

# Designing, Visualizing and Understanding Deep Neural Networks

## Lecture 5: Convolutional Networks I

---

CS 194/294-129 Spring 2018  
John Canny

Slides originated from Efros, Karpathy, Ransato, Seitz, and Palmer

# Last Time - Gradient Descent

To reach a minimum of loss, we should follow the negative gradient.

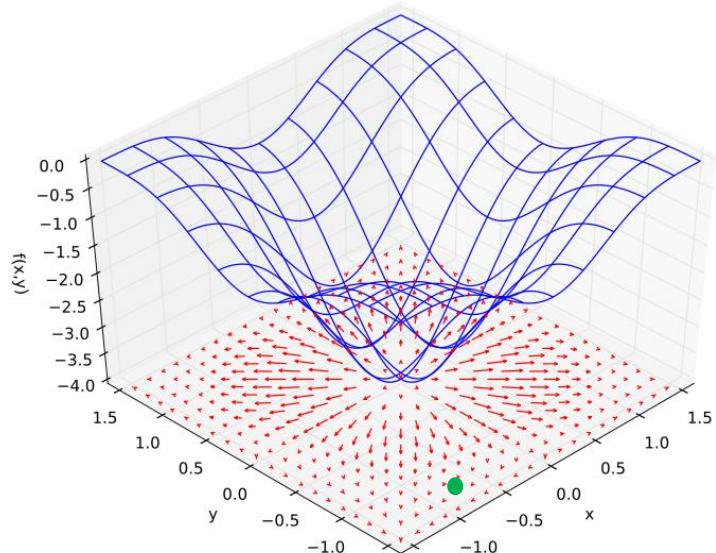
i.e. we should take small steps in direction

$$-\nabla_W L(W)$$

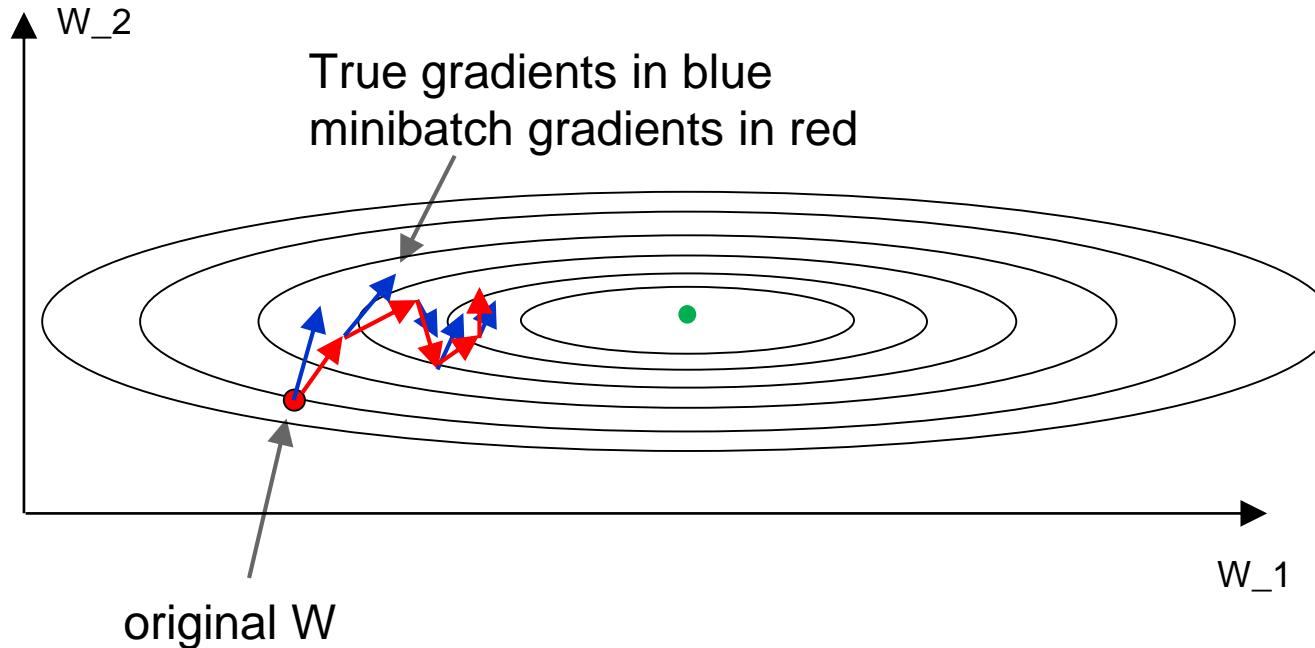
Let  $W^t$  denote the weights at step  $t$  of gradient descent. Then

$$W^{t+1} = W^t - \alpha \nabla_W L(W)$$

Where  $\alpha$  is called the **learning rate**.

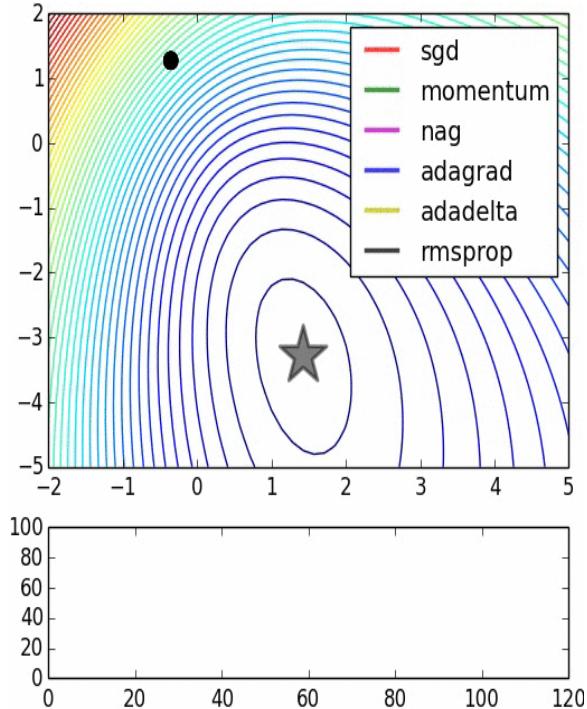


# Last time: Minibatches and SGD



Gradients are noisy but still make good progress on average

# Last Time: SGD refinements: Momentum, Nesterov Momentum, RMSprop, ADAGRAD



# This Time: Backpropagation

So far we have been using gradient methods to minimize a loss over some parameters.

Our loss is of the form  $L(f(x, W), y)$ .

$x$  is an input,  $y$  is a target, and  $W$  are the parameters.

To compute the gradient of  $L$  wrt  $W$ , we need the **chain rule**. If  $f$  is single-valued,  $W$  a single parameter the chain rule is just:

$$\frac{dL}{dW} = \frac{dL}{df} \frac{df}{dW}$$

If  $W$  is a vector of parameters, then we have:

$$\nabla_W L = \frac{dL}{df} \nabla_W f$$

Which is really just the first rule applied to all the partial derivatives wrt elements of  $W$ .

# Updates

- Project Proposals are due in two weeks.
- If you need help putting a team together, we suggest you post to Piazza in the “project\_teams” folder.
- You can either post your own ideas on project topics, or respond to other groups looking to complete their team.
- Assignment 1 due on Feb 12, make sure at least that you’ve started by now...

# The Chain Rule

If  $f$  is also vector-valued with  $k$  values is a vector of  $m$  parameters, then we can apply the chain rule parameter-wise and then sum over the contributions:

$$\frac{\partial L}{\partial W_j} = \sum_{i=1}^k \frac{\partial L}{\partial f_i} \frac{\partial f_i}{\partial W_j}$$

For  $i = 1, \dots, m$ . This can be written as a matrix multiply

$$J_L(W) = J_L(f) J_f(W)$$

Where  $J_f(W)$  is a Jacobian matrix.

$$J_f(W)_{ij} = \frac{\partial f_i}{\partial W_j}$$

# Jacobians

The Jacobian generalizes the gradient of a scalar-valued function  $f$  to a  $k$ -valued function. Here we think of the function as a neural layer with  $m$  inputs and  $k$  outputs.

$$J_f(W) = \begin{bmatrix} \frac{\partial f_1}{\partial W_1} & \frac{\partial f_1}{\partial W_2} & \dots & \frac{\partial f_1}{\partial W_m} \\ \frac{\partial f_2}{\partial W_1} & \frac{\partial f_2}{\partial W_2} & \dots & \frac{\partial f_2}{\partial W_m} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_k}{\partial W_1} & \frac{\partial f_k}{\partial W_2} & \dots & \frac{\partial f_k}{\partial W_m} \end{bmatrix}$$

The Jacobian has dimensions  $k \times m$  which is  $n_{outputs} \times n_{inputs}$ .

# N-step Chain Rule

Now suppose we have several vector-valued functions  $A(\cdot)$ ,  $B(\cdot)$ ,  $C(\cdot)$ , ... composed in a chain (e.g. a deep network):

$$W \rightarrow A \rightarrow B \rightarrow C \rightarrow \cdots L$$

Algebraically, that looks like:

$$L(W) = L(K \cdots C(B(A(W))) \cdots)$$

Then we just multiply Jacobians (matrix multiply  $*$ ) to get the gradient:

$$J_L(W) = J_L(K) * \cdots * J_C(B) * J_B(A) * J_A(W)$$

And  $J_L(W) = (\nabla_W L)^T$  which is the gradient we need to minimize loss over  $W$ .

# Backpropagation

Now all the Jacobians are matrices, and matrix multiply is associative.

$$J_L(W) = J_L(K) * \cdots J_C(B) * J_B(A) * J_A(W)$$

So we could actually evaluate the product of Jacobians in any order, including left-to-right and right-to-left.

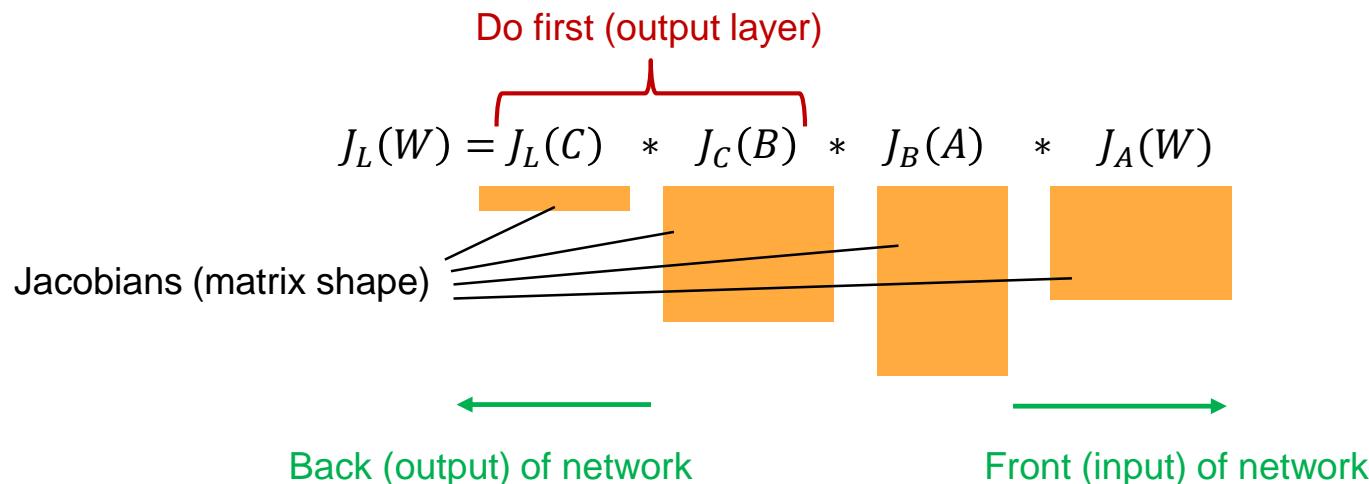
**Backpropagation:** evaluate the Jacobian product (loss gradient wrt params) left-to-right, i.e. from the output of the neural network toward its input.

Why not some other order?

# Backpropagation

## Reason 1: Efficiency

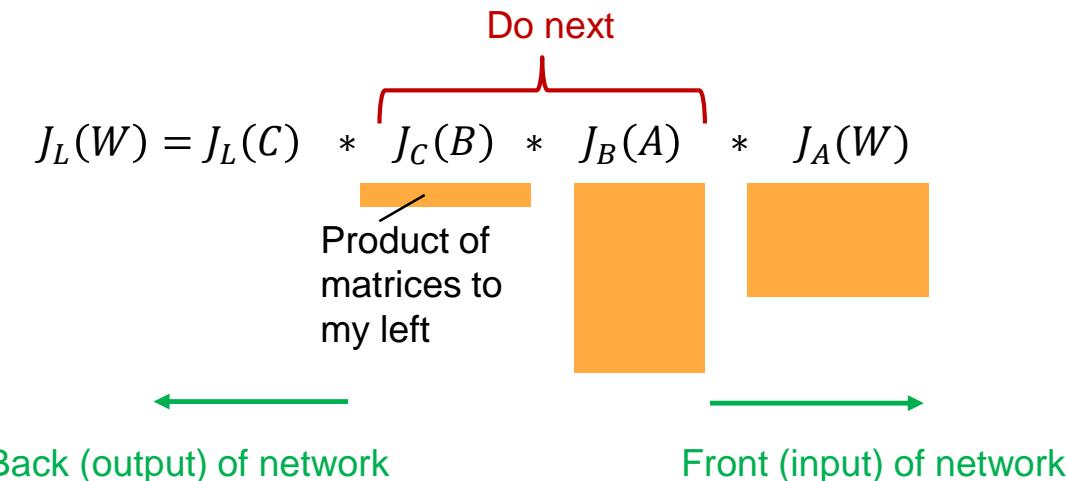
**Output Jacobian is always a row vector** (because loss is a scalar). Matrix-vector multiply is much less expensive than matrix-matrix multiply.



# Backpropagation

## Reason 1: Efficiency

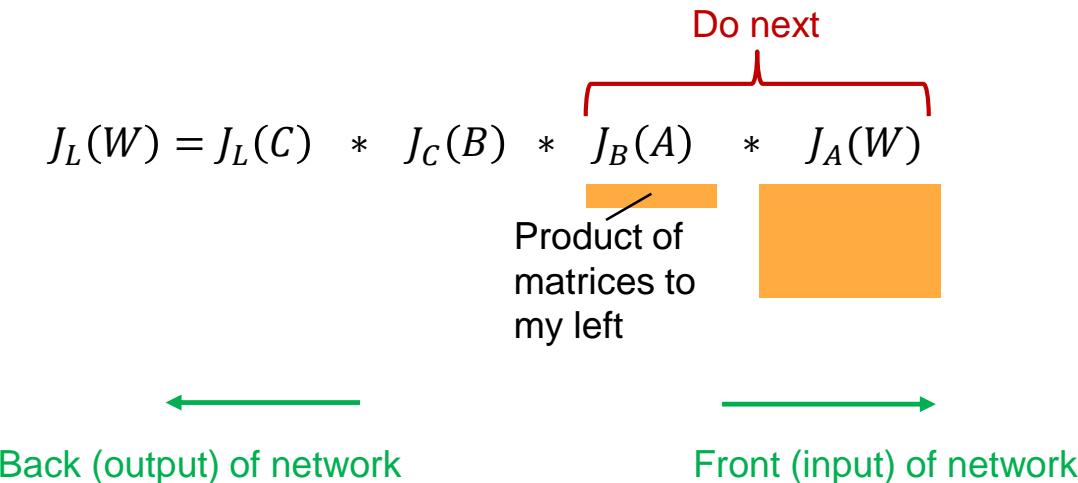
Output Jacobian is always a row vector. Matrix-vector multiply is much less expensive than matrix-matrix multiply.



# Backpropagation

## Reason 1: Efficiency

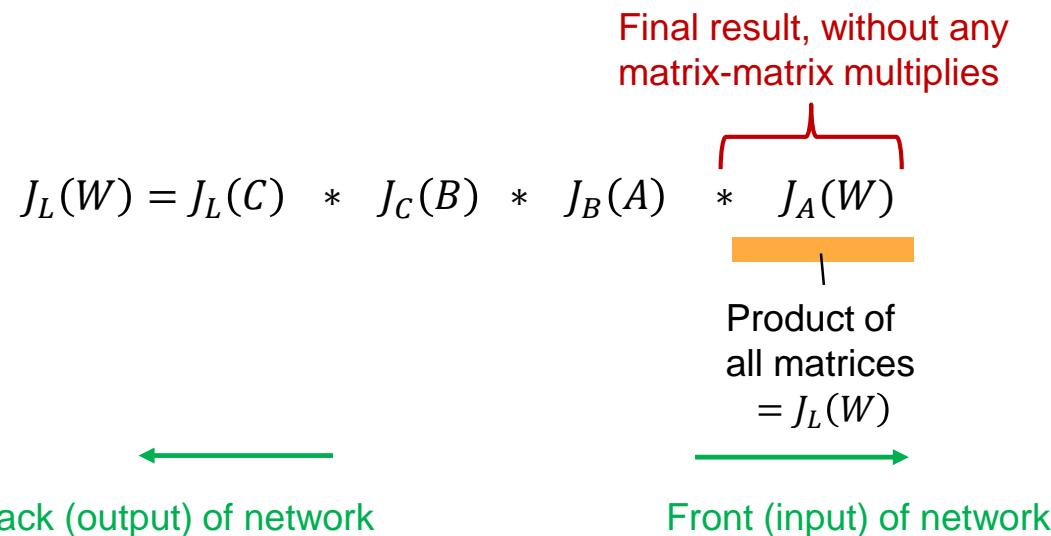
Output Jacobian is always a row vector. Matrix-vector multiply is much less expensive than matrix-matrix multiply.



# Backpropagation

## Reason 1: Efficiency

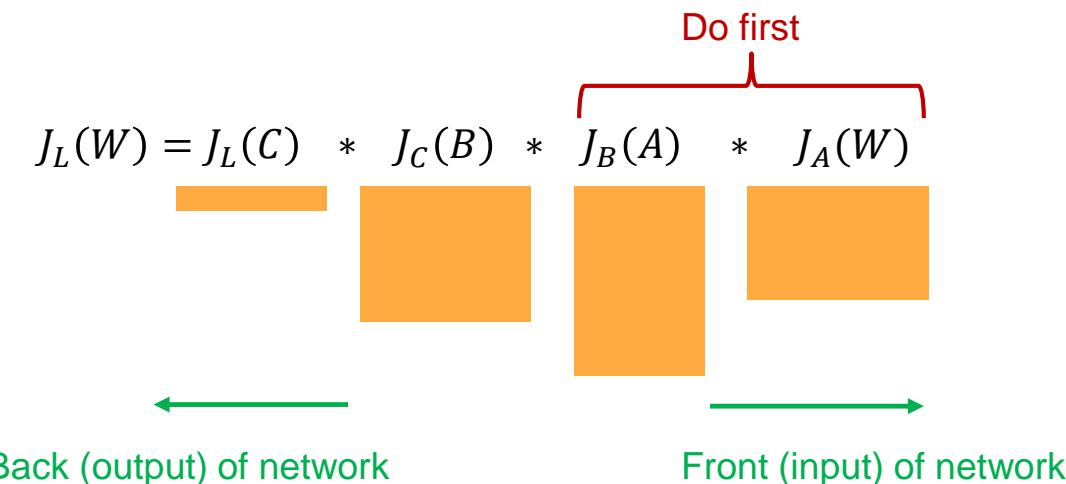
Output Jacobian is always a row vector. Matrix-vector multiply is much less expensive than matrix-matrix multiply.



# For-propagation

## Reason 1: Efficiency

By comparison, we could invent “for-propagation”:

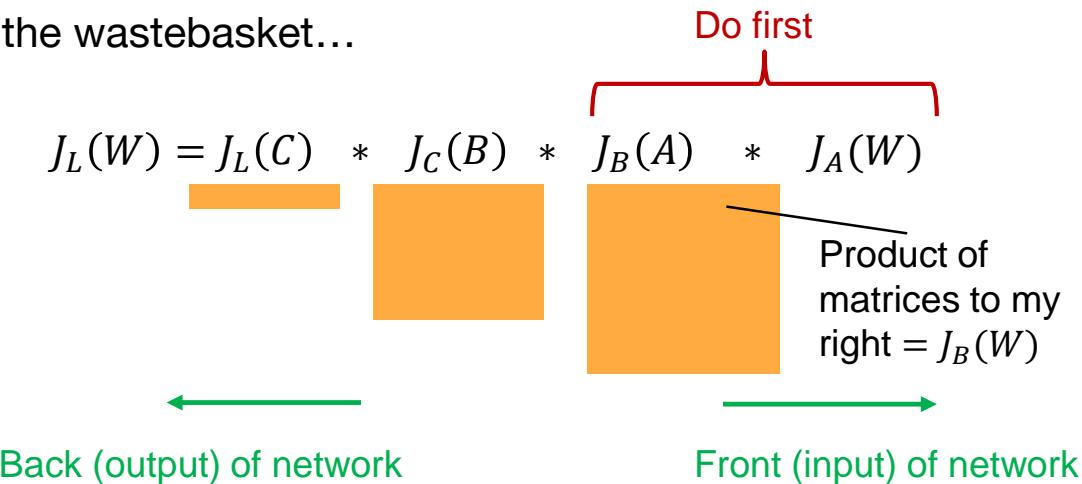


# For-propagation

# Reason 1: Efficiency

We could invent “for-propagation”: Oops, cost  $O(n^3)$  instead of  $O(n^2)$  for the first multiply and we still have another matrix-matrix multiply to do.

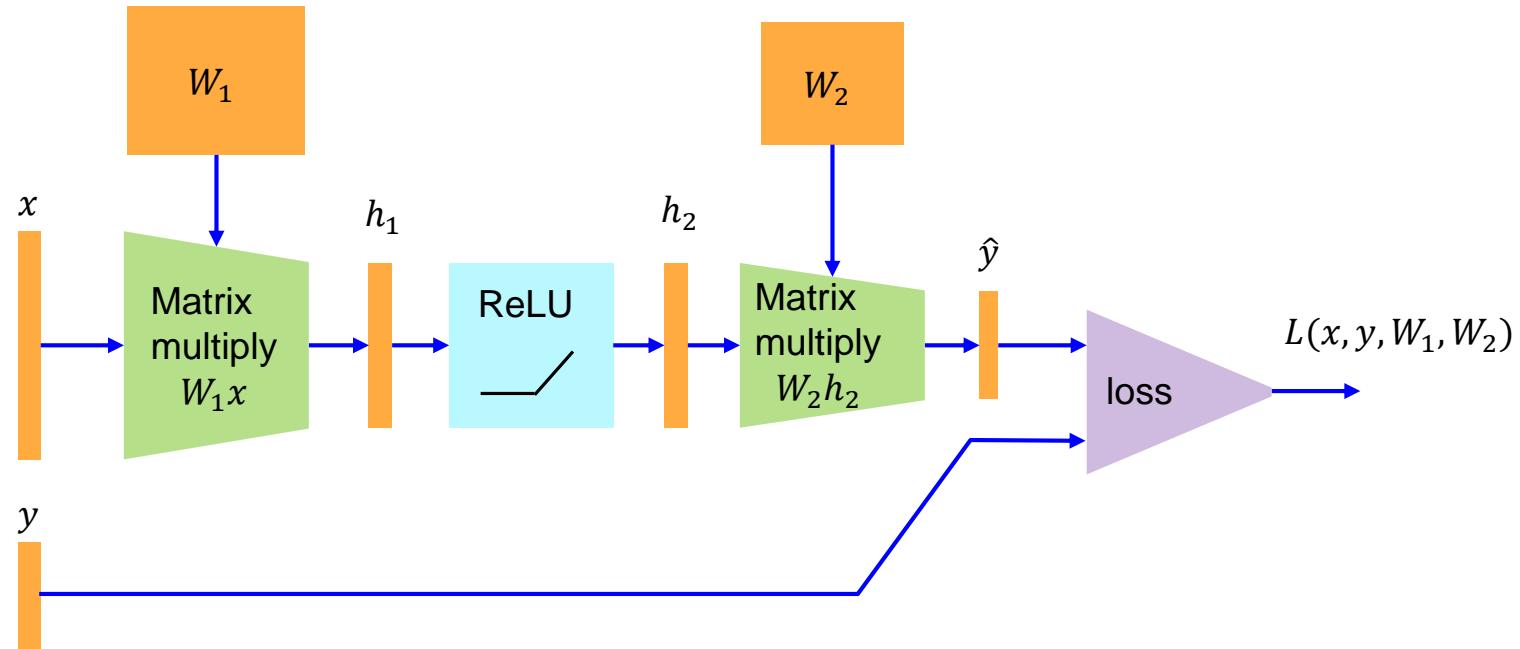
A good invention for the wastebasket....



# Backpropagation

**Reason 2 to use backpropagation:** Common subexpressions.

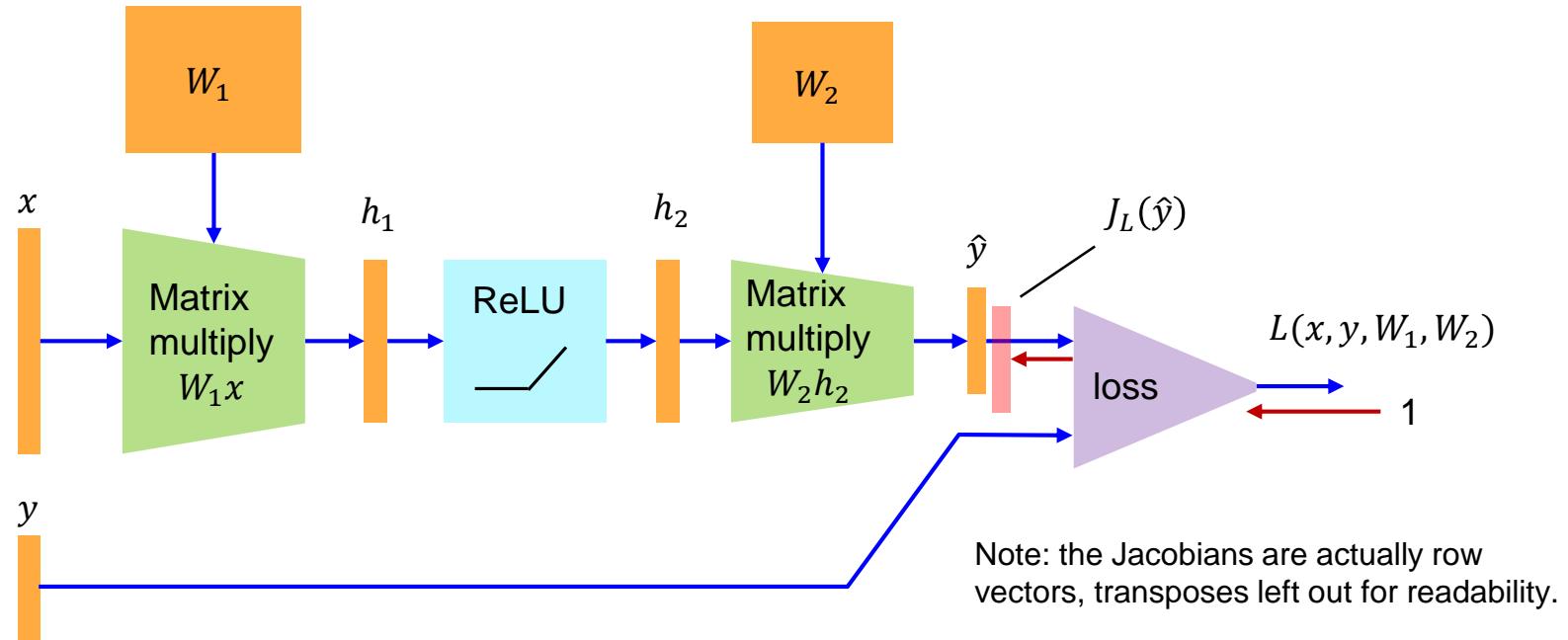
Let's build a real neural network (Tensorflow style):



# Backpropagation

**Reason 2 to use backpropagation:** Common subexpressions.

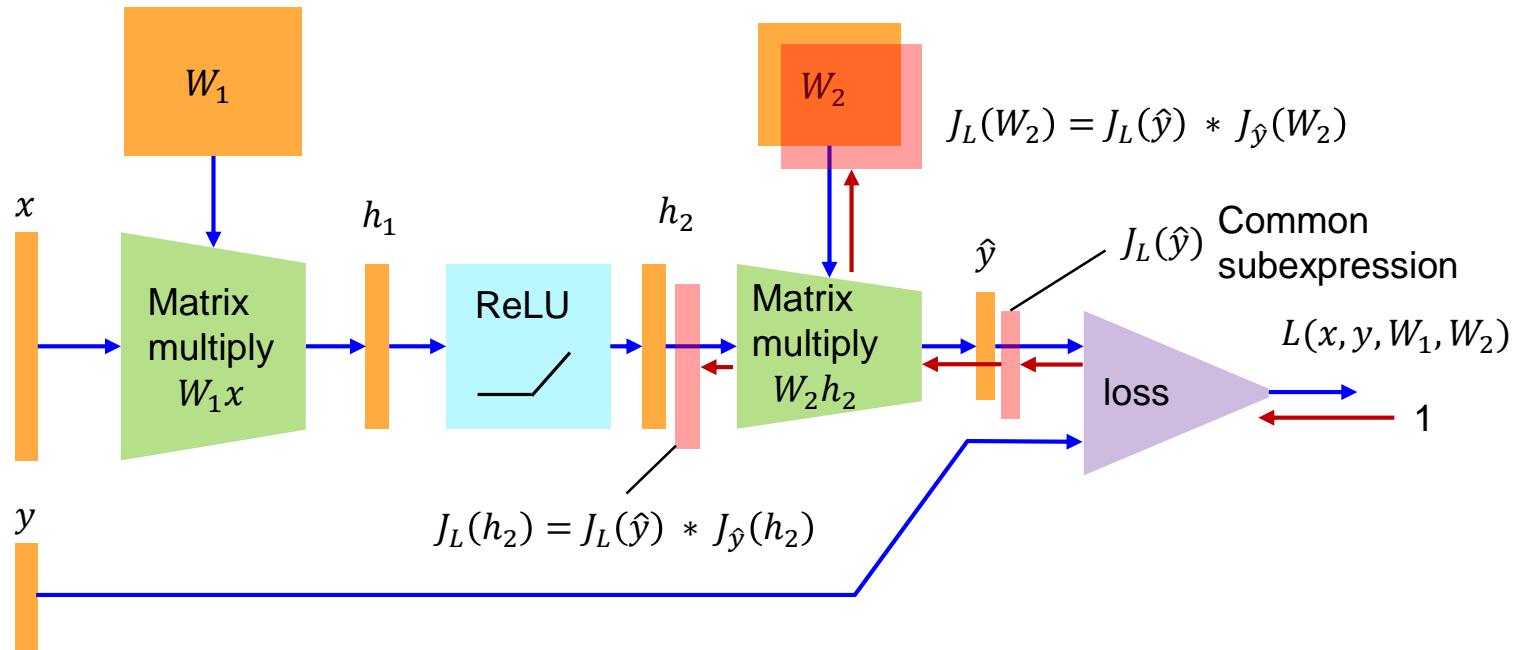
Let's build a real neural network (Tensorflow style):



# Backpropagation

**Reason 2 to use backpropagation:** Common subexpressions.

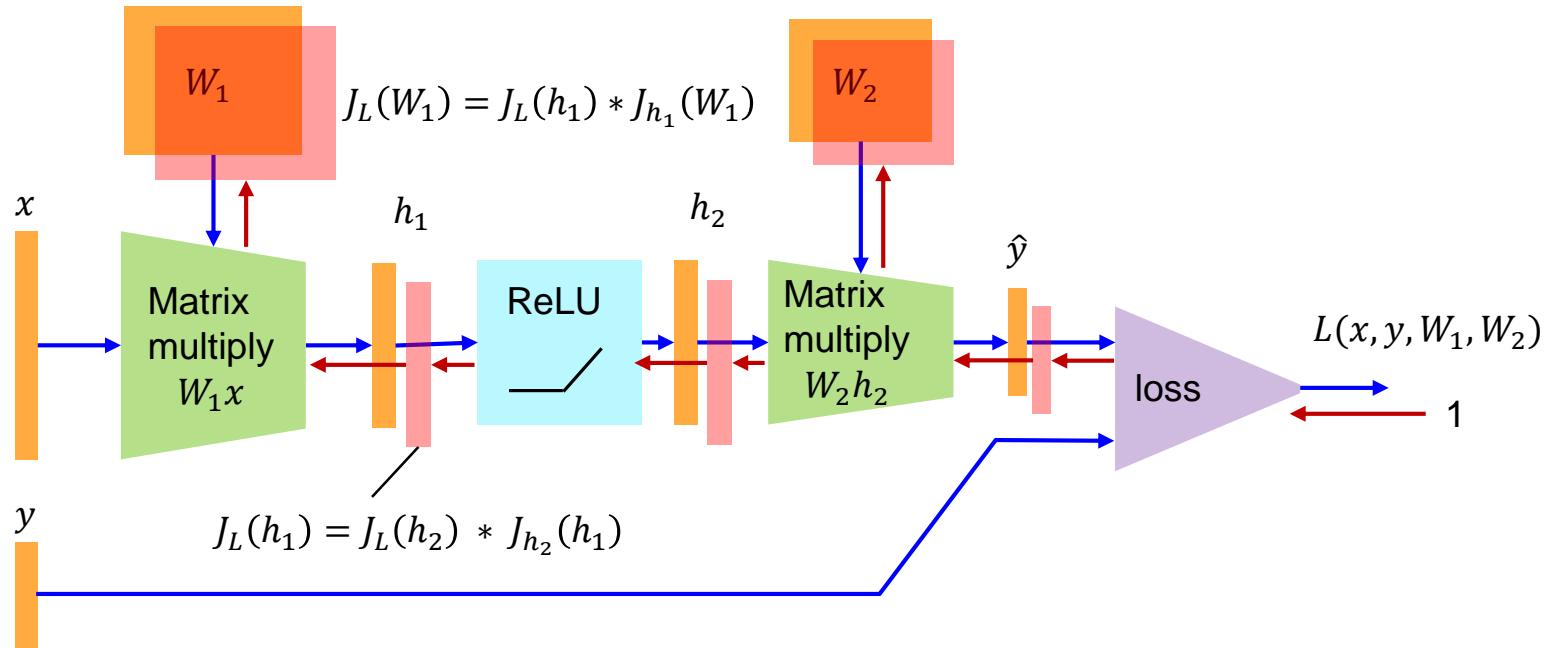
Let's build a real neural network (Tensorflow style):



# Backpropagation

**Reason 2 to use backpropagation:** Common subexpressions.

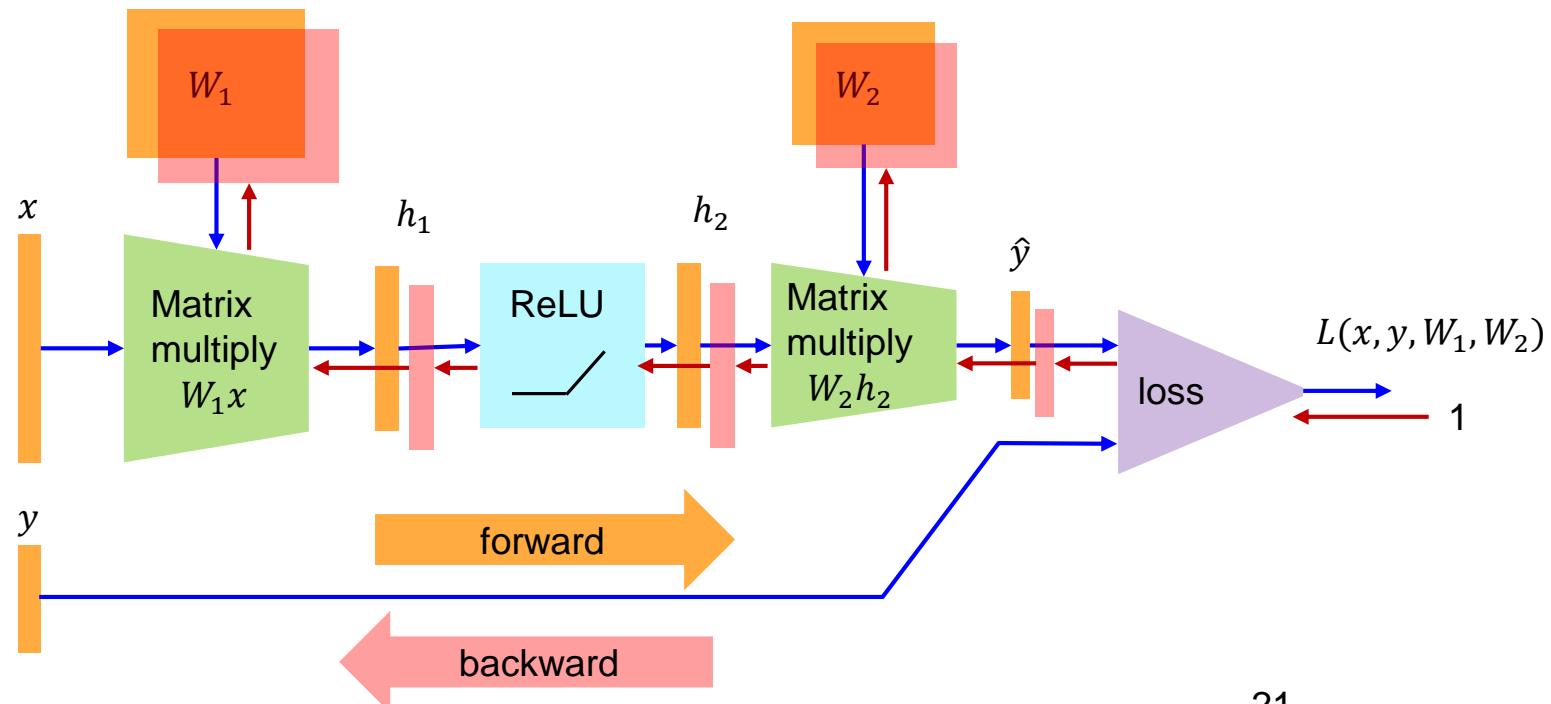
Let's build a real neural network (Tensorflow style):



# Last Time: Backpropagation

**Backprop Efficiency:** matrix-vector multiply only and common subexpressions.

Example: a real neural network (Tensorflow style):



# Backpropagation Recipes: Cross Entropy Loss

- This is just calculus: always check yourself:
- Cross-entropy loss (correct class label  $y \in \{1, \dots, n\}$  known exactly, input vector  $s$ ):

$$L = -\log s_y$$

So

$$J_L(s)_i = \nabla_s L_i^T = \begin{cases} -\frac{1}{s_i} & \text{for } i = y \\ 0 & \text{otherwise} \end{cases}$$

In general if  $y$  is a distribution,  $L = -\sum_{i=1}^n y_i \log s_i$ , then

$$J_L(s)_i = \nabla_s L_i^T = -\frac{y_i}{s_i}$$

# Backprop Recipes: Multiclass SVM Loss

- Assume label  $y \in \{1, \dots, n\}$  known exactly, input vector  $s$ :

$$L = \sum_{i=1}^n \max(0, 1 - s_y + s_i)$$

So

$$J_L(s)_i = \nabla_s L_i^T = \begin{cases} 1 & \text{if } s_y - s_i < 1 \text{ and } i \neq y \\ 0 & \text{otherwise} \end{cases}$$

# Backprop Recipes: Multiclass Logistic Fn

Assume input vector  $s$ , output  $y$  of same dimension  $n$ .

$$y_i = \frac{\exp s_i}{\exp s_1 + \dots + \exp s_n} = \frac{f}{g}$$

Then

$$\frac{\partial y_i}{\partial s_j} = \frac{f'g - g'f}{g^2} = \frac{\exp s_i g \delta_{ij} - \exp s_j \exp s_i}{g^2}$$

Simplifying:

$$\frac{\partial y_i}{\partial s_j} = y_i \delta_{ij} - y_i y_j$$

Where:

$$\delta_{ij} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}$$

# Backprop Recipes: Matrix Multiply (FC) Layer

Assume input vector  $s$  of dimension  $n$ , output  $y$  of dimension  $m$ , layer  $y = Ws$

$$y_i = \sum_{j=1}^n W_{ij}s_j$$

Then

$$\frac{\partial y_i}{\partial s_j} = W_{ij} \quad \text{and so} \quad J_L(s) = J_L(y)W$$

Also we want to backprop to the weight matrix and compute:

$$\frac{\partial y_i}{\partial W_{jk}}$$

Which is a tensor (ouch), but notice that its zero unless  $i = j$

# Backprop Recipes: Matrix Multiply (FC) Layer

We want:  $\frac{\partial y_i}{\partial W_{jk}}$

Which is a tensor, but notice that its zero unless  $i = j$

So it simplifies to

$$\frac{\partial y_i}{\partial W_{ik}} = s_k$$

And if we already know the output loss Jacobian (from backprop):  $J_L(y)$ , then

$$\frac{\partial L}{\partial W_{ik}} = \frac{\partial L}{\partial y_i} \frac{\partial y_i}{\partial W_{ik}} = J_L(y)_i s_k$$

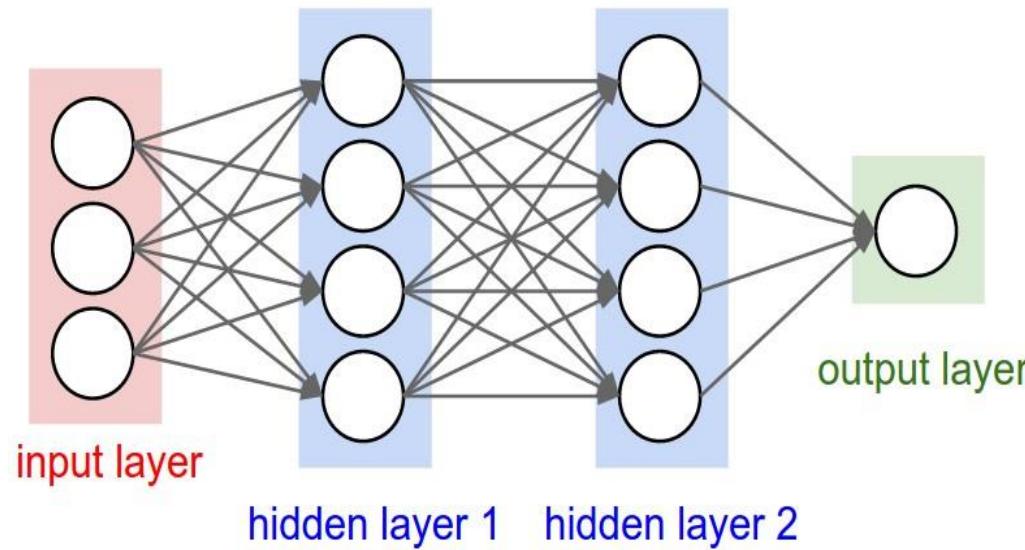
Or more compactly:

$$\frac{\partial L}{\partial W_{ik}} = \underbrace{J_L(y)^T}_{\text{Col vector}} \underbrace{s^T}_{\text{Row vector}}$$

# Backpropagation

- Compute function values (activations) from the first layer to the last.
- Compute derivatives of the loss wrt other layers from the last layer to the first (backpropagation).
- This only requires matrix-vector multiplies.
- Paths from the loss layer to inner layers are re-used.

# This time: Neural Networks for Visual Data

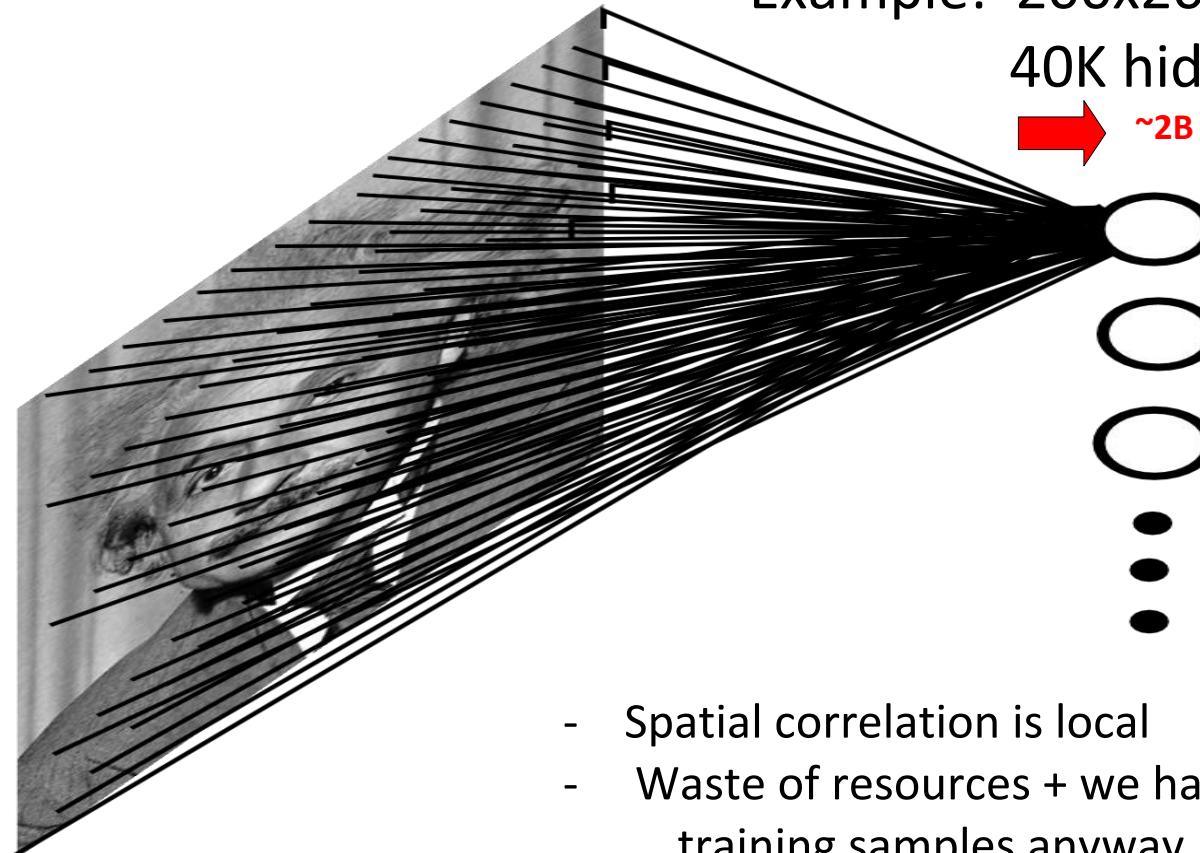


# Fully Connected Layer

Example: 200x200 image

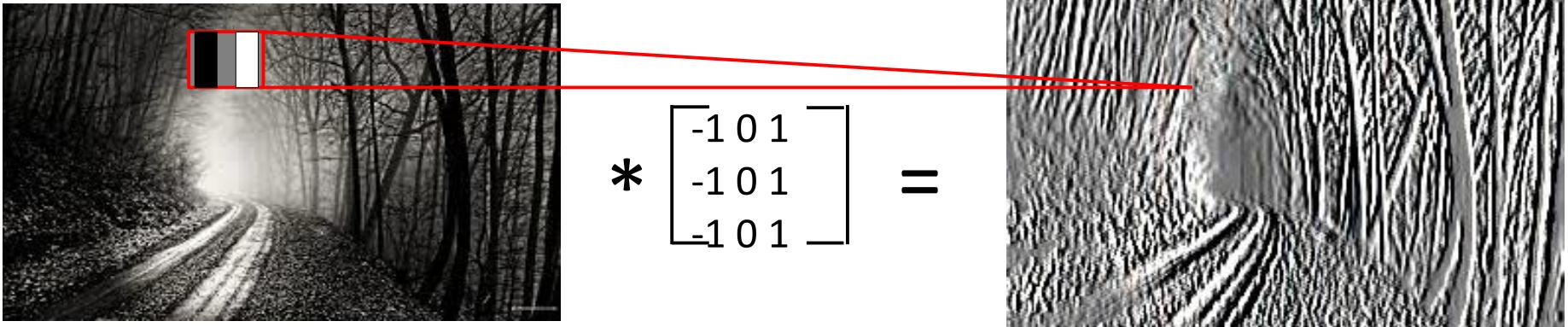
40K hidden units

→ **~2B parameters!!!**



- Spatial correlation is local
- Waste of resources + we have not enough training samples anyway..

# Convolutional of Two Signals



$$g(x, y) = (h * f)(x, y) = \sum_{i=-\infty}^{\infty} \sum_{j=-\infty}^{\infty} h(i, j)f(x - i, y - j)$$

- elementwise multiplication and sum of a filter and the signal (image)

# Convolutional of Two Signals

The convolution formula is:

$$(h * f)(x, y) = \sum_{i=-\infty}^{\infty} \sum_{j=-\infty}^{\infty} h(i, j)f(x - i, y - j)$$

It follows (by redefining indices) that convolution is commutative:

$$h * f = f * h$$

**Finite filters:** Typically one function, the “filter” (here its  $h$ ) is defined over a finite support  $[-w/2, w/2]^2$ , then

$$(h * f)(x, y) = \sum_{i=-w/2}^{w/2} \sum_{j=-w/2}^{w/2} h(i, j)f(x - i, y - j)$$

# Image filtering

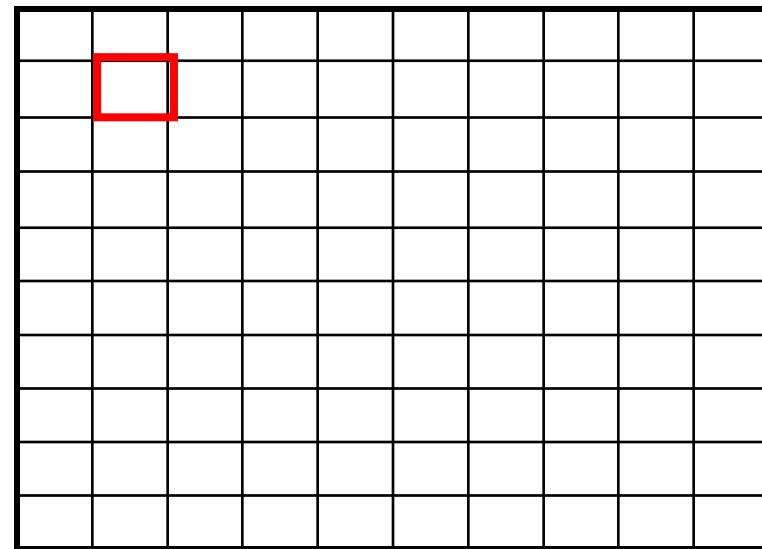
$$g = h * f$$

$$h[\cdot, \cdot] \frac{1}{9} \begin{matrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{matrix}$$

$f[.,.]$

$g[.,.]$

0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0	0
0	0	0	90	90	90	90	90	0	0	0
0	0	0	90	90	90	90	90	0	0	0
0	0	0	90	0	90	90	90	0	0	0
0	0	0	90	90	90	90	90	0	0	0
0	0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0



# Image filtering

$$g = h * f$$

$$h[\cdot, \cdot] \frac{1}{9} \begin{matrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{matrix}$$

$f[\cdot, \cdot]$

0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0	0
0	0	0	90	90	90	90	90	0	0	0
0	0	0	90	90	90	90	90	0	0	0
0	0	0	90	0	90	90	90	0	0	0
0	0	0	90	90	90	90	90	0	0	0
0	0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0

$g[\cdot, \cdot]$


# Image filtering

$$g = h * f$$

$$h[\cdot, \cdot] \frac{1}{9} \begin{matrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{matrix}$$

$f[\cdot, \cdot]$

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

$g[\cdot, \cdot]$

			0	10	20				

# Image filtering

$$g = h * f$$

$$h[\cdot, \cdot] \frac{1}{9} \begin{matrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{matrix}$$

$f[\cdot, \cdot]$

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

$g[\cdot, \cdot]$

			0	10	20	30			

# Image filtering

$$g = h * f$$

$$h[\cdot, \cdot] \frac{1}{9} \begin{matrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{matrix}$$

$f[\cdot, \cdot]$

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

$g[\cdot, \cdot]$


# Image filtering

$$g = h * f$$

$$h[\cdot, \cdot] \frac{1}{9} \begin{matrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{matrix}$$

$f[\cdot, \cdot]$

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

$g[\cdot, \cdot]$

			0	10	20	30	30		

# Image filtering

$$g = h * f$$

$$h[\cdot, \cdot] \frac{1}{9} \begin{matrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{matrix}$$

$f[\cdot, \cdot]$

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

$g[\cdot, \cdot]$

			0	10	20	30	30		

# Image filtering

$$g = h * f \quad h[\cdot, \cdot] \quad \frac{1}{9} \begin{array}{|c|c|c|} \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline \end{array}$$

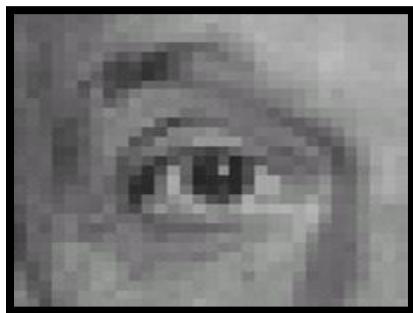
$f[\cdot, \cdot]$

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	90	0	90	90	90	0	0
0	0	0	90	90	90	90	90	0	0
0	0	0	0	0	0	0	0	0	0
0	0	90	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

$g[\cdot, \cdot]$

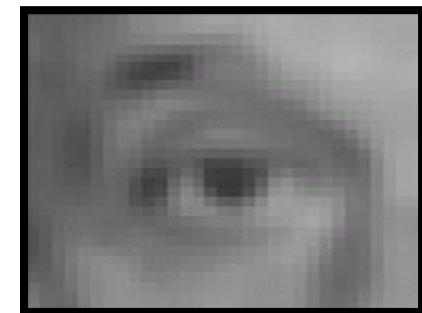
	0	10	20	30	30	30	20	10	
	0	20	40	60	60	60	40	20	
	0	30	60	90	90	90	60	30	
	0	30	50	80	80	90	60	30	
	0	30	50	80	80	90	60	30	
	0	20	30	50	50	60	40	20	
	10	20	30	30	30	30	20	10	
	10	10	10	0	0	0	0	0	

# Linear filters: examples



Original

$$\frac{1}{9} \begin{array}{|c|c|c|} \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline 1 & 1 & 1 \\ \hline \end{array} =$$



Blur (with a mean filter)

# Practice with linear filters



Original

0	0	0
0	1	0
0	0	0

?

# Practice with linear filters



Original

0	0	0
0	1	0
0	0	0



Filtered  
(no change)

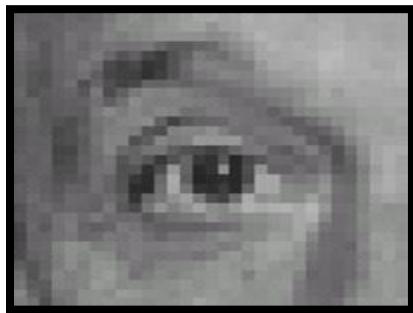
# Practice with linear filters



0	0	0
0	0	1
0	0	0

?

# Practice with linear filters



Original

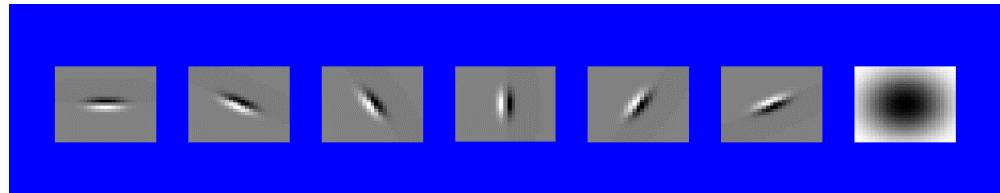
0	0	0
0	0	1
0	0	0



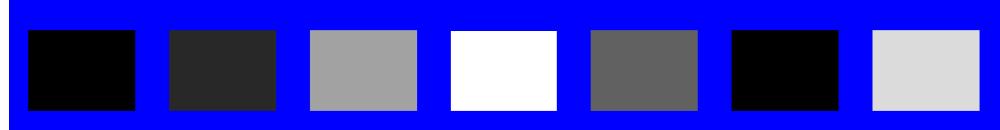
Shifted left  
By 1 pixel

# Can you match the texture to the response?

Filters



1



2

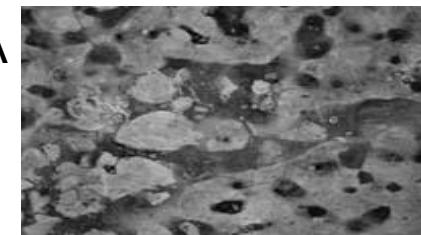


3



Mean abs responses

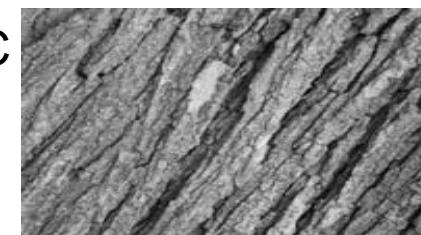
A



B



C

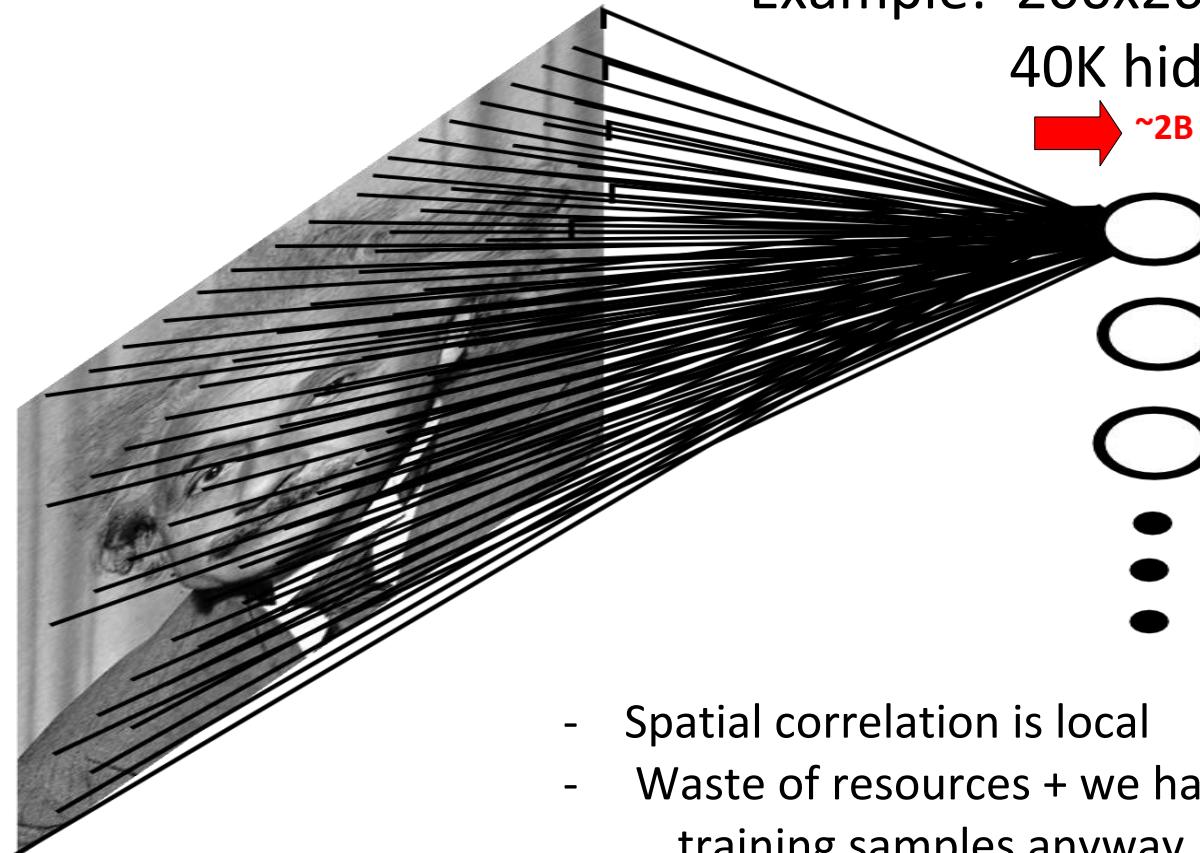


# Fully Connected Layer

Example: 200x200 image

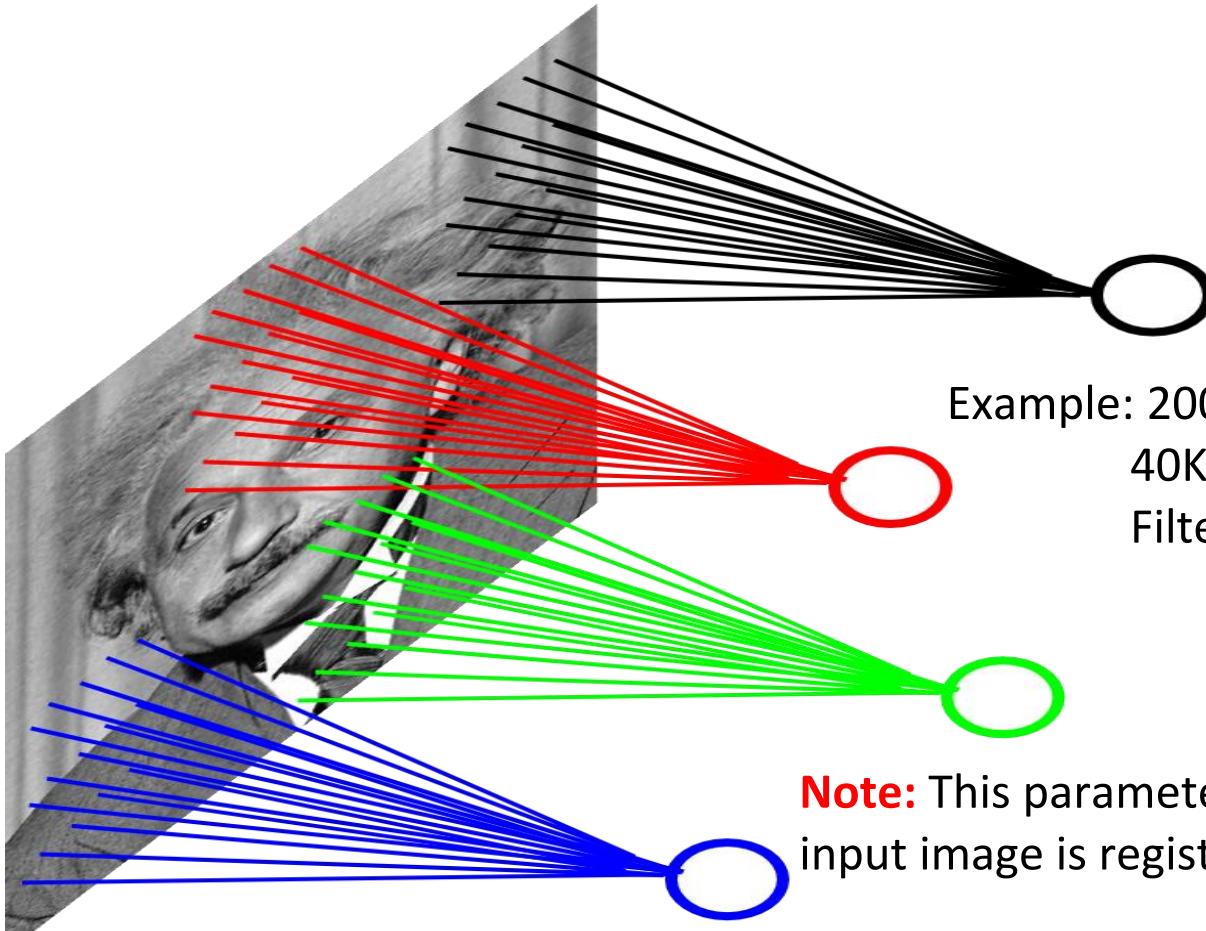
40K hidden units

→ **~2B parameters!!!**



- Spatial correlation is local
- Waste of resources + we have not enough training samples anyway..

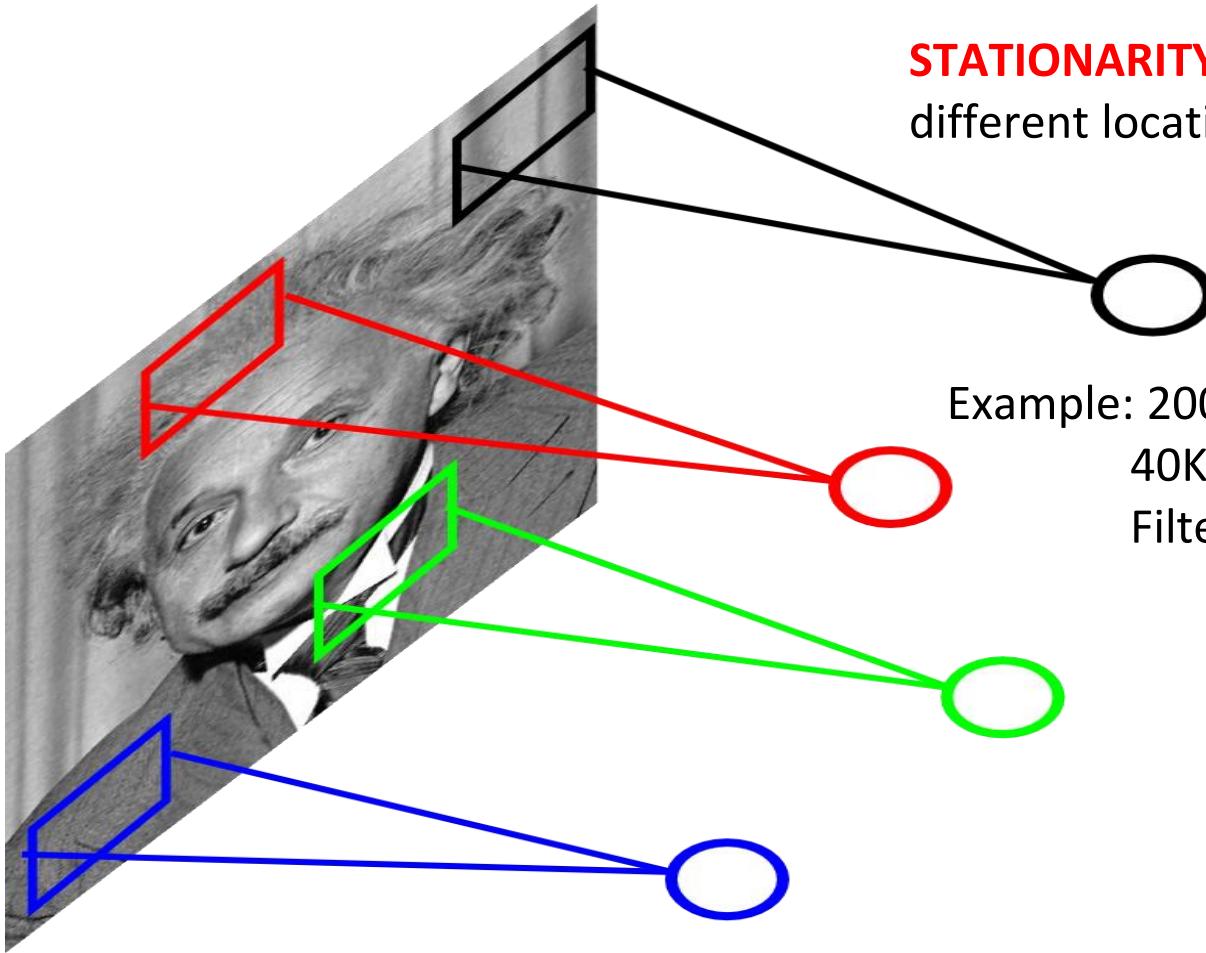
# Locally Connected Layer



Example: 200x200 image  
40K hidden units  
Filter size: 10x10  
4M parameters

**Note:** This parameterization is good when input image is registered (e.g., face recognition).

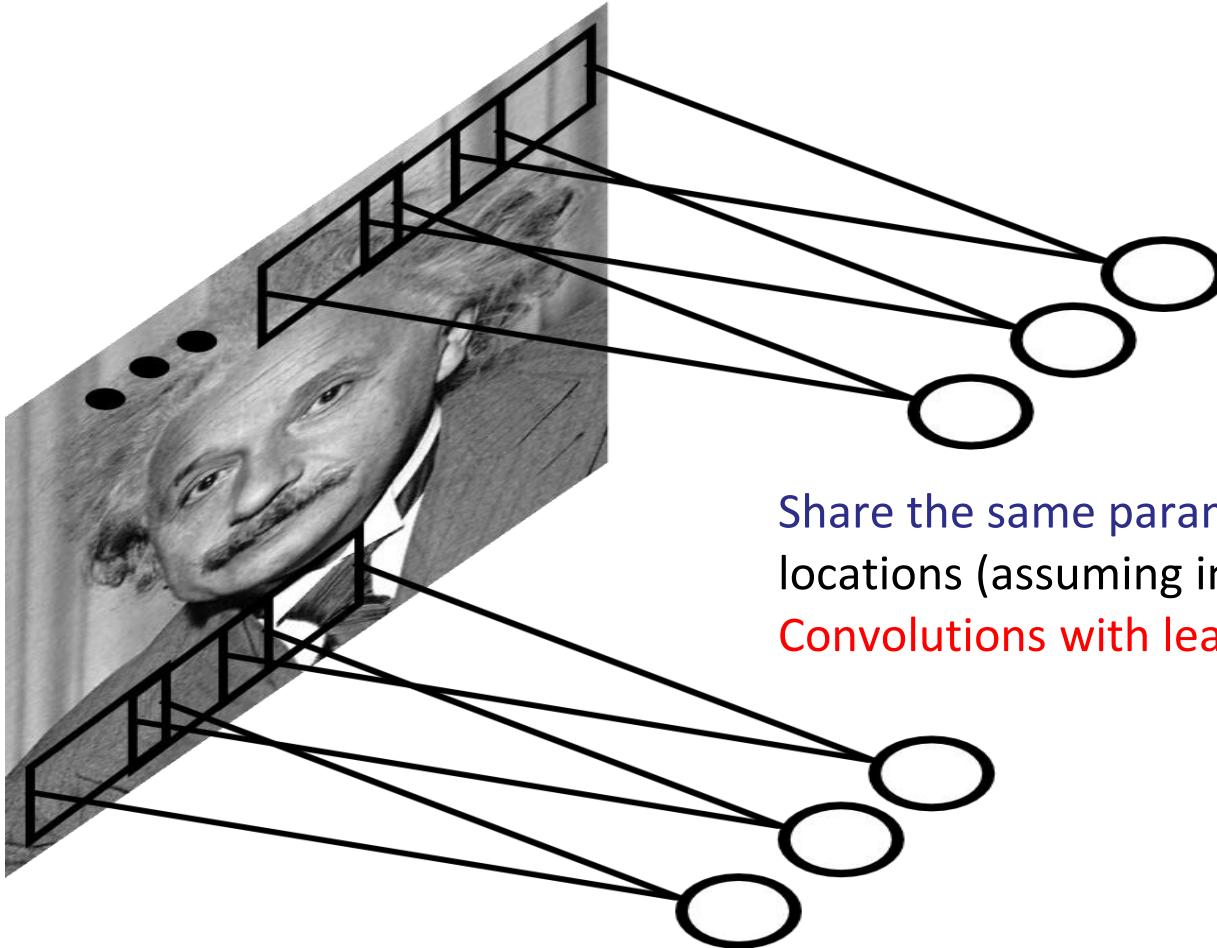
# Locally Connected Layer



**STATIONARITY?** Statistics is similar at different locations

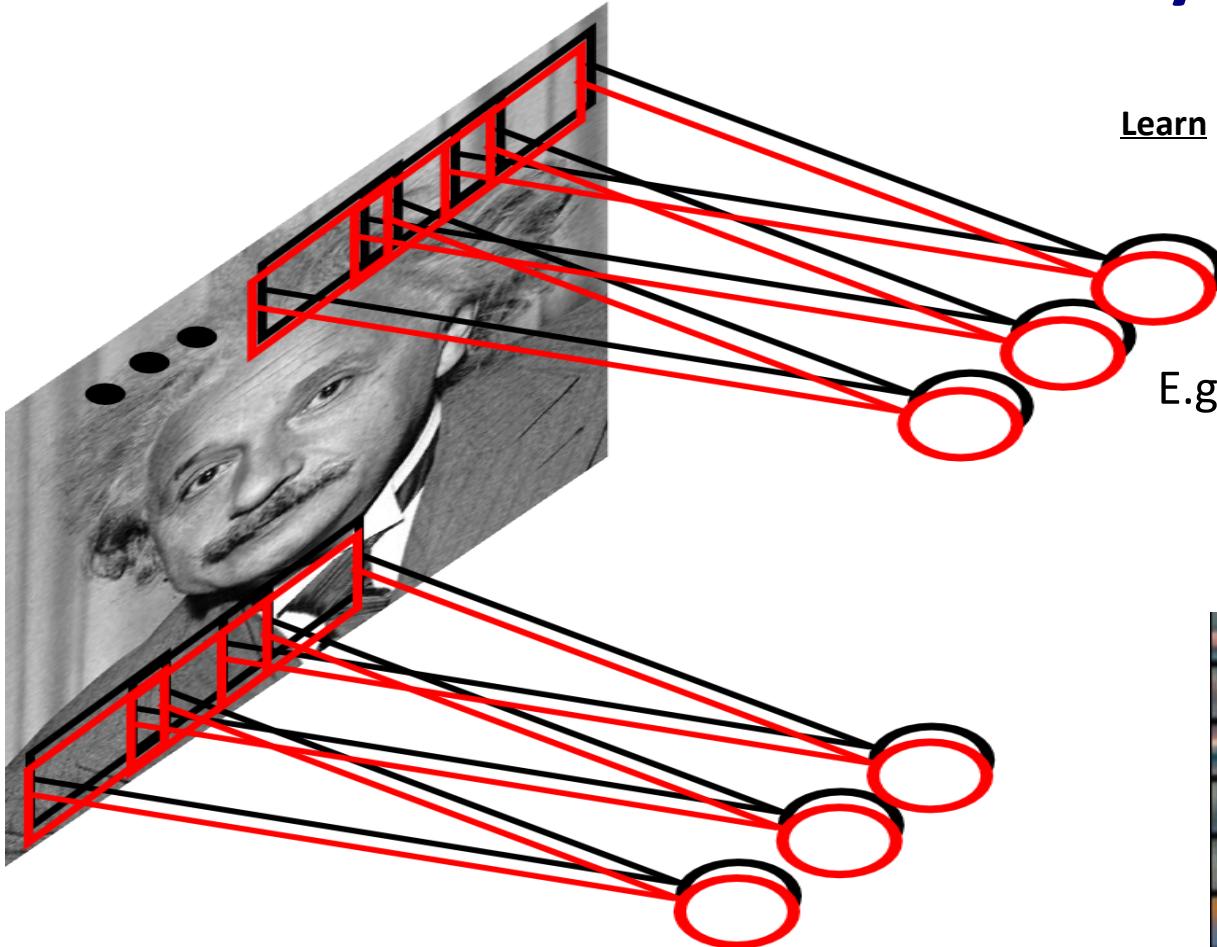
Example: 200x200 image  
40K hidden units  
Filter size: 10x10  
4M parameters

# Convolutional Layer



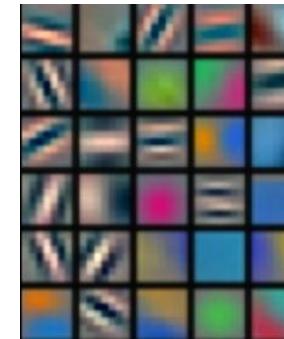
Share the same parameters across different locations (assuming input is stationary):  
**Convolutions with learned kernels**

# Convolutional Layer

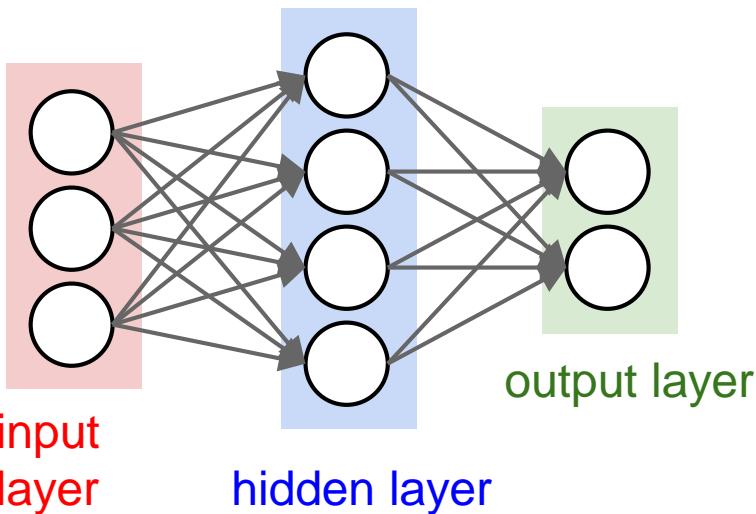


Learn multiple filters.

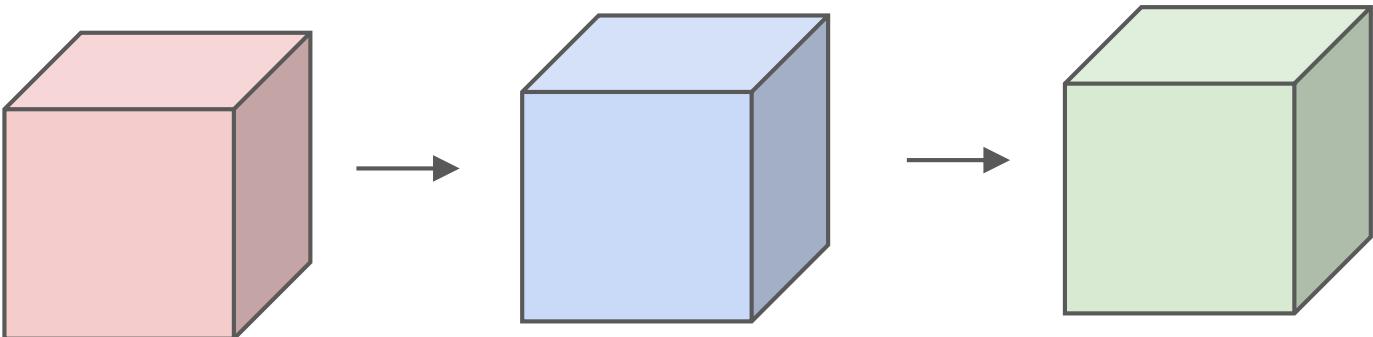
E.g.: 200x200 image  
100 Filters  
Filter size: 10x10  
10K parameters



before:

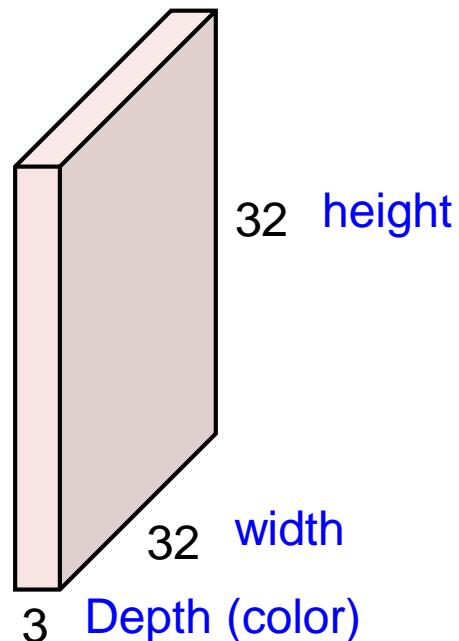


now:



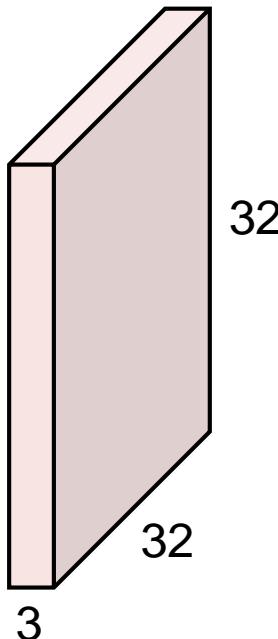
# Convolution Layer

32x32x3 image



# Convolution Layer

32x32x3 image



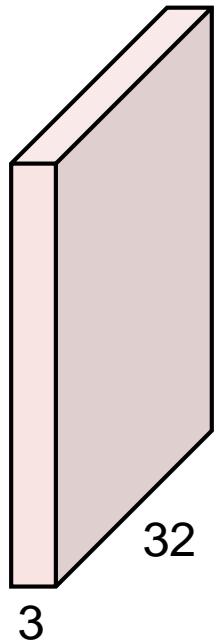
5x5x3 filter



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

# Convolution Layer

32x32x3 image



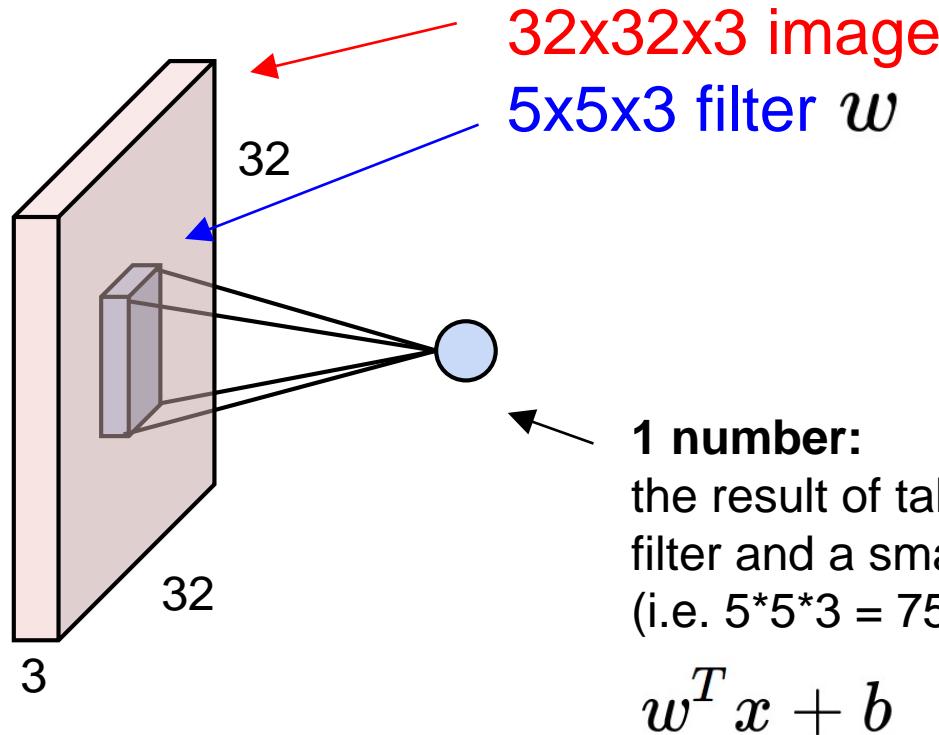
5x5x3 filter



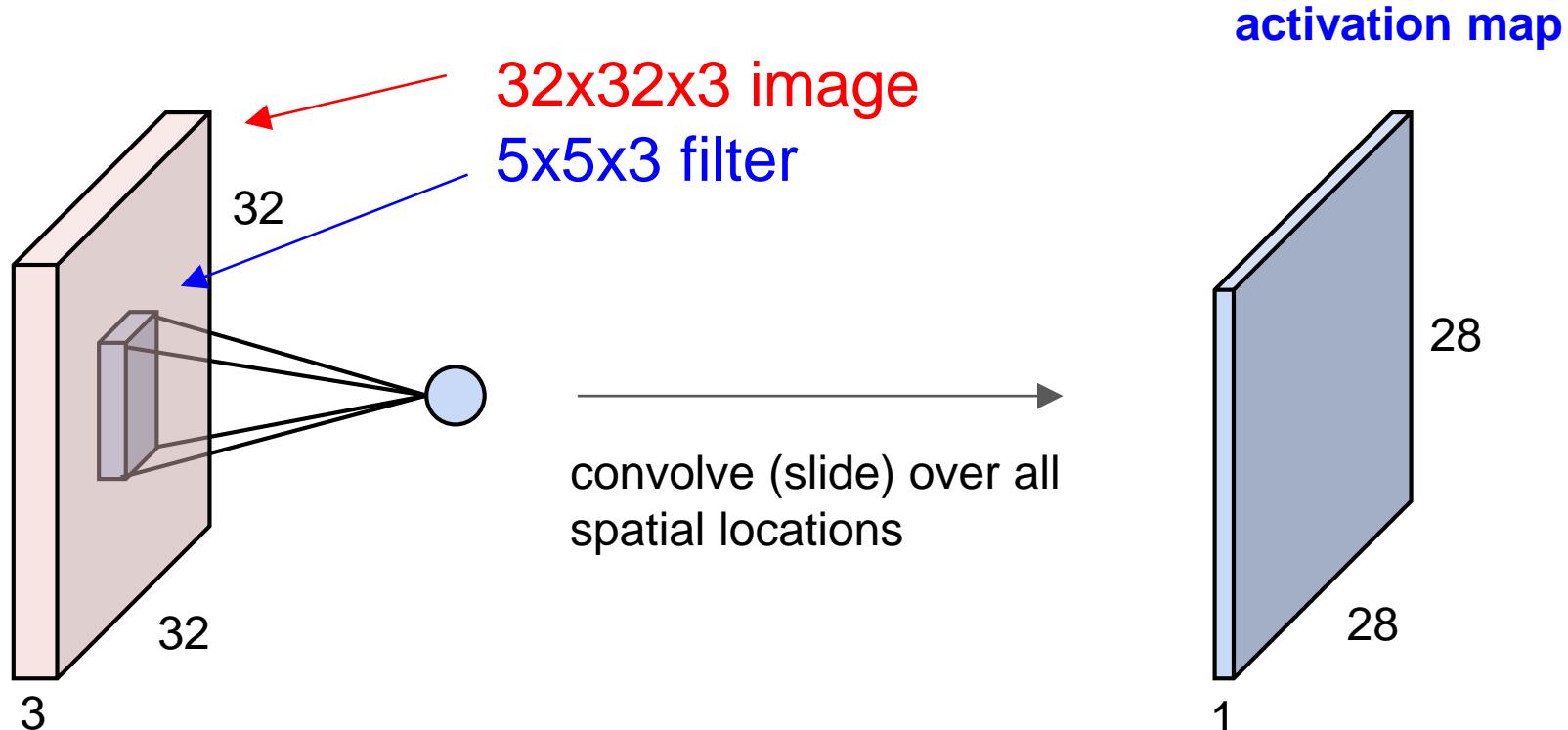
Filters always extend to the full depth of the input volume

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

# Convolution Layer



# Convolution Layer

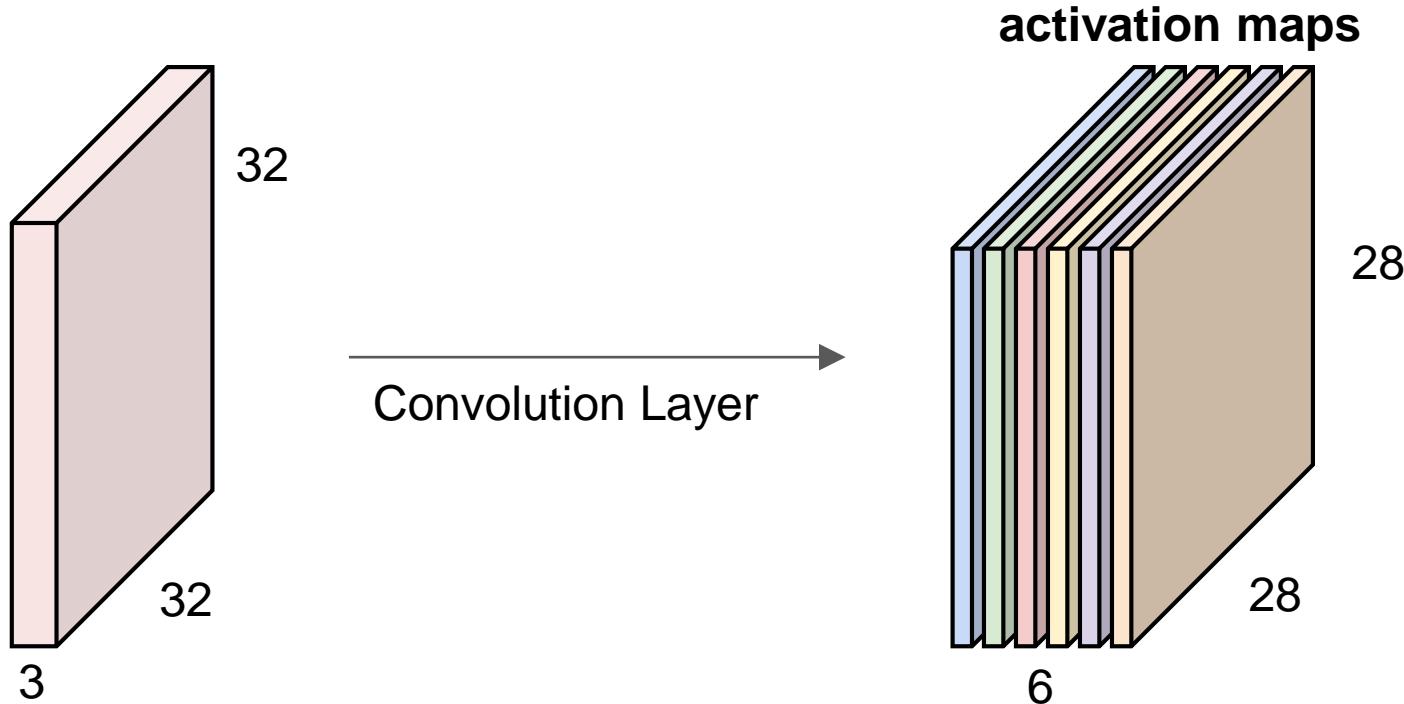


# Convolution Layer

consider a second, green filter

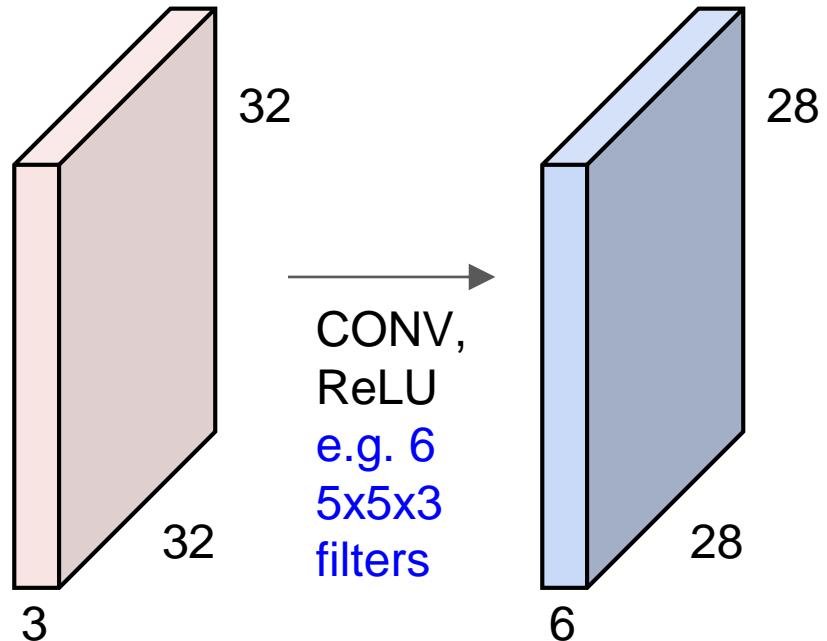


For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

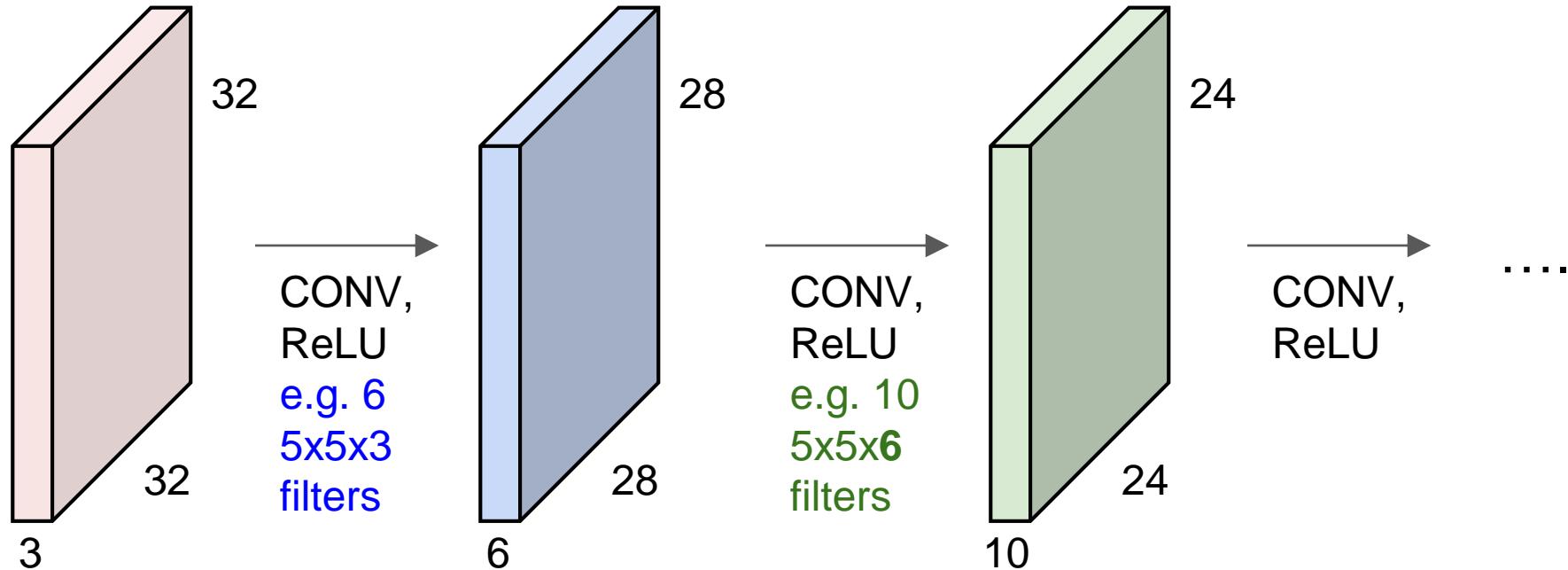


We stack these up to get a “new image” of size  $28 \times 28 \times 6$ !

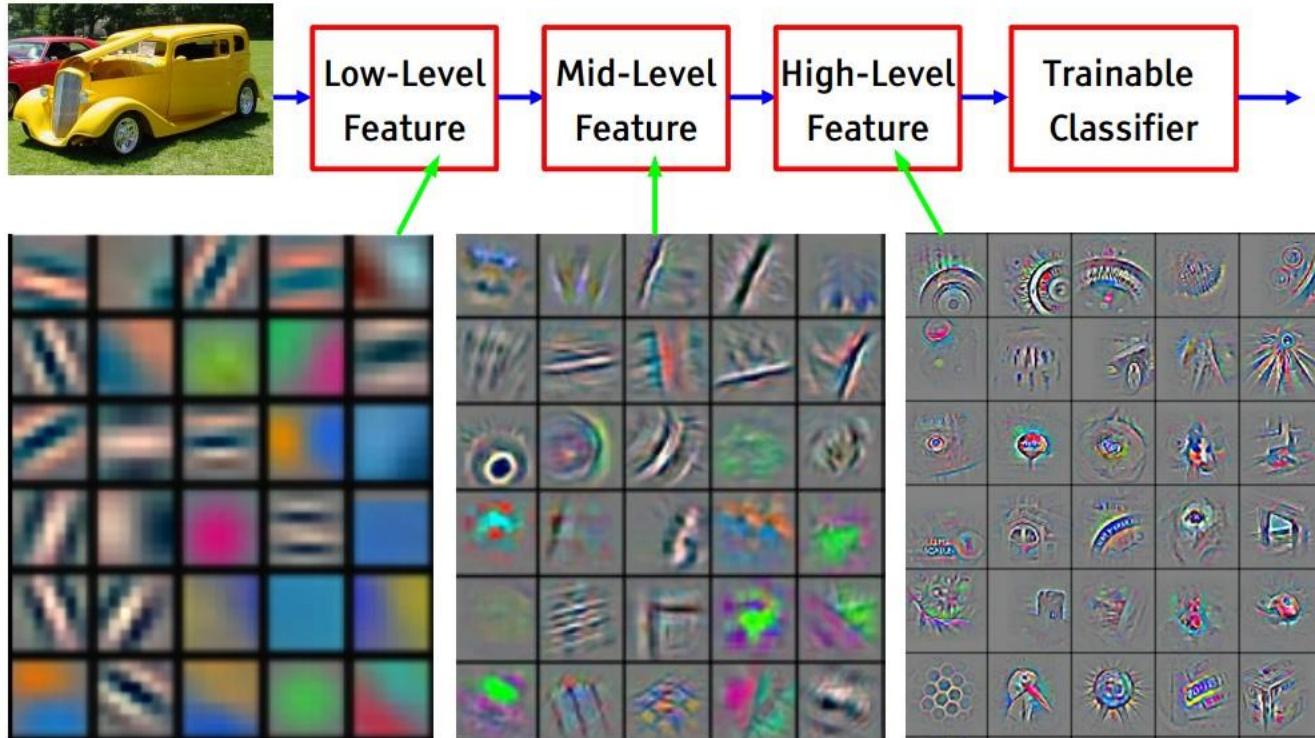
**Preview:** ConvNet is a sequence of Convolution Layers, interspersed with non-linear activation functions



**Preview:** ConvNet is a sequence of Convolution Layers, interspersed with non-linear activation functions



# Recall



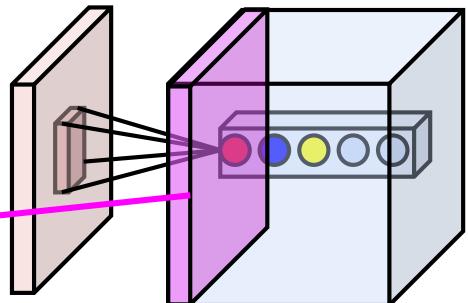
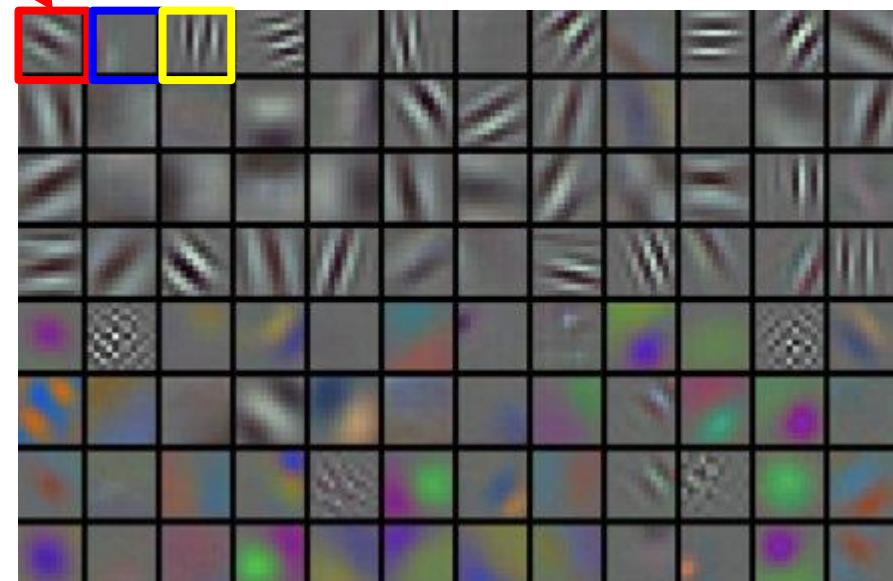
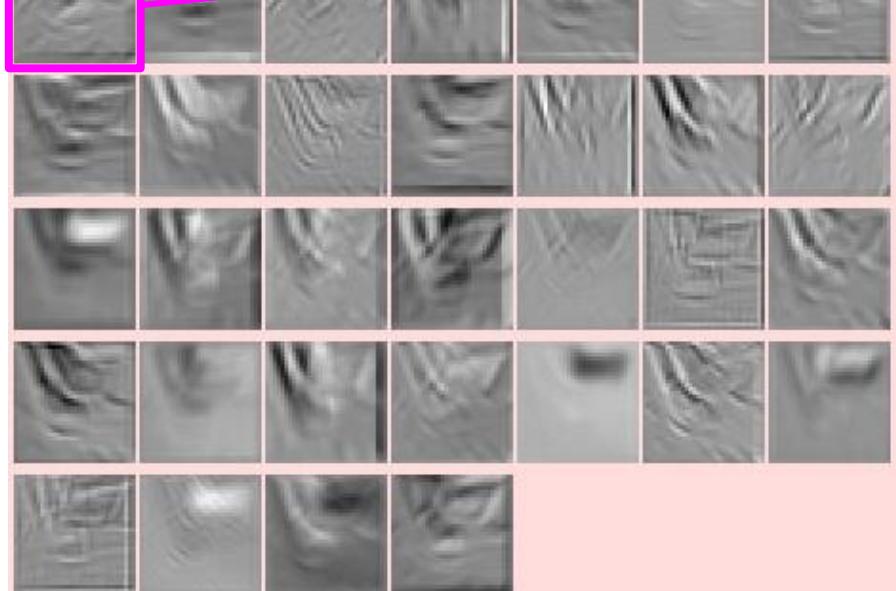
Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

Activations:

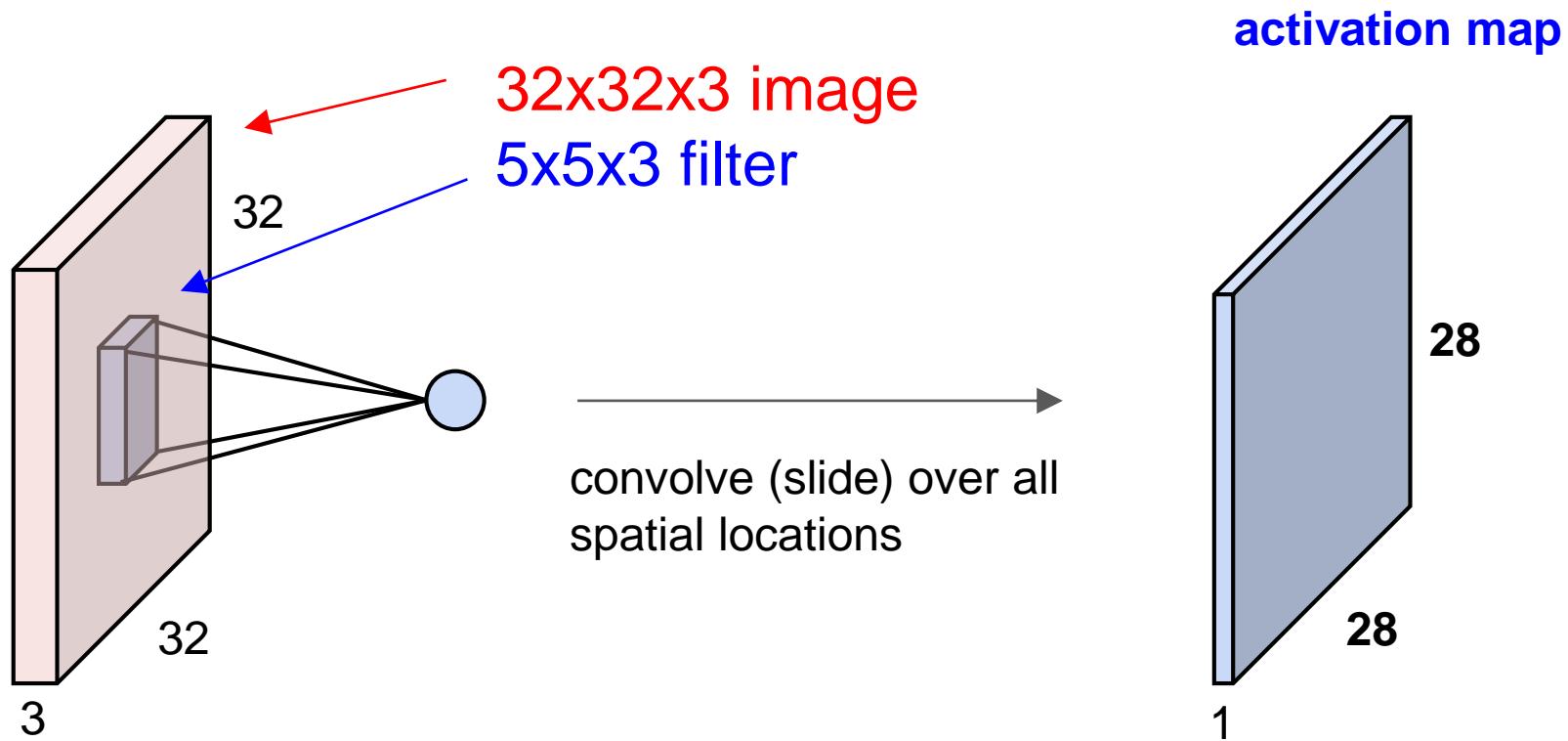


convolving the first filter in the input gives  
the first slice of depth in output volume

Activations:

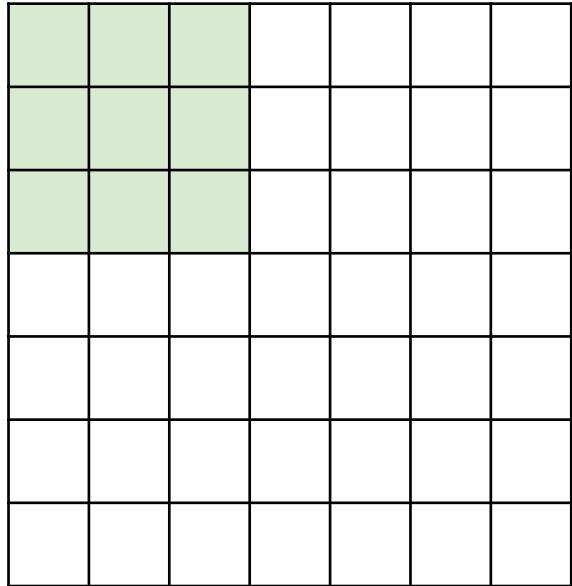


## A closer look at spatial dimensions:



## A closer look at spatial dimensions:

7

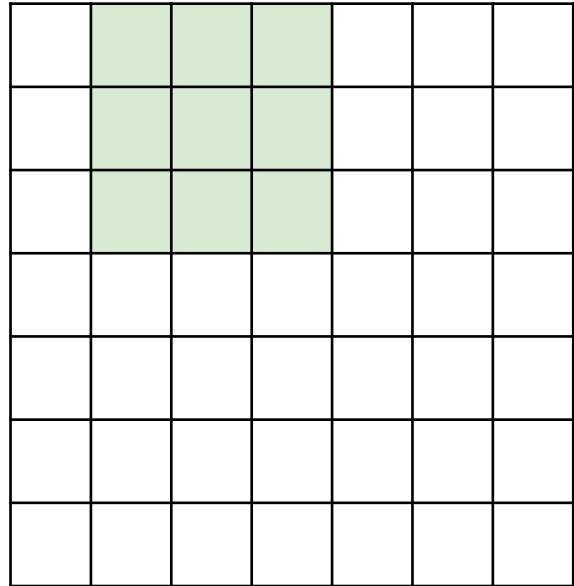


7x7 input (spatially)  
assume 3x3 filter

7

## A closer look at spatial dimensions:

7

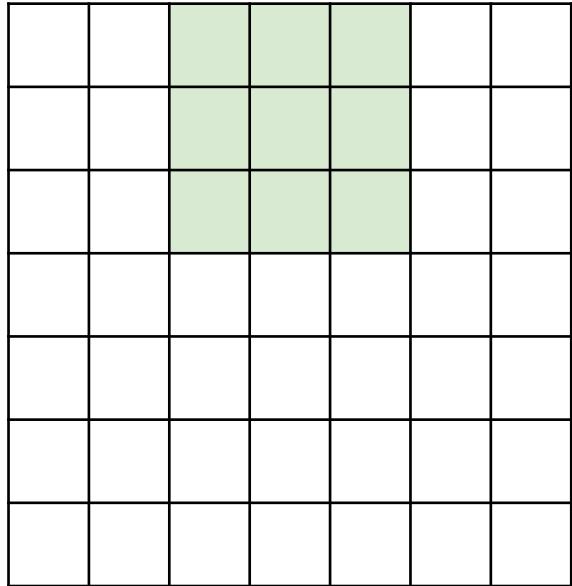


7x7 input (spatially)  
assume 3x3 filter

7

## A closer look at spatial dimensions:

7

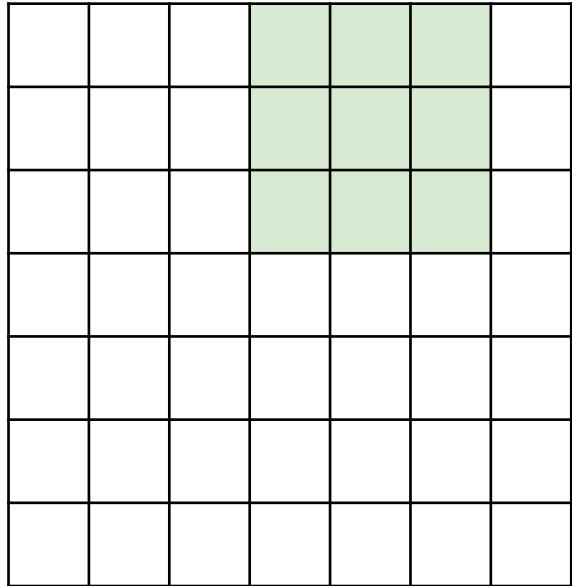


7x7 input (spatially)  
assume 3x3 filter

7

## A closer look at spatial dimensions:

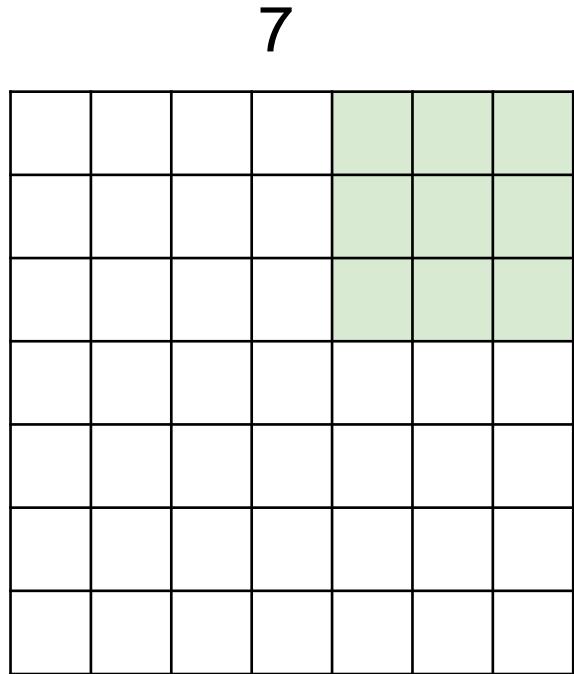
7



7x7 input (spatially)  
assume 3x3 filter

7

## A closer look at spatial dimensions:



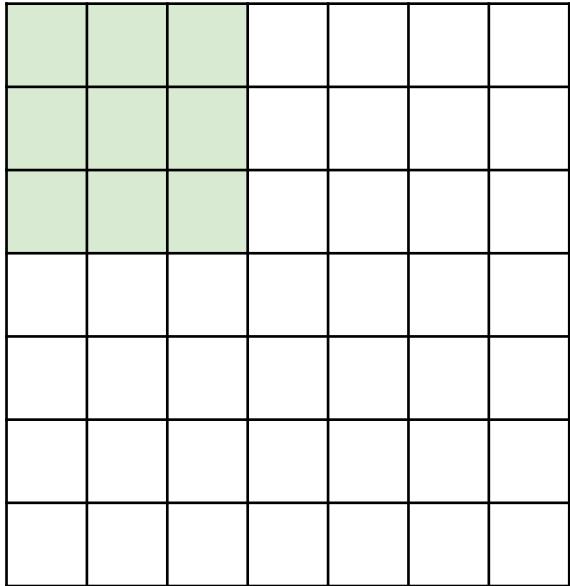
7x7 input (spatially)  
assume 3x3 filter

**=> 5x5 output**

7

## A closer look at spatial dimensions:

7

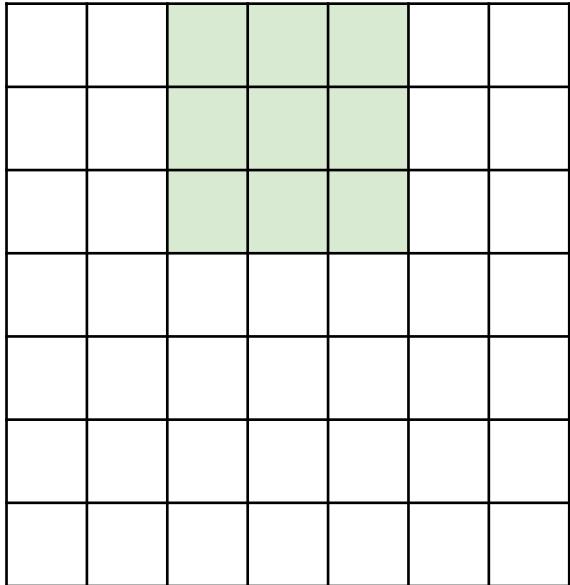


7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 2**

7

## A closer look at spatial dimensions:

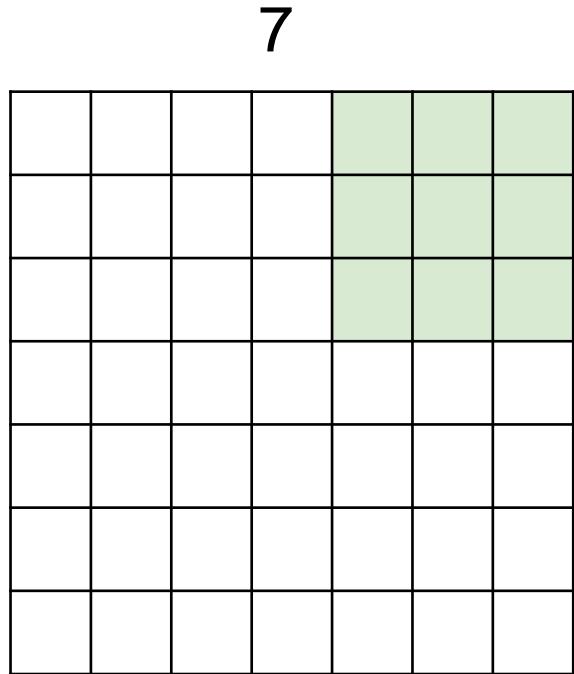
7



7

7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 2**

## A closer look at spatial dimensions:

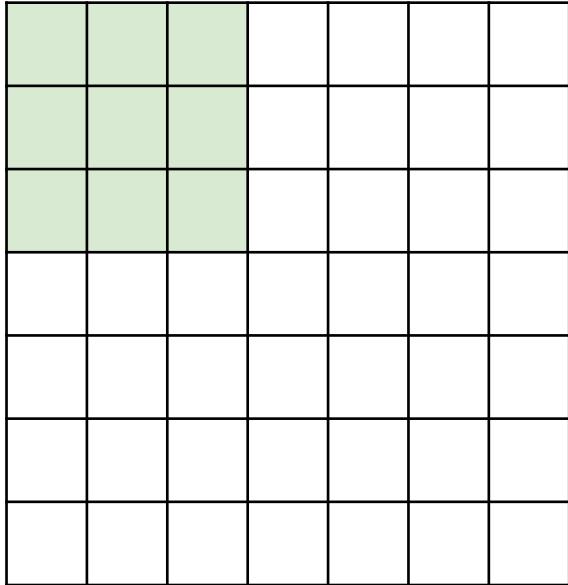


7

7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 2**  
**=> 3x3 output!**

## A closer look at spatial dimensions:

7

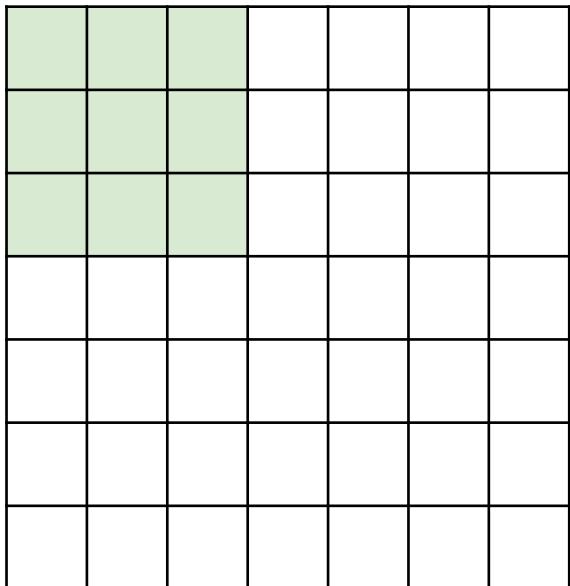


7

7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 3?**

## A closer look at spatial dimensions:

7

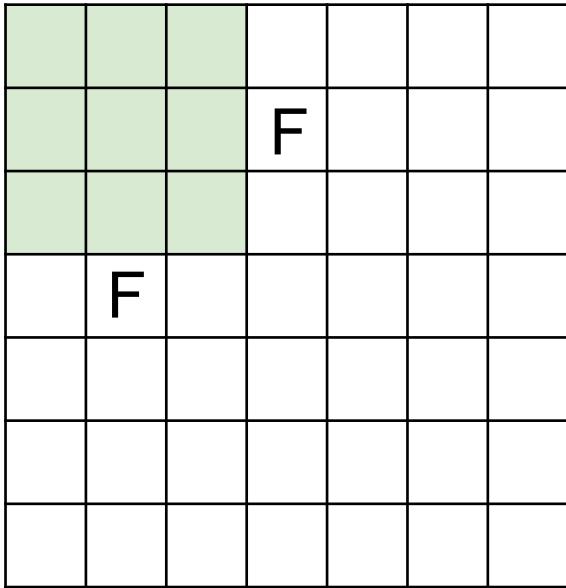


7

7x7 input (spatially)  
assume 3x3 filter  
applied **with stride 3?**

**doesn't fit!**  
Wastes some input pixels.

N



N

Output size:

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

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

$$\text{stride } 1 \Rightarrow (7 - 3)/1 + 1 = 5$$

$$\text{stride } 2 \Rightarrow (7 - 3)/2 + 1 = 3$$

$$\text{stride } 3 \Rightarrow (7 - 3)/3 + 1 = 2.33 :\backslash$$

# In practice: Common to zero pad the border

0	0	0	0	0	0		
0							
0							
0							
0							

Aside: remove image mean first!

e.g. input 7x7

3x3 filter, applied with **stride 1**

**pad with 1 pixel border => what is the output?**

(recall:)

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

# In practice: Common to zero pad the border

0	0	0	0	0	0		
0							
0							
0							
0							

e.g. input 7x7

**3x3 filter, applied with stride 1**

**pad with 1 pixel border => what is the output?**

**7x7 output!**

# In practice: Common to zero pad the border

0	0	0	0	0	0		
0							
0							
0							
0							

e.g. input 7x7

3x3 filter, applied with **stride 1**

**pad with 1 pixel border => what is the output?**

**7x7 output!**

in general, common to see CONV layers with stride 1, filters of size FxF, and zero-padding with  $(F-1)/2$ . (will preserve size spatially)

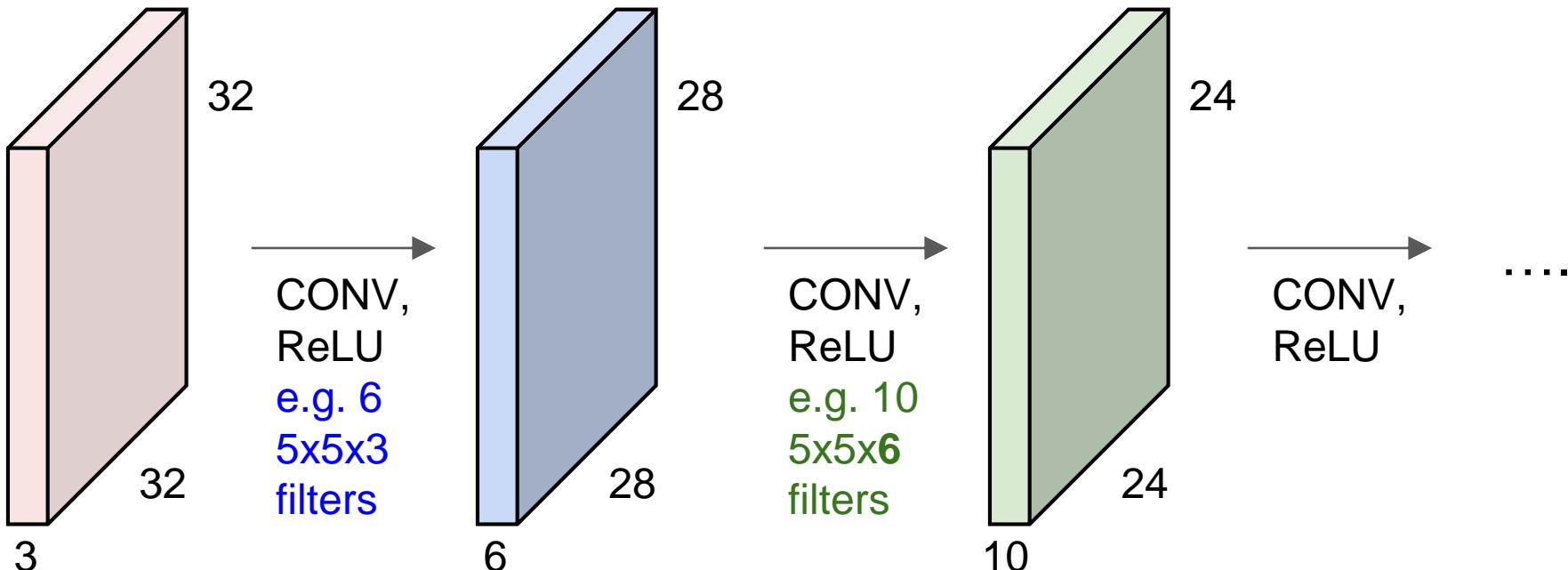
e.g.  $F = 3 \Rightarrow$  zero pad with 1

$F = 5 \Rightarrow$  zero pad with 2

$F = 7 \Rightarrow$  zero pad with 3

## Remember back to...

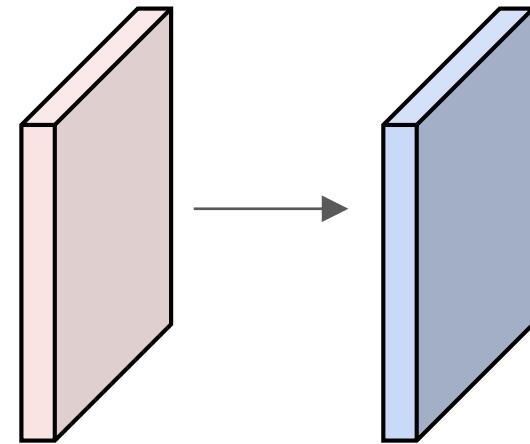
E.g. 32x32 input convolved repeatedly with 5x5 filters shrinks volumes spatially!  
(32 -> 28 -> 24 ...). Shrinking too fast is not good, doesn't work well.



# Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2

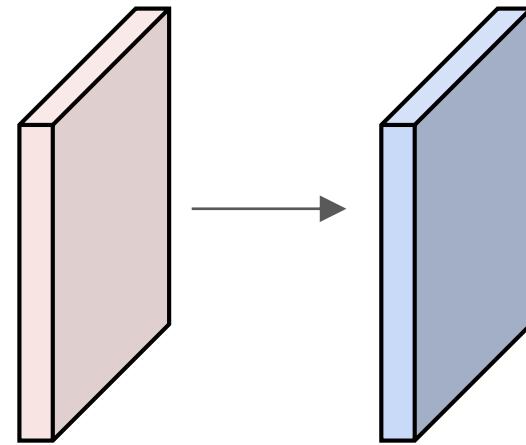


Output volume size: ?

Examples time:

Input volume: **32x32x3**

**10 5x5** filters with stride 1, pad 2



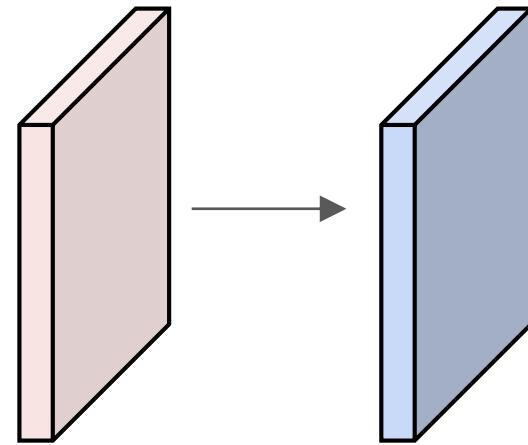
Output volume size:

$(32+2*2-5)/1+1 = 32$  spatially, so  
**32x32x10**

# Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2

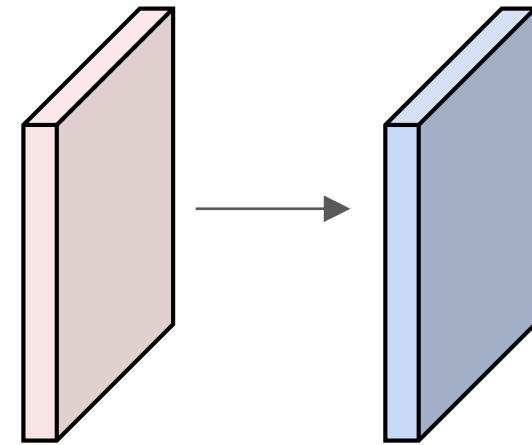


Number of parameters in this layer?

# Examples time:

Input volume: **32x32x3**

**10 5x5** filters with stride 1, pad 2



Number of parameters in this layer?

each filter has **5\*5\*3 + 1 = 76** params

(+1 for bias)

$$\Rightarrow \text{76} * \text{10} = \text{760}$$

**Summary.** To summarize, the Conv Layer:

- Accepts a volume of size  $W_1 \times H_1 \times D_1$
- Requires four hyperparameters:
  - Number of filters  $K$ ,
  - their spatial extent  $F$ ,
  - the stride  $S$ ,
  - the amount of zero padding  $P$ .
- Produces a volume of size  $W_2 \times H_2 \times D_2$  where:
  - $W_2 = (W_1 - F + 2P)/S + 1$
  - $H_2 = (H_1 - F + 2P)/S + 1$  (i.e. width and height are computed equally by symmetry)
  - $D_2 = K$
- With parameter sharing, it introduces  $F \cdot F \cdot D_1$  weights per filter, for a total of  $(F \cdot F \cdot D_1) \cdot K$  weights and  $K$  biases.
- In the output volume, the  $d$ -th depth slice (of size  $W_2 \times H_2$ ) is the result of performing a valid convolution of the  $d$ -th filter over the input volume with a stride of  $S$ , and then offset by  $d$ -th bias.

## Common settings:

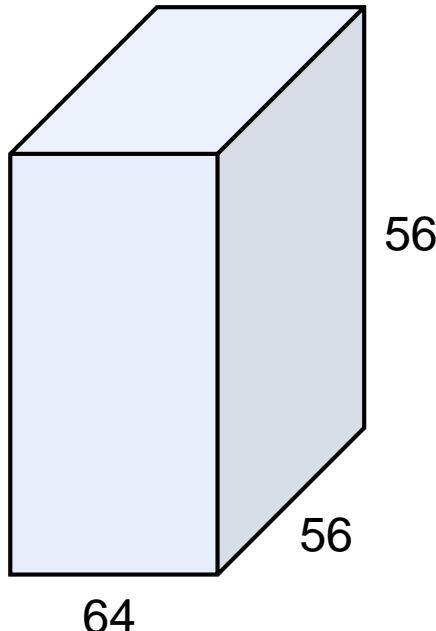
**Summary.** To summarize, the Conv Layer:

- Accepts a volume of size  $W_1 \times H_1 \times D_1$
- Requires four hyperparameters:
  - Number of filters  $K$ ,
  - their spatial extent  $F$ ,
  - the stride  $S$ ,
  - the amount of zero padding  $P$ .
- Produces a volume of size  $W_2 \times H_2 \times D_2$  where:
  - $W_2 = (W_1 - F + 2P)/S + 1$
  - $H_2 = (H_1 - F + 2P)/S + 1$  (i.e. width and height are computed equally by symmetry)
  - $D_2 = K$
- With parameter sharing, it introduces  $F \cdot F \cdot D_1$  weights per filter, for a total of  $(F \cdot F \cdot D_1) \cdot K$  weights and  $K$  biases.
- In the output volume, the  $d$ -th depth slice (of size  $W_2 \times H_2$ ) is the result of performing a valid convolution of the  $d$ -th filter over the input volume with a stride of  $S$ , and then offset by  $d$ -th bias.

$K = (\text{powers of 2, e.g. } 32, 64, 128, 512)$

- $F = 3, S = 1, P = 1$
- $F = 5, S = 1, P = 2$
- $F = 5, S = 2, P = ? \text{ (whatever fits)}$
- $F = 1, S = 1, P = 0$

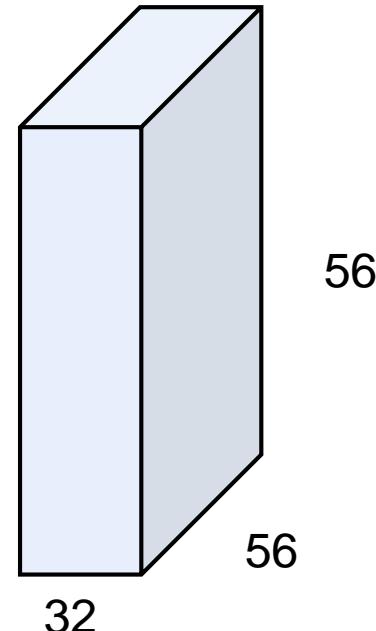
(btw, 1x1 convolution layers make perfect sense)



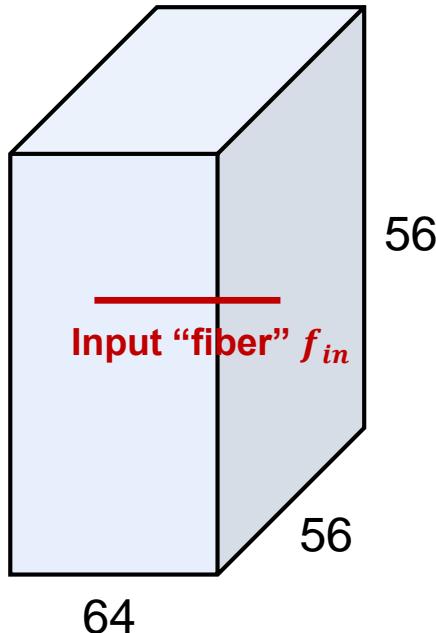
1x1 CONV  
with 32 filters

---

(each filter has size  
1x1x64, and performs a  
64-dimensional dot  
product)

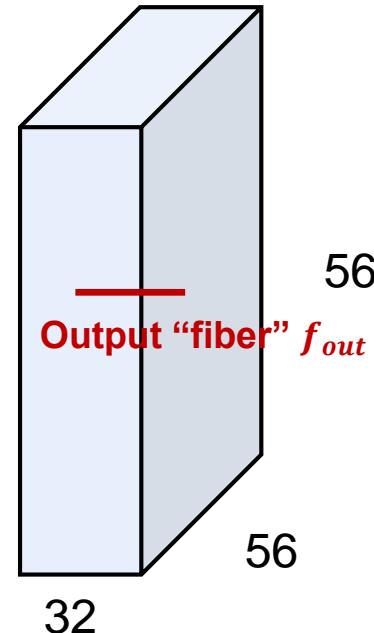


(btw, 1x1 convolution layers make perfect sense)

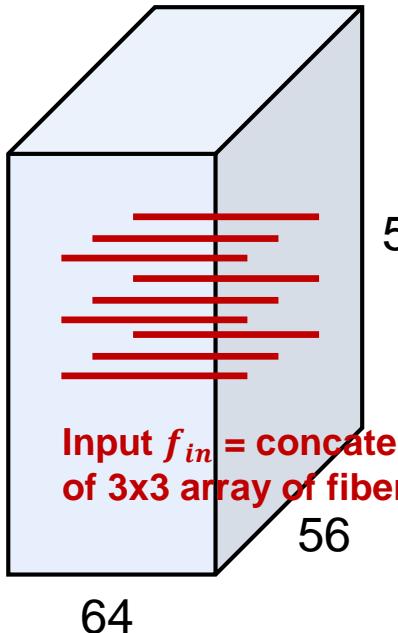


1x1 CONV  
with 32 filters

$f_{out} = Mf_{in}$  where  $M$   
is a  $32 \times 64$  matrix

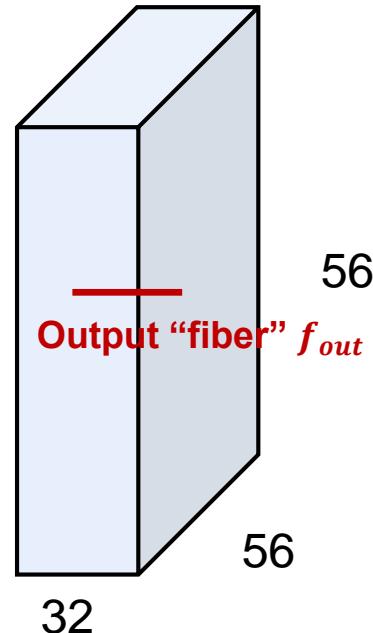


## Aside: convolution via matrix multiply: im2col

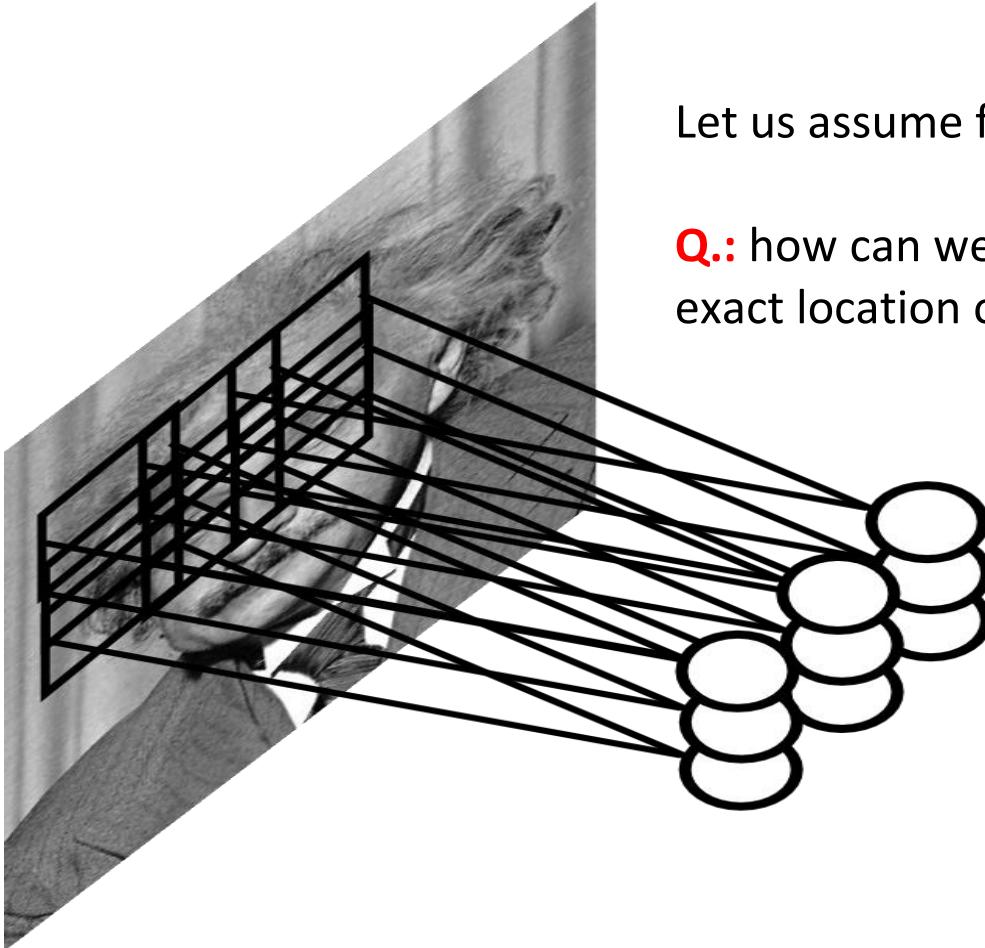


1x1 CONV  
with 32 filters

$f_{out} = Mf_{in}$  where  $M$   
is a  $32 \times (64 \times 9)$  matrix



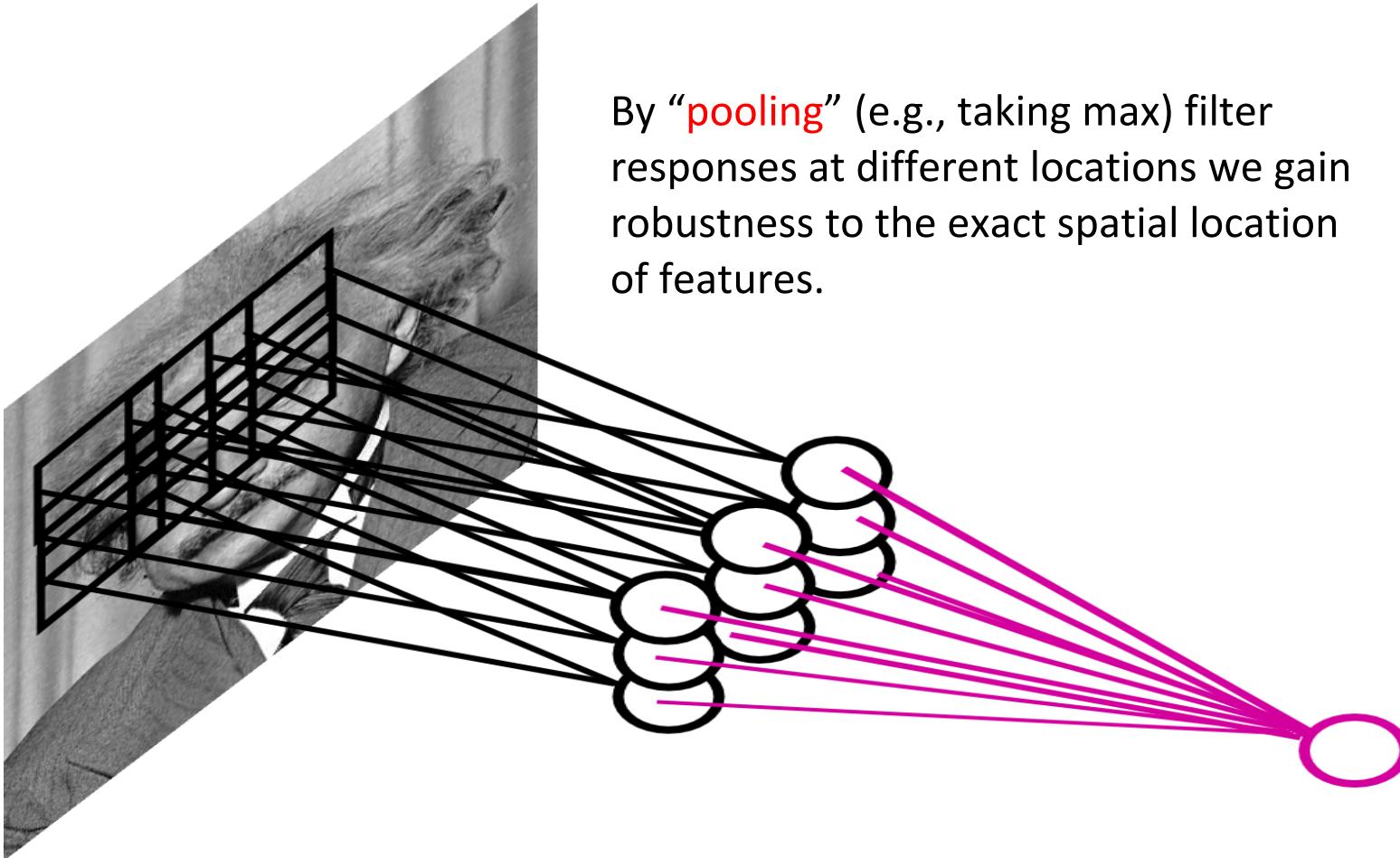
# Pooling Layer



Let us assume filter is an “eye” detector.

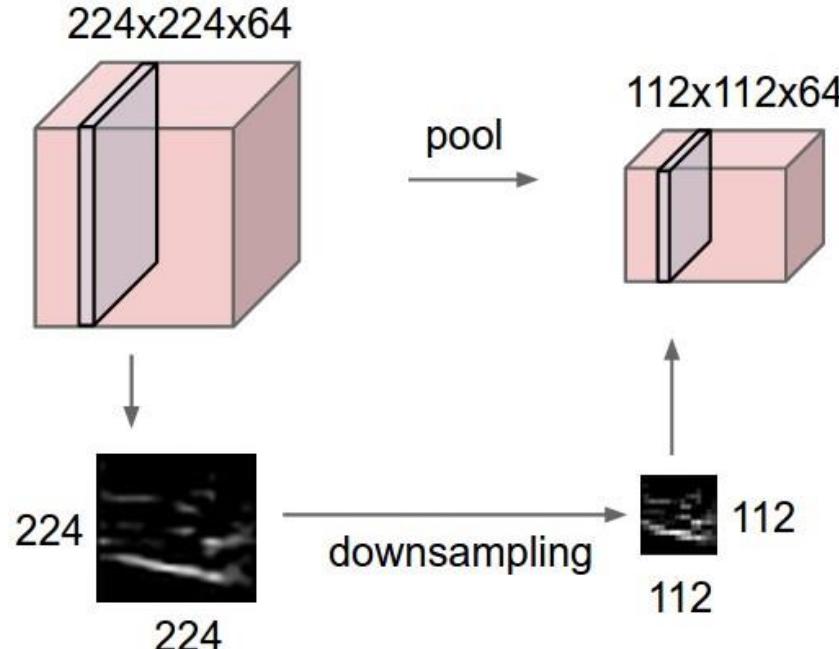
**Q.:** how can we make the detection robust to the exact location of the eye?

# Pooling Layer

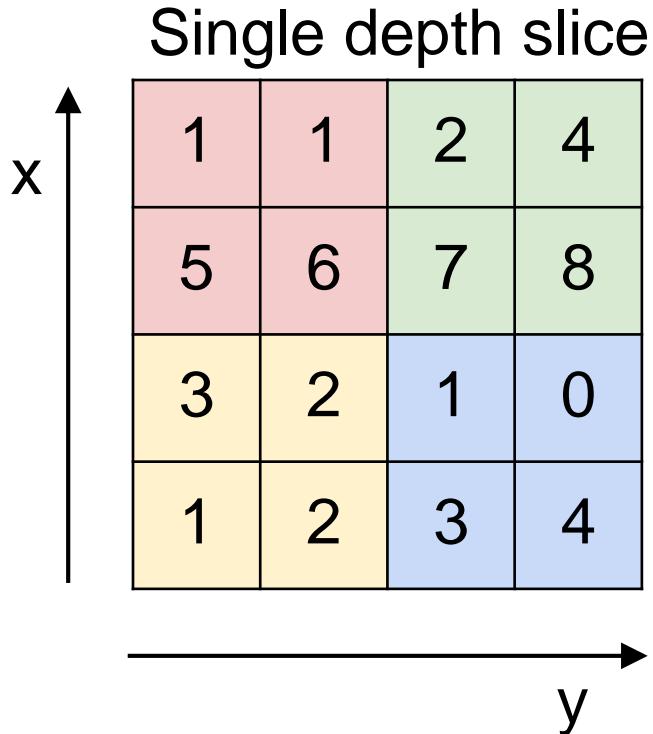


# Pooling layer

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



# MAX POOLING

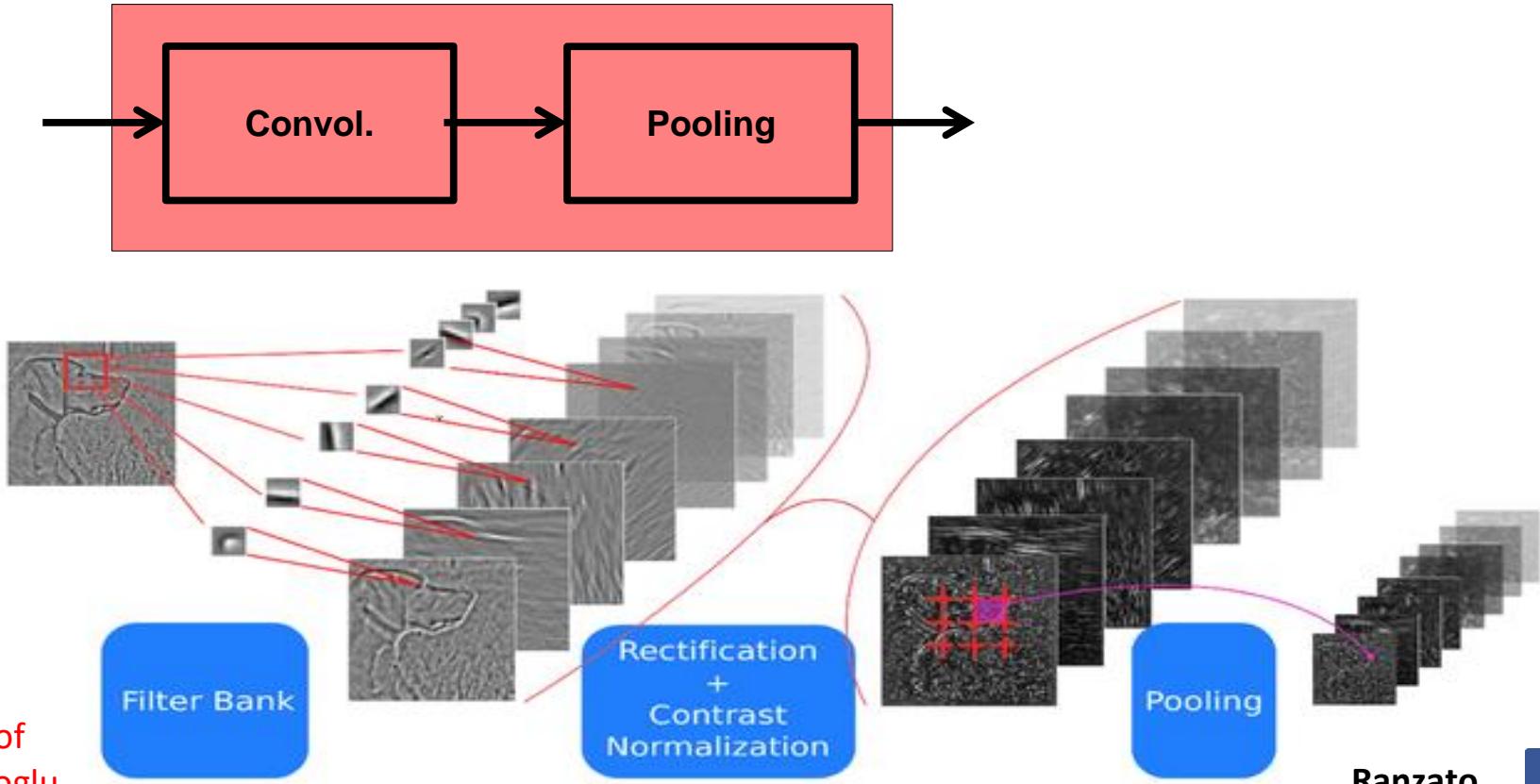


max pool with 2x2 filters  
and stride 2

6	8
3	4

# ConvNets: Typical Stage

One stage (zoom)

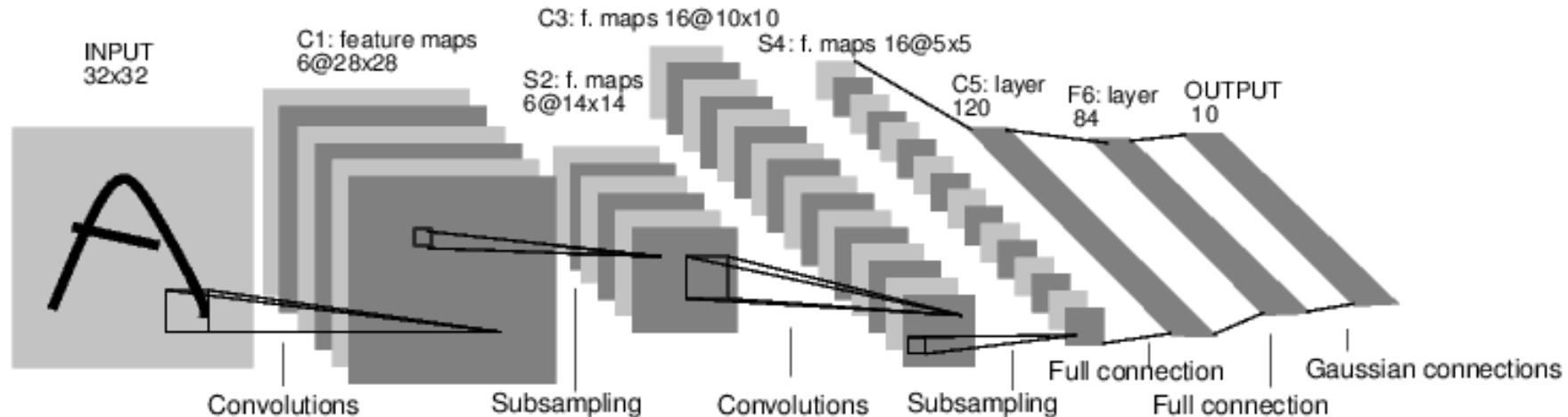


## [ConvNetJS demo: training on CIFAR-10]

<http://cs.stanford.edu/people/karpathy/convnetjs/demo/cifar10.html>

# Case Study: LeNet-5

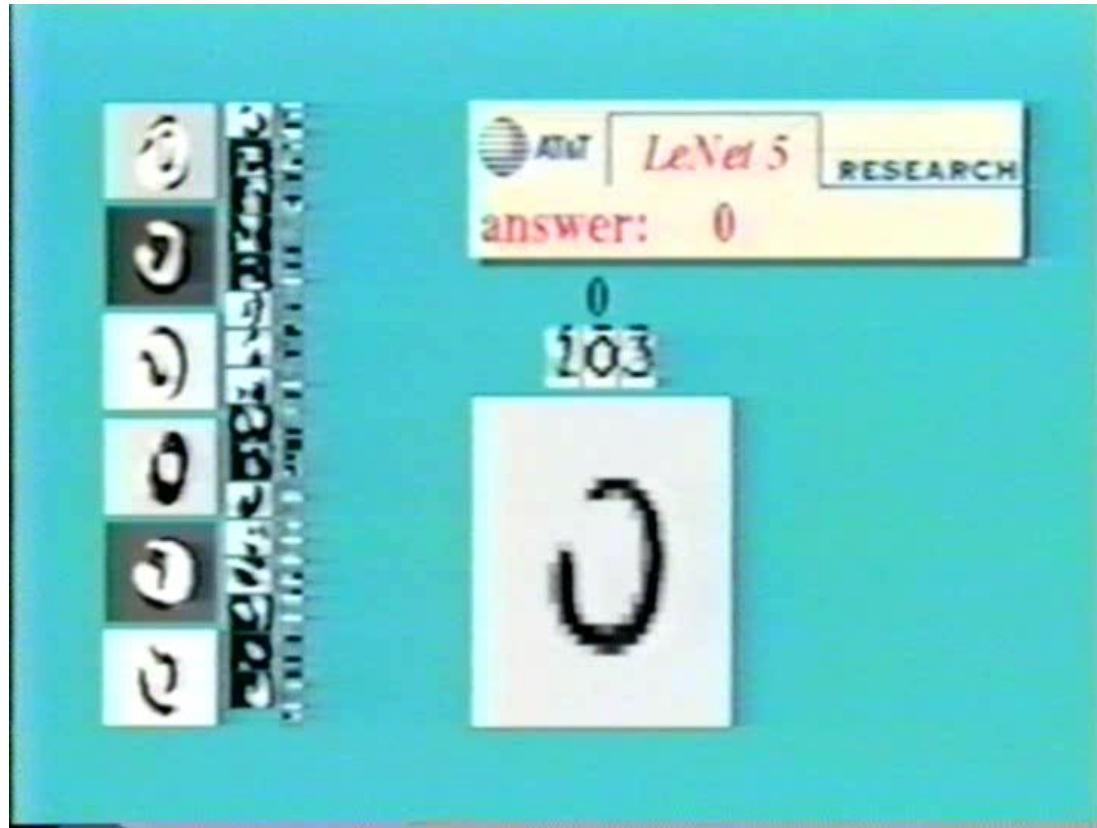
[LeCun et al., 1998]



Conv filters were  $5 \times 5$ , applied at stride 1

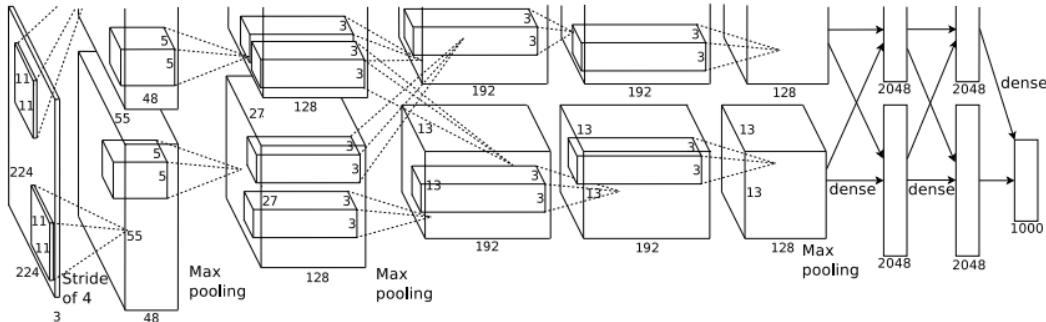
Subsampling (Pooling) layers were  $2 \times 2$  applied at stride 2  
i.e. architecture is [CONV-POOL-CONV-POOL-CONV-FC]

# Handwritten digit classification



# Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

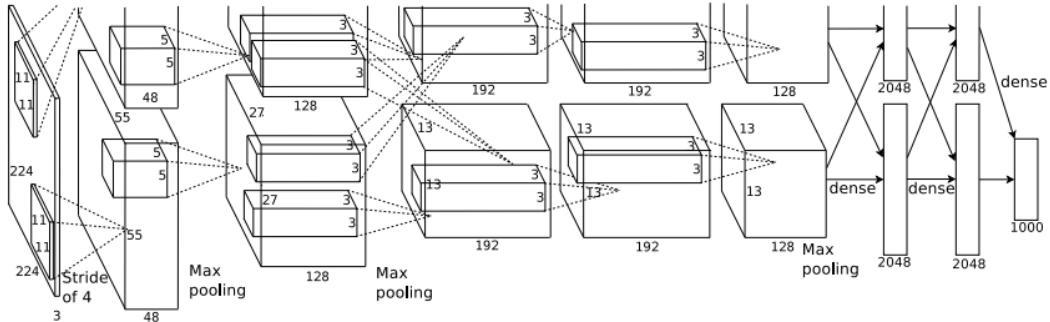
**First layer (CONV1):** 96 11x11 filters applied at stride 4

=>

Q: what is the output volume size? Hint:  $(227-11)/4+1 = 55$

# Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

**First layer (CONV1):** 96 11x11 filters applied at stride 4

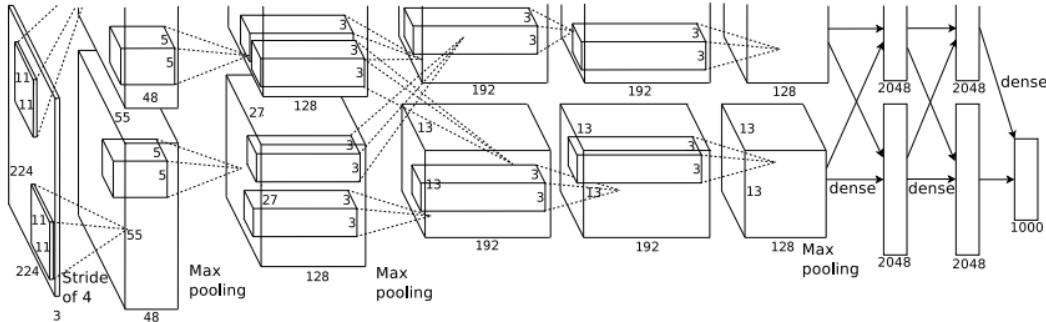
=>

Output volume **[55x55x96]**

Q: What is the total number of parameters in this layer?

# Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

**First layer (CONV1):** 96 11x11 filters applied at stride 4

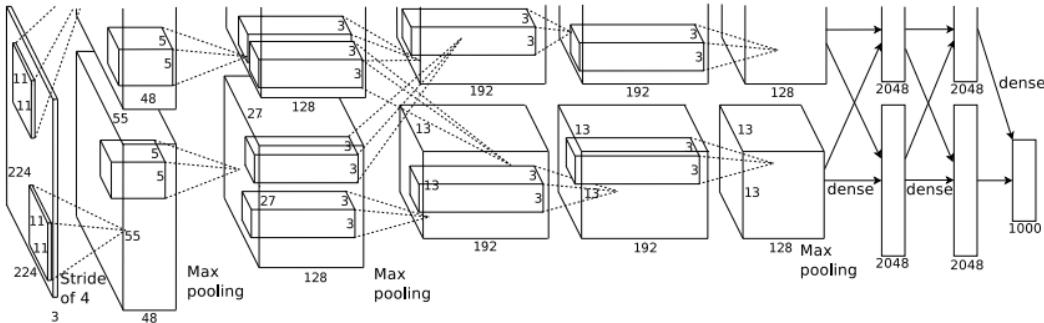
=>

Output volume **[55x55x96]**

Parameters:  $(11 \times 11 \times 3) \times 96 = 35\text{K}$

# Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

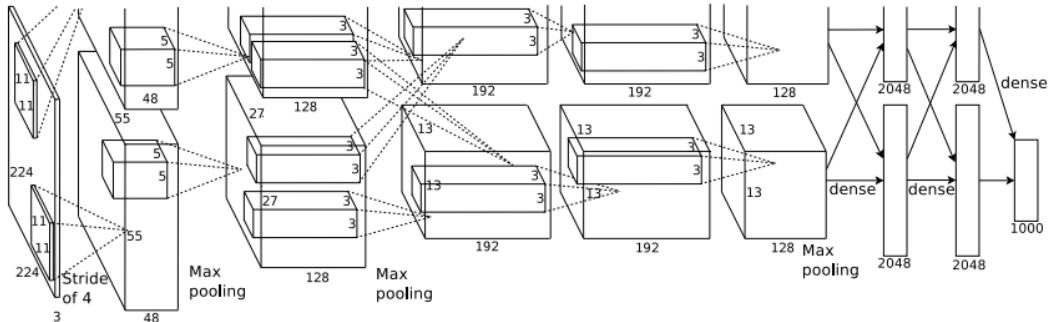
After CONV1: 55x55x96

**Second layer (POOL1): 3x3 filters applied at stride 2**

Q: what is the output volume size? Hint:  $(55-3)/2+1 = 27$

# Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

After CONV1: 55x55x96

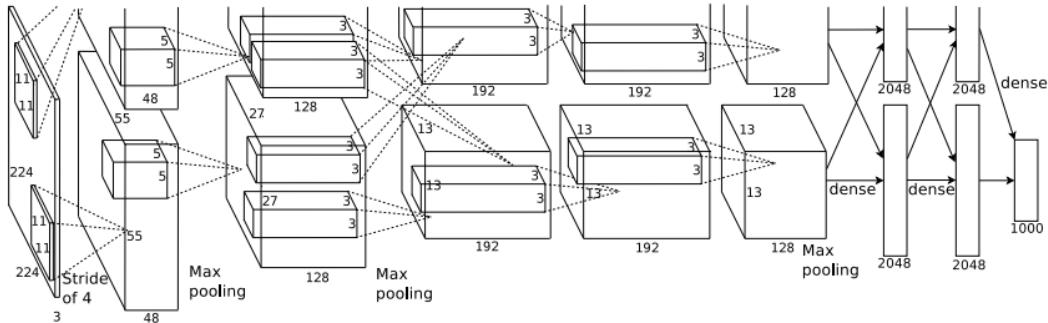
**Second layer (POOL1):** 3x3 filters applied at stride 2

Output volume: 27x27x96

Q: what is the number of parameters in this layer?

# Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

After CONV1: 55x55x96

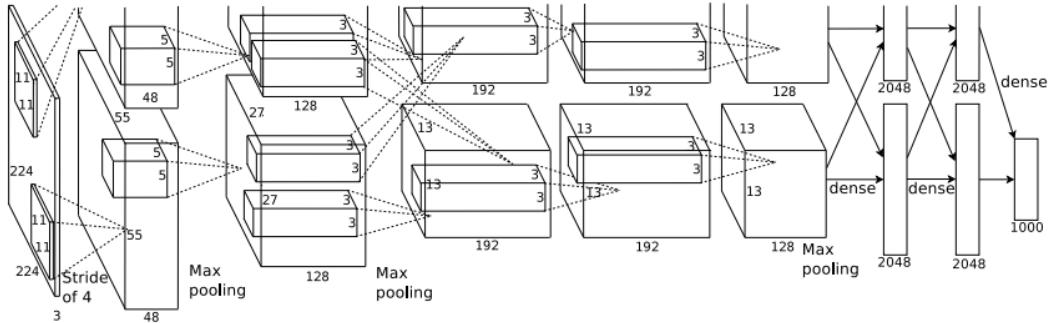
**Second layer (POOL1):** 3x3 filters applied at stride 2

Output volume: 27x27x96

Parameters: 0!

# Case Study: AlexNet

[Krizhevsky et al. 2012]



Input: 227x227x3 images

After CONV1: 55x55x96

After POOL1: 27x27x96

...

# Case Study: AlexNet

[Krizhevsky et al. 2012]

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

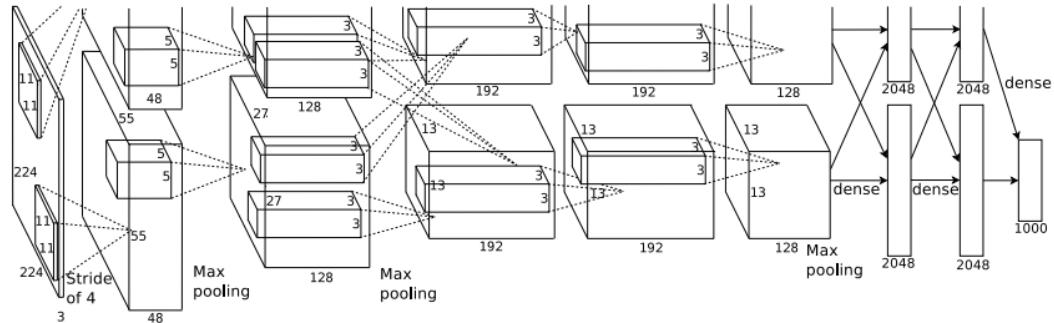
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1

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

[4096] FC6: 4096 neurons

[4096] FC7: 4096 neurons

[1000] FC8: 1000 neurons (class scores)



# Case Study: AlexNet

[Krizhevsky et al. 2012]

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

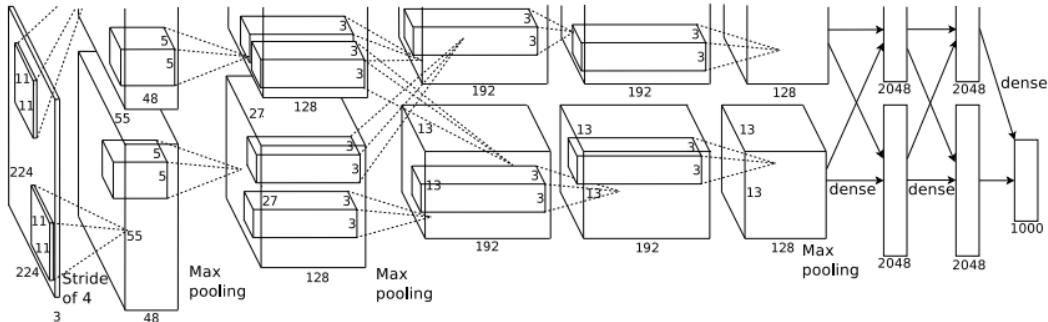
[13x13x256] CONV5: 256 3x3 filters at stride 1, pad 1

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

[4096] FC6: 4096 neurons

[4096] FC7: 4096 neurons

[1000] FC8: 1000 neurons (class scores)



## 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%

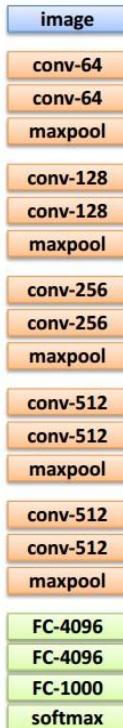
“You need a lot of data if you want to  
train/use CNNs”

# Transfer Learning

“You need a lot of data if you want to  
train the CNN”

**NOT  
ALWAYS**

# Transfer Learning with CNNs

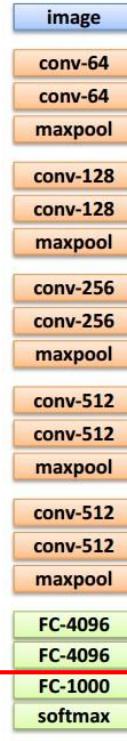


1. Train on  
Imagenet

# Transfer Learning with CNNs



1. Train on  
Imagenet



2. If small dataset: fix  
all weights (treat CNN  
as fixed feature  
extractor), retrain only  
the classifier

i.e. swap the Softmax  
layer at the end

# Transfer Learning with CNNs

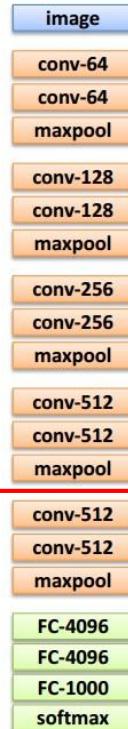


1. Train on  
Imagenet



2. If small dataset: fix  
all weights (treat CNN  
as fixed feature  
extractor), retrain only  
the classifier

i.e. swap the Softmax  
layer at the end



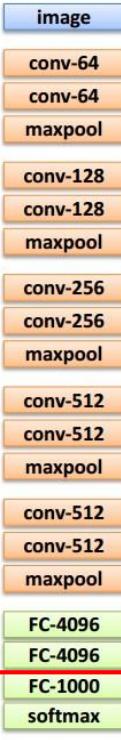
3. If you have medium sized  
dataset, “**finetune**”  
instead: use the old weights  
as initialization, train the full  
network or only some of the  
higher layers

retrain bigger portion of the  
network, or even all of it.

# Transfer Learning with CNNs



1. Train on  
Imagenet



2. If small dataset: fix  
all weights (treat CNN  
as fixed feature  
extractor), retrain only  
the classifier

i.e. swap the Softmax  
layer at the end



3. If you have medium sized  
dataset, “**finetune**”  
instead: use the old weights  
as initialization, train the full  
network or only some of the  
higher layers

retrain bigger portion of the  
network, or even all of it.

tip: use only ~1/10th of  
the original learning rate  
in finetuning to player,  
and ~1/100th on  
intermediate layers

# Case Study: VGGNet

[Simonyan and Zisserman, 2014]

Only 3x3 CONV stride 1, pad 1  
and 2x2 MAX POOL stride 2

best model

11.2% top 5 error in ILSVRC 2013

->

7.3% top 5 error

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224 × 224 RGB image)					
conv3-64	conv3-64 LRN	conv3-64 <b>conv3-64</b>	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 <b>conv3-128</b>	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 <b>conv3-256</b> <b>conv1-256</b>	conv3-256 conv3-256 <b>conv3-256</b>	conv3-256 conv3-256 conv3-256 <b>conv3-256</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

Table 2: Number of parameters (in millions).

Network	A,A-LRN	B	C	D	E
Number of parameters	133	133	134	138	144

INPUT: [224x224x3] memory: 224\*224\*3=150K params: 0 (not counting biases)

CONV3-64: [224x224x64] memory: 224\*224\*64=3.2M params:  $(3*3*3)*64 = 1,728$

CONV3-64: [224x224x64] memory: 224\*224\*64=3.2M params:  $(3*3*64)*64 = 36,864$

POOL2: [112x112x64] memory: 112\*112\*64=800K params: 0

CONV3-128: [112x112x128] memory: 112\*112\*128=1.6M params:  $(3*3*64)*128 = 73,728$

CONV3-128: [112x112x128] memory: 112\*112\*128=1.6M params:  $(3*3*128)*128 = 147,456$

POOL2: [56x56x128] memory: 56\*56\*128=400K params: 0

CONV3-256: [56x56x256] memory: 56\*56\*256=800K params:  $(3*3*128)*256 = 294,912$

CONV3-256: [56x56x256] memory: 56\*56\*256=800K params:  $(3*3*256)*256 = 589,824$

CONV3-256: [56x56x256] memory: 56\*56\*256=800K params:  $(3*3*256)*256 = 589,824$

POOL2: [28x28x256] memory: 28\*28\*256=200K params: 0

CONV3-512: [28x28x512] memory: 28\*28\*512=400K params:  $(3*3*256)*512 = 1,179,648$

CONV3-512: [28x28x512] memory: 28\*28\*512=400K params:  $(3*3*512)*512 = 2,359,296$

CONV3-512: [28x28x512] memory: 28\*28\*512=400K params:  $(3*3*512)*512 = 2,359,296$

POOL2: [14x14x512] memory: 14\*14\*512=100K params: 0

CONV3-512: [14x14x512] memory: 14\*14\*512=100K params:  $(3*3*512)*512 = 2,359,296$

CONV3-512: [14x14x512] memory: 14\*14\*512=100K params:  $(3*3*512)*512 = 2,359,296$

CONV3-512: [14x14x512] memory: 14\*14\*512=100K params:  $(3*3*512)*512 = 2,359,296$

POOL2: [7x7x512] memory: 7\*7\*512=25K params: 0

FC: [1x1x4096] memory: 4096 params:  $7*7*512*4096 = 102,760,448$

FC: [1x1x4096] memory: 4096 params:  $4096*4096 = 16,777,216$

FC: [1x1x1000] memory: 1000 params:  $4096*1000 = 4,096,000$

ConvNet Configuration			
B	C	D	E
13 weight layers	16 weight layers	16 weight layers	19 weight layers
put (224 × 224 RGB image)			
conv3-64	conv3-64	conv3-64	conv3-64
<b>conv3-64</b>	conv3-64	conv3-64	conv3-64
maxpool			
conv3-128	conv3-128	conv3-128	conv3-128
<b>conv3-128</b>	conv3-128	conv3-128	conv3-128
maxpool			
conv3-256	conv3-256	conv3-256	conv3-256
conv3-256	conv3-256	conv3-256	conv3-256
<b>conv1-256</b>		<b>conv3-256</b>	conv3-256
maxpool			
conv3-512	conv3-512	conv3-512	conv3-512
conv3-512	conv3-512	conv3-512	conv3-512
<b>conv1-512</b>		<b>conv3-512</b>	conv3-512
maxpool			
conv3-512	conv3-512	conv3-512	conv3-512
conv3-512	conv3-512	conv3-512	conv3-512
<b>conv1-512</b>		<b>conv3-512</b>	conv3-512
maxpool			
FC-4096			
FC-4096			
FC-1000			
soft-max			

INPUT: [224x224x3] memory: 224\*224\*3=150K params: 0 (not counting biases)

CONV3-64: [224x224x64] memory: 224\*224\*64=3.2M params:  $(3*3*3)*64 = 1,728$

CONV3-64: [224x224x64] memory: 224\*224\*64=3.2M params:  $(3*3*64)*64 = 36,864$

POOL2: [112x112x64] memory: 112\*112\*64=800K params: 0

CONV3-128: [112x112x128] memory: 112\*112\*128=1.6M params:  $(3*3*64)*128 = 73,728$

CONV3-128: [112x112x128] memory: 112\*112\*128=1.6M params:  $(3*3*128)*128 = 147,456$

POOL2: [56x56x128] memory: 56\*56\*128=400K params: 0

CONV3-256: [56x56x256] memory: 56\*56\*256=800K params:  $(3*3*128)*256 = 294,912$

CONV3-256: [56x56x256] memory: 56\*56\*256=800K params:  $(3*3*256)*256 = 589,824$

CONV3-256: [56x56x256] memory: 56\*56\*256=800K params:  $(3*3*256)*256 = 589,824$

POOL2: [28x28x256] memory: 28\*28\*256=200K params: 0

CONV3-512: [28x28x512] memory: 28\*28\*512=400K params:  $(3*3*256)*512 = 1,179,648$

CONV3-512: [28x28x512] memory: 28\*28\*512=400K params:  $(3*3*512)*512 = 2,359,296$

CONV3-512: [28x28x512] memory: 28\*28\*512=400K params:  $(3*3*512)*512 = 2,359,296$

POOL2: [14x14x512] memory: 14\*14\*512=100K params: 0

CONV3-512: [14x14x512] memory: 14\*14\*512=100K params:  $(3*3*512)*512 = 2,359,296$

CONV3-512: [14x14x512] memory: 14\*14\*512=100K params:  $(3*3*512)*512 = 2,359,296$

CONV3-512: [14x14x512] memory: 14\*14\*512=100K params:  $(3*3*512)*512 = 2,359,296$

POOL2: [7x7x512] memory: 7\*7\*512=25K params: 0

FC: [1x1x4096] memory: 4096 params:  $7*7*512*4096 = 102,760,448$

FC: [1x1x4096] memory: 4096 params:  $4096*4096 = 16,777,216$

FC: [1x1x1000] memory: 1000 params:  $4096*1000 = 4,096,000$

**TOTAL memory:** 24M \* 4 bytes  $\approx$  93MB / image (only forward!  $\sim 2$  for bwd)

**TOTAL params:** 138M parameters

ConvNet Configuration			
B	C	D	E
13 weight layers	16 weight layers	16 weight layers	19 weight layers
put (224 × 224 RGB image)			
conv3-64	conv3-64	conv3-64	conv3-64
<b>conv3-64</b>	conv3-64	conv3-64	conv3-64
maxpool			
conv3-128	conv3-128	conv3-128	conv3-128
<b>conv3-128</b>	conv3-128	conv3-128	conv3-128
maxpool			
conv3-256	conv3-256	conv3-256	conv3-256
conv3-256	conv3-256	conv3-256	conv3-256
<b>conv1-256</b>		<b>conv3-256</b>	conv3-256
maxpool			
conv3-512	conv3-512	conv3-512	conv3-512
conv3-512	conv3-512	conv3-512	conv3-512
<b>conv1-512</b>		<b>conv3-512</b>	conv3-512
maxpool			
conv3-512	conv3-512	conv3-512	conv3-512
conv3-512	conv3-512	conv3-512	conv3-512
<b>conv1-512</b>		<b>conv3-512</b>	conv3-512
maxpool			
FC-4096			
FC-4096			
FC-1000			
soft-max			

INPUT: [224x224x3] memory: 224\*224\*3=150K params: 0 (not counting biases)

CONV3-64: [224x224x64] memory: **224\*224\*64=3.2M** params:  $(3*3*3)*64 = 1,728$

CONV3-64: [224x224x64] memory: **224\*224\*64=3.2M** params:  $(3*3*64)*64 = 36,864$

POOL2: [112x112x64] memory: 112\*112\*64=800K params: 0

CONV3-128: [112x112x128] memory: 112\*112\*128=1.6M params:  $(3*3*64)*128 = 73,728$

CONV3-128: [112x112x128] memory: 112\*112\*128=1.6M params:  $(3*3*128)*128 = 147,456$

POOL2: [56x56x128] memory: 56\*56\*128=400K params: 0

CONV3-256: [56x56x256] memory: 56\*56\*256=800K params:  $(3*3*128)*256 = 294,912$

CONV3-256: [56x56x256] memory: 56\*56\*256=800K params:  $(3*3*256)*256 = 589,824$

CONV3-256: [56x56x256] memory: 56\*56\*256=800K params:  $(3*3*256)*256 = 589,824$

POOL2: [28x28x256] memory: 28\*28\*256=200K params: 0

CONV3-512: [28x28x512] memory: 28\*28\*512=400K params:  $(3*3*256)*512 = 1,179,648$

CONV3-512: [28x28x512] memory: 28\*28\*512=400K params:  $(3*3*512)*512 = 2,359,296$

CONV3-512: [28x28x512] memory: 28\*28\*512=400K params:  $(3*3*512)*512 = 2,359,296$

POOL2: [14x14x512] memory: 14\*14\*512=100K params: 0

CONV3-512: [14x14x512] memory: 14\*14\*512=100K params:  $(3*3*512)*512 = 2,359,296$

CONV3-512: [14x14x512] memory: 14\*14\*512=100K params:  $(3*3*512)*512 = 2,359,296$

CONV3-512: [14x14x512] memory: 14\*14\*512=100K params:  $(3*3*512)*512 = 2,359,296$

POOL2: [7x7x512] memory: 7\*7\*512=25K params: 0

FC: [1x1x4096] memory: 4096 params:  $7*7*512*4096 = \text{102,760,448}$

FC: [1x1x4096] memory: 4096 params:  $4096*4096 = 16,777,216$

FC: [1x1x1000] memory: 1000 params:  $4096*1000 = 4,096,000$

Note:

Most memory is in early CONV

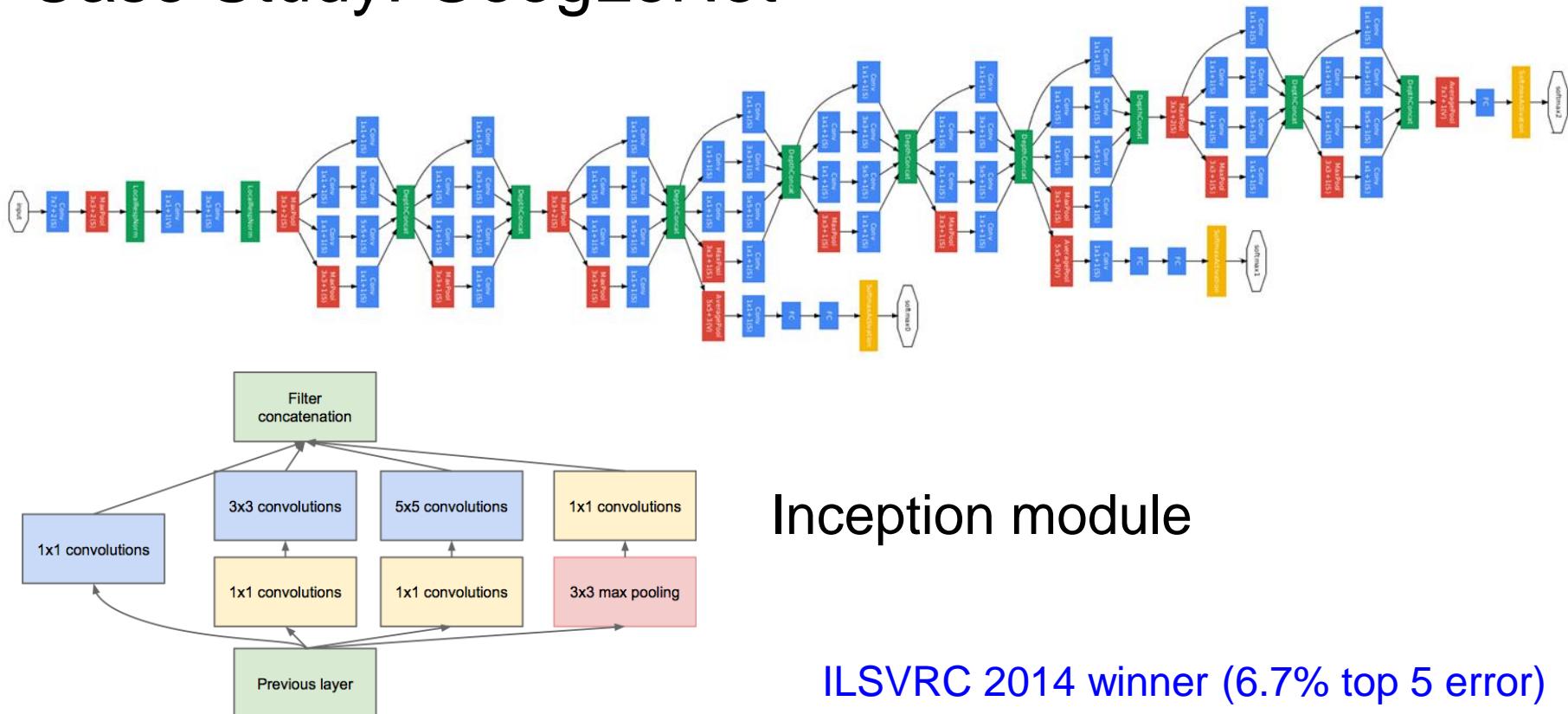
Most params are in late FC

**TOTAL memory: 24M \* 4 bytes ≈ 93MB / image (only forward! ~\*2 for bwd)**

**TOTAL params: 138M parameters**

# Case Study: GoogLeNet

[Szegedy et al., 2014]



# Case Study: GoogLeNet

type	patch size/ stride	output size	depth	#1×1	#3×3 reduce	#3×3	#5×5 reduce	#5×5	pool proj	params	ops
convolution	7×7/2	112×112×64	1							2.7K	34M
max pool	3×3/2	56×56×64	0								
convolution	3×3/1	56×56×192	2		64	192				112K	360M
max pool	3×3/2	28×28×192	0								
inception (3a)		28×28×256	2	64	96	128	16	32	32	159K	128M
inception (3b)		28×28×480	2	128	128	192	32	96	64	380K	304M
max pool	3×3/2	14×14×480	0								
inception (4a)		14×14×512	2	192	96	208	16	48	64	364K	73M
inception (4b)		14×14×512	2	160	112	224	24	64	64	437K	88M
inception (4c)		14×14×512	2	128	128	256	24	64	64	463K	100M
inception (4d)		14×14×528	2	112	144	288	32	64	64	580K	119M
inception (4e)		14×14×832	2	256	160	320	32	128	128	840K	170M
max pool	3×3/2	7×7×832	0								
inception (5a)		7×7×832	2	256	160	320	32	128	128	1072K	54M
inception (5b)		7×7×1024	2	384	192	384	48	128	128	1388K	71M
avg pool	7×7/1	1×1×1024	0								
dropout (40%)		1×1×1024	0								
linear		1×1×1000	1							1000K	1M
softmax		1×1×1000	0								

Fun features:

- Only 5 million params!  
(Removes FC layers completely)

Compared to AlexNet:

- 12X less params
- 2x more compute
- 6.67% (vs. 16.4%)

# Case Study: ResNet [He et al., 2015]

ILSVRC 2015 winner (3.6% top 5 error)



## MSRA @ ILSVRC & COCO 2015 Competitions

- **1st places in all five main tracks**

- ImageNet Classification: “Ultra-deep” (quote Yann) **152-layer nets**
- ImageNet Detection: **16%** better than 2nd
- ImageNet Localization: **27%** better than 2nd
- COCO Detection: **11%** better than 2nd
- COCO Segmentation: **12%** better than 2nd

\*improvements are relative numbers

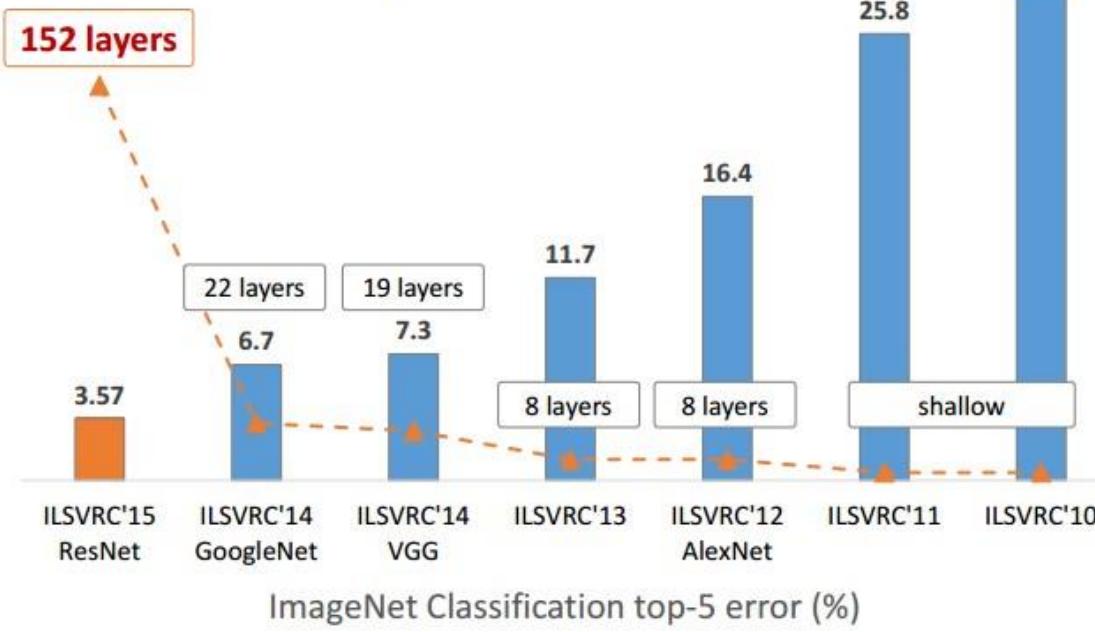


Kaiming He, Xiangyu Zhang, Shaoqing Ren, & Jian Sun. “Deep Residual Learning for Image Recognition”. arXiv 2015.

Slide from Kaiming He’s recent presentation <https://www.youtube.com/watch?v=1PGLj-uKT1w>

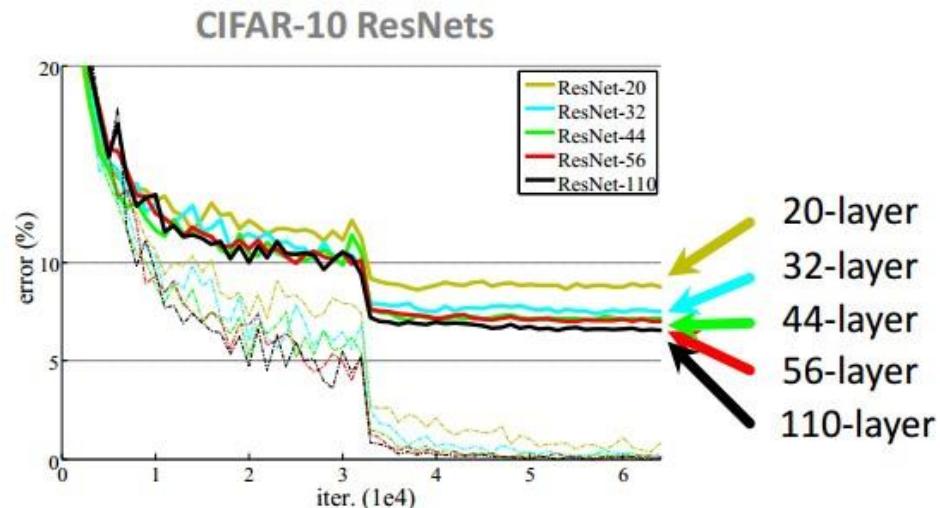
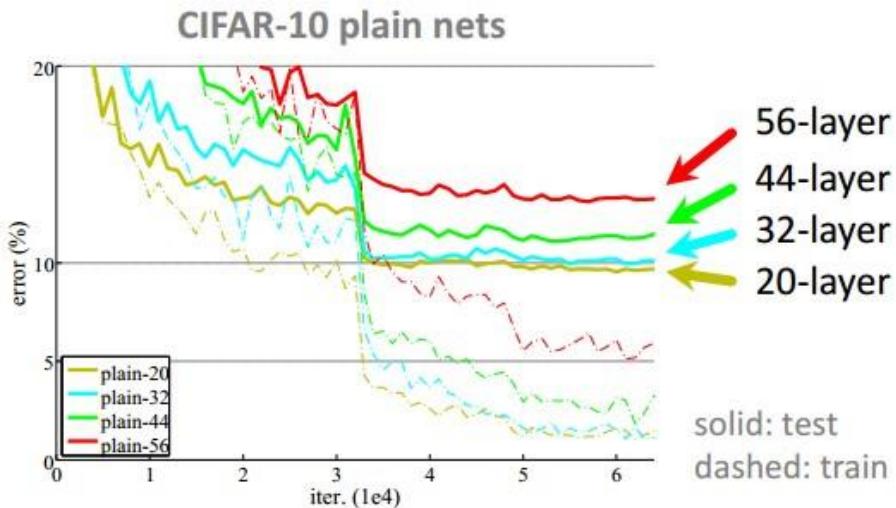
Slide based on cs231n by Fei-Fei Li & Andrej Karpathy & Justin Johnson

# Revolution of Depth



(slide from Kaiming He's recent presentation)

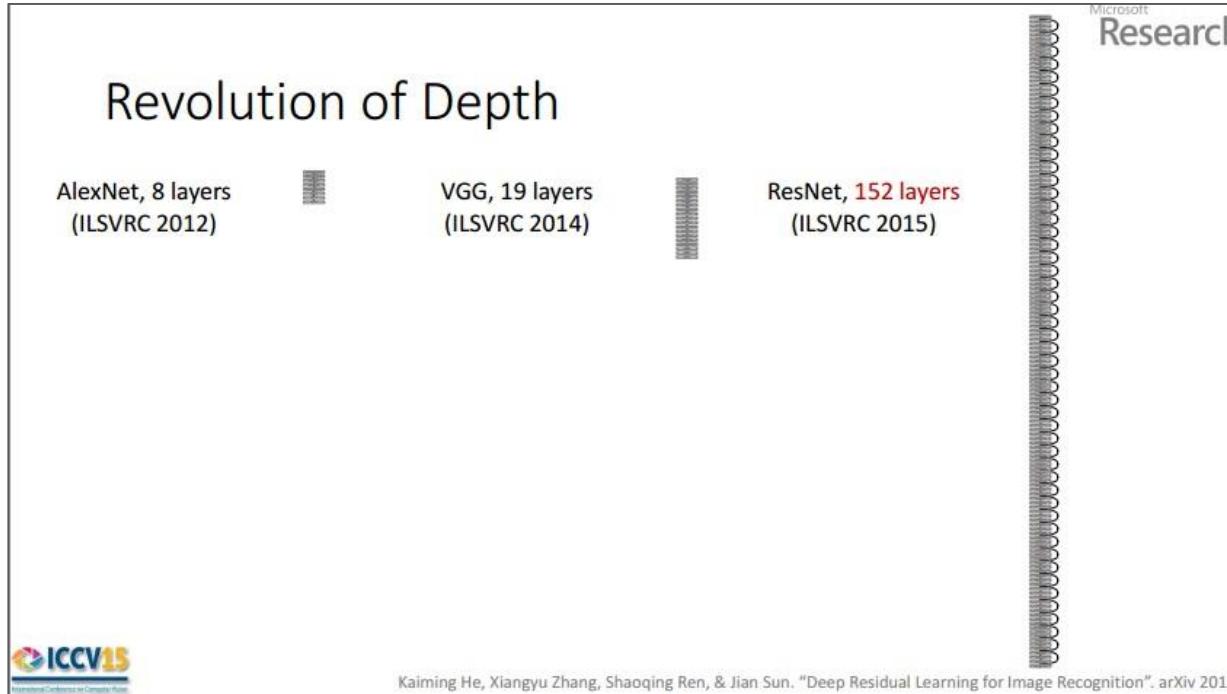
# CIFAR-10 experiments



# Case Study: ResNet

[He et al., 2015]

ILSVRC 2015 winner (3.6% top 5 error)



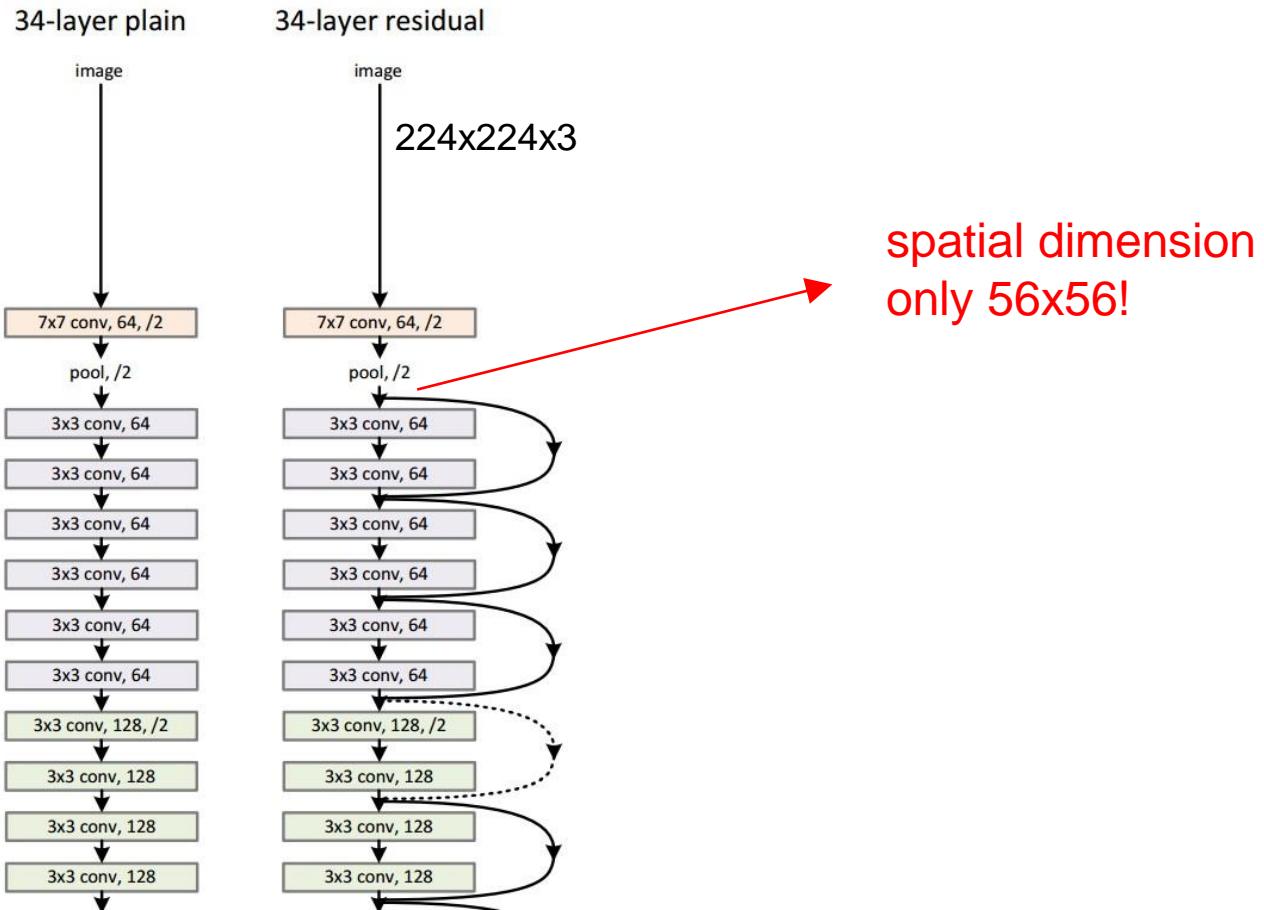
(slide from Kaiming He's recent presentation)

2-3 weeks of training  
on 8 GPU machine

at runtime: faster  
than a VGGNet!  
(even though it has  
8x more layers)

# Case Study: ResNet

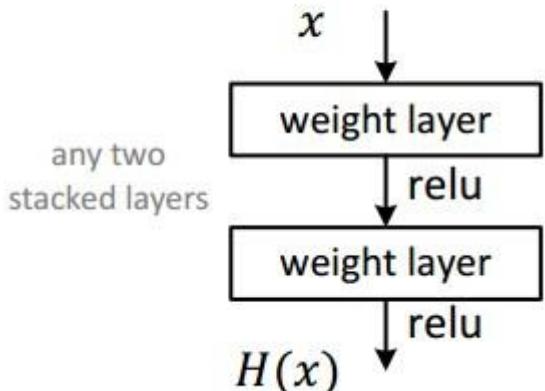
[He et al., 2015]



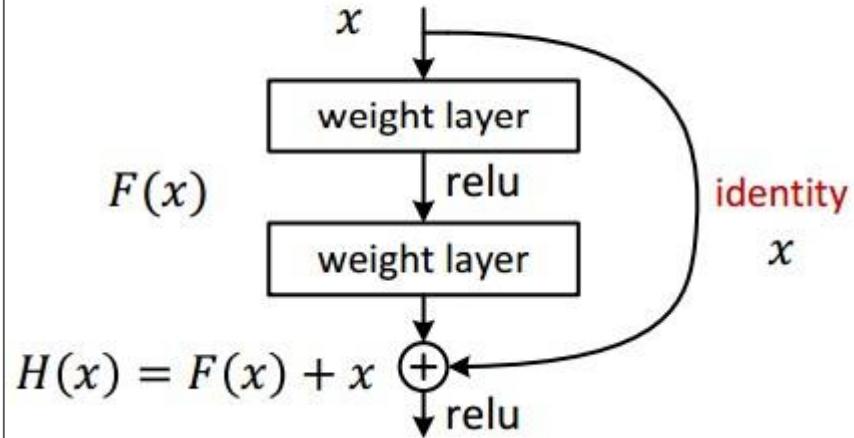
# Case Study: ResNet

[He et al., 2015]

- Plain net



- Residual net



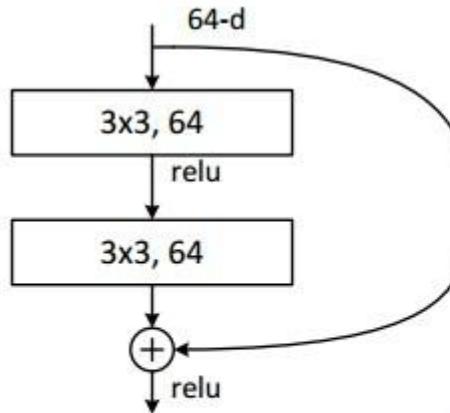
# Case Study: ResNet

[He et al., 2015]

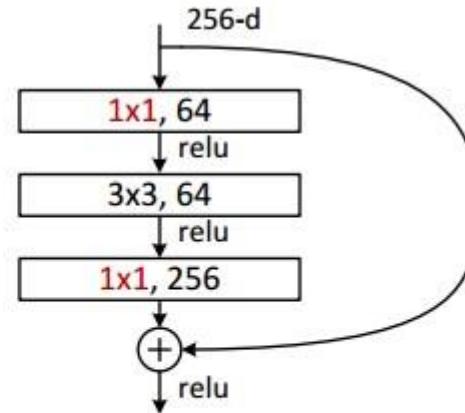
- Batch Normalization after every CONV layer
- Xavier/2 initialization from He et al.
- SGD + Momentum (0.9)
- Learning rate: 0.1, divided by 10 when validation error plateaus
- Mini-batch size 256
- Weight decay of 1e-5
- No dropout used

# Case Study: ResNet

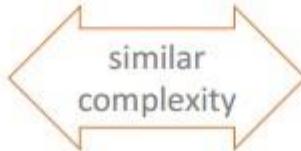
[He et al., 2015]



all-3x3

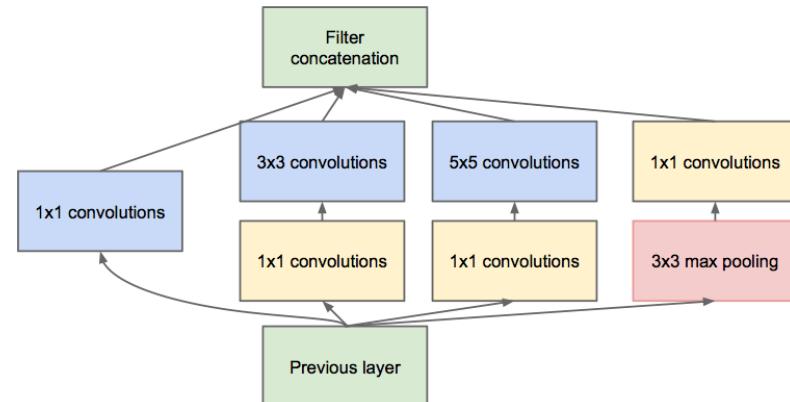
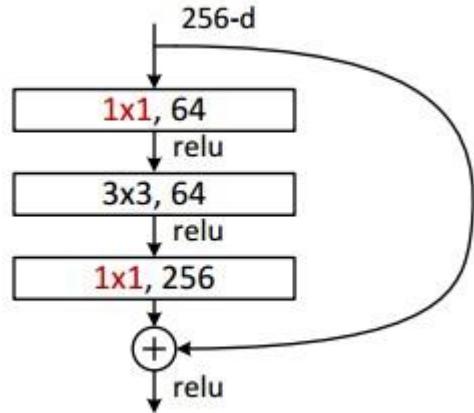


bottleneck  
(for ResNet-50/101/152)

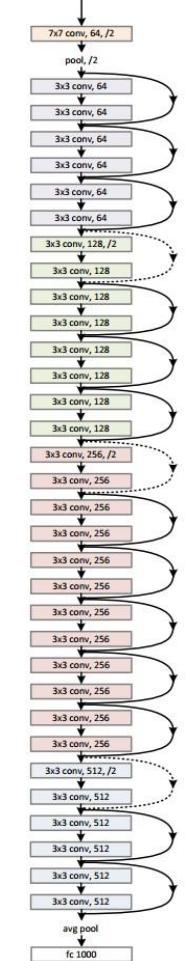


# Case Study: ResNet

[He et al., 2015]



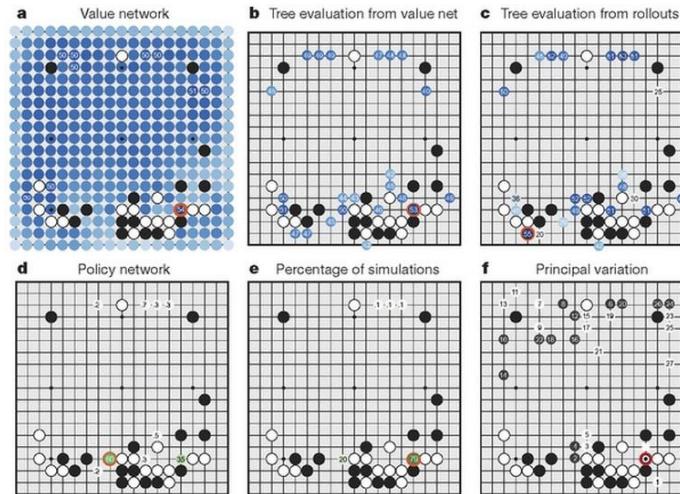
(this trick is also used in GoogLeNet)



# Case Study: ResNet [He et al., 2015]

layer name	output size	18-layer	34-layer	50-layer	101-layer	152-layer
conv1	112×112			7×7, 64, stride 2		
conv2_x	56×56	$\left[ \begin{array}{c} 3 \times 3, 64 \\ 3 \times 3, 64 \end{array} \right] \times 2$	$\left[ \begin{array}{c} 3 \times 3, 64 \\ 3 \times 3, 64 \end{array} \right] \times 3$	$\left[ \begin{array}{c} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{array} \right] \times 3$	$\left[ \begin{array}{c} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{array} \right] \times 3$	$\left[ \begin{array}{c} 1 \times 1, 64 \\ 3 \times 3, 64 \\ 1 \times 1, 256 \end{array} \right] \times 3$
conv3_x	28×28	$\left[ \begin{array}{c} 3 \times 3, 128 \\ 3 \times 3, 128 \end{array} \right] \times 2$	$\left[ \begin{array}{c} 3 \times 3, 128 \\ 3 \times 3, 128 \end{array} \right] \times 4$	$\left[ \begin{array}{c} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{array} \right] \times 4$	$\left[ \begin{array}{c} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{array} \right] \times 4$	$\left[ \begin{array}{c} 1 \times 1, 128 \\ 3 \times 3, 128 \\ 1 \times 1, 512 \end{array} \right] \times 8$
conv4_x	14×14	$\left[ \begin{array}{c} 3 \times 3, 256 \\ 3 \times 3, 256 \end{array} \right] \times 2$	$\left[ \begin{array}{c} 3 \times 3, 256 \\ 3 \times 3, 256 \end{array} \right] \times 6$	$\left[ \begin{array}{c} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{array} \right] \times 6$	$\left[ \begin{array}{c} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{array} \right] \times 23$	$\left[ \begin{array}{c} 1 \times 1, 256 \\ 3 \times 3, 256 \\ 1 \times 1, 1024 \end{array} \right] \times 36$
conv5_x	7×7	$\left[ \begin{array}{c} 3 \times 3, 512 \\ 3 \times 3, 512 \end{array} \right] \times 2$	$\left[ \begin{array}{c} 3 \times 3, 512 \\ 3 \times 3, 512 \end{array} \right] \times 3$	$\left[ \begin{array}{c} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{array} \right] \times 3$	$\left[ \begin{array}{c} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{array} \right] \times 3$	$\left[ \begin{array}{c} 1 \times 1, 512 \\ 3 \times 3, 512 \\ 1 \times 1, 2048 \end{array} \right] \times 3$
	1×1			average pool, 1000-d fc, softmax		
FLOPs		$1.8 \times 10^9$	$3.6 \times 10^9$	$3.8 \times 10^9$	$7.6 \times 10^9$	$11.3 \times 10^9$

# Case Study Bonus: DeepMind's AlphaGo



Slide based on cs231n by Fei-Fei Li & Andrej Karpathy & Justin Johnson

The input to the policy network is a  $19 \times 19 \times 48$  image stack consisting of 48 feature planes. The first hidden layer zero pads the input into a  $23 \times 23$  image, then convolves  $k$  filters of kernel size  $5 \times 5$  with stride 1 with the input image and applies a rectifier nonlinearity. Each of the subsequent hidden layers 2 to 12 zero pads the respective previous hidden layer into a  $21 \times 21$  image, then convolves  $k$  filters of kernel size  $3 \times 3$  with stride 1, again followed by a rectifier nonlinearity. The final layer convolves 1 filter of kernel size  $1 \times 1$  with stride 1, with a different bias for each position, and applies a softmax function. The match version of AlphaGo used  $k = 192$  filters; [Fig. 2b](#) and [Extended Data Table 3](#) additionally show the results of training with  $k = 128, 256$  and  $384$  filters.

## policy network:

[ $19 \times 19 \times 48$ ] Input

CONV1: 192  $5 \times 5$  filters , stride 1, pad 2 => [ $19 \times 19 \times 192$ ]

CONV2..12: 192  $3 \times 3$  filters, stride 1, pad 1 => [ $19 \times 19 \times 192$ ]

CONV: 1  $1 \times 1$  filter, stride 1, pad 0 => [ $19 \times 19$ ] (*probability map of promising moves*)

# The Chain Rule

- ConvNets stack CONV,ReLU,POOL,FC layers
- Trend towards smaller filters and deeper architectures
- Trend towards getting rid of POOL/FC layers (just CONV)
- Early architectures look like  
 **$[(CONV-RELU)*N-POOL?] * M - (FC-RELU)*K, SOFTMAX$**
- but recent advances such as ResNet/GoogLeNet use only Conv-ReLU, 1x1 convolutions and Softmax