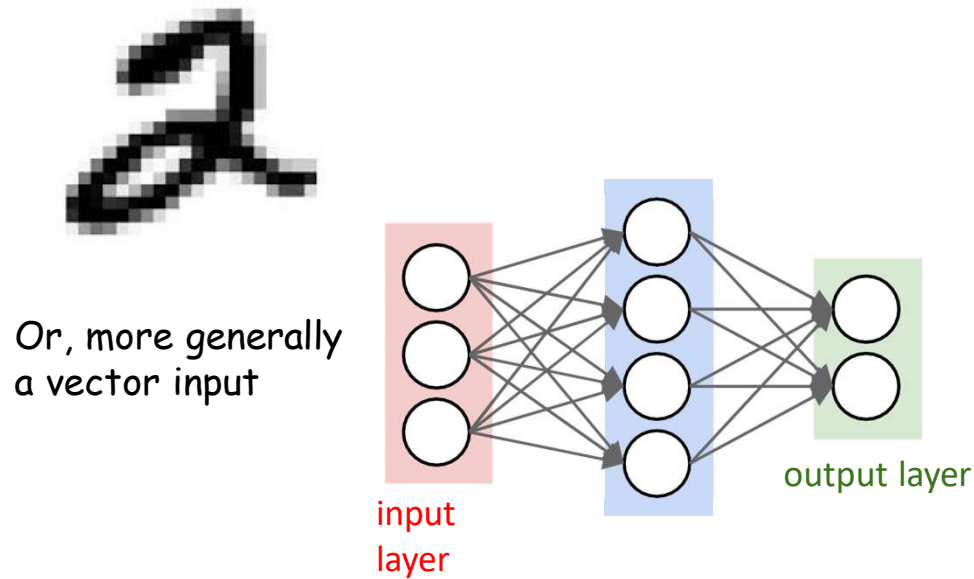


# Introduction to Deep Learning

Presentation Title

H M Sabbir Ahmad  
Phd, Systems Engineering.  
09/12/2024

# Multi-Layer Perception Issue



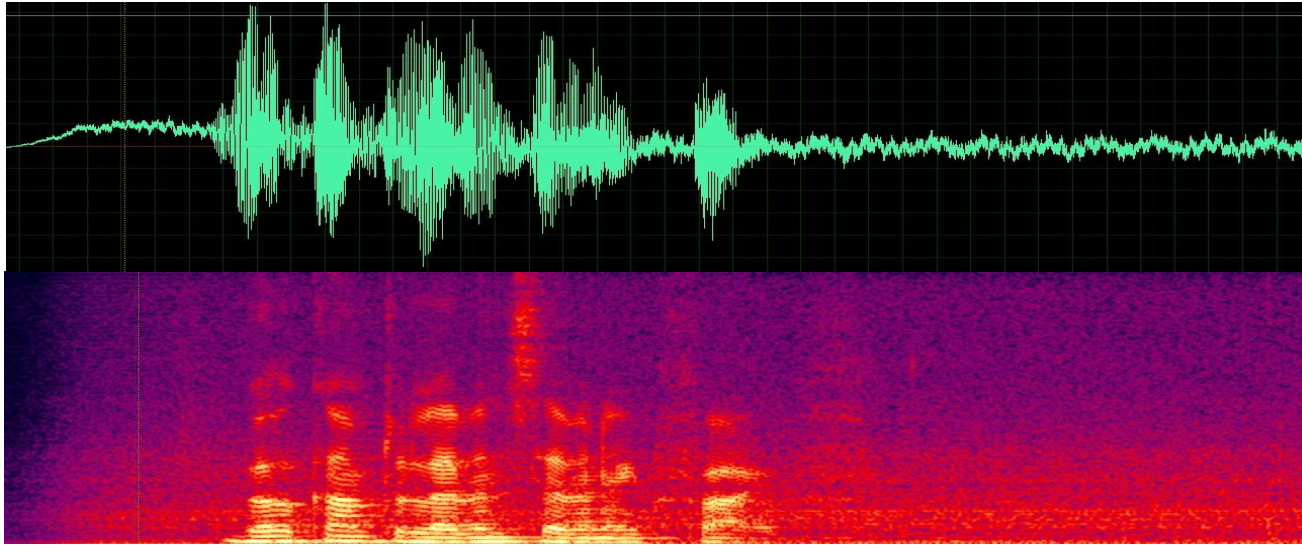
- Can recognize patterns in data
  - E.g. digits
  - Or any other vector data

Boston University School/college name here



Slide credit: Bhiksha Raj

# A new problem



- Does this signal contain the word “Welcome”?
- Compose an MLP for this problem.
  - Assuming all recordings are exactly the same length..

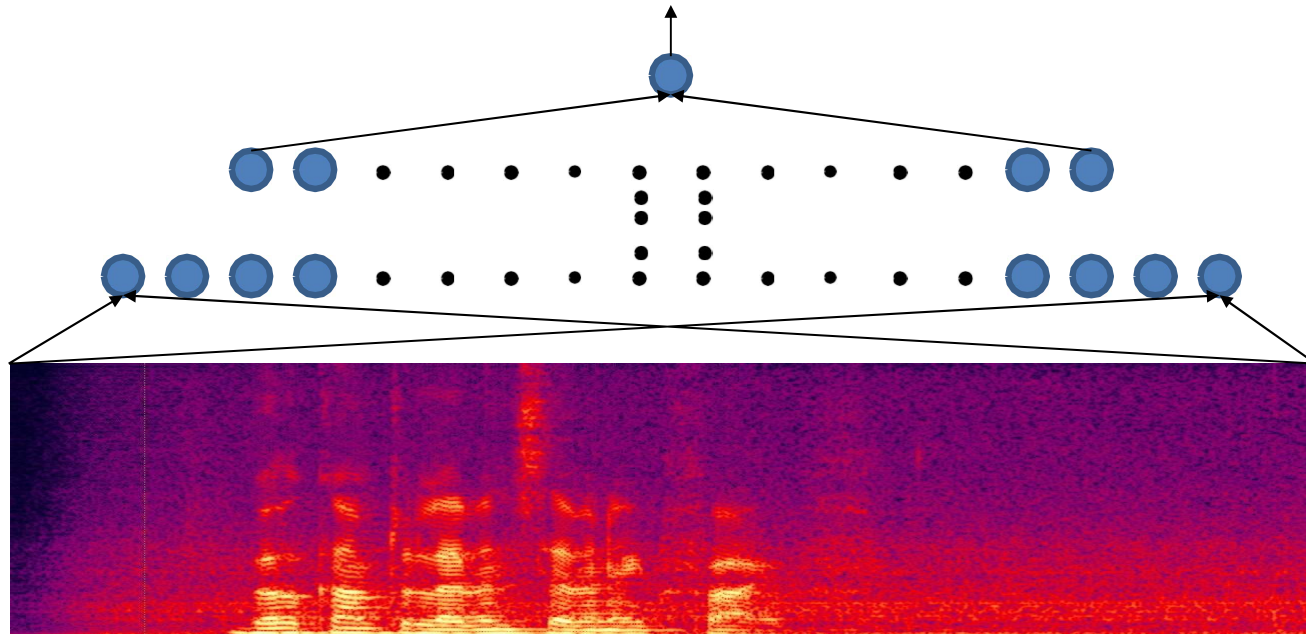
Boston University School/college name here

4



Slide credit: Bhiksha Raj

# Finding a Welcome



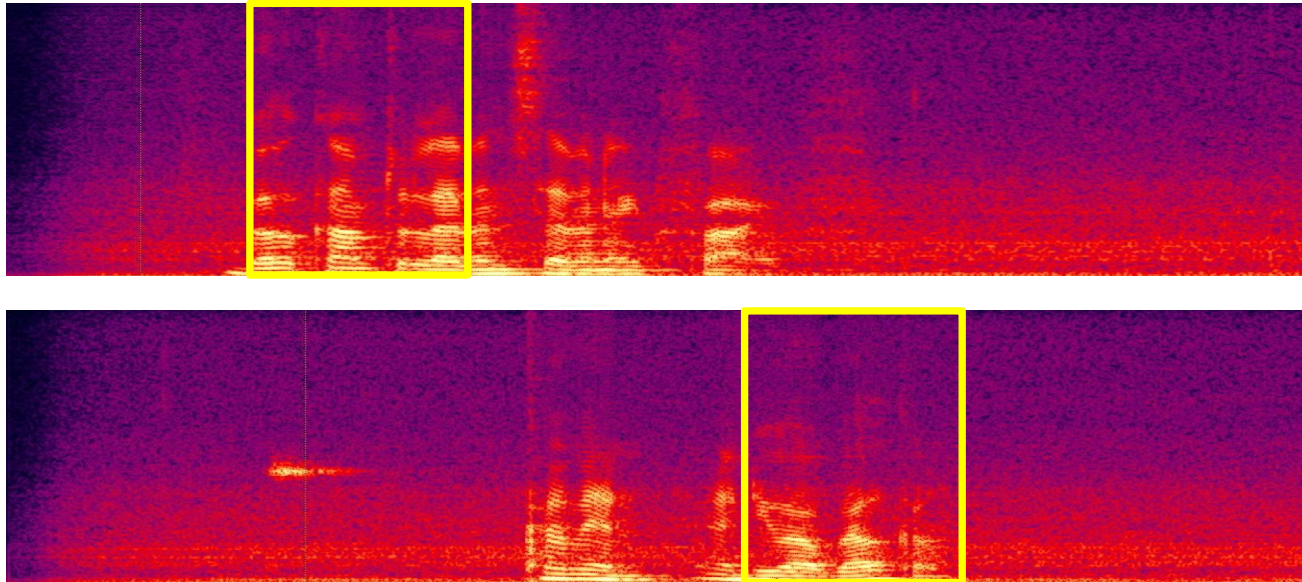
- Trivial solution: Train an MLP for the entire recording

Boston University School/college name here



Slide credit: Bhiksha Raj

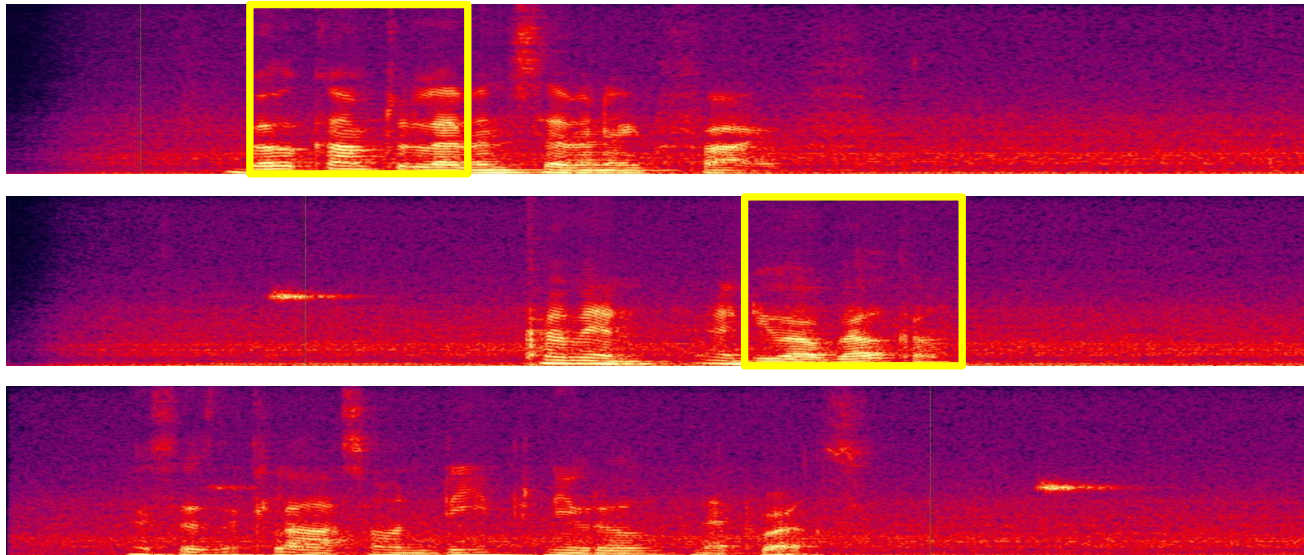
# Finding a Welcome



- Problem with trivial solution: Network that finds a “welcome” in the top recording will not find it in the lower one
  - Unless trained with both
  - Will require a very large network and a large amount of training data to cover every case



# Finding a Welcome



- Need a *simple* network that will fire regardless of the location of “Welcome”
  - and not fire when there is none

# Flowers



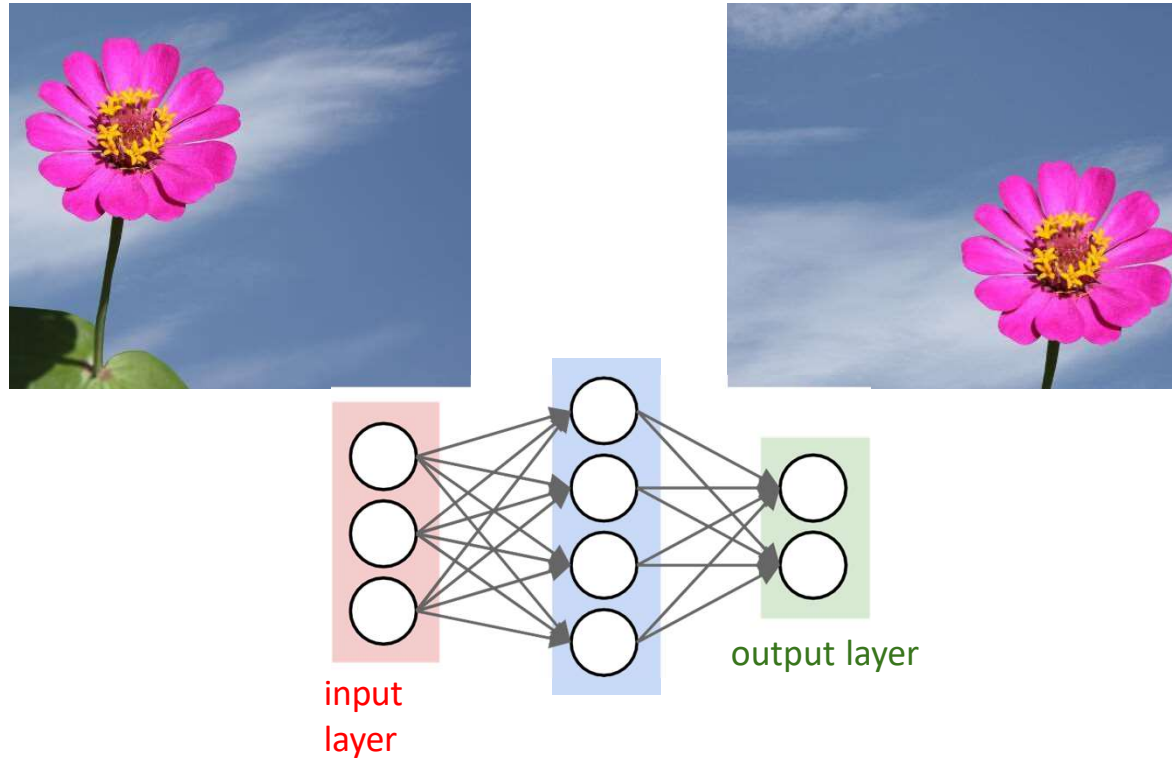
- Is there a flower in any of these images

Boston University School/college name here



Slide credit: Bhiksha Raj

# A problem



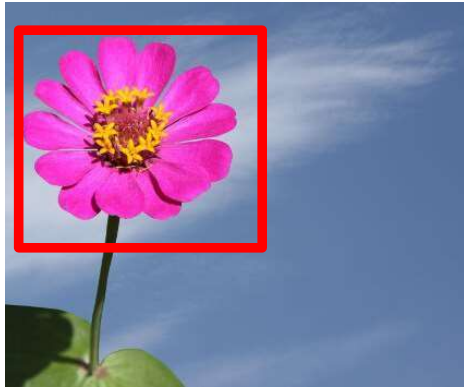
- Will an MLP that recognizes the left image as a flower also recognize the one on the right as a flower?

Boston University School/college name here



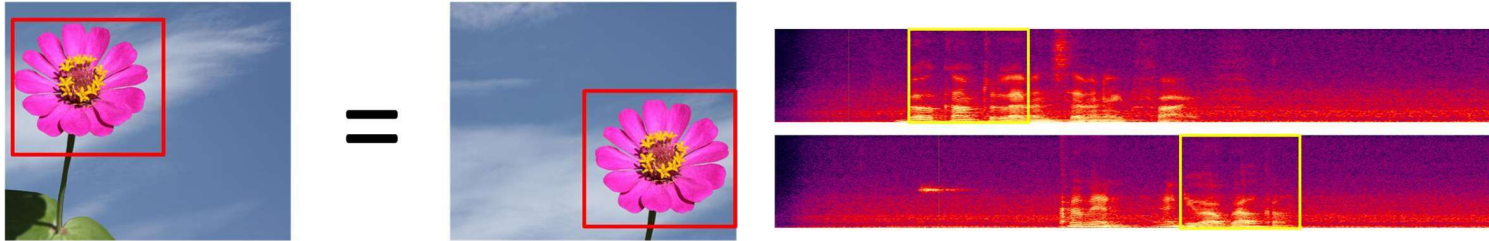


# A problem



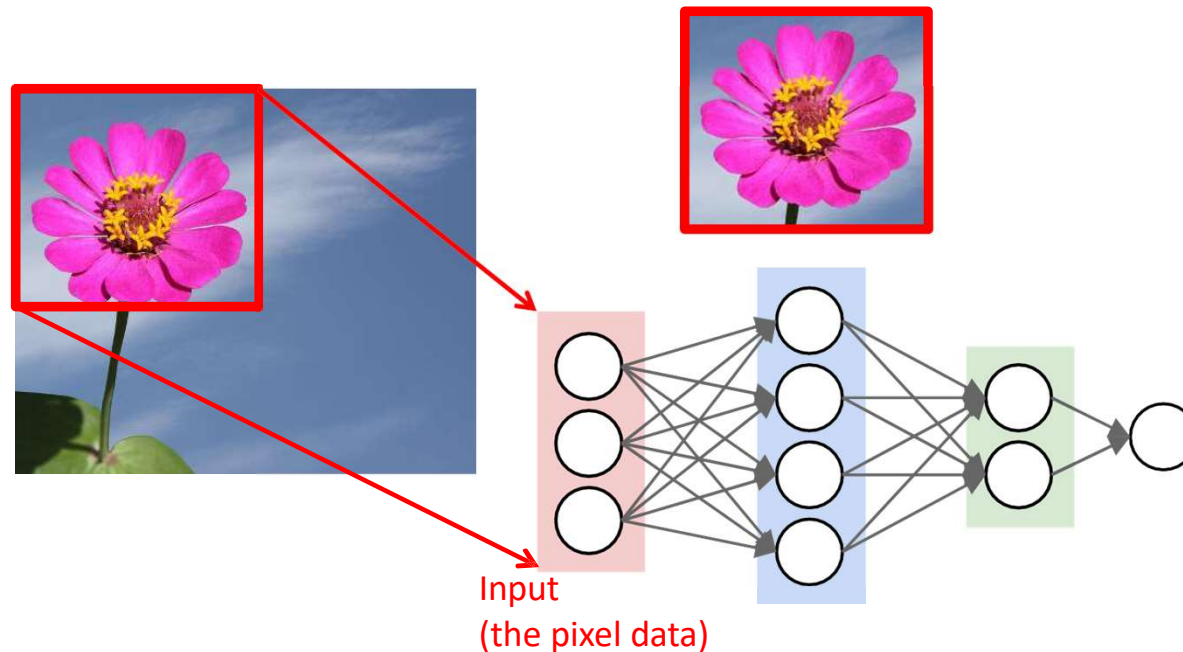
- Need a network that will “fire” regardless of the precise location of the target object

# The need for *shift invariance*



- In many problems the *location* of a pattern is not important
  - Only the presence of the pattern
- Conventional MLPs are sensitive to the location of the pattern
  - Moving it by one component results in an entirely different input that the MLP won't recognize
- Requirement: Network must be *shift invariant*

# The 2-d analogue: Does this picture have a flower?



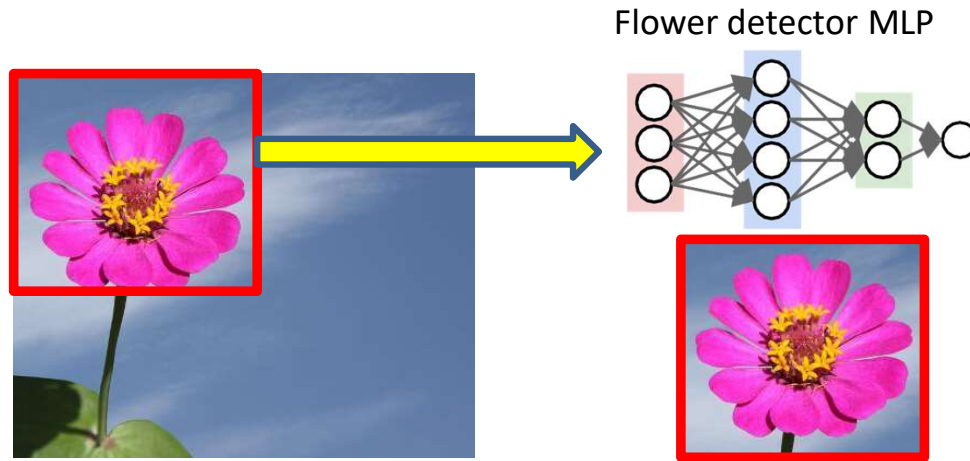
- *Scan* for the desired object
  - “Look” for the target object at each position

Boston University School/college name here



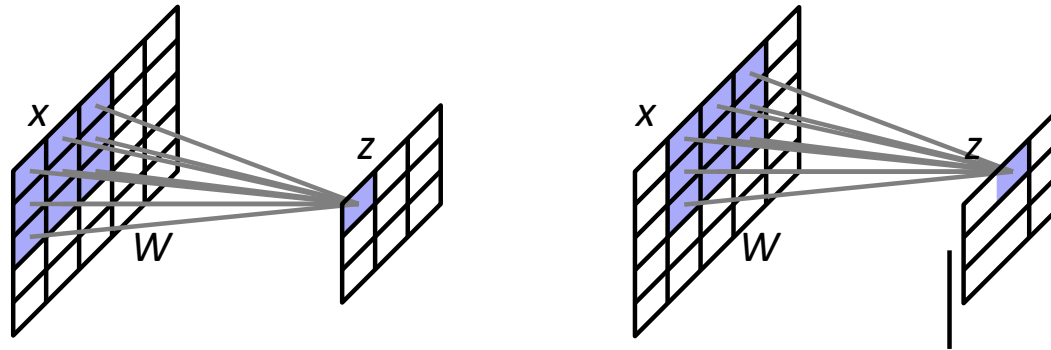
Slide credit: Bhiksha Raj

# Solution: Scan



- *Scan* for the desired object
- At each location, the entire region is sent through the MLP

# Advantages of convolutions



*Drastically* reduces the parameter count

- 256x256 grayscale image  $\Rightarrow$  256x256 single-channel hidden layer: 4 *billion parameters* in fully connected network to 9 *parameters* in 3x3 convolution

Captures (some) “natural” invariances

- Shifting input image one pixel to the right shifts the hidden unit “image”

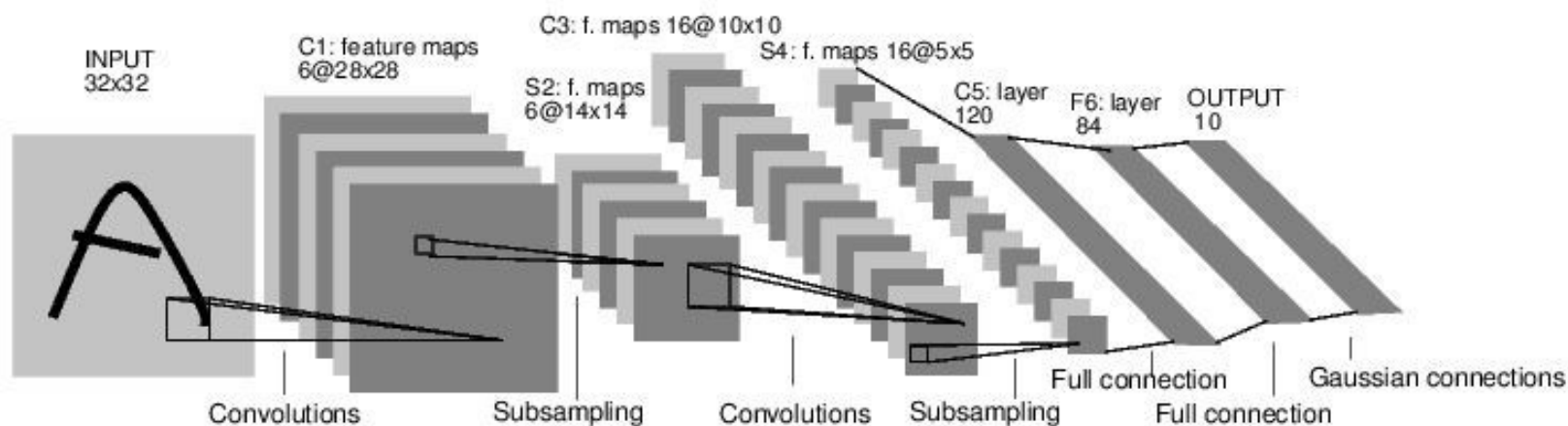
Boston University School/college name here



Slide credit: J. Zico Kolter and Tianqi Chen



# Convolutional Neural Networks



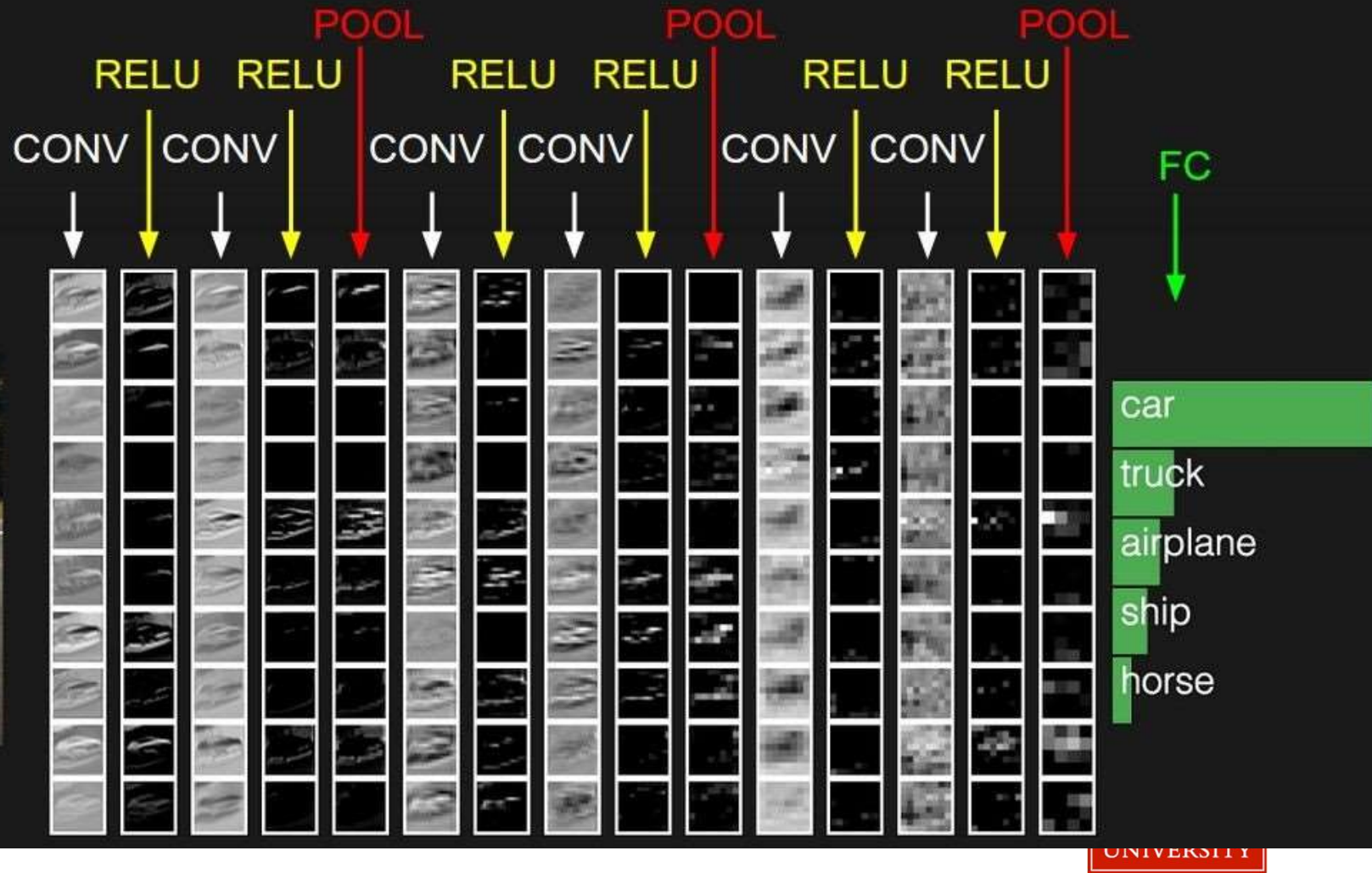
*[LeNet-5, LeCun 1980]*

Boston University School/college name here



Slide credit: Fei-Fei Li, Andrej Karpathy and Justin Johnson

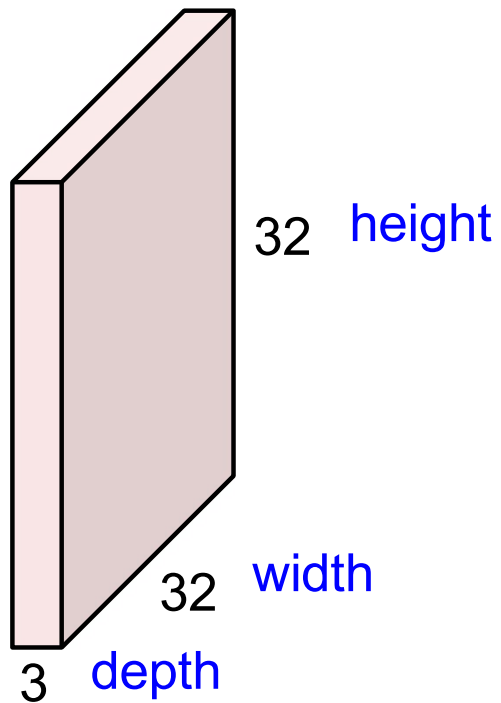
preview:



Slide credit: Fei-Fei Li, Andrej Karpathy and Justin Johnson

# Convolution Layer

32x32x3 image



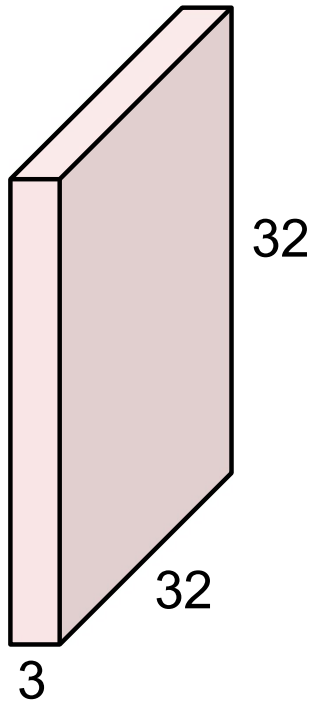
Boston University School/college name here



Slide credit: Fei-Fei Li, Andrej Karpathy and Justin Johnson

# Convolution Layer

32x32x3 image



5x5x3 filter



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

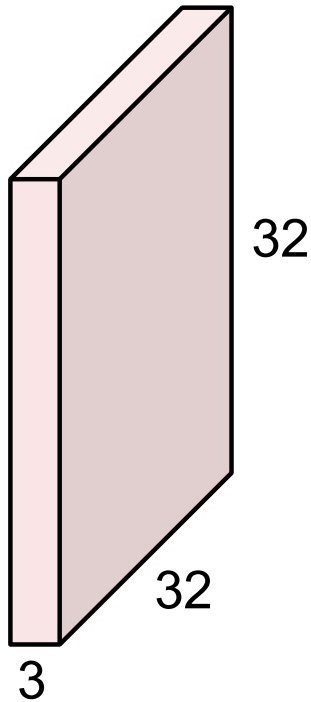
Boston University School/college name here



Slide credit: Fei-Fei Li, Andrej Karpathy and Justin Johnson

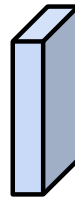
# Convolution Layer

32x32x3 image



Filters always extend the full depth of the input volume

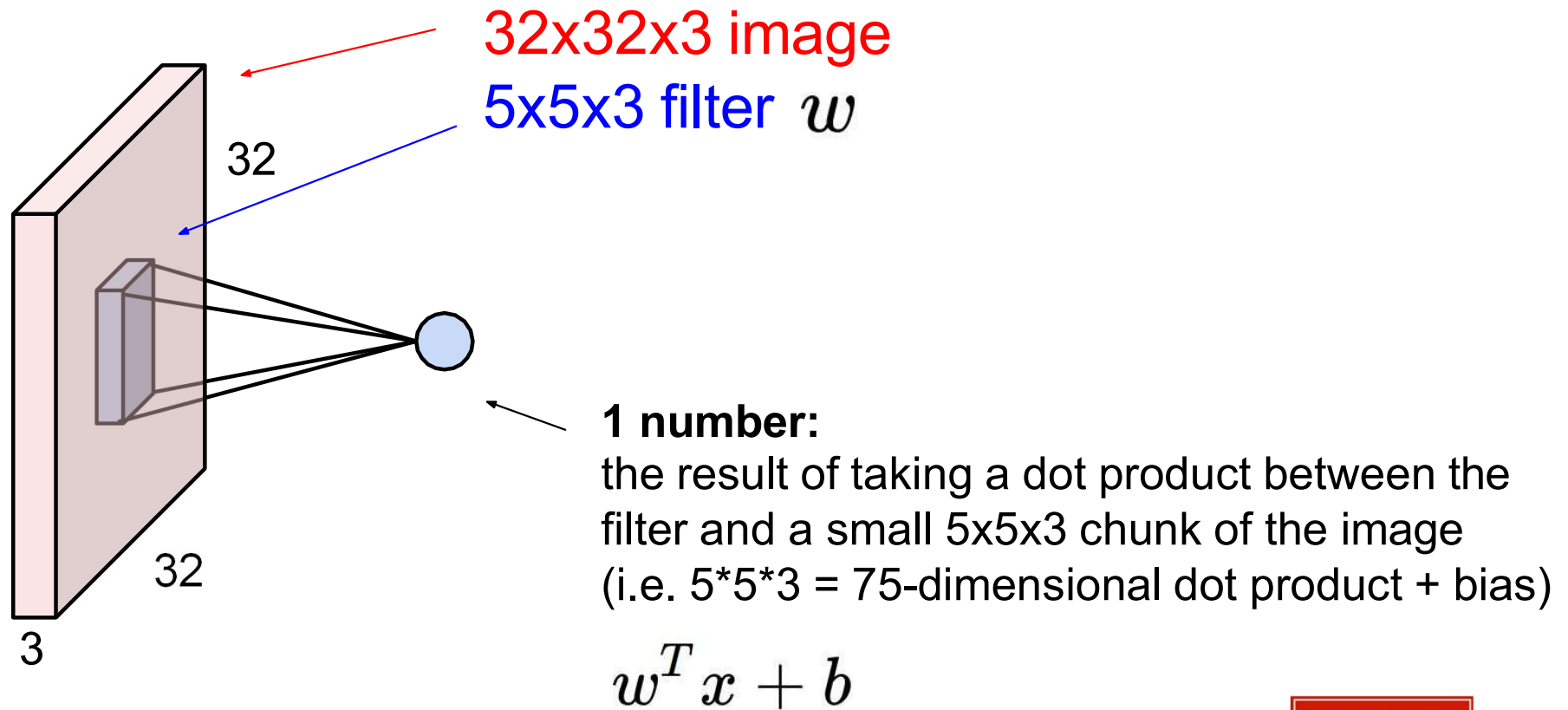
5x5x3 filter



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



# Convolution Layer

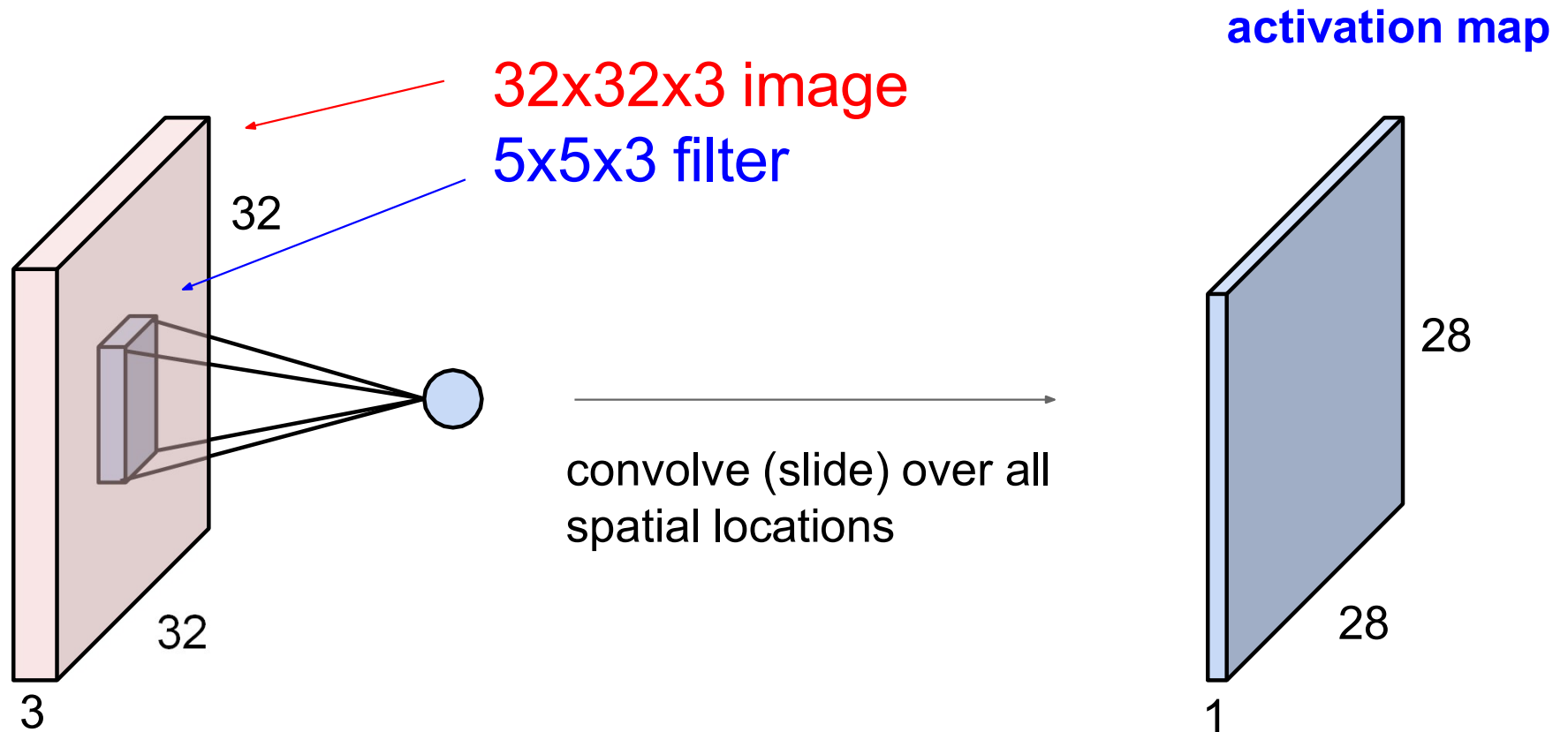


Boston University School/college name here



Slide credit: Fei-Fei Li, Andrej Karpathy and Justin Johnson

# Convolution Layer



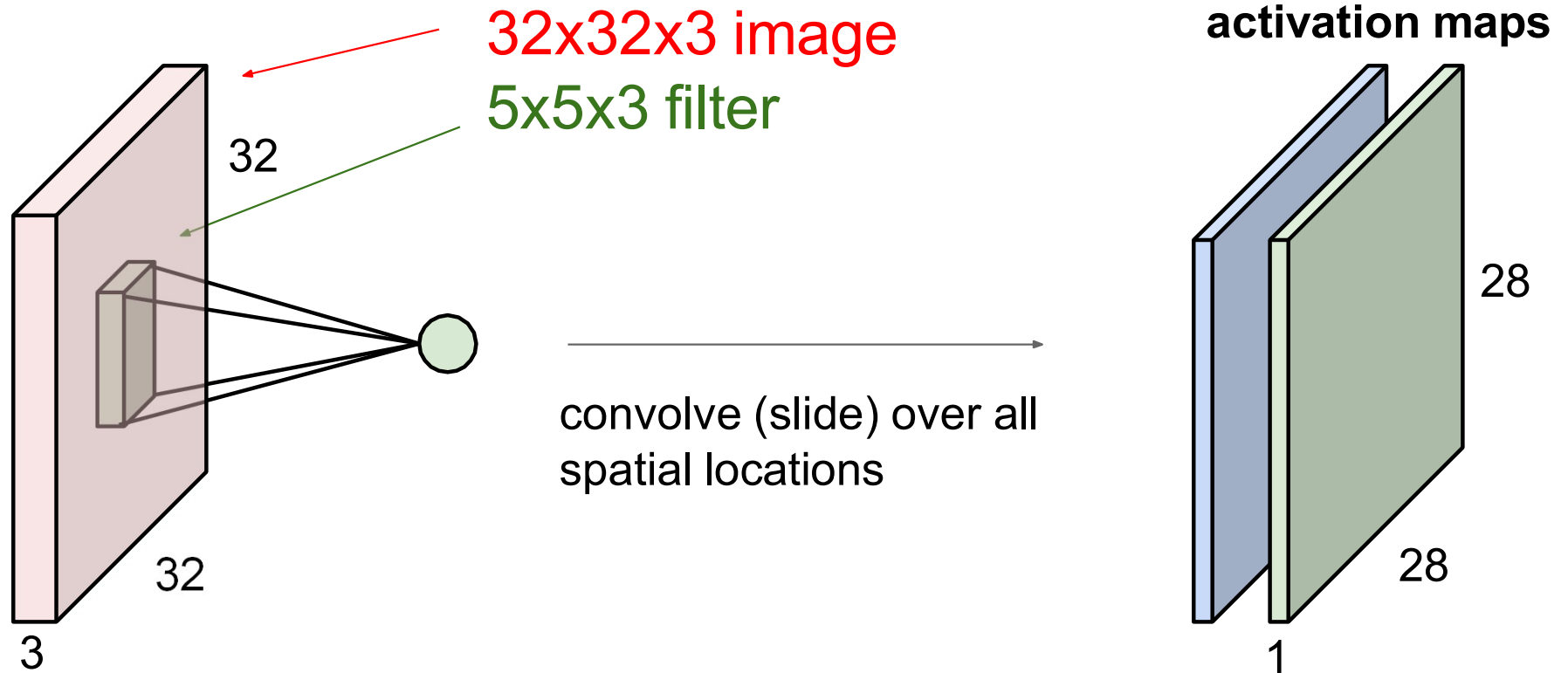
Boston University School/college name here



Slide credit: Fei-Fei Li, Andrej Karpathy and Justin Johnson

# Convolution Layer

consider a second, **green** filter



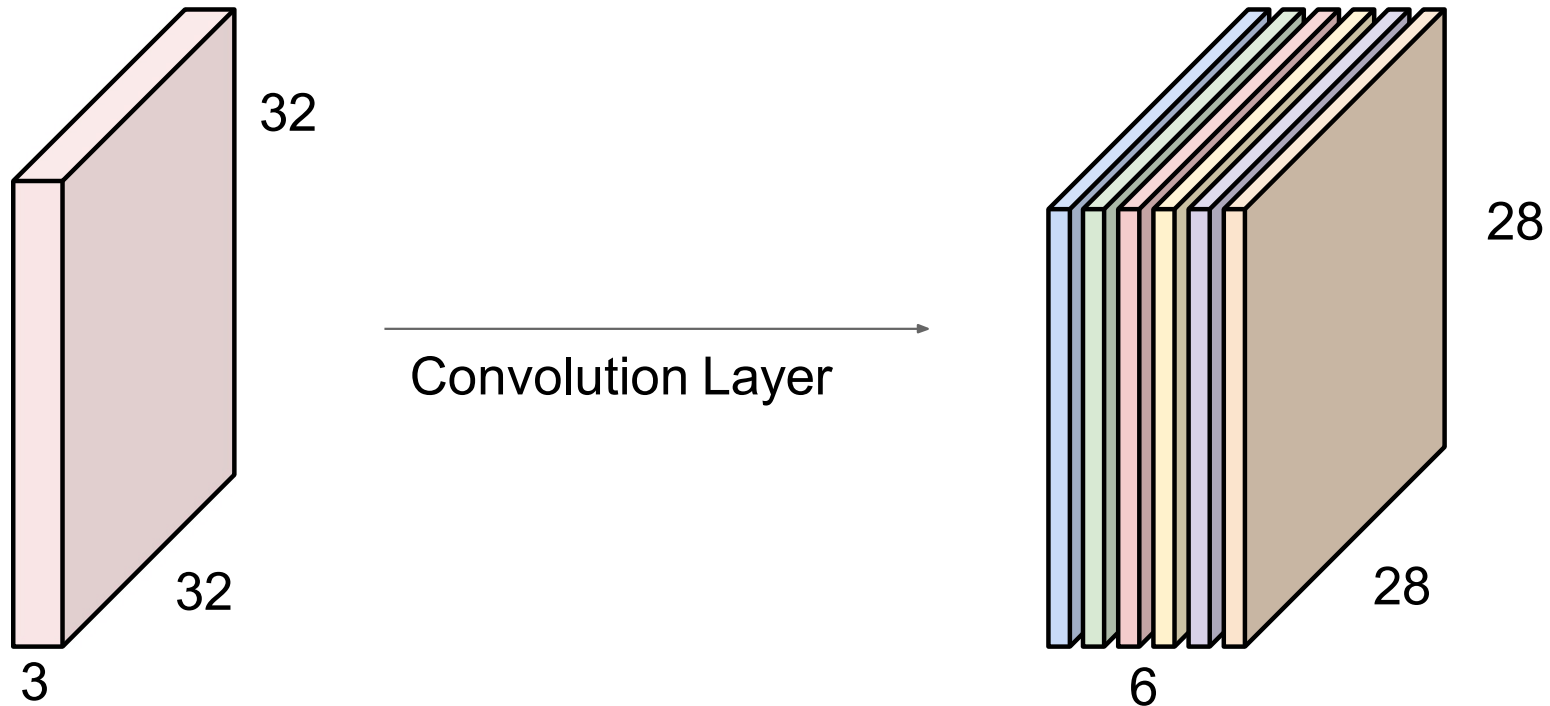
Boston University School/college name here



Slide credit: Fei-Fei Li, Andrej Karpathy and Justin Johnson

# Convolution Layer

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



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

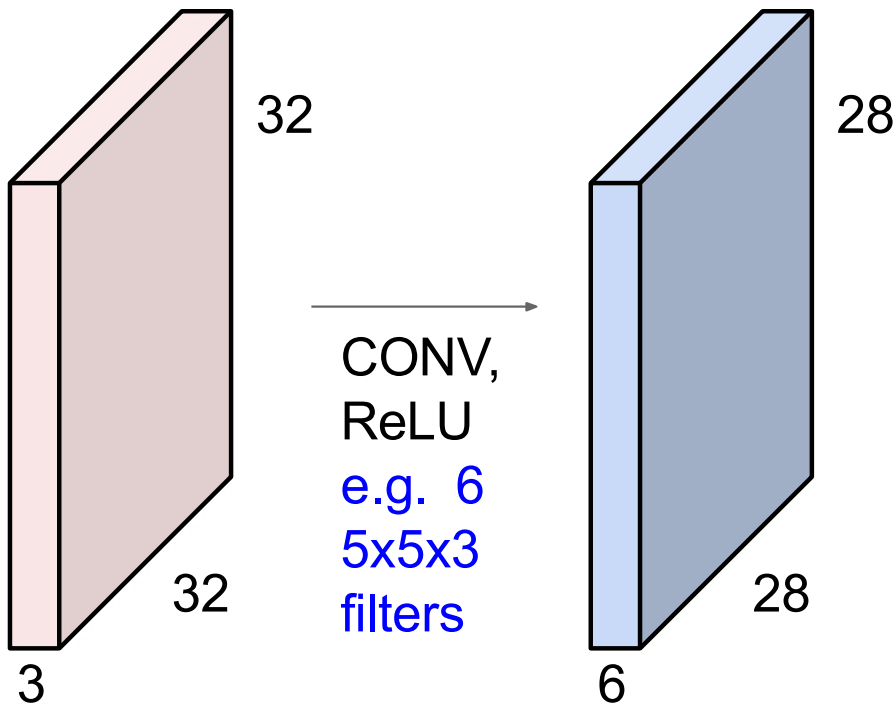
Boston University School/college name here



Slide credit: Fei-Fei Li, Andrej Karpathy and Justin Johnson

# Convolution Layer

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



Boston University School/college name here

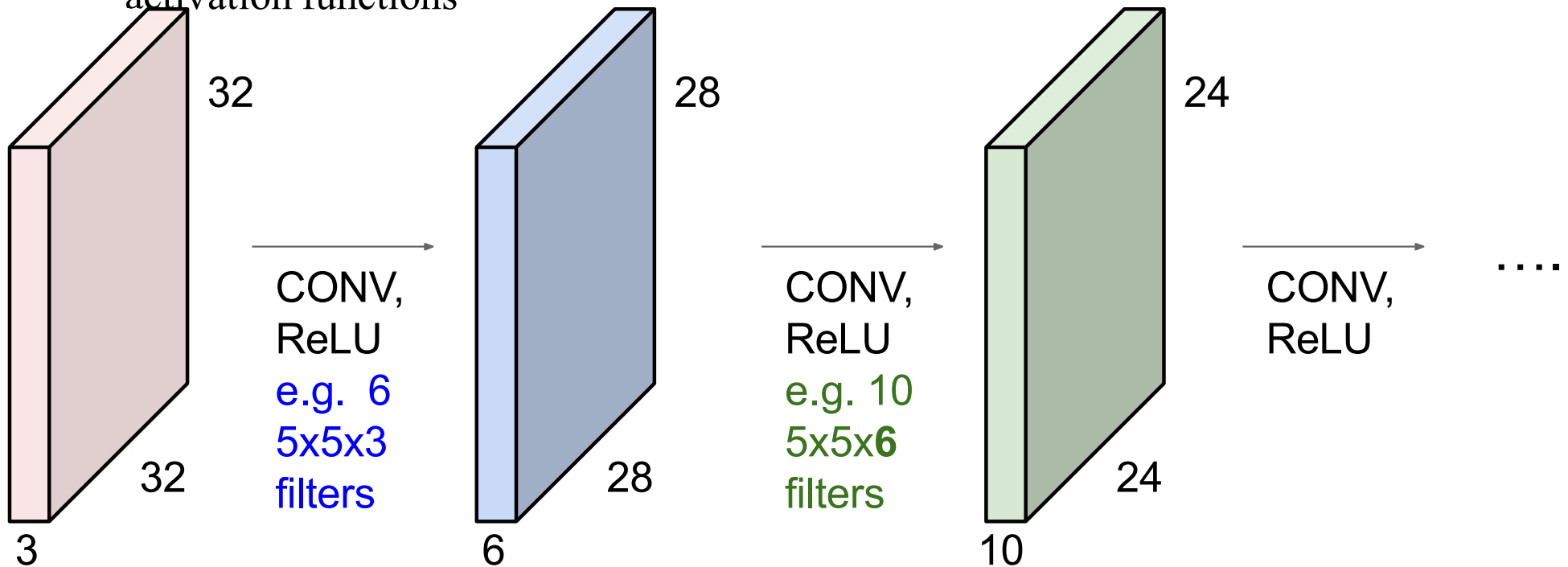


Slide credit: Fei-Fei Li, Andrej Karpathy and Justin Johnson



# Convolution Layer

**Preview:** ConvNet is a sequence of Convolutional Layers, interspersed with activation functions



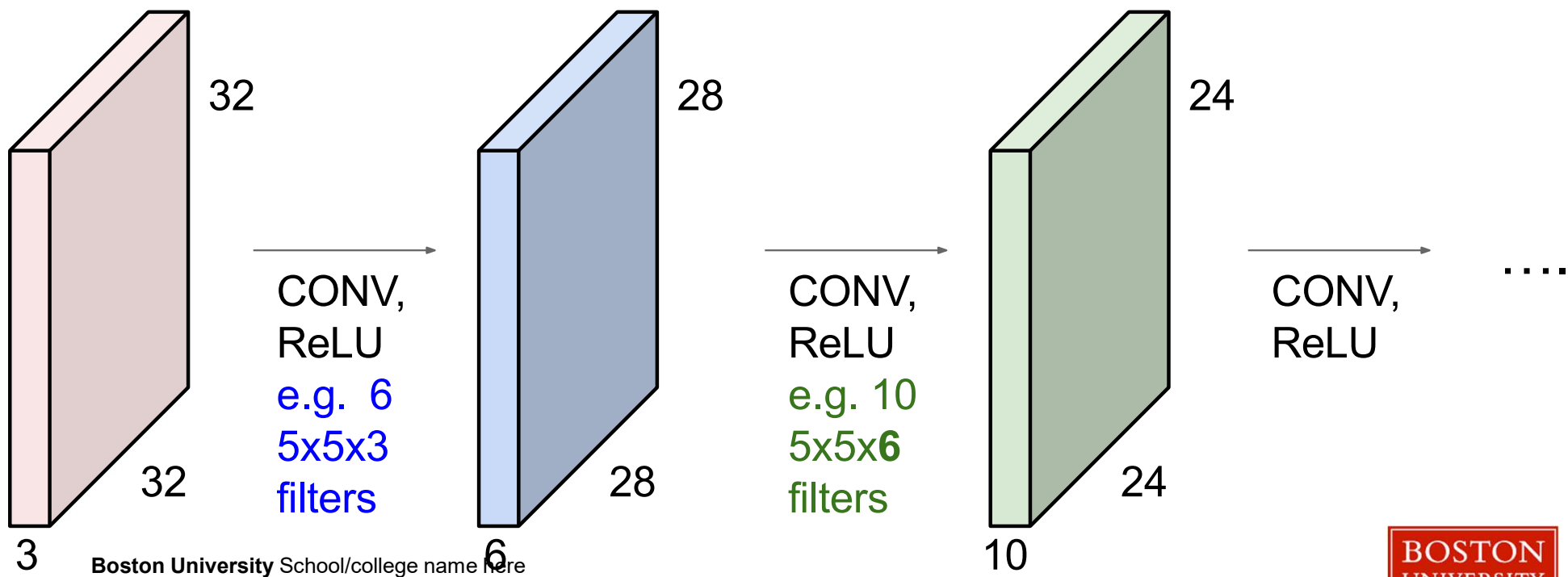
Boston University School/college name here



Slide credit: Fei-Fei Li, Andrej Karpathy and Justin Johnson

## Remember back to...

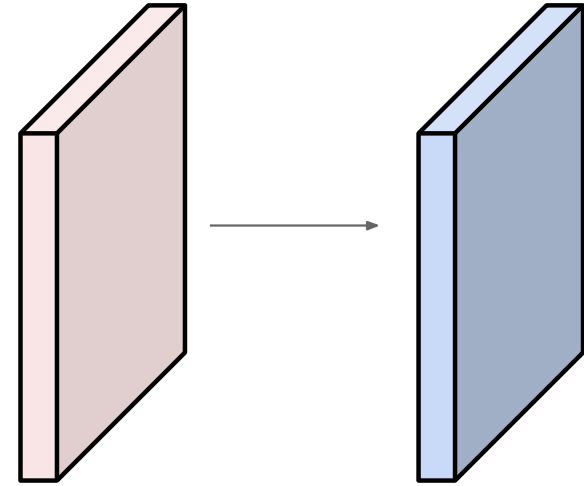
E.g. 32x32 input convolved repeatedly with 5x5 filters shrinks volumes spatially! (32 → 28 → 24 ...). Shrinking too fast is not good, doesn't work well.



Examples time:

Input volume: **32x32x3**  
10 5x5 filters with stride 1, pad 2

Output volume size: ?



Examples time:

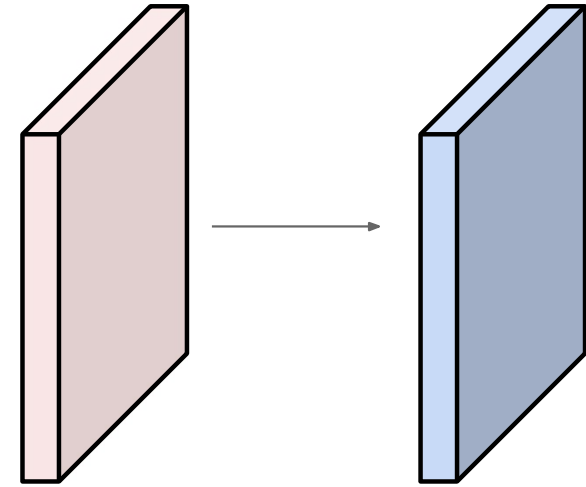
Input volume: **32x32x3**

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

Output volume size:

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

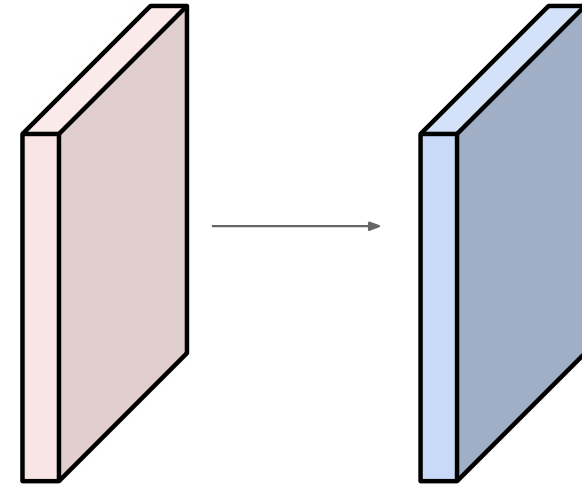
**32x32x10**



Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2



Number of parameters in this layer?

Boston University School/college name here



Slide credit: Fei-Fei Li, Andrej Karpathy and Justin Johnson



Examples time:

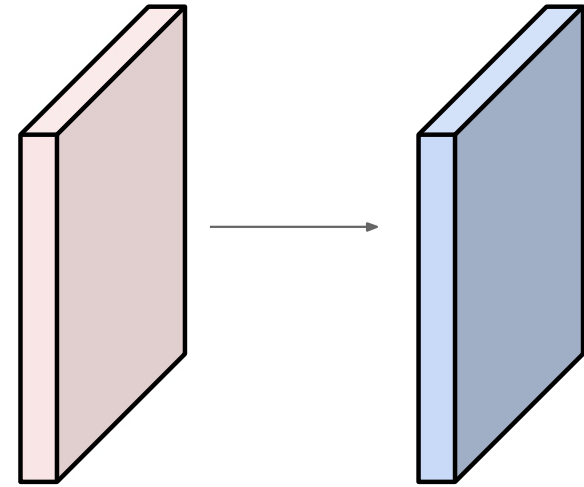
Input volume: **32x32x3**

**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)

=>  $76*10 = 760$



# Convolutions in image processing

Convolutions (typically with *prespecified* filters) are a common operation in many computer vision applications: convolution networks just move to *learned* filters



Original image  $x$



Gaussian blur



Image gradient

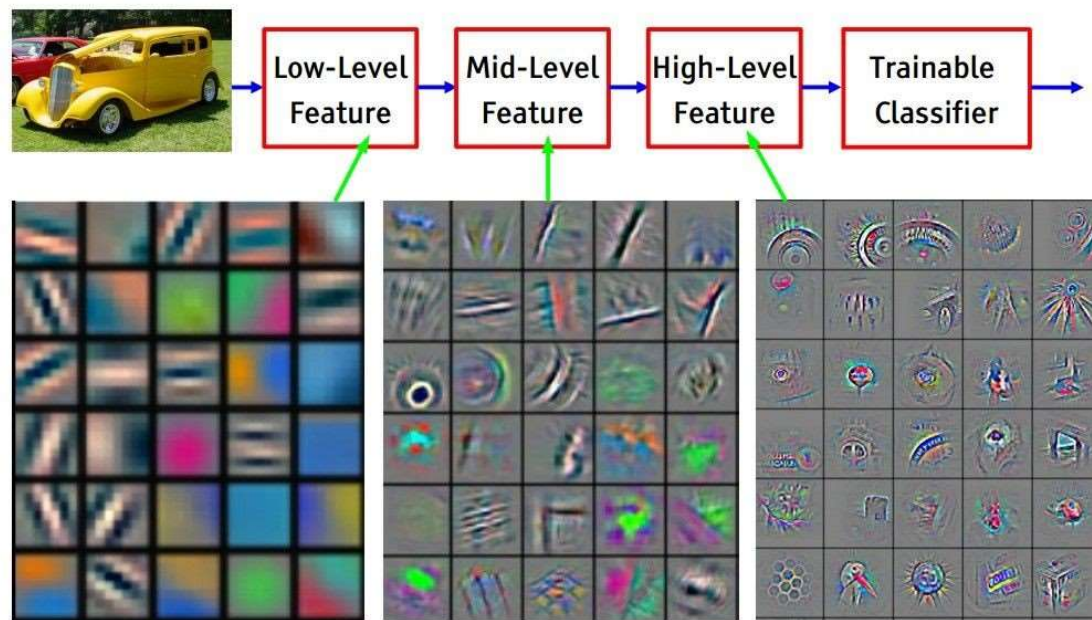
$$x * \begin{bmatrix} 1 & 4 & 7 & 4 & 1 \\ 4 & 16 & 26 & 16 & 4 \\ 7 & 26 & 41 & 26 & 7 \\ 4 & 16 & 26 & 16 & 4 \\ 1 & 4 & 4 & 4 & 1 \end{bmatrix} / 273$$



$$\left( \left( x * \begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix} \right)^2 + \left( x * \begin{bmatrix} -1 & -2 & -1 \\ 0 & 0 & 0 \\ 1 & 2 & 1 \end{bmatrix} \right)^2 \right)^{\frac{1}{2}}$$



# Preview



*[From recent Yann LeCun slides]*

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

**Boston University** School/college name here



Slide credit: Fei-Fei Li, Andrej Karpathy and Justin Johnson

# Convolutions in detail

- Convolutions are a basic primitive in many computer vision and image processing algorithms
- Idea is to “slide” the weights  $h \times h$  weight  $m$  (called a filter, with kernel size  $h$ ) over the image to produce a new image, written  $y = x * m$

$z_{11}$	$z_{12}$	$z_{13}$	$z_{14}$	$z_{15}$
$z_{21}$	$z_{22}$	$z_{23}$	$z_{24}$	$z_{25}$
$z_{31}$	$z_{32}$	$z_{33}$	$z_{34}$	$z_{35}$
$z_{41}$	$z_{42}$	$z_{43}$	$z_{44}$	$z_{45}$
$z_{51}$	$z_{52}$	$z_{53}$	$z_{54}$	$z_{55}$

\*

$w_{11}$	$w_{12}$	$w_{13}$
$w_{21}$	$w_{22}$	$w_{23}$
$w_{31}$	$w_{32}$	$w_{33}$

=

$y_{11}$	$y_{12}$	$y_{13}$
$y_{21}$	$y_{22}$	$y_{23}$
$y_{31}$	$y_{32}$	$y_{33}$

# Convolutions in detail

- Convolutions are a basic primitive in many computer vision and image processing algorithms
- Idea is to “slide” the weights  $h \times h$  weight  $m$  (called a filter, with kernel size  $h$ ) over the image to produce a new image, written  $y = x * m$

$z_{11}$	$z_{12}$	$z_{13}$	$z_{14}$	$z_{15}$
$z_{21}$	$z_{22}$	$z_{23}$	$z_{24}$	$z_{25}$
$z_{31}$	$z_{32}$	$z_{33}$	$z_{34}$	$z_{35}$
$z_{41}$	$z_{42}$	$z_{43}$	$z_{44}$	$z_{45}$
$z_{51}$	$z_{52}$	$z_{53}$	$z_{54}$	$z_{55}$

 $*$ 

$w_{11}$	$w_{12}$	$w_{13}$
$w_{21}$	$w_{22}$	$w_{23}$
$w_{31}$	$w_{32}$	$w_{33}$

 $=$ 

$y_{11}$	$y_{12}$	$y_{13}$
$y_{21}$	$y_{22}$	$y_{23}$
$y_{31}$	$y_{32}$	$y_{33}$

  
$$y_{11} = z_{11}w_{11} + z_{12}w_{12} + z_{13}w_{13} + z_{21}w_{21} + \dots$$

# Convolutions in detail

- Convolutions are a basic primitive in many computer vision and image processing algorithms
- Idea is to “slide” the weights  $h \times h$  weight  $m$  (called a filter, with kernel size  $h$ ) over the image to produce a new image, written  $y = x * m$

$z_{11}$	$z_{12}$	$z_{13}$	$z_{14}$	$z_{15}$
$z_{21}$	$z_{22}$	$z_{23}$	$z_{24}$	$z_{25}$
$z_{31}$	$z_{32}$	$z_{33}$	$z_{34}$	$z_{35}$
$z_{41}$	$z_{42}$	$z_{43}$	$z_{44}$	$z_{45}$
$z_{51}$	$z_{52}$	$z_{53}$	$z_{54}$	$z_{55}$

 $*$ 

$w_{11}$	$w_{12}$	$w_{13}$
$w_{21}$	$w_{22}$	$w_{23}$
$w_{31}$	$w_{32}$	$w_{33}$

 $=$ 

$y_{11}$	$y_{12}$	$y_{13}$
$y_{21}$	$y_{22}$	$y_{23}$
$y_{31}$	$y_{32}$	$y_{33}$

  
$$y_{12} = z_{12}w_{11} + z_{13}w_{12} + z_{14}w_{13} + z_{22}w_{21} + \dots$$

# Convolutions in detail

- Convolutions are a basic primitive in many computer vision and image processing algorithms
- Idea is to “slide” the weights  $h \times h$  weight  $m$  (called a filter, with kernel size  $h$ ) over the image to produce a new image, written  $y = x * m$

$z_{11}$	$z_{12}$	$z_{13}$	$z_{14}$	$z_{15}$
$z_{21}$	$z_{22}$	$z_{23}$	$z_{24}$	$z_{25}$
$z_{31}$	$z_{32}$	$z_{33}$	$z_{34}$	$z_{35}$
$z_{41}$	$z_{42}$	$z_{43}$	$z_{44}$	$z_{45}$
$z_{51}$	$z_{52}$	$z_{53}$	$z_{54}$	$z_{55}$

\*

$w_{11}$	$w_{12}$	$w_{13}$
$w_{21}$	$w_{22}$	$w_{23}$
$w_{31}$	$w_{32}$	$w_{33}$

=

$y_{11}$	$y_{12}$	$y_{13}$
$y_{21}$	$y_{22}$	$y_{23}$
$y_{31}$	$y_{32}$	$y_{33}$

$$y_{13} = z_{13}w_{11} + z_{14}w_{12} + z_{15}w_{13} + z_{23}w_{21} + \dots$$

Boston University School of College name here



# Convolutions in detail

- Convolutions are a basic primitive in many computer vision and image processing algorithms
- Idea is to “slide” the weights  $h \times h$  weight  $m$  (called a filter, with kernel size  $h$ ) over the image to produce a new image, written  $y = x * m$

$z_{11}$	$z_{12}$	$z_{13}$	$z_{14}$	$z_{15}$
$z_{21}$	$z_{22}$	$z_{23}$	$z_{24}$	$z_{25}$
$z_{31}$	$z_{32}$	$z_{33}$	$z_{34}$	$z_{35}$
$z_{41}$	$z_{42}$	$z_{43}$	$z_{44}$	$z_{45}$
$z_{51}$	$z_{52}$	$z_{53}$	$z_{54}$	$z_{55}$

 $*$ 

$w_{11}$	$w_{12}$	$w_{13}$
$w_{21}$	$w_{22}$	$w_{23}$
$w_{31}$	$w_{32}$	$w_{33}$

 $=$ 

$y_{11}$	$y_{12}$	$y_{13}$
$y_{21}$	$y_{22}$	$y_{23}$
$y_{31}$	$y_{32}$	$y_{33}$

$$y_{21} = z_{21}w_{11} + z_{22}w_{12} + z_{23}w_{13} + z_{31}w_{21} + \dots$$



# Convolutions in detail

- Convolutions are a basic primitive in many computer vision and image processing algorithms
- Idea is to “slide” the weights  $h \times h$  weight  $m$  (called a filter, with kernel size  $h$ ) over the image to produce a new image, written  $y = x * m$

$z_{11}$	$z_{12}$	$z_{13}$	$z_{14}$	$z_{15}$
$z_{21}$	$z_{22}$	$z_{23}$	$z_{24}$	$z_{25}$
$z_{31}$	$z_{32}$	$z_{33}$	$z_{34}$	$z_{35}$
$z_{41}$	$z_{42}$	$z_{43}$	$z_{44}$	$z_{45}$
$z_{51}$	$z_{52}$	$z_{53}$	$z_{54}$	$z_{55}$

 $*$ 

$w_{11}$	$w_{12}$	$w_{13}$
$w_{21}$	$w_{22}$	$w_{23}$
$w_{31}$	$w_{32}$	$w_{33}$

 $=$ 

$y_{11}$	$y_{12}$	$y_{13}$
$y_{21}$	$y_{22}$	$y_{23}$
$y_{31}$	$y_{32}$	$y_{33}$

$$y_{22} = z_{22}w_{11} + z_{23}w_{12} + z_{24}w_{13} + z_{32}w_{21} + \dots$$

# Convolutions in detail

- Convolutions are a basic primitive in many computer vision and image processing algorithms
- Idea is to “slide” the weights  $h \times h$  weight  $m$  (called a filter, with kernel size  $h$ ) over the image to produce a new image, written  $y = x * m$

$z_{11}$	$z_{12}$	$z_{13}$	$z_{14}$	$z_{15}$
$z_{21}$	$z_{22}$	$z_{23}$	$z_{24}$	$z_{25}$
$z_{31}$	$z_{32}$	$z_{33}$	$z_{34}$	$z_{35}$
$z_{41}$	$z_{42}$	$z_{43}$	$z_{44}$	$z_{45}$
$z_{51}$	$z_{52}$	$z_{53}$	$z_{54}$	$z_{55}$

 $*$ 

$w_{11}$	$w_{12}$	$w_{13}$
$w_{21}$	$w_{22}$	$w_{23}$
$w_{31}$	$w_{32}$	$w_{33}$

 $=$ 

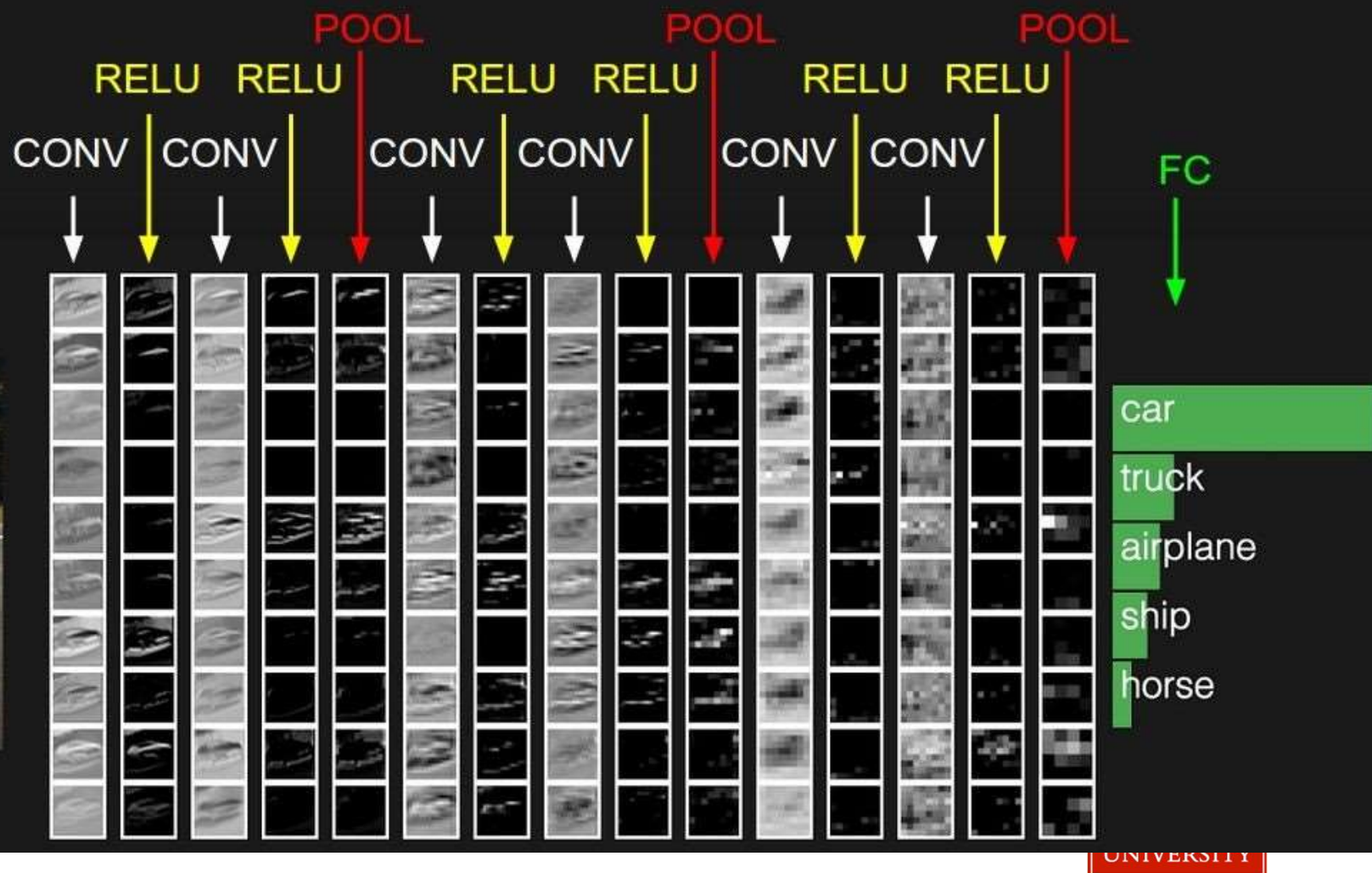
$y_{11}$	$y_{12}$	$y_{13}$
$y_{21}$	$y_{22}$	$y_{23}$
$y_{31}$	$y_{32}$	$y_{33}$

  
$$y_{23} = z_{23}w_{11} + z_{24}w_{12} + z_{25}w_{13} + z_{33}w_{21} + \dots$$

Boston University School/college name here



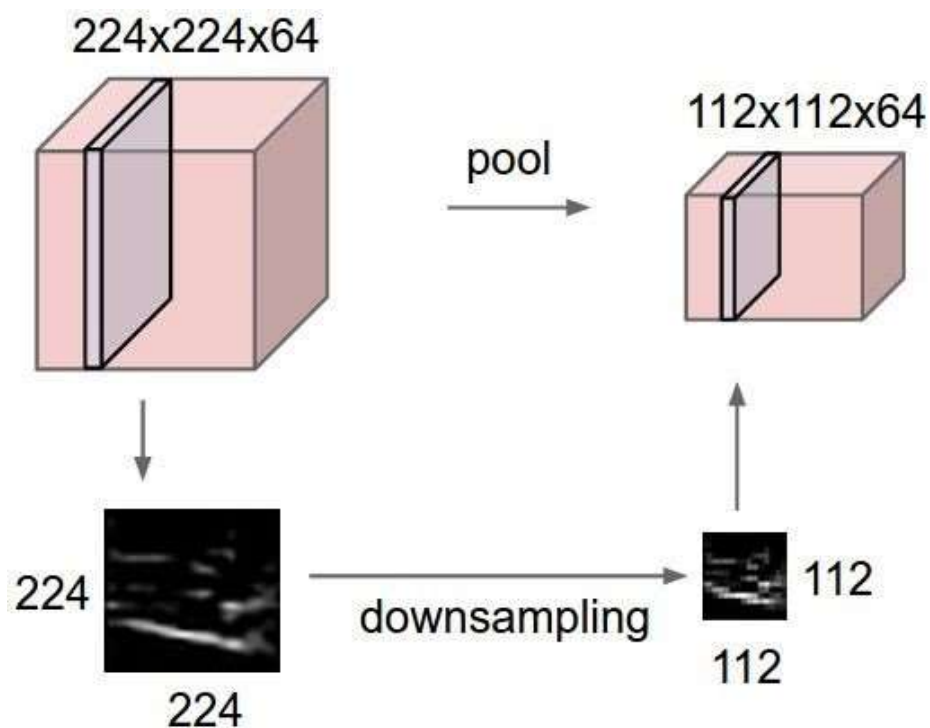
two more layers to go: POOL/FC



Slide credit: Fei-Fei Li, Andrej Karpathy and Justin Johnson

# Pooling layer

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

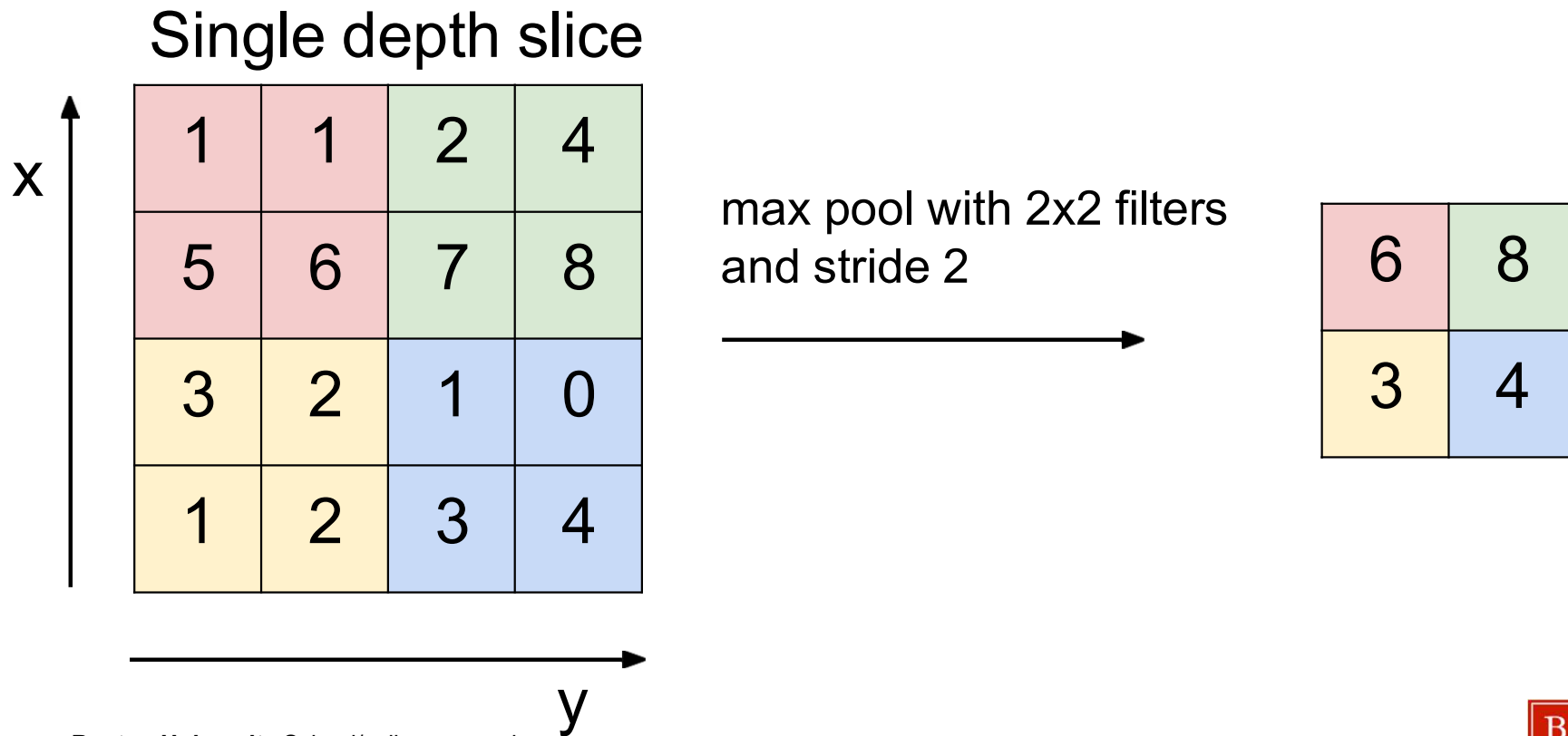


Boston University School/college name here



Slide credit: Fei-Fei Li, Andrej Karpathy and Justin Johnson

# MAX POOLING

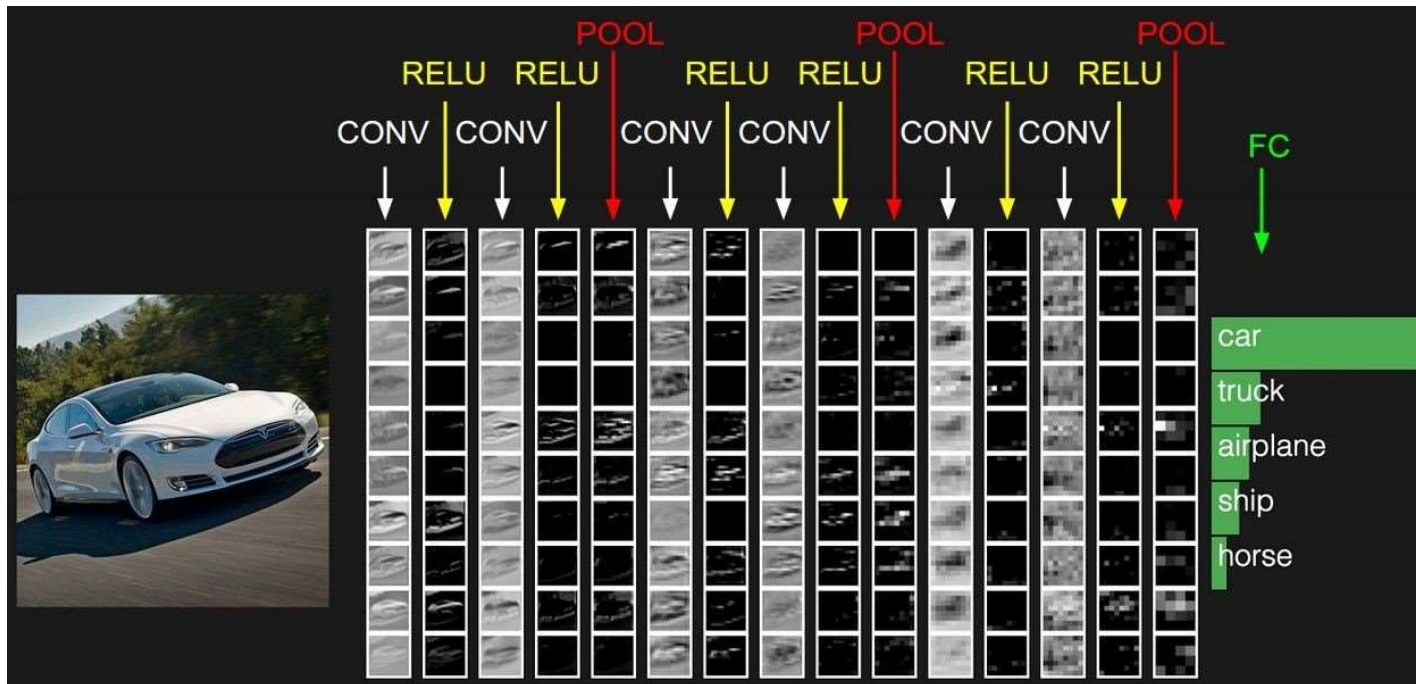


Boston University School/college name here



# Fully Connected Layer (FC layer)

Contains neurons that connect to the entire input volume, as in ordinary Neural Networks



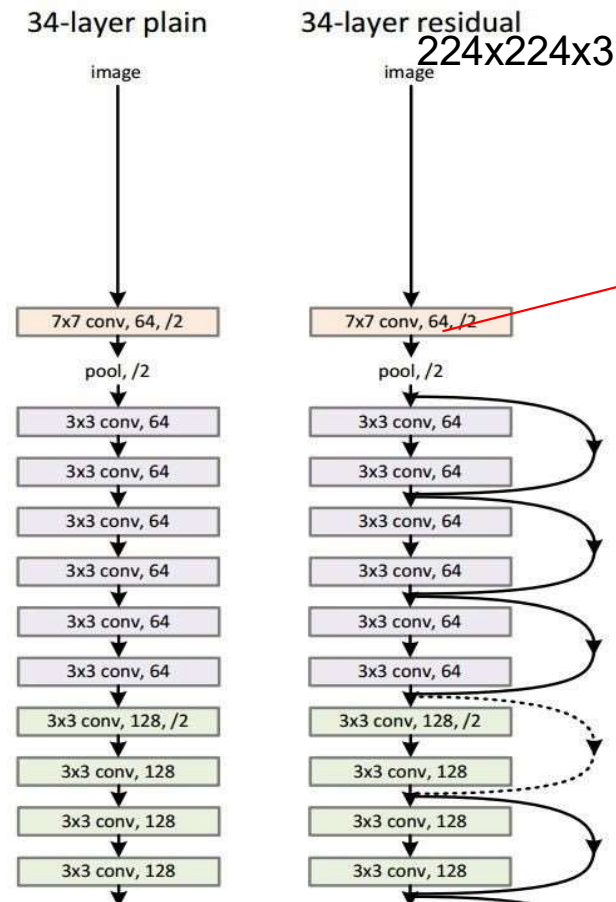
Boston University School/college name here



Slide credit: Fei-Fei Li, Andrej Karpathy and Justin Johnson

# Case Study: ResNet

[He et al., 2015]



spatial dimension  
only 56x56!

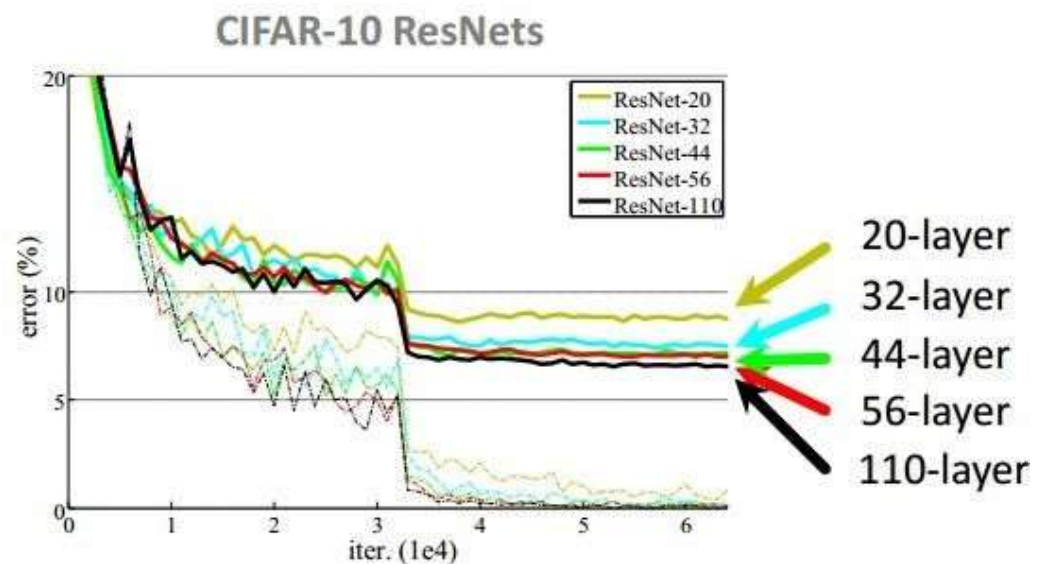
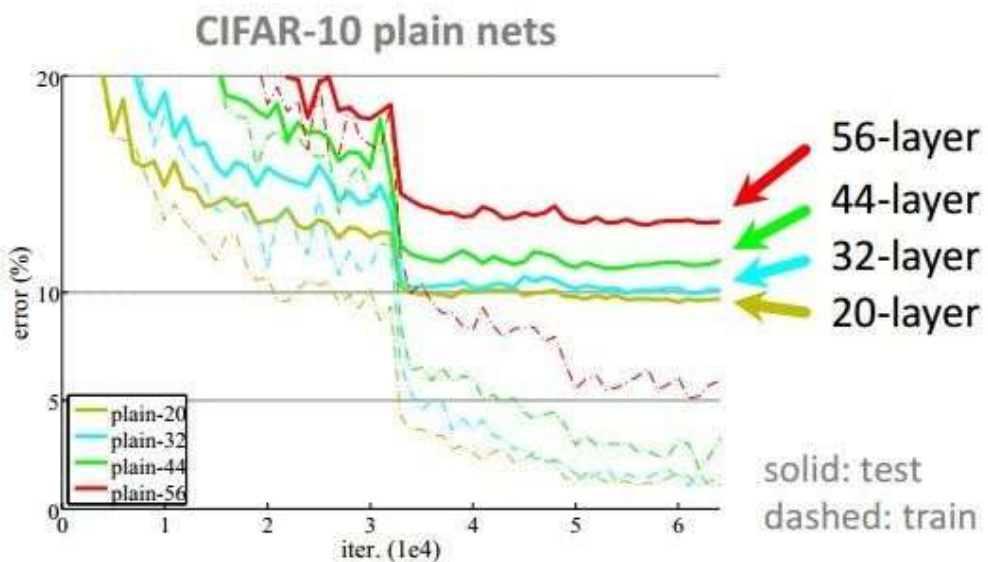
Boston University School/college name here



Slide credit: Fei-Fei Li, Andrej Karpathy and Justin Johnson



# CIFAR-10 experiments



Boston University School/college name here



Slide credit: Fei-Fei Li, Andrej Karpathy and Justin Johnson



# The End

Thanks for your attention.

I would be glad if you have any question.