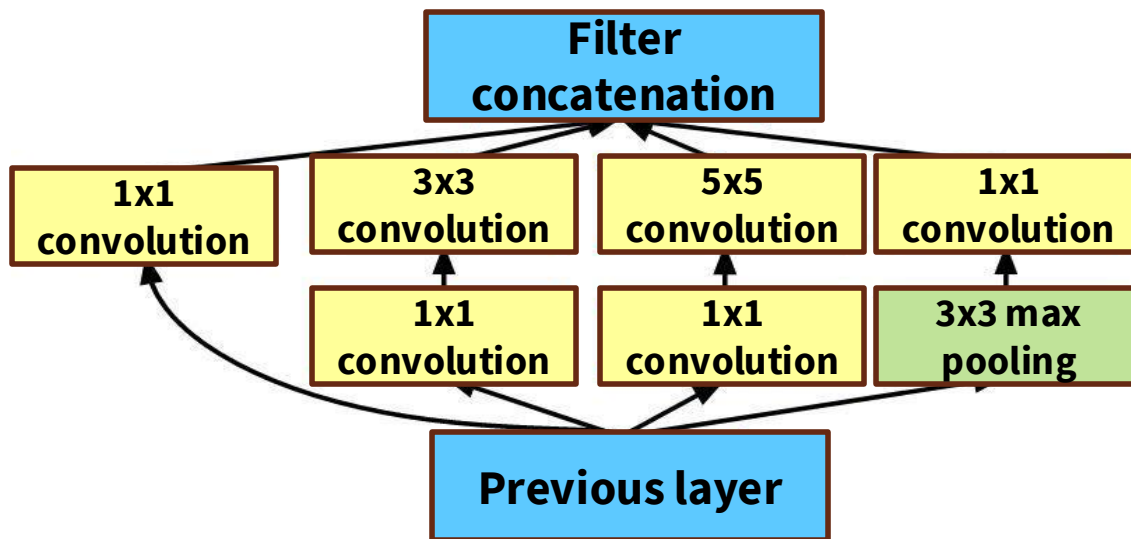
ONLY 5 MILLION PARAMETERS!

ILSVRC'14 CLASSIFICATION WINNER (6.7% TOP 5
ERROR)



GoogleNet

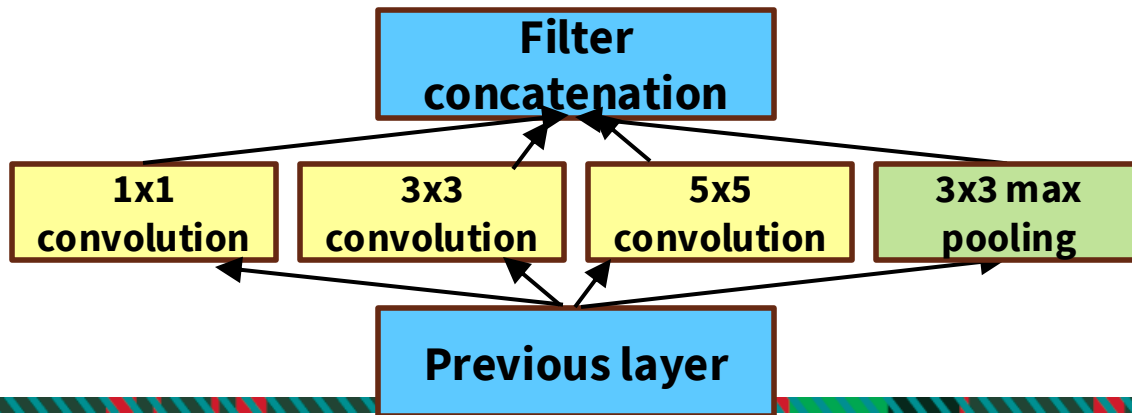
“INCEPTION MODULE”: DESIGN A GOOD LOCAL NETWORK TOPOLOGY (NETWORK WITHIN A NETWORK) AND THEN STACK THESE MODULES ON TOP OF EACH OTHER



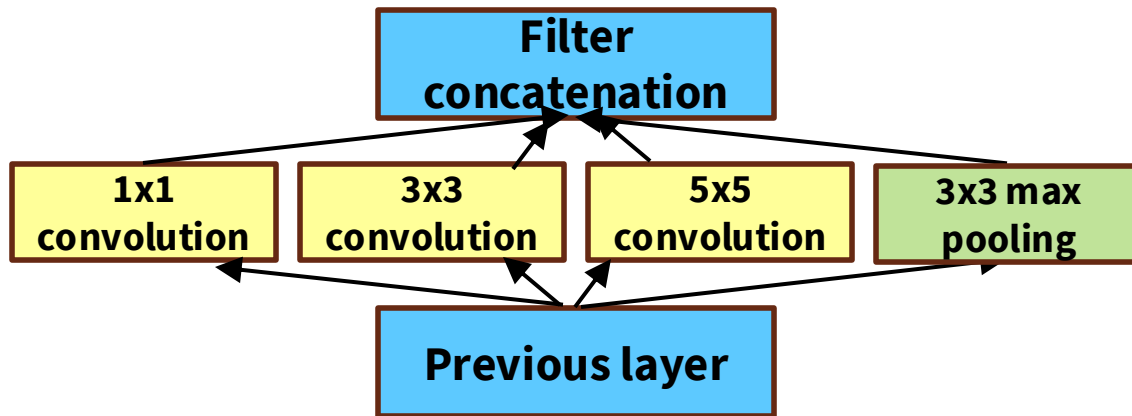
GoogleNet

Naïve Inception Model

- Apply parallel filter operations on the input :
 - › Multiple receptive field sizes for convolution (1x1, 3x3, 5x5)
 - › Pooling operation (3x3)
- Concatenate all filter outputs together depth-wise



- What's the problem with this?
High computational complexity



OUTPUT VOLUME SIZES:

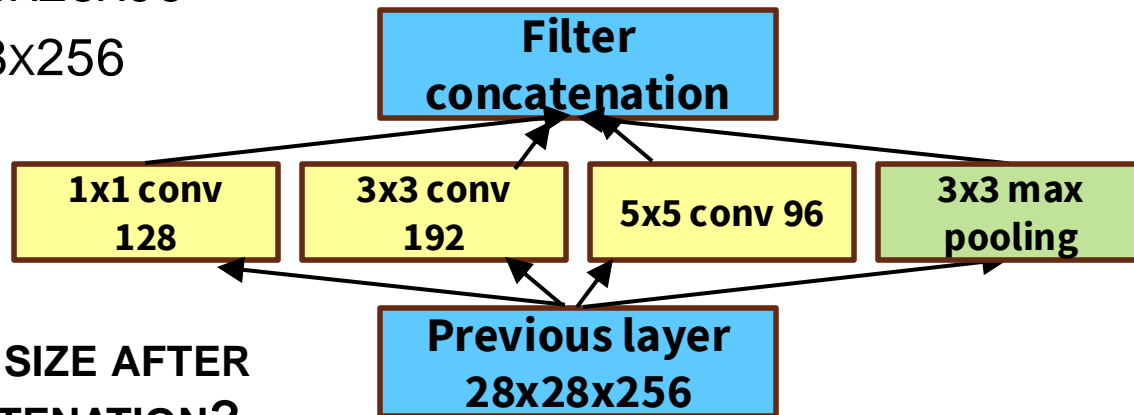
1X1 CONV, 128: 28X28X128

3X3 CONV, 192: 28X28X192

5X5 CONV, 96: 28X28X96

3X3 POOL: 28X28X256

Example:



WHAT IS OUTPUT SIZE AFTER
FILTER CONCATENATION?

$$28 \times 28 \times (128 + 192 + 96 + 256) = 28 \times 28 \times 672$$

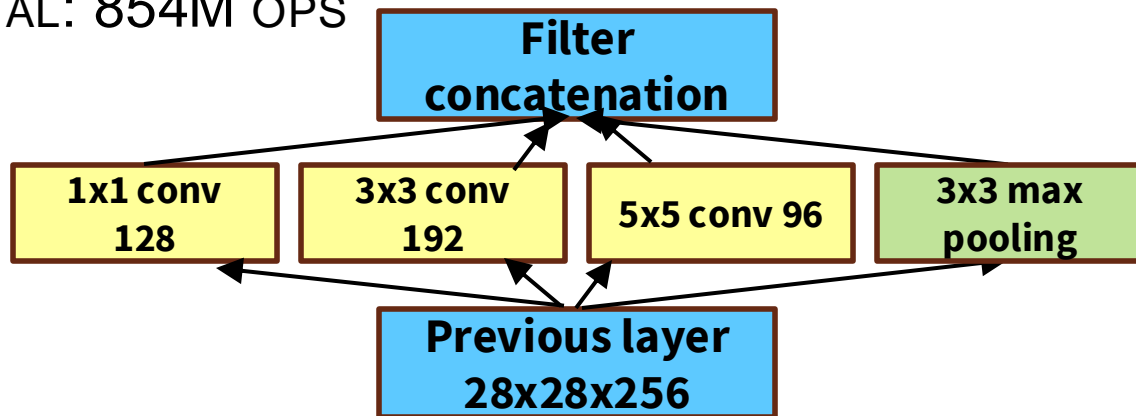
NUMBER OF CONVOLUTION OPERATIONS:

1x1 CONV, 128: $28 \times 28 \times 128 \times 1 \times 1 \times 256$

3x3 CONV, 192: $28 \times 28 \times 192 \times 3 \times 3 \times 256$

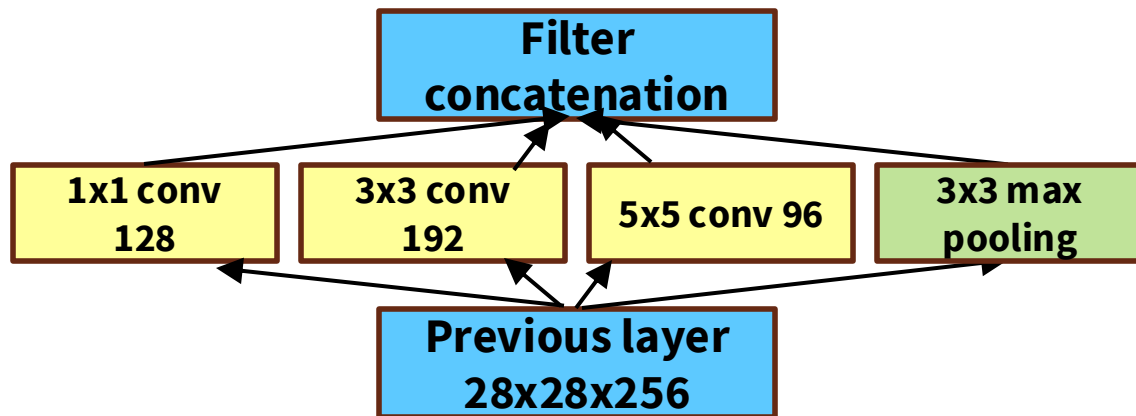
5x5 CONV, 96: $28 \times 28 \times 96 \times 5 \times 5 \times 256$

TOTAL: 854M OPS

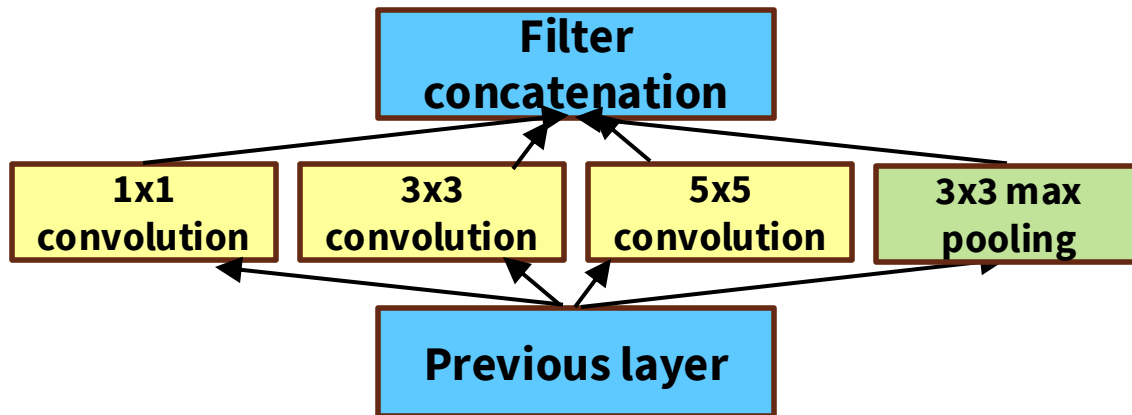


VERY EXPENSIVE COMPUTE!

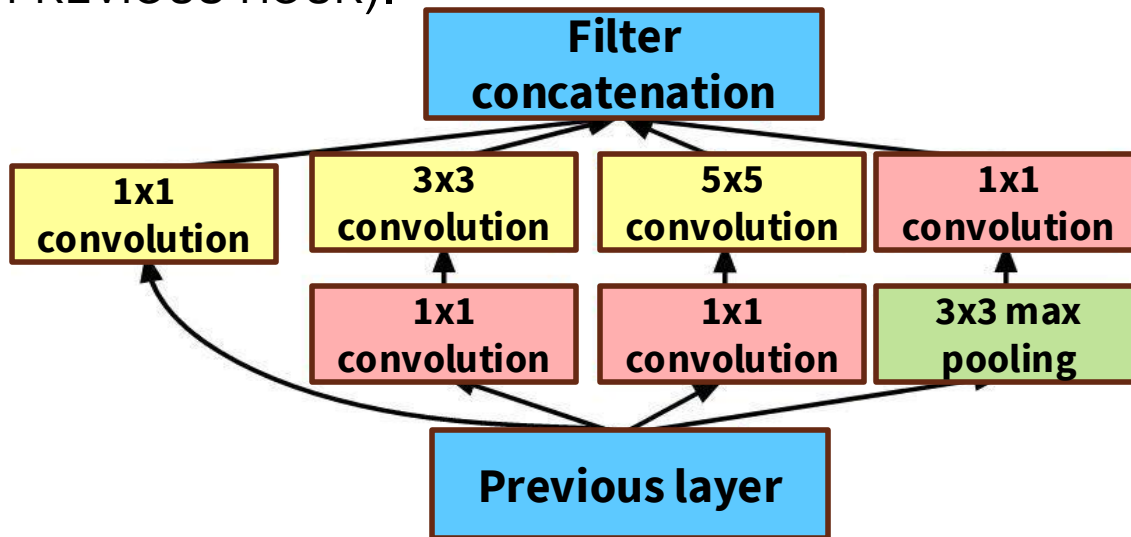
POOLING LAYER ALSO PRESERVES FEATURE DEPTH, WHICH MEANS TOTAL DEPTH AFTER CONCATENATION CAN ONLY GROW AT EVERY LAYER.



SOLUTION: “BOTTLENECK” LAYERS THAT USE 1X1 CONVOLUTIONS TO REDUCE FEATURE DEPTH (FROM PREVIOUS HOUR).



SOLUTION: “BOTTLENECK” LAYERS THAT USE 1X1 CONVOLUTIONS TO REDUCE FEATURE DEPTH (FROM PREVIOUS HOUR).



NUMBER OF CONVOLUTION OPERATIONS:

1x1 CONV, 64: $28 \times 28 \times 64 \times 1 \times 1 \times 256$

1x1 CONV, 64: $28 \times 28 \times 64 \times 1 \times 1 \times 256$

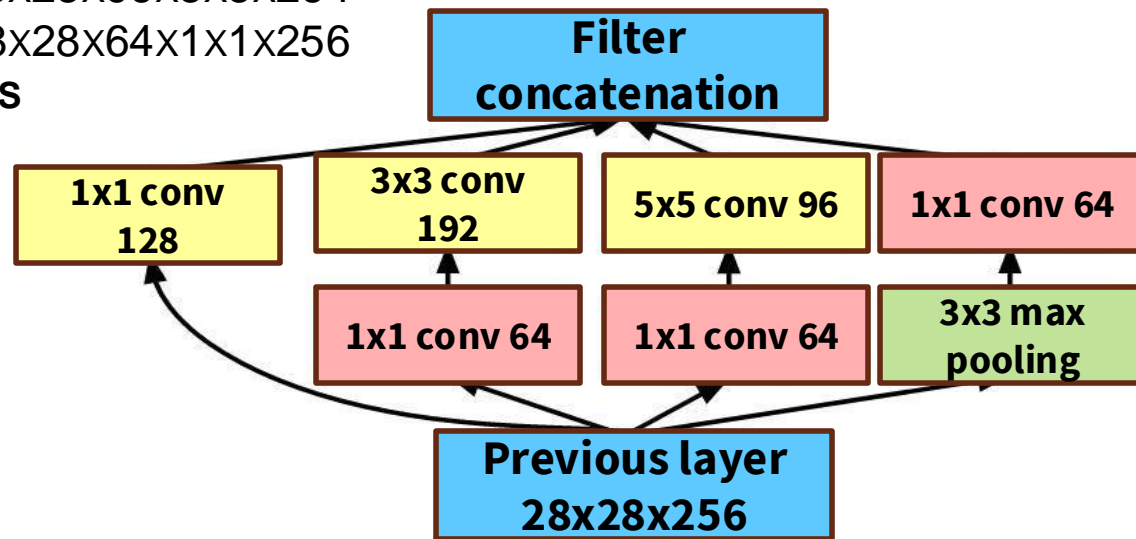
1x1 CONV, 128: $28 \times 28 \times 128 \times 1 \times 1 \times 256$

3x3 CONV, 192: $28 \times 28 \times 192 \times 3 \times 3 \times 64$

5x5 CONV, 96: $28 \times 28 \times 96 \times 5 \times 5 \times 264$

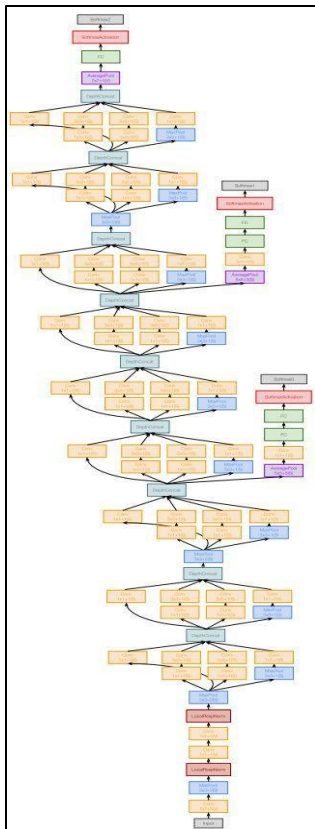
1x1 CONV, 64: $28 \times 28 \times 64 \times 1 \times 1 \times 256$

TOTAL: 353M OPS



COMPARED TO 854M OPS FOR NAIVE VERSION

GoogleNet

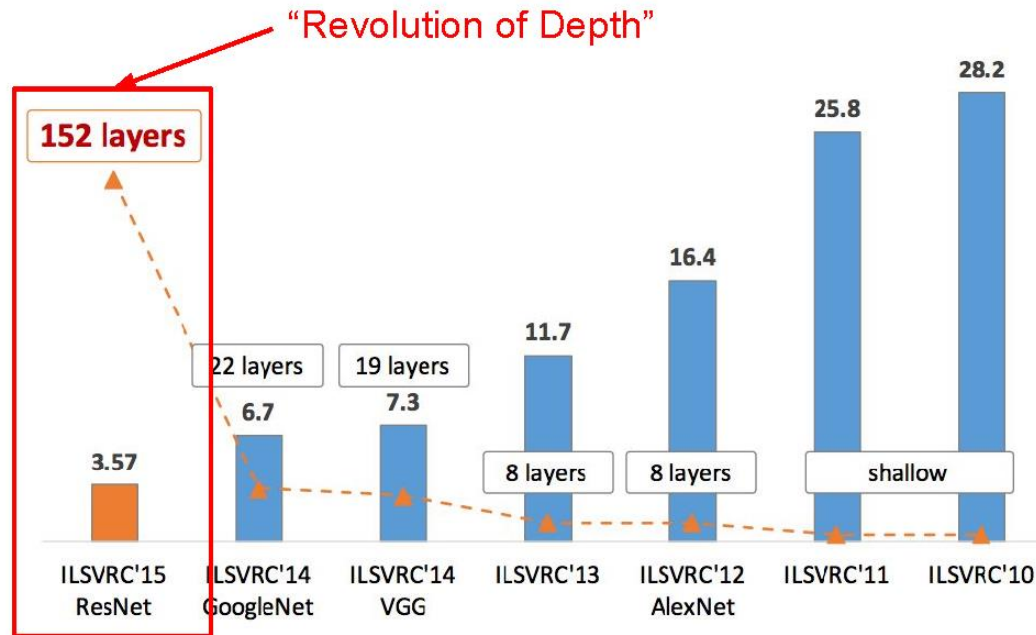


Details/Retrospectives :

- Deeper networks, with computational efficiency
- 22 layers
- Efficient “Inception” module
- No FC layers
- 12x less params than AlexNet
- ILSVRC’14 classification winner (6.7% top 5 error)

INTRODUCED THE IDEA THAT CNN LAYERS **DIDN'T ALWAYS HAVE TO BE STACKED UP SEQUENTIALLY**. COMING UP WITH THE INCEPTION MODULE, THE AUTHORS SHOWED THAT A CREATIVE STRUCTURING OF LAYERS CAN LEAD TO IMPROVED PERFORMANCE AND **COMPUTATIONALLY EFFICIENCY**.

ImageNet Large Scale Visual Recognition Challenge (ILSVRC) winners



ResNet

DEEP RESIDUAL LEARNING FOR IMAGE RECOGNITION
- KAIMING HE, XIANGYU ZHANG, SHAOQING REN,
JIAN SUN; 2015

EXTREMELY DEEP NETWORK – 152 LAYERS

DEEPER NEURAL NETWORKS ARE MORE DIFFICULT
TO TRAIN.

DEEP NETWORKS SUFFER FROM VANISHING AND
EXPLODING GRADIENTS.

PRESENT A RESIDUAL LEARNING FRAMEWORK TO
EASE THE TRAINING OF NETWORKS THAT ARE
SUBSTANTIALLY DEEPER THAN THOSE USED
PREVIOUSLY.

The diagram illustrates the VGG-16 architecture, showing the sequence of layers from input to output. The layers are color-coded and connected by arrows indicating the data flow. The architecture consists of the following layers (from bottom to top):

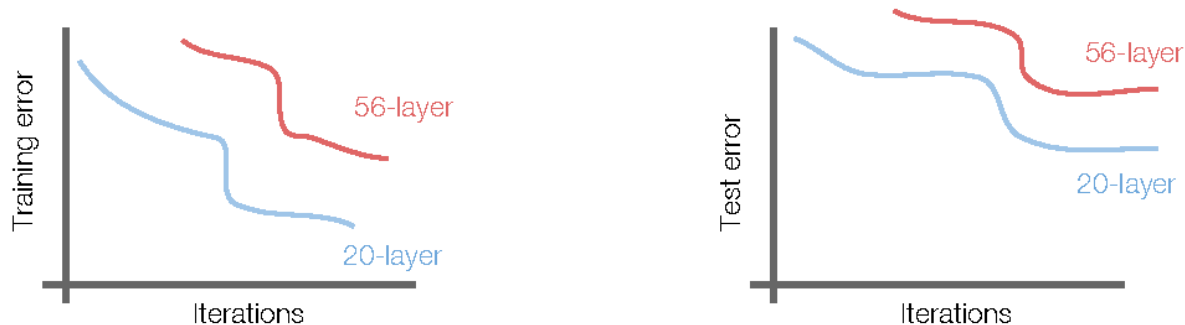
- Input** (Grey box)
- 7x7 conv. 64 / 2** (Orange box)
- Pool** (Blue box)
- 3x3 conv. 64** (Maroon box)
- 3x3 conv. 64** (Maroon box)
- 3x3 conv. 64** (Maroon box)
- 3x3 conv. 64** (Maroon box)
- 3x3 conv. 128 / 2** (Purple box)
- 3x3 conv. 128** (Purple box)
- 3x3 conv. 128** (Purple box)
- 3x3 conv. 128** (Purple box)
- 3x3 conv. 128** (Purple box)
- ...** (Ellipsis indicating more layers)
- 3x3 conv. 64** (Green box)
- 3x3 conv. 64** (Green box)
- 3x3 conv. 64** (Green box)
- 3x3 conv. 64** (Green box)
- Pool** (Blue box)
- FC 1000** (Green box)
- Softmax** (Red box)

The diagram shows that the architecture is composed of two main stages of convolutional layers, each followed by a pooling layer. The first stage consists of two blocks of three 3x3 convolutional layers (64 filters) followed by a pooling layer. The second stage consists of two blocks of three 3x3 convolutional layers (128 filters) followed by a pooling layer. The final output is a fully connected layer (FC 1000) followed by a Softmax layer for classification.

-

ResNet

- What happens when we continue stacking deeper layers on a convolutional neural network?



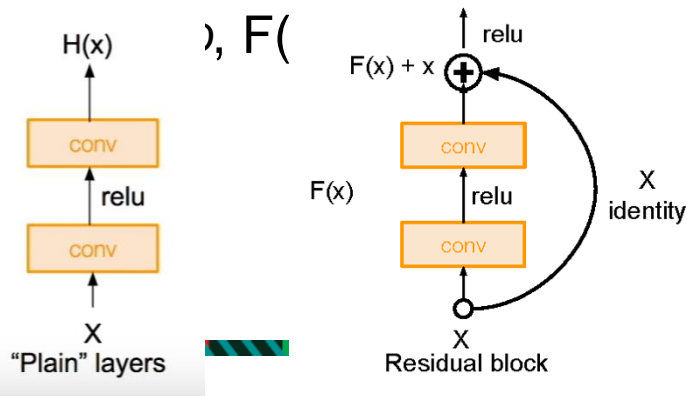
- 56-layer model performs worse on both training and test error
-> The deeper model performs worse (not caused by overfitting)!

ResNet

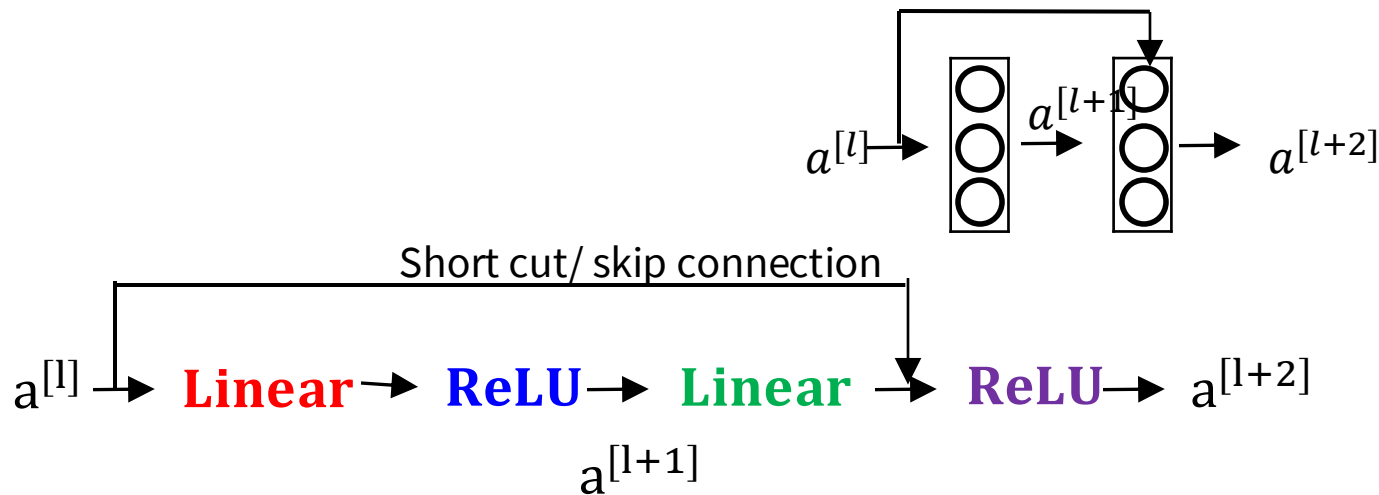
- **Hypothesis:** The problem is an optimization problem. Very deep networks are harder to optimize.
- **Solution:** Use network layers to fit residual mapping instead of directly trying to fit a desired underlying mapping.
- We will use **skip connections** allowing us to take the activation from one layer and feed it into another layer, much deeper into the network.
- Use layers to fit residual $F(x) = H(x) - x$ instead of $H(x)$ directly

RESIDUAL BLOCK

INPUT X GOES THROUGH CONV-RELU-CONV SERIES AND GIVES US $F(X)$. THAT RESULT IS THEN ADDED TO THE ORIGINAL INPUT X . LET'S CALL THAT $H(X) = F(X) + X$. IN TRADITIONAL CNNs, $H(X)$ WOULD JUST BE EQUAL TO $F(X)$. SO, INSTEAD OF JUST COMPUTING THAT TRANSFORMATION (STRAIGHT FROM X TO $F(X)$), WE'RE COMPUTING THE TERM THAT WE H/



ResNet



$$\mathbf{z}^{[l+1]} = \mathbf{W}^{[l+1]} \mathbf{a}^{[l]} + \mathbf{b}^{[l+1]}$$

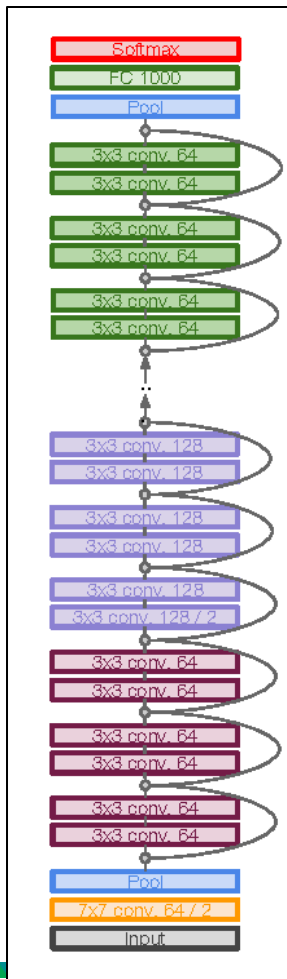
$$\mathbf{z}^{[l+2]} = \mathbf{W}^{[l+2]} \mathbf{a}^{[l+1]} + \mathbf{b}^{[l+2]}$$

$$\mathbf{a}^{[l+1]} = \mathbf{g}(\mathbf{z}^{[l+1]})$$

$$\mathbf{a}^{[l+2]} = \mathbf{g}(\mathbf{z}^{[l+2]})$$

$$\mathbf{a}^{[l+2]} = \mathbf{g}(\mathbf{z}^{[l+2]} + \mathbf{a}^{[l]}) = \mathbf{g}(\mathbf{W}^{[l+2]} \mathbf{a}^{[l+1]} + \mathbf{b}^{[l+2]} + \mathbf{a}^{[l]})$$

ResNet



FULL RESNET ARCHITECTURE:

STACK RESIDUAL BLOCKS

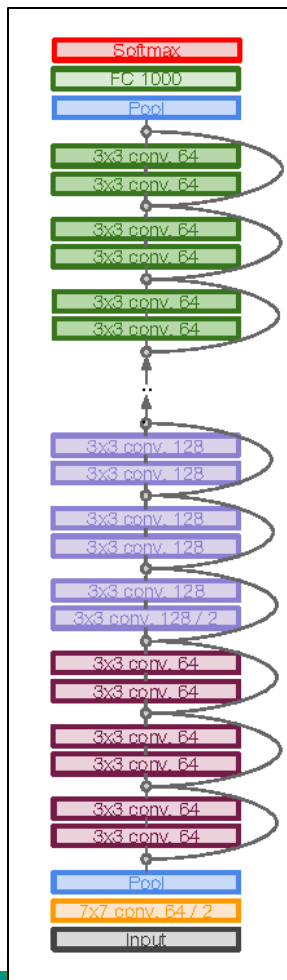
EVERY RESIDUAL BLOCK HAS TWO 3X3 CONV LAYERS

PERIODICALLY, DOUBLE # OF FILTERS AND
DOWNSAMPLE SPATIALLY USING STRIDE 2 (IN EACH
DIMENSION)

ADDITIONAL CONV LAYER AT THE BEGINNING

NO FC LAYERS AT THE END (ONLY FC 1000 TO
OUTPUT CLASSES)

ResNet



TOTAL DEPTHS OF 34, 50, 101, OR 152
LAYERS FOR IMAGENET

FOR DEEPER NETWORKS (RESNET-50+),
USE “BOTTLENECK” LAYER TO IMPROVE
EFFICIENCY (SIMILAR TO GOOGLNET)

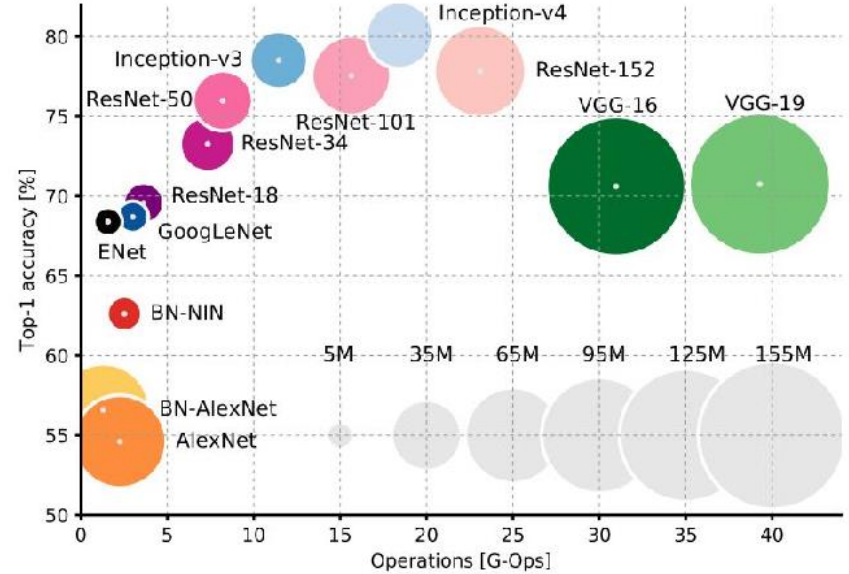
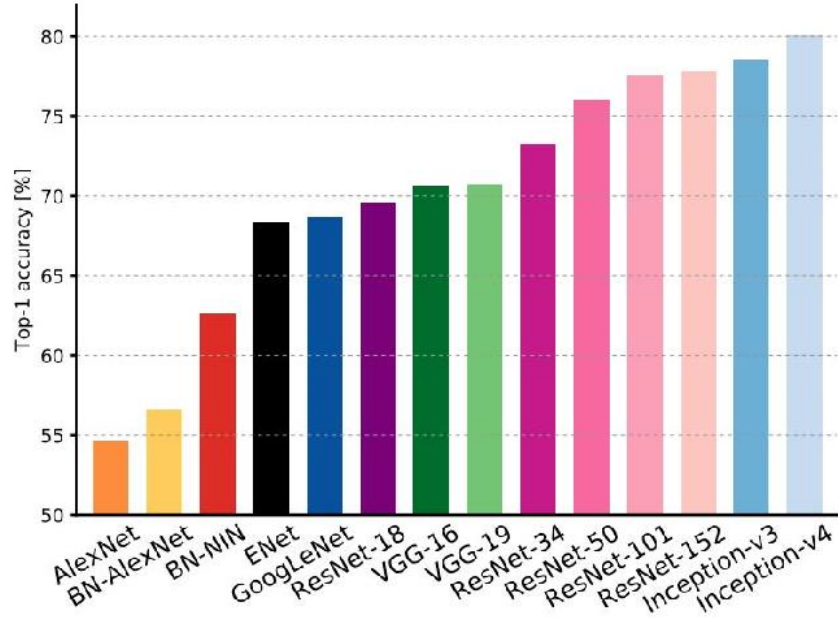
EXPERIMENTAL RESULTS:

ABLE TO TRAIN VERY DEEP NETWORKS WITHOUT DEGRADING
DEEPER NETWORKS NOW ACHIEVE LOWER TRAINING ERRORS
AS EXPECTED

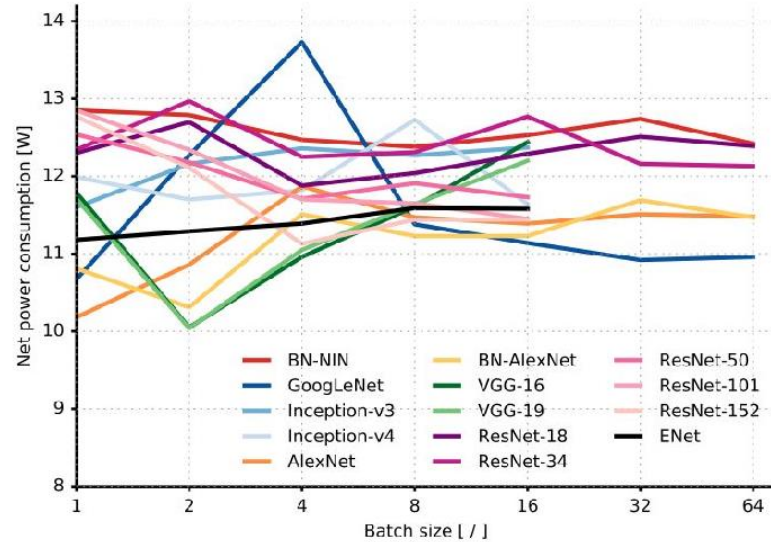
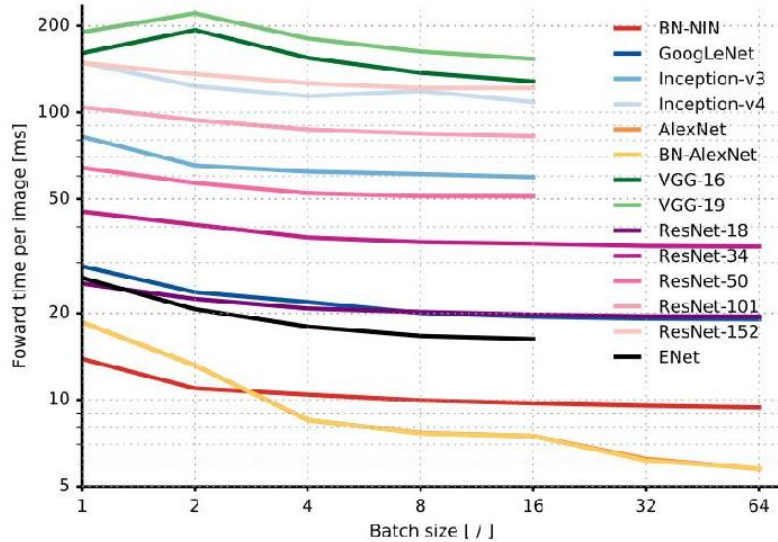
THE **BEST** CNN ARCHITECTURE THAT WE CURRENTLY HAVE
AND IS A GREAT INNOVATION FOR THE IDEA OF RESIDUAL
LEARNING.

EVEN BETTER THAN HUMAN PERFORMANCE!

Accuracy comparison



Forward pass time and power consumption



Summary

LENET-5

ALEXNET

VGG

GOOGLENET – INCEPTION MODULE

RESNET – RESIDUAL BLOCK



References

GRADIENT-BASED LEARNING APPLIED TO DOCUMENT RECOGNITION; ANN LECUN, LÉON BOTTOU, YOSHUA BENGIO, PATRICK HAFFNER; 1998

IMAGENET CLASSIFICATION WITH DEEP CONVOLUTIONAL NEURAL NETWORKS - ALEX KRIZHEVSKY, ILYA SUTSKEVER, GEOFFREY E. HINTON; 2012

VERY DEEP CONVOLUTIONAL NETWORKS FOR LARGE SCALE IMAGE RECOGNITION - KAREN SIMONYAN AND ANDREW ZISSERMAN; 2015

GOING DEEPER WITH CONVOLUTIONS - CHRISTIAN SZEGEDY ET AL.; 2015

DEEP RESIDUAL LEARNING FOR IMAGE RECOGNITION - KAIMING HE, XIANGYU ZHANG, SHAOQING REN, JIAN SUN; 2015

*STANFORD CS231- FEI-FEI & JUSTIN JOHNSON & SERENA YEUNG. LECTURE 9
COURSERA, MACHINE LEARNING COURSE BY ANDREW NG.*

References

THE 9 DEEP LEARNING PAPERS YOU NEED TO KNOW ABOUT (UNDERSTANDING CNNs PART 3) BY ADIT DESHPANDE

[HTTPS://ADESHPANDE3.GITHUB.IO/ADESHPANDE3.GITHUB.IO/THE-9-DEEP-LEARNING-PAPERS-YOU-NEED-TO-KNOW-ABOUT.HTML](https://adeshpande3.github.io/adeshpande3.github.io/The-9-Deep-Learning-Papers-You-Need-To-Know-About.html)

CNNs ARCHITECTURES: LENET, ALEXNET, VGG, GOOGLNET, RESNET AND MORE ...

BY SIDDHARTH DAS [HTTPS://MEDIUM.COM/@SIDDHARTHDAS_32104/CNNS-ARCHITECTURES-LENET-ALEXNET-VGG-GOOGLENET-RESNET-AND-MORE-666091488DF5](https://medium.com/@siddharthdas_32104/cnns-architectures-leetnet-alexnet-vgg-googlenet-resnet-and-more-666091488df5)

SLIDE TAKEN FROM FORWARD AND BACKPROPAGATION IN CONVOLUTIONAL NEURAL NETWORK. – MEDIUM , BY SUJIT RAI

[HTTPS://MEDIUM.COM/@2017CSM1006/FORWARD-AND-BACKPROPAGATION-IN-CONVOLUTIONAL-NEURAL-NETWORK-4DFA96D7B37E](https://medium.com/@2017CSM1006/forward-and-backpropagation-in-convolutional-neural-network-4dfa96d7b37e)