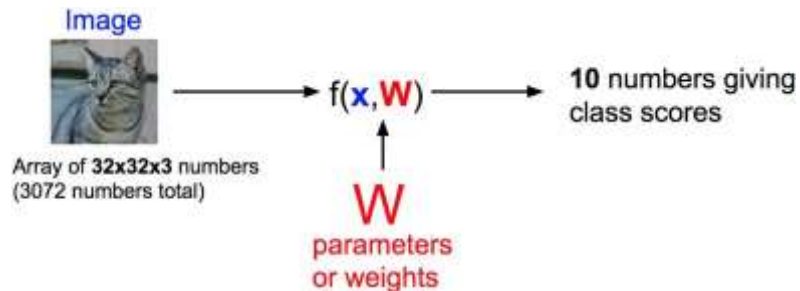# Image Classification with CNNs

Fei-Fei Li, Yunzhu Li, Ruohan Gao
2023
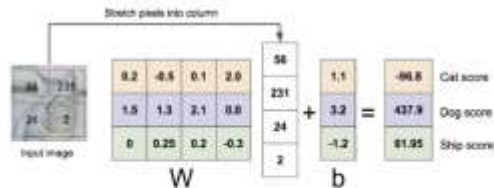
# Recap: Image Classification with Linear Classifier



$$f(x,W) = Wx + b$$
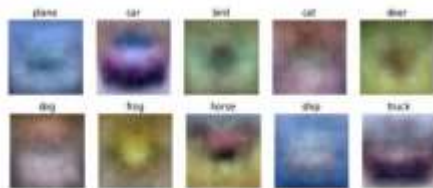
# Recap: Loss Function

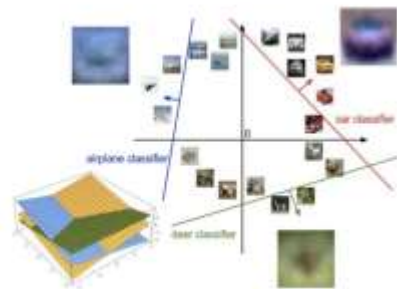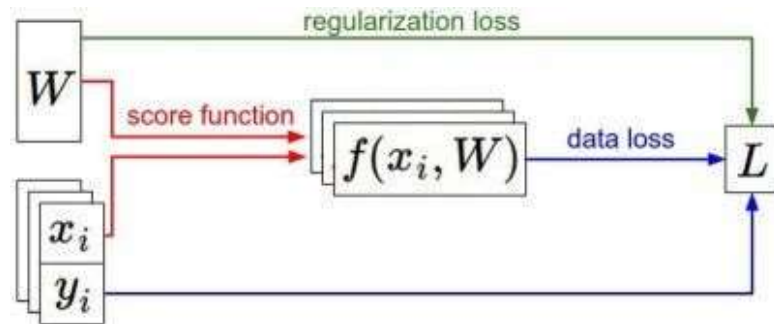- We have some dataset of (x,y)
- We have a **score function:** $s = f(x; W) = Wx$
- We have a **loss function**:

$$L_i = -\log\left(\frac{e^{s_{y_i}}}{\sum_j e^{s_j}}\right)$$ Softmax

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

$$L = \frac{1}{N} \sum_{i=1}^{N} L_i + R(W)$$ Full loss

# Recap: Optimization

# Problem: Linear Classifiers are not very powerful

## Visual Viewpoint



Linear classifiers learn
one template per class

## Geometric Viewpoint



Linear classifiers
can only draw linear
decision boundaries

# Last time: Neural Networks

Linear score function:

$$f = Wx$$

2-layer Neural Network

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

# Last time: Computation Graph

$$f = Wx$$

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

x

W

\* → **s** (scores) → hinge loss → + → L

R

$$R(W)$$

# Last time: Backpropagation
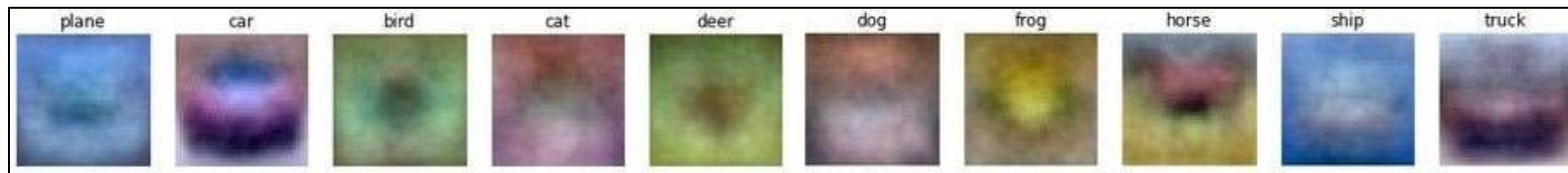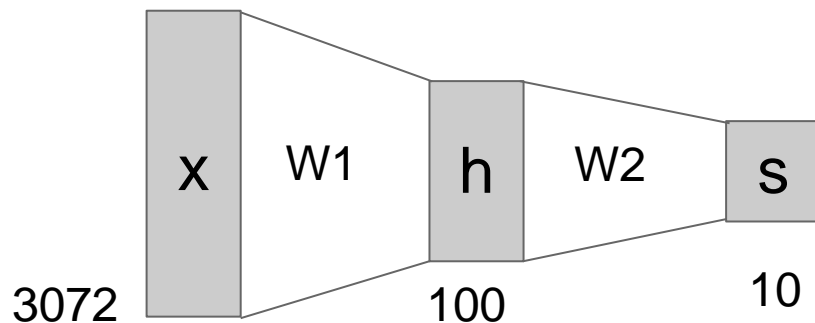


$x$

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z}\frac{\partial z}{\partial x}$$

"local gradient"

$$\frac{\partial z}{\partial x}$$

$$\frac{\partial z}{\partial y}$$

f

"Downstream gradients"

$y$

$$\frac{\partial L}{\partial y} = \frac{\partial L}{\partial z}\frac{\partial z}{\partial y}$$

$z$

$$\frac{\partial L}{\partial z}$$

"Upstream gradient"

# Backprop with Matrices (or Tensors)

Loss L still a scalar!

dL/dx always has the same shape as x!

$[D_x \times M_x]$  $x$

$[D_x \times M_x]$  $\dfrac{\partial L}{\partial x} = \dfrac{\partial z}{\partial x}\dfrac{\partial L}{\partial z}$

"Downstream gradients"

Matrix-vector multiply

$[D_y \times M_y]$  $y$

$\dfrac{\partial L}{\partial y} = \dfrac{\partial z}{\partial y}\dfrac{\partial L}{\partial z}$

$[D_y \times M_y]$

"local gradients"

$\dfrac{\partial z}{\partial x}$  $[(D_x \times M_x) \times (D_z \times M_z)]$

$\dfrac{\partial z}{\partial y}$  $[(D_y \times M_y) \times (D_z \times M_z)]$

Jacobian matrices

For each element of y, how much does it influence each element of z?

$z$  $[D_z \times M_z]$

$\dfrac{\partial L}{\partial z}$  $[D_z \times M_z]$

"Upstream gradient"
For each element of z, how much does it influence L?

# **Image Classification**: A core task in Computer Vision

(assume given a set of labels)
{dog, cat, truck, plane, ...}

⟶

cat
dog
bird
deer
truck

# Pixel space



f(x) = Wx

plane car bird cat deer

dog frog horse ship truck
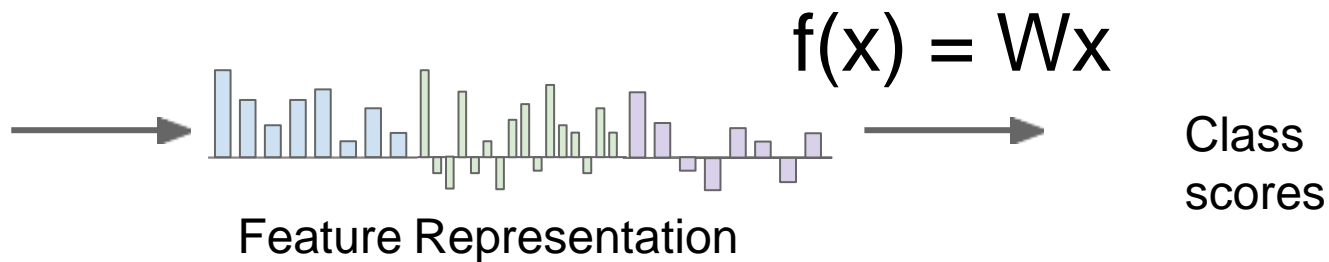
Class scores
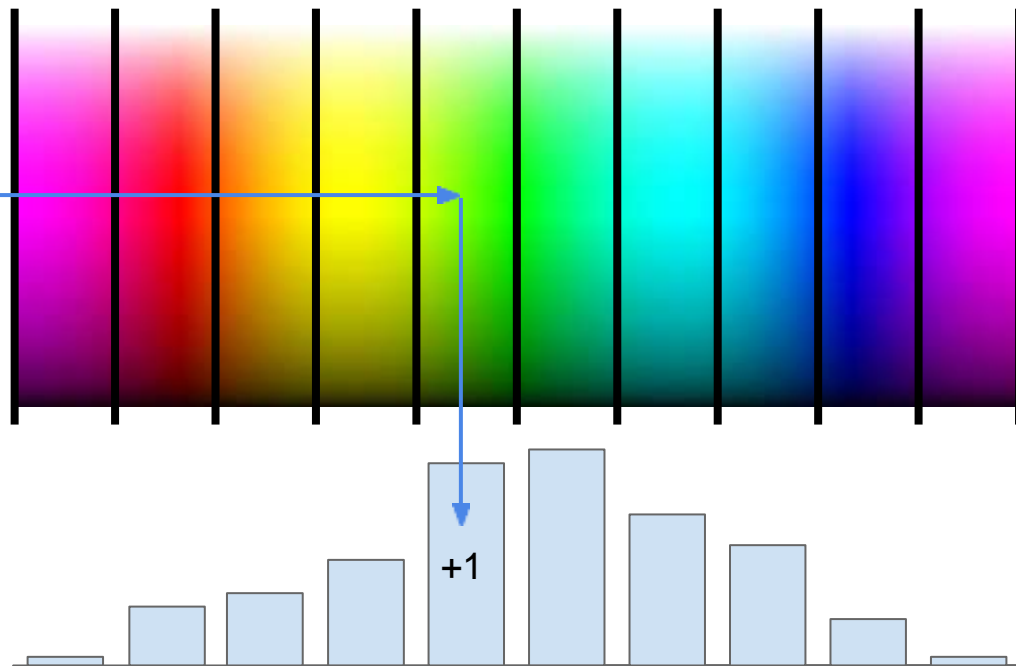
# Image features



$$f(x) = Wx$$

Feature Representation
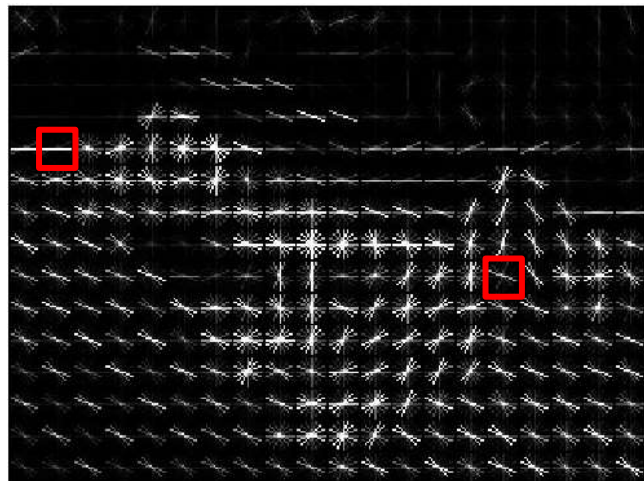
Class scores

# Example: Color Histogram



+1

# Example: Histogram of Oriented Gradients (HoG)



Divide image into 8x8 pixel regions
Within each region quantize edge
direction into 9 bins

Example: 320x240 image gets divided
into 40x30 bins; in each bin there are
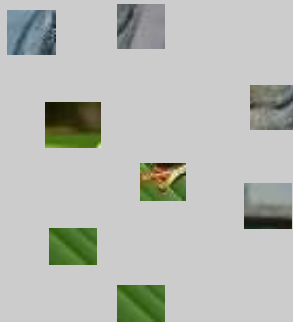9 numbers so feature vector has
30*40*9 = 10,800 numbers

Lowe, "Object recognition from local scale-invariant features", ICCV 1999
Dalal and Triggs, "Histograms of oriented gradients for human detection," CVPR 2005
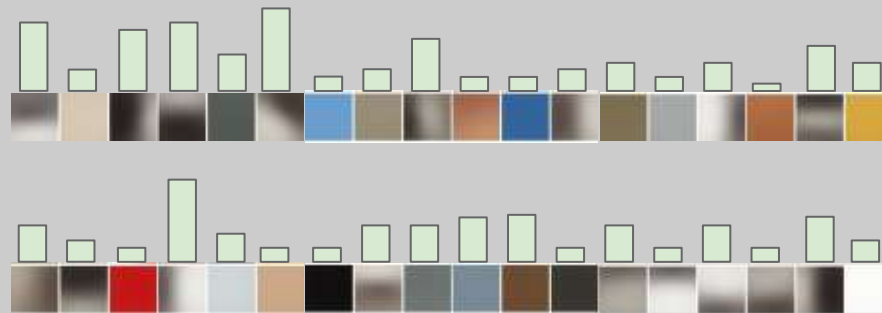
# Example: Bag of Words

**Step 1: Build codebook**



Extract random patches

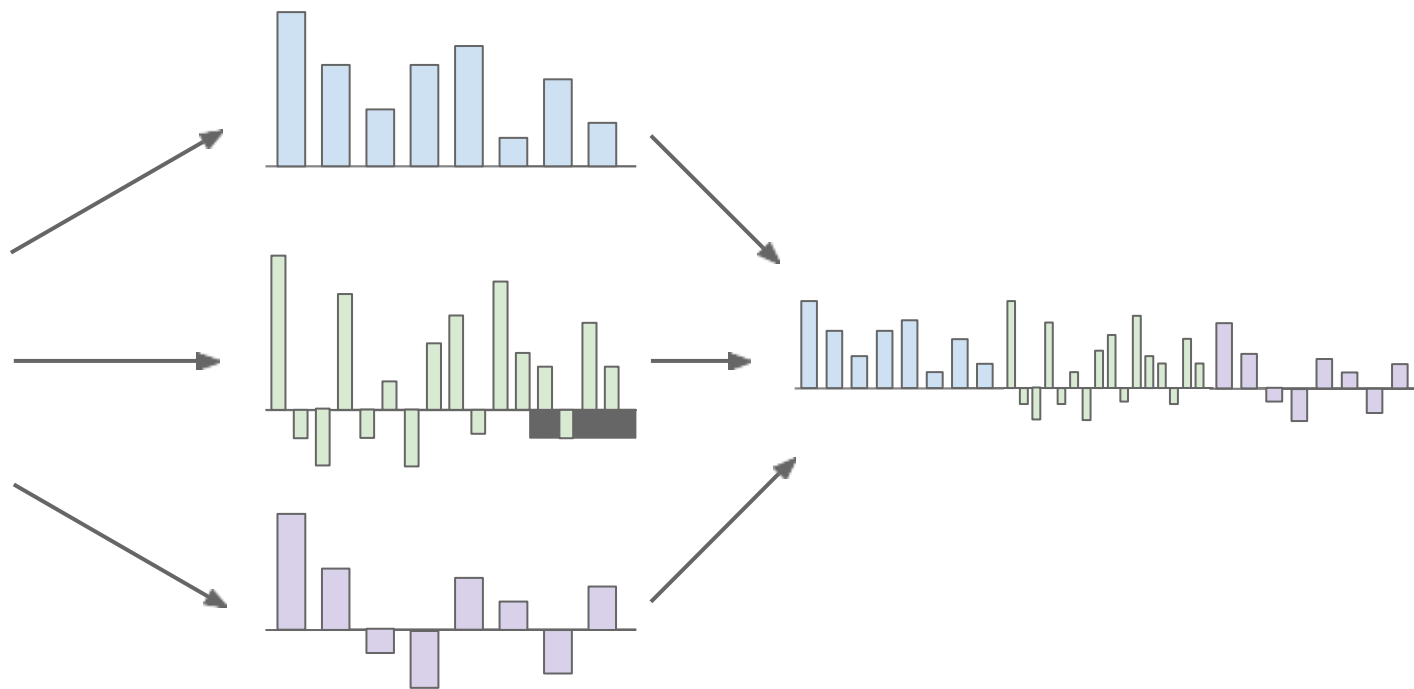Cluster patches to form "codebook" of "visual words"
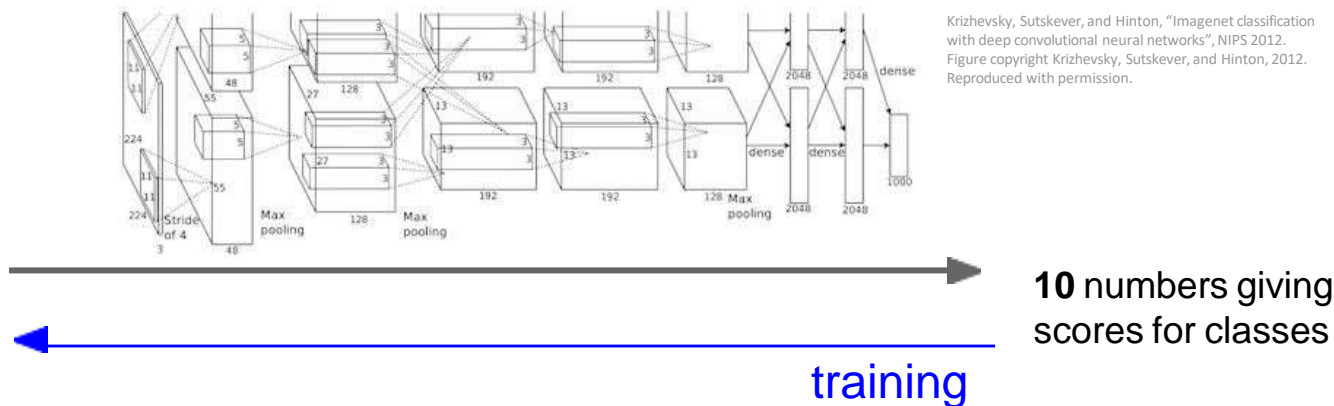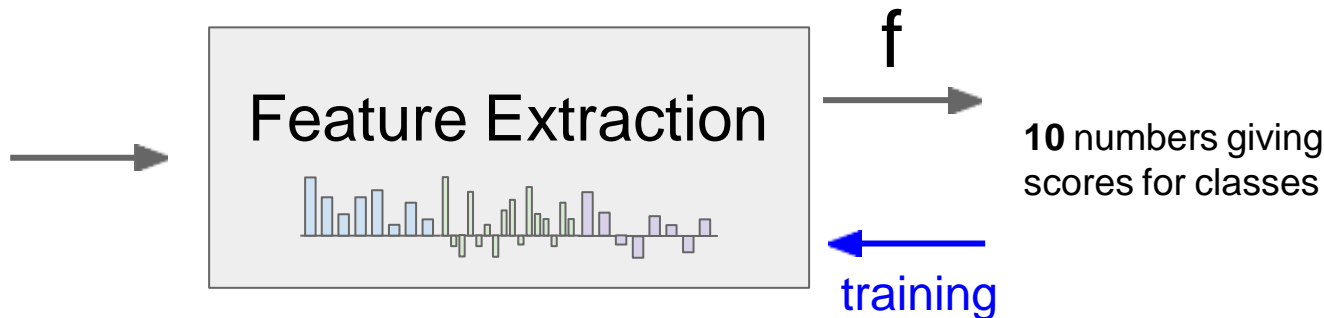
**Step 2: Encode images**



Fei-Fei and Perona, "A bayesian hierarchical model for learning natural scene categories", CVPR 2005
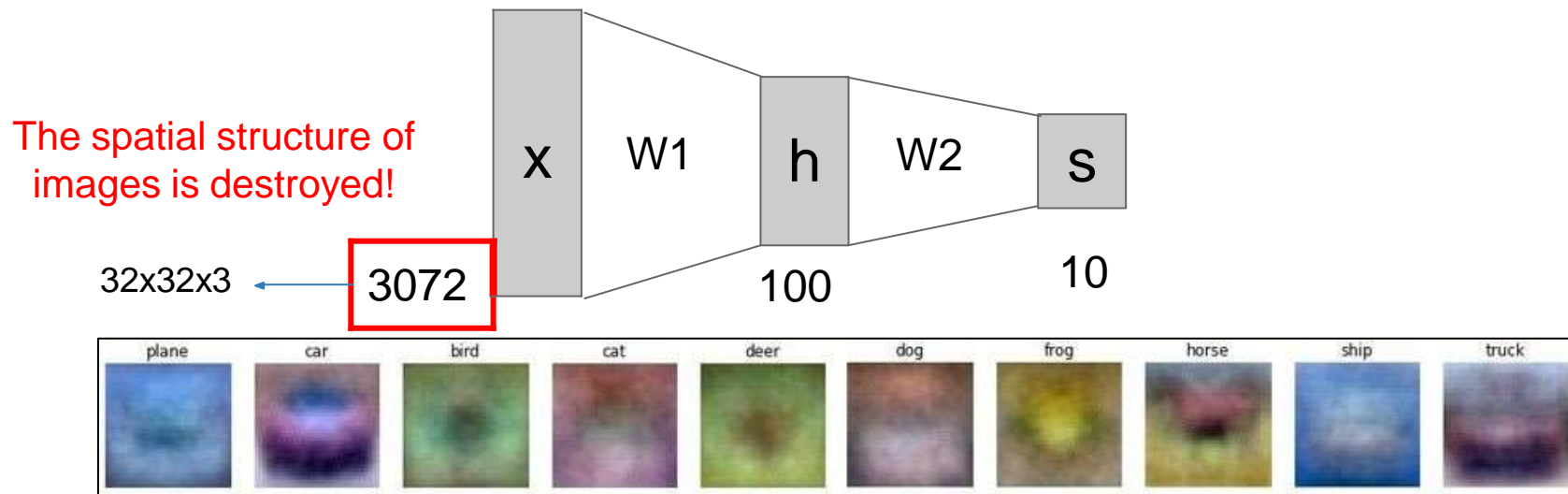
# Image Features

# Image features vs. ConvNets



f

**10** numbers giving scores for classes

training

Krizhevsky, Sutskever, and Hinton, "Imagenet classification with deep convolutional neural networks", NIPS 2012. Figure copyright Krizhevsky, Sutskever, and Hinton, 2012. Reproduced with permission.

**10** numbers giving scores for classes

training

# Last Time: Neural Networks

Linear score function:

2-layer Neural Network

$$f = Wx$$

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

The spatial structure of images is destroyed!

x    W1    h    W2    s

32x32x3    3072    100    10

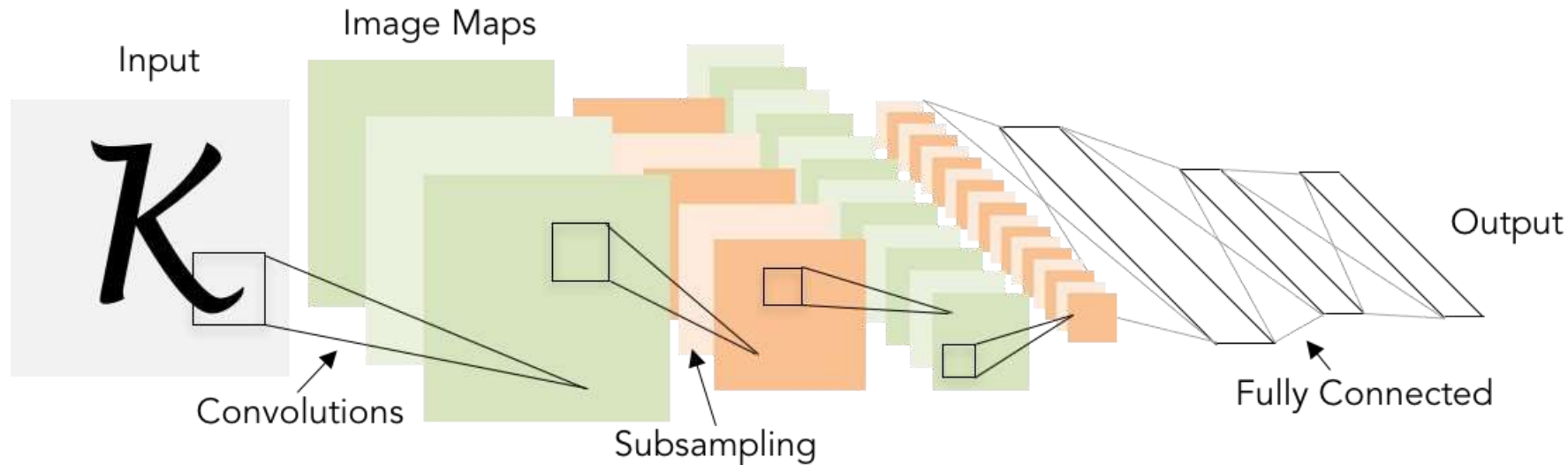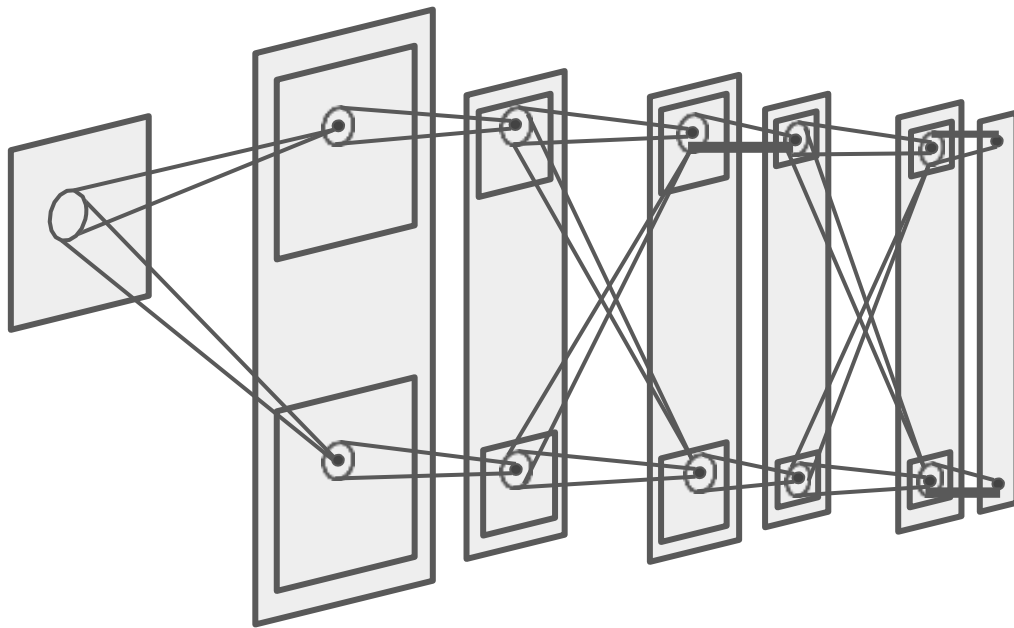# Next: Convolutional Neural Networks



Illustration of LeCun et al. 1998 from CS231n 2017 Lecture 1

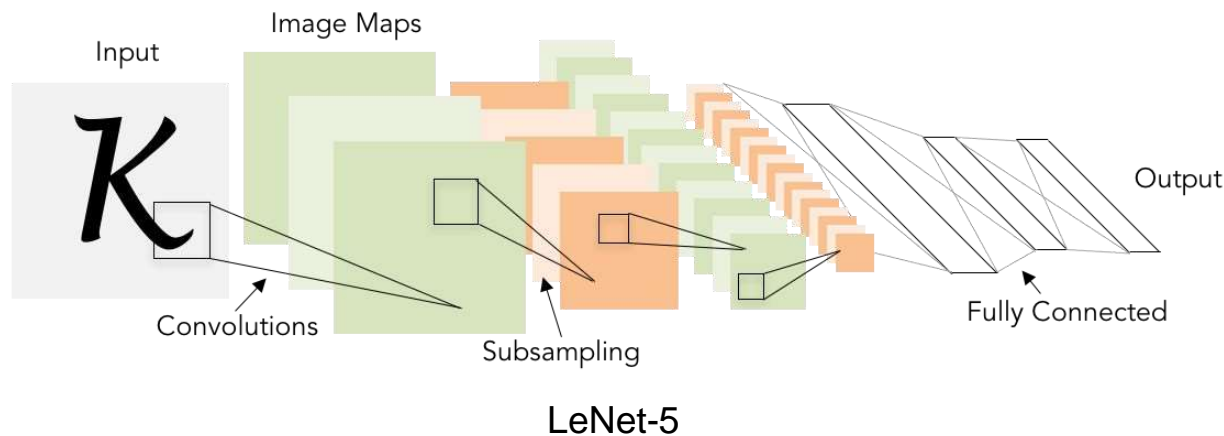# A bit of history:

## **Neocognitron**
## *[Fukushima 1980]*

"sandwich" architecture (SCSCSC...)
simple cells: modifiable parameters
complex cells: perform pooling

# A bit of history:

**Gradient-based learning applied to document recognition**
*[LeCun, Bottou, Bengio, Haffner 1998]*



LeNet-5

# A bit of history:
## ImageNet Classification with Deep Convolutional Neural Networks
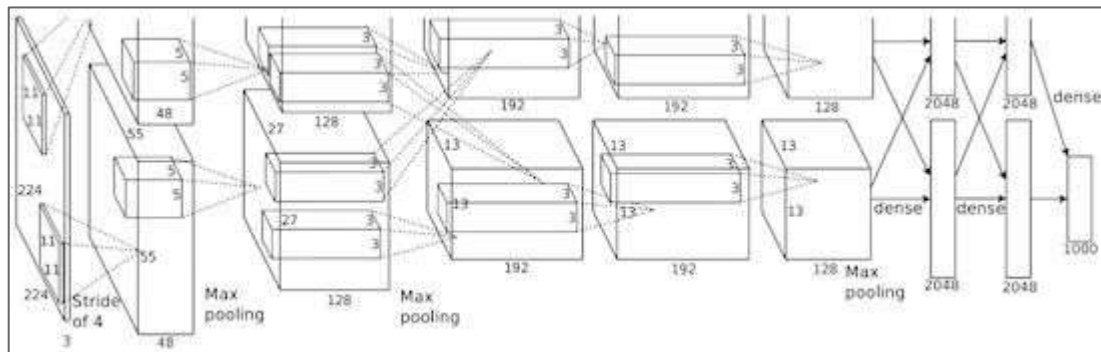*[Krizhevsky, Sutskever, Hinton, 2012]*



Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

"AlexNet"

# Fast-forward to today: ConvNets are everywhere
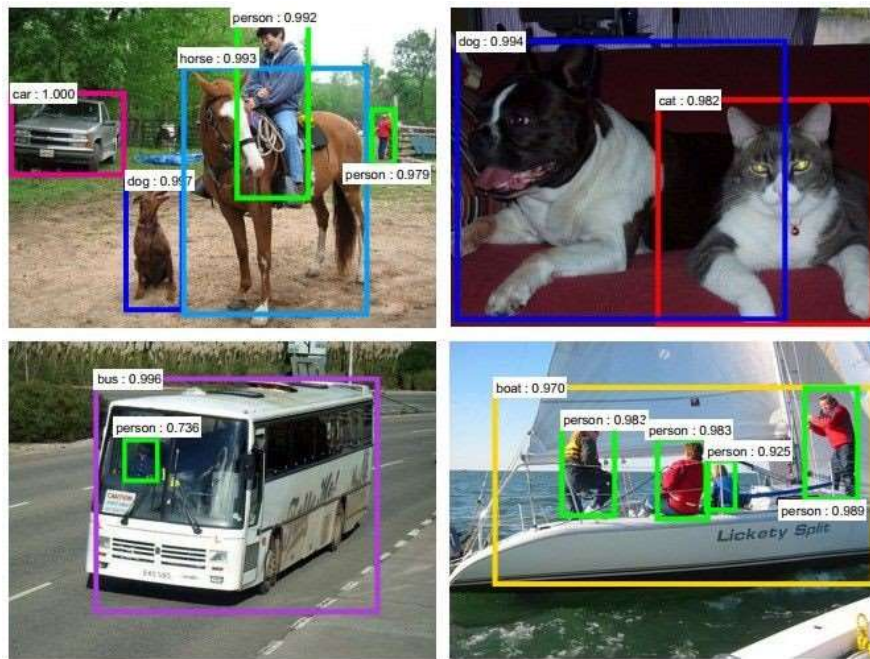
Classification

Retrieval



Figures copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.
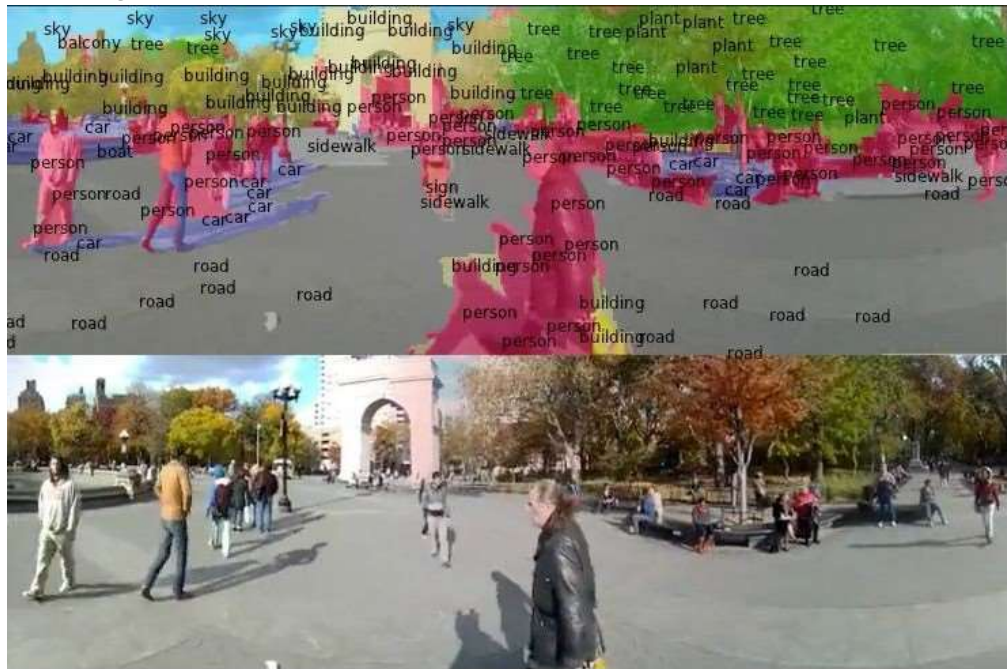
# Fast-forward to today: ConvNets are everywhere

Detection



Segmentation



Figures copyright Shaoqing Ren, Kaiming He, Ross Girschick, Jian Sun, 2015. Reproduced with permission.

*[Faster R-CNN: Ren, He, Girshick, Sun 2015]*

Figures copyright Clement Farabet, 2012.
Reproduced with permission.

*[Farabet et al., 2012]*

# Fast-forward to today: ConvNets are everywhere



self-driving cars

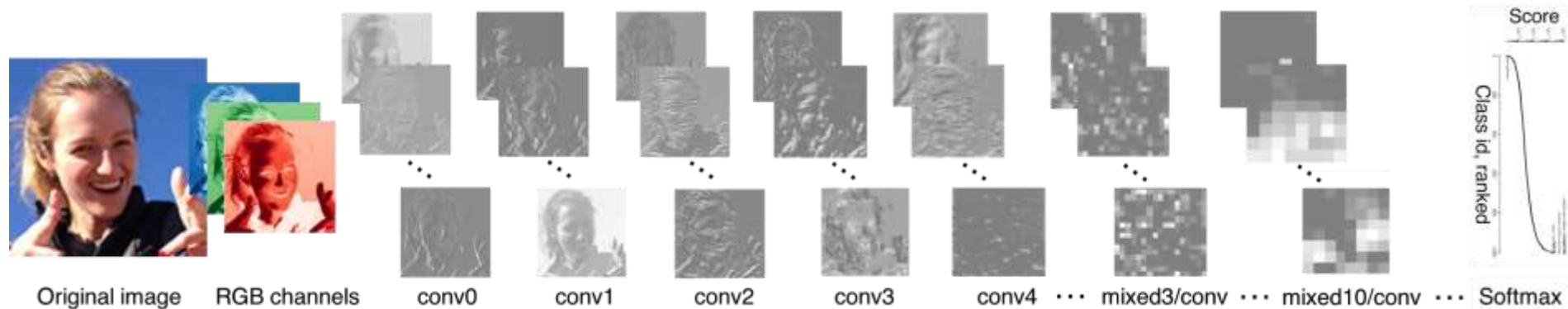Photo by Lane McIntosh. Copyright CS231n 2017.



This image by GBPublic_PR is licensed under CC-BY 2.0

NVIDIA Tesla line
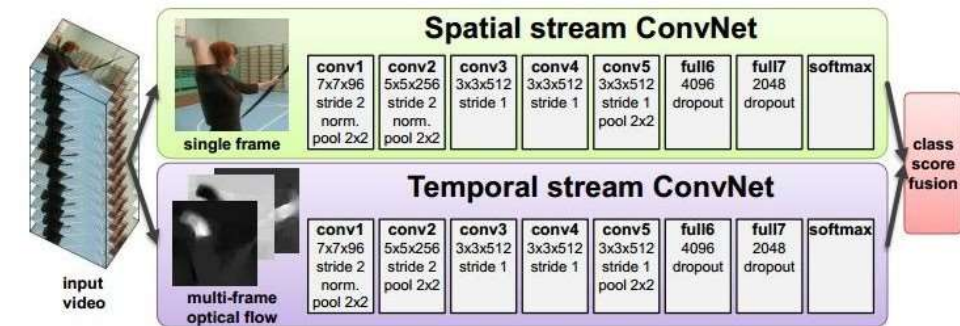(these are the GPUs on rye01.stanford.edu)

Note that for embedded systems a typical setup would involve NVIDIA Tegras, with integrated GPU and ARM-based CPU cores.

# Fast-forward to today: ConvNets are everywhere



Original image    RGB channels    conv0    conv1    conv2    conv3    conv4    ⋯    mixed3/conv    ⋯    mixed10/conv    ⋯    Softmax

*[Taigman et al. 2014]*

Activations of inception-v3 architecture [Szegedy et al. 2015] to image of Emma McIntosh, used with permission. Figure and architecture not from Taigman et al. 2014.



*[Simonyan et al. 2014]*

Figures copyright Simonyan et al., 2014. Reproduced with permission.
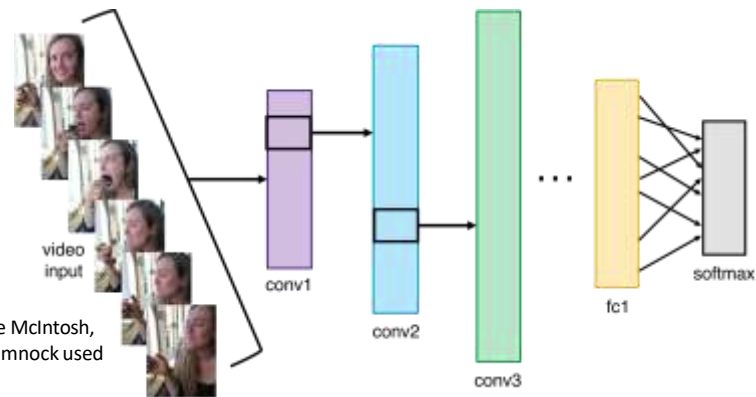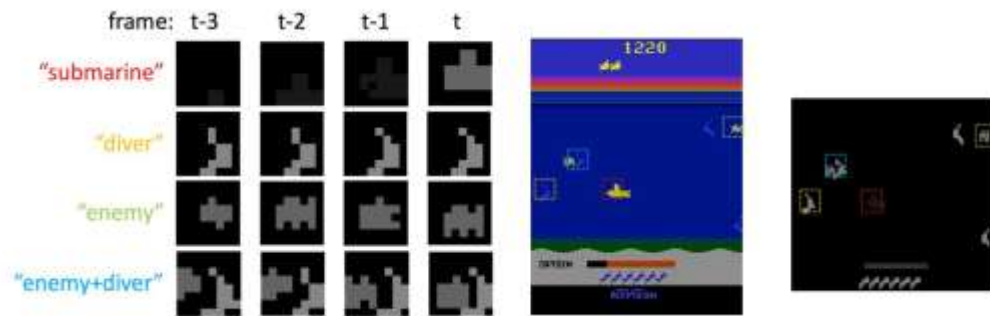
Illustration by Lane McIntosh, photos of Katie Cumnock used with permission.

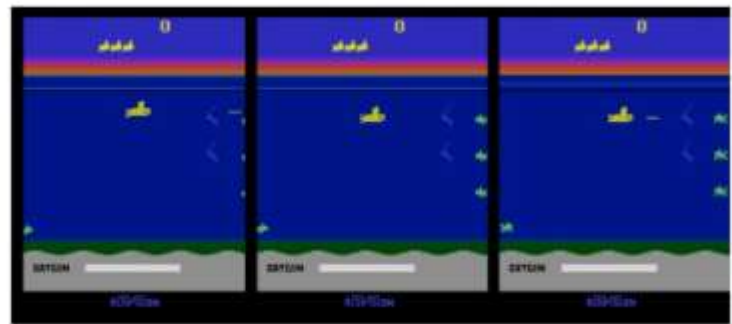# Fast-forward to today: ConvNets are everywhere



Images are examples of pose estimation, not actually from Toshev & Szegedy 2014. Copyright Lane McIntosh.
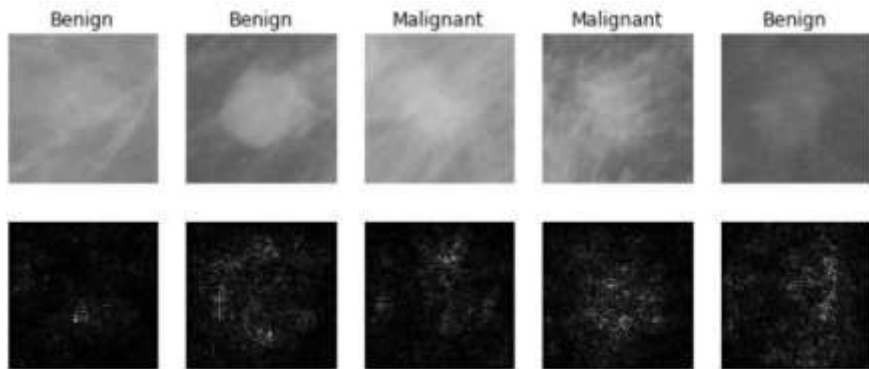
*[Toshev, Szegedy 2014]*



*[Guo et al. 2014]*

Figures copyright Xiaoxiao Guo, Satinder Singh, Honglak Lee, Richard Lewis, and Xiaoshi Wang, 2014. Reproduced with permission.

# Fast-forward to today: ConvNets are everywhere



*[Levy et al. 2016]*

Figure copyright Levy et al. 2016.
Reproduced with permission.



*[Dieleman et al. 2014]*

From left to right: public domain by NASA, usage permitted by
ESA/Hubble, public domain by NASA, and public domain.



*[Sermanet et al. 2011]*
*[Ciresan et al.]*

Photos by Lane McIntosh.
Copyright CS231n 2017.

*Whale recognition, Kaggle Challenge*

*Mnih and Hinton, 2010*

**No errors**



*A white teddy bear sitting in the grass*

**Minor errors**



*A man in a baseball uniform throwing a ball*

**Somewhat related**



*A woman is holding a cat in her hand*

# Image Captioning

*[Vinyals et al., 2015]*
*[Karpathy and Fei-Fei, 2015]*



*A man riding a wave on top of a surfboard*
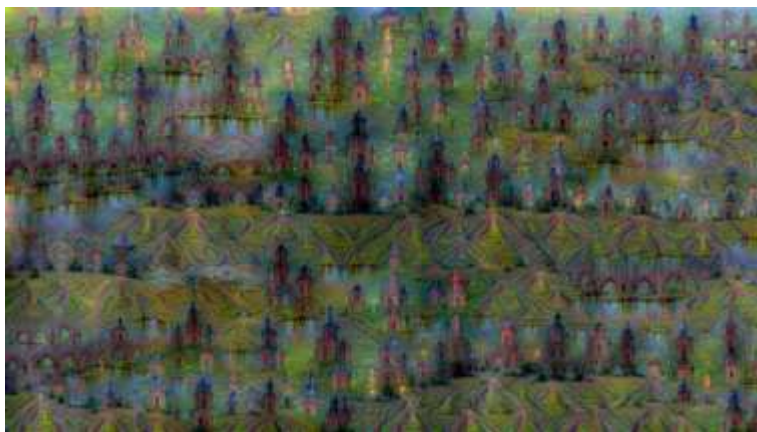


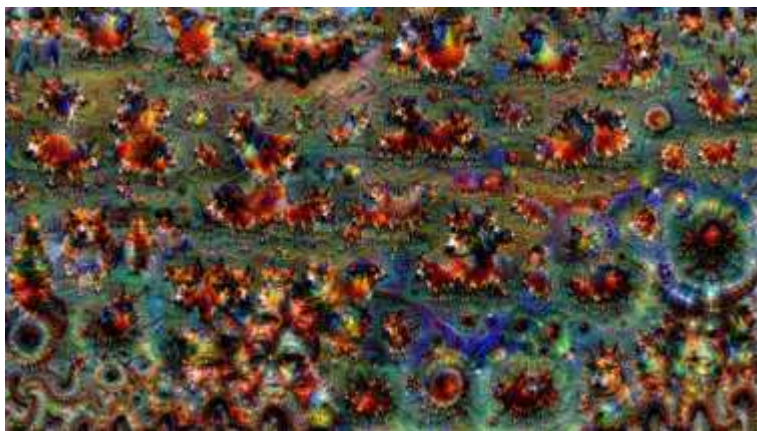*A cat sitting on a suitcase on the floor*



*A woman standing on a beach holding a surfboard*

All images are CC0 Public domain:
https://pixabay.com/en/luggage-antique-cat-1643010/
https://pixabay.com/en/teddy-plush-bears-cute-teddy-bear-1623436/
https://pixabay.com/en/surf-wave-summer-sport-litoral-1668716/
https://pixabay.com/en/woman-female-model-portrait-adult-983967/
https://pixabay.com/en/handstand-lake-meditation-496008/
https://pixabay.com/en/baseball-player-shortstop-infield-1045263/

Captions generated by Justin Johnson using Neuraltalk2

Figures copyright Justin Johnson, 2015. Reproduced with permission. Generated using the Inceptionism approach from a blog post by Google Research.

Original image is CC0 public domain
Starry Night and Tree Roots by Van Gogh are in the public domain
Bokeh image is in the public domain
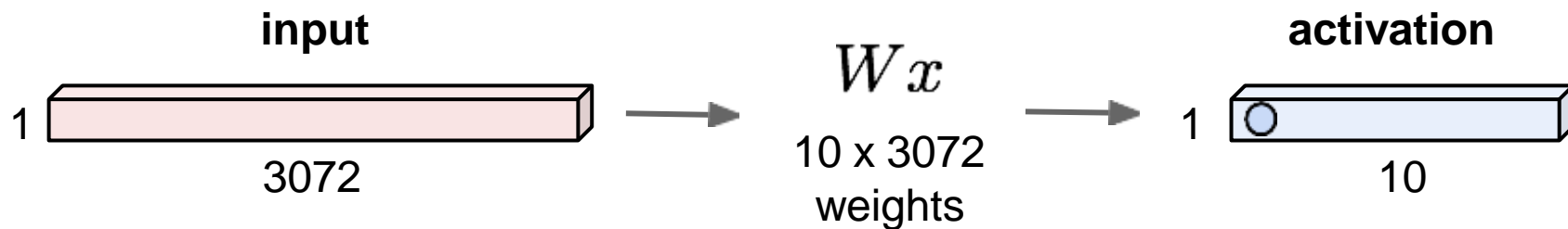Stylized images copyright Justin Johnson, 2017;
reproduced with permission

Gatys et al, "Image Style Transfer using Convolutional Neural Networks", CVPR 2016
Gatys et al, "Controlling Perceptual Factors in Neural Style Transfer", CVPR 2017

# Convolutional Neural Networks

# Recap: Fully Connected Layer
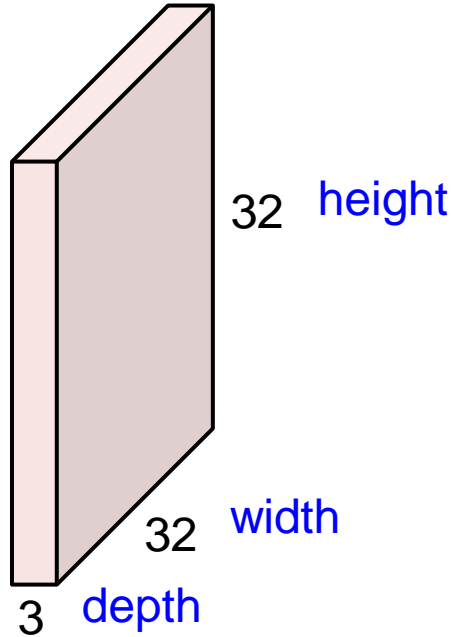
32x32x3 image -> stretch to 3072 x 1

**input**

1

3072

$Wx$

10 x 3072
weights

**activation**

1

10

# Fully Connected Layer

32x32x3 image -> stretch to 3072 x 1

**input**

1

3072

$Wx$

10 x 3072
weights

**activation**

1

10

**1 number:**
the result of taking a dot product
between a row of W and the input
(a 3072-dimensional dot product)

# Convolution Layer

32x32x3 image -> preserve spatial structure

32 height
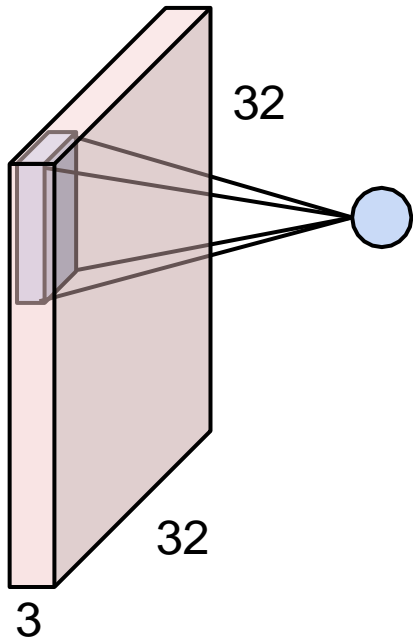
32 width

3 depth

# Convolution Layer

32x32x3 image



5x5x3 filter

32

32

3

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

# Convolution Layer

32x32x3 image



32

32

3

5x5x3 filter



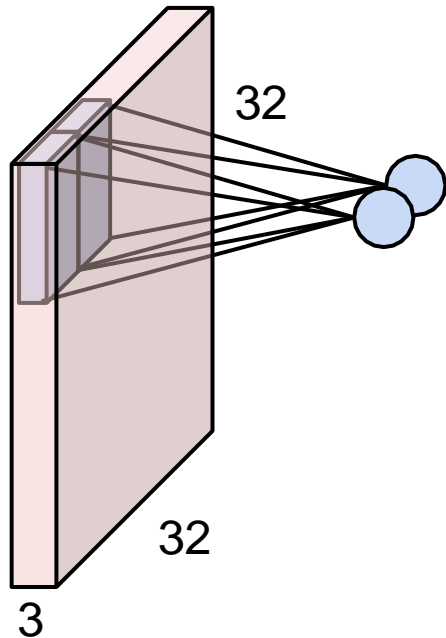**Convolve** the filter with the image i.e. "slide over the image spatially, computing dot products"
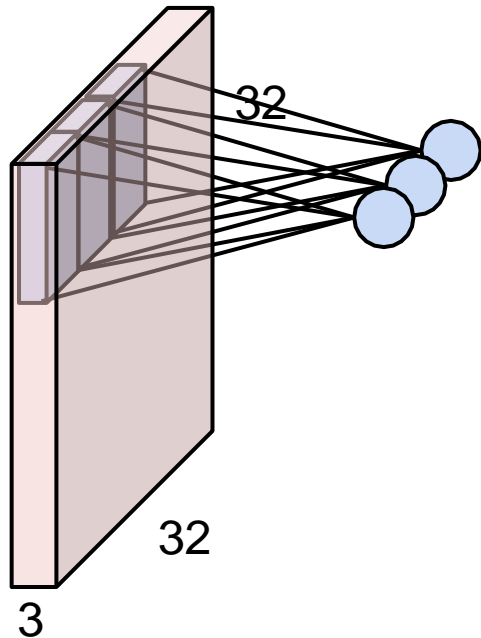
# Convolution Layer



32x32x3 image

5x5x3 filter $w$

**1 number:**
the result of taking a dot product between the filter and a small 5x5x3 chunk of the image (i.e. 5*5*3 = 75-dimensional dot product + bias)

$$w^T x + b$$

# Convolution Layer

32

32

3

# Convolution Layer



32

32

3

# Convolution Layer



32

32

3

# Convolution Layer



32

32

3

# Convolution Layer



**activation map**

32x32x3 image

5x5x3 filter

32

32

3

convolve (slide) over all spatial locations

28

28

1

# Convolution Layer

consider a second, green filter

32x32x3 image
5x5x3 filter

32

32

3

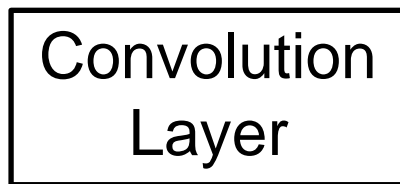convolve (slide) over all
spatial locations

**activation maps**

28

28

1

# Convolution Layer

3x32x32 image

Consider 6 filters,
each 3x5x5

6 activation maps,
each 1x28x28



32

32

3

6x3x5x5
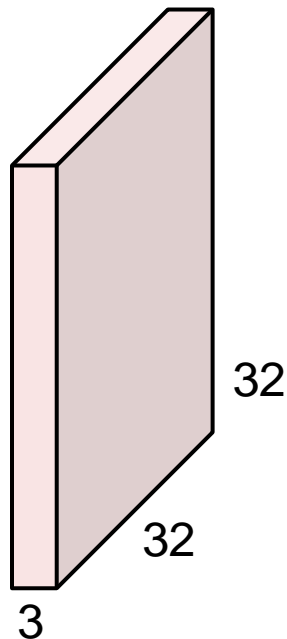filters

Convolution
Layer

Stack activations to get a
6x28x28 output image!

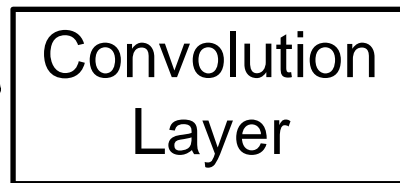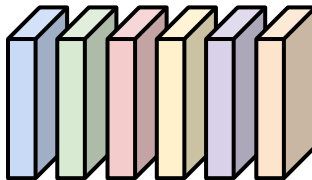Slide inspiration: Justin Johnson

# Convolution Layer

3x32x32 image

Also 6-dim bias vector:

6 activation maps,
each 1x28x28



Convolution
Layer

32

32

3

6x3x5x5
filters

Stack activations to get a
6x28x28 output image!

# Convolution Layer

**3x32x32 image**

Also 6-dim bias vector:

28x28 grid, at each point a 6-dim vector



Convolution Layer

32

32

3

6x3x5x5 filters

Stack activations to get a 6x28x28 output image!

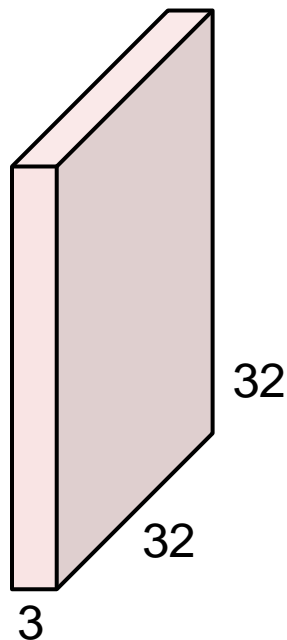Slide inspiration: Justin Johnson

# Convolution Layer



2x3x32x32
Batch of images

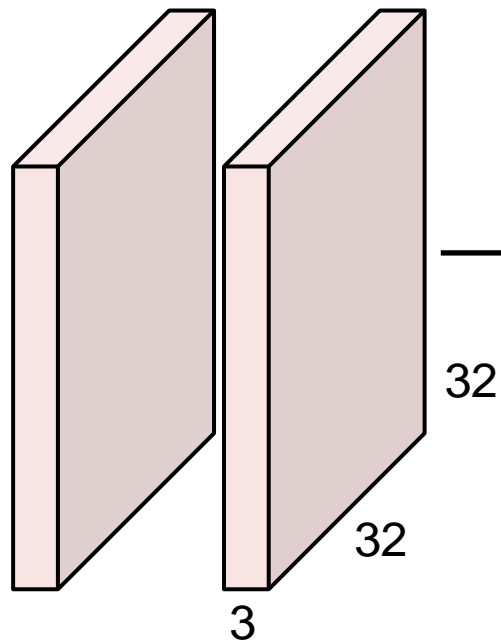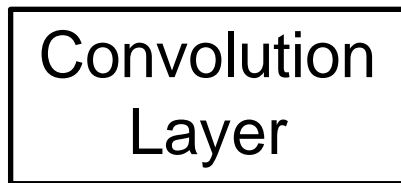Also 6-dim bias vector:

Convolution Layer

6x3x5x5 filters

2x6x28x28
Batch of outputs

32

32

3

Slide inspiration: Justin Johnson

# Convolution Layer

$N \times C_{in} \times H \times W$
Batch of images

Also $C_{out}$-dim bias vector:

$N \times C_{out} \times H' \times W'$
Batch of outputs

Convolution Layer

H

W

$C_{in}$

$C_{out} \times C_{in} \times K_w \times K_h$
filters

$C_{out}$

Slide inspiration: Justin Johnson

**Preview:** ConvNet is a sequence of Convolution Layers



32

28

CONV
e.g. 6
5x5x3
filters

32

28

3

6

**Preview:** ConvNet is a sequence of Convolution Layers



32
32
3

CONV

e.g. 6
5x5x3
filters

28
28
6

CONV

e.g. 10
5x5x**6**
filters

24
24
10

CONV

....

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

**Preview:** What do convolutional filters learn?



32

32

3

Conv → ReLU

28

28

Linear classifier: One template per class



| plane | car | bird | cat | deer |
| dog | frog | horse | ship | truck |

**Preview:** What do convolutional filters learn?

MLP: Bank of whole-image templates

**Preview:** What do convolutional filters learn?

32

32

3

Conv → ReLU →

28

28

First-layer conv filters: local image templates
(Often learns oriented edges, opposing colors)



AlexNet: 64 filters, each 3x11x11

one filter =>
one activation map

example 5x5 filters
(32 total)

Activations:

We call the layer convolutional because it is related to convolution of two signals:

$$f[x,y] * g[x,y] \;=\; \sum_{n_1=-\infty}^{\infty} \sum_{n_2=-\infty}^{\infty} f[n_1,n_2] \cdot g[x-n_1, y-n_2]$$

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

Figure copyright Andrej Karpathy.

preview:

# A closer look at spatial dimensions:

**activation map**

32x32x3 image

5x5x3 filter

32

32

3

convolve (slide) over all spatial locations

28

28

1

A closer look at spatial dimensions:

7



7

7x7 input (spatially)
assume 3x3 filter

A closer look at spatial dimensions:

7

7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

7



7

7x7 input (spatially)
assume 3x3 filter

A closer look at spatial dimensions:

7

7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

7



7

7x7 input (spatially)
assume 3x3 filter

**=> 5x5 output**

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**

A closer look at spatial dimensions:

7

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?**

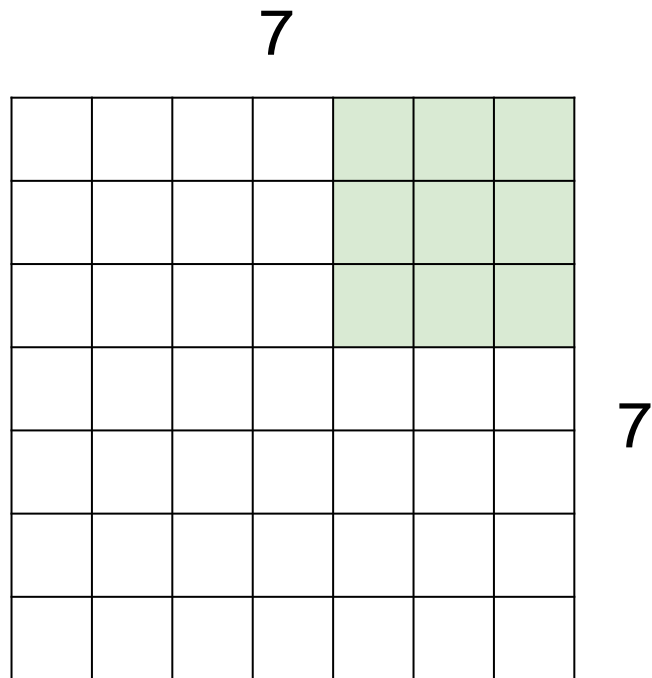A closer look at spatial dimensions:

7



7

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

**doesn't fit!**
cannot apply 3x3 filter on
7x7 input with stride 3.

Output size:
**(N - F) / stride + 1**

e.g. N = 7, F = 3:
stride 1 => (7 - 3)/1 + 1 = 5
stride 2 => (7 - 3)/2 + 1 = 3
stride 3 => (7 - 3)/3 + 1 = 2.33 :\

# 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?

(recall:)
(N - F) / 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!**

(recall:)
(N + 2P - F) / 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 general, common to see CONV layers with stride 1, filters of size FxF, and zero-padding with (F-1)/2. (will preserve size spatially)
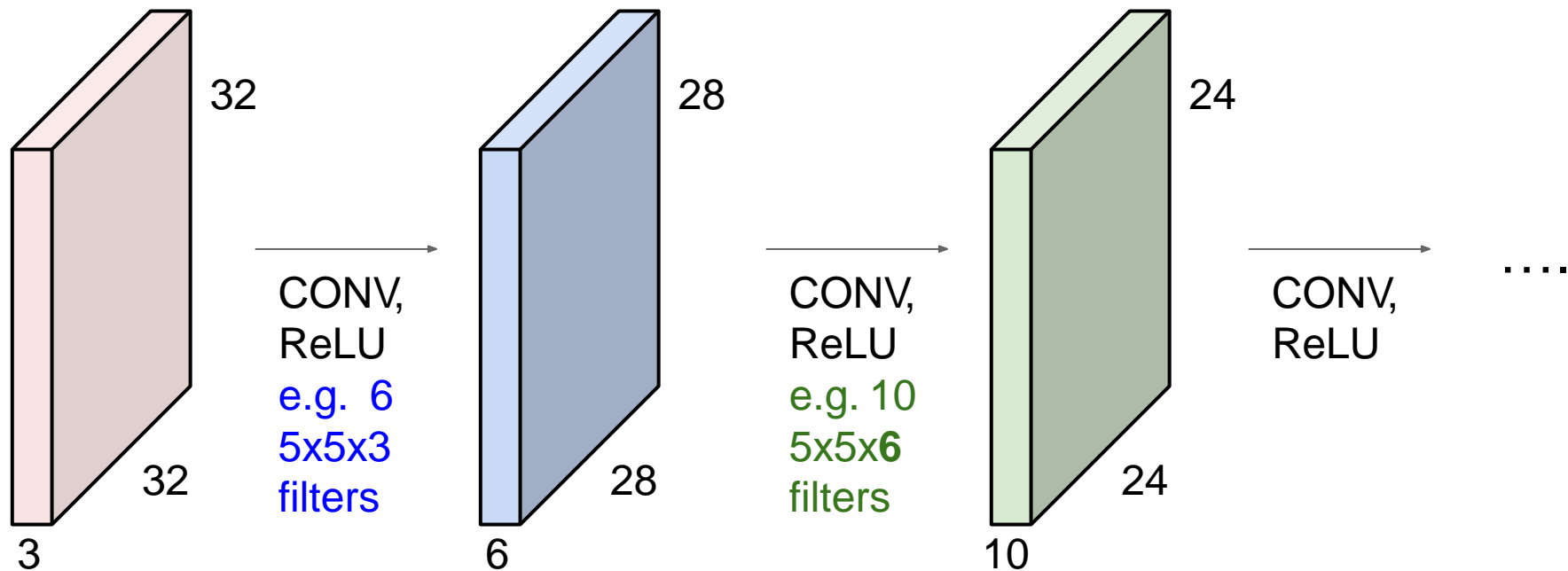e.g. F = 3 => zero pad with 1
     F = 5 => zero pad with 2
     F = 7 => zero pad with 3

**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.



32
32
3

CONV,
ReLU
e.g. 6
5x5x3
filters

28
28
6

CONV,
ReLU
e.g. 10
5x5x**6**
filters

24
24
10

CONV,
ReLU

....
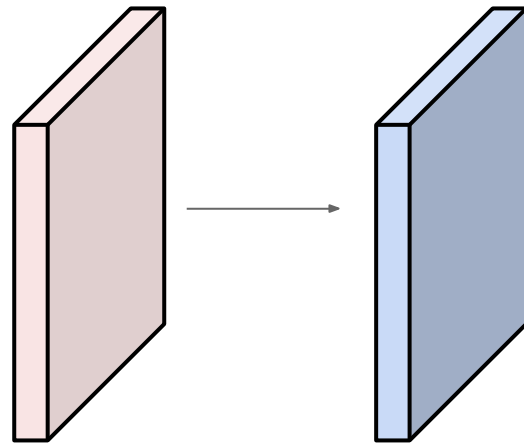
Examples time:

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

Output volume size: ?

Examples time:

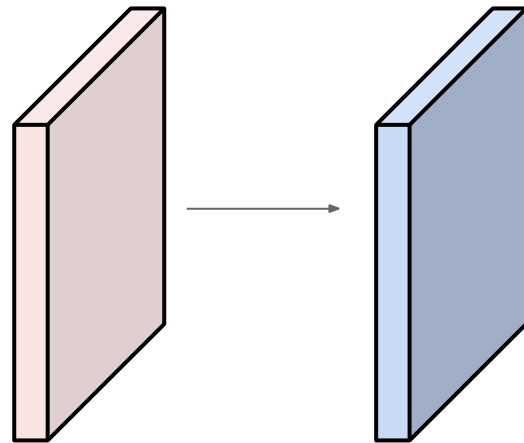Input volume: **32x32x3**
10 5x5 filters with stride 1, pad 2

Output volume size:
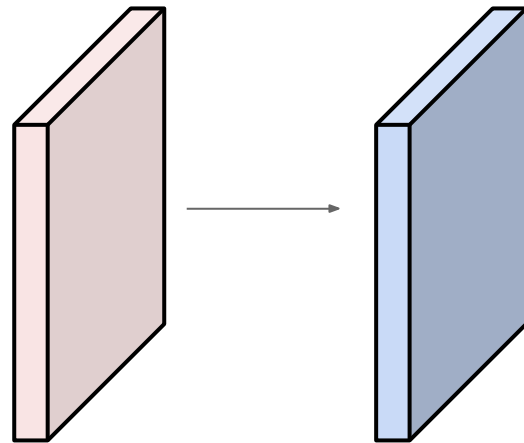(32+2*2-5)/1+1 = 32 spatially, so
**32x32x10**

Examples time:

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



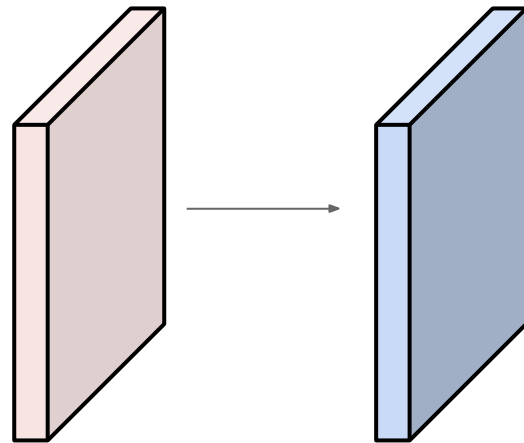Number of parameters in this layer?

Examples time:

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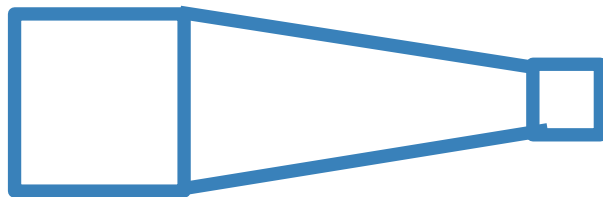10 5x5 filters with stride 1, pad 2

Number of parameters in this layer?
each filter has 5*5*3 + 1 = 76 params        (+1 for bias)
=> 76*10 = **760**

# Receptive Fields

For convolution with kernel size K, each element in the output depends on a K x K **receptive field** in the input



Input                    Output

# Receptive Fields

Each successive convolution adds K – 1 to the receptive field size
With L layers the receptive field size is 1 + L * (K – 1)



Input

Output

Be careful – "receptive field in the input" vs. "receptive field in the previous layer"
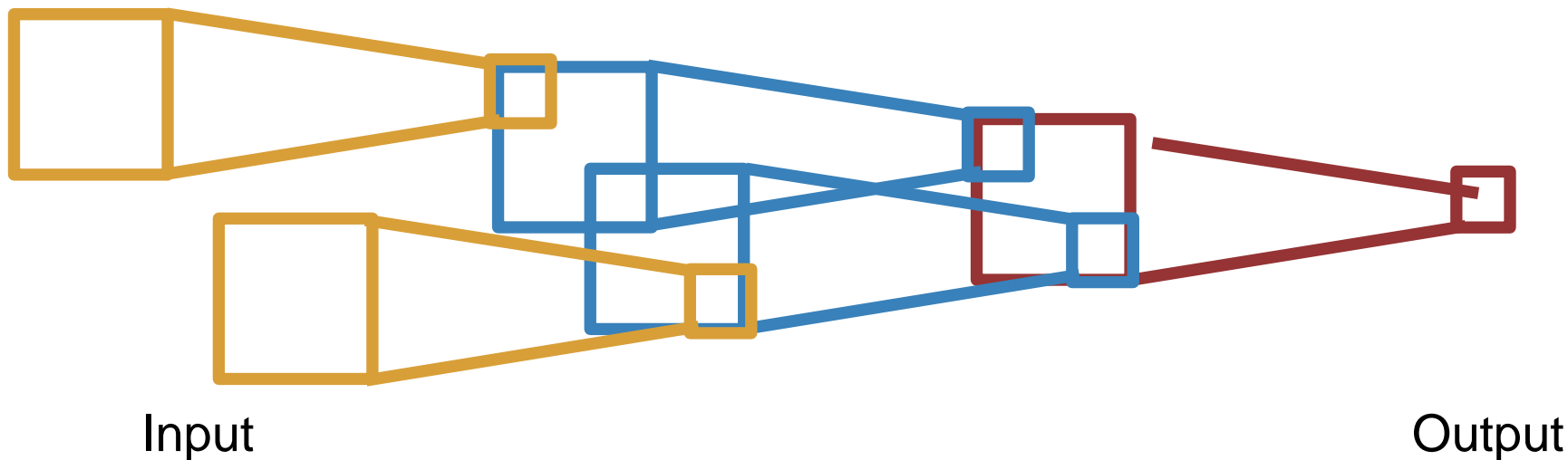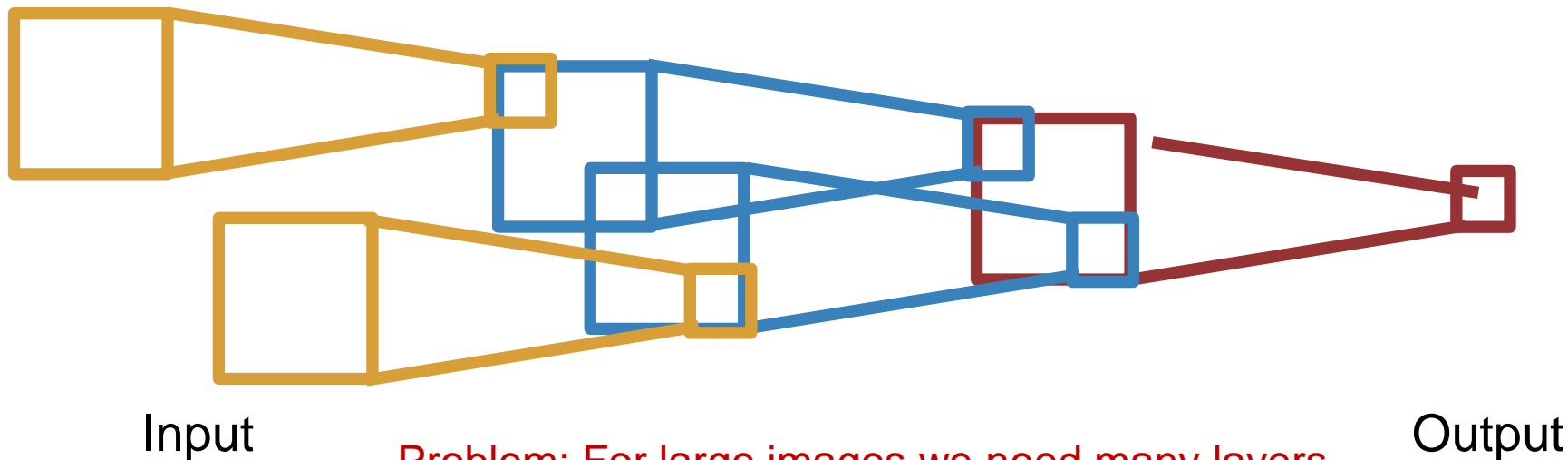
Slide inspiration: Justin Johnson

# Receptive Fields

Each successive convolution adds K − 1 to the receptive field size
With L layers the receptive field size is 1 + L * (K − 1)



Input

Output

Problem: For large images we need many layers
for each output to "see" the whole image image

# Receptive Fields

Each successive convolution adds K – 1 to the receptive field size
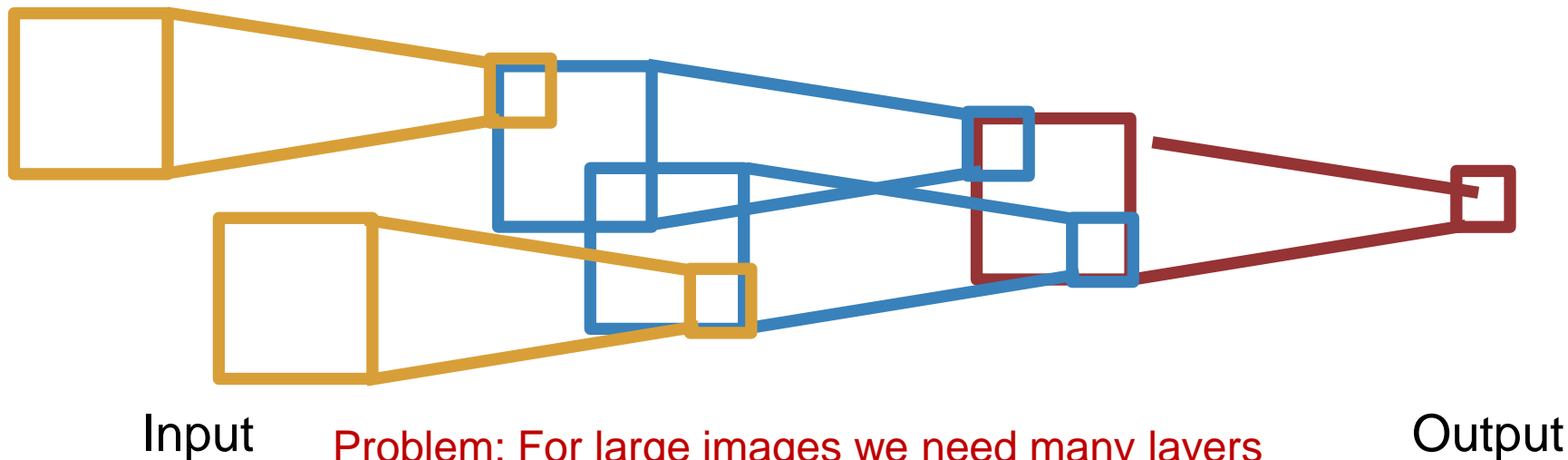With L layers the receptive field size is 1 + L * (K – 1)



Input

Output
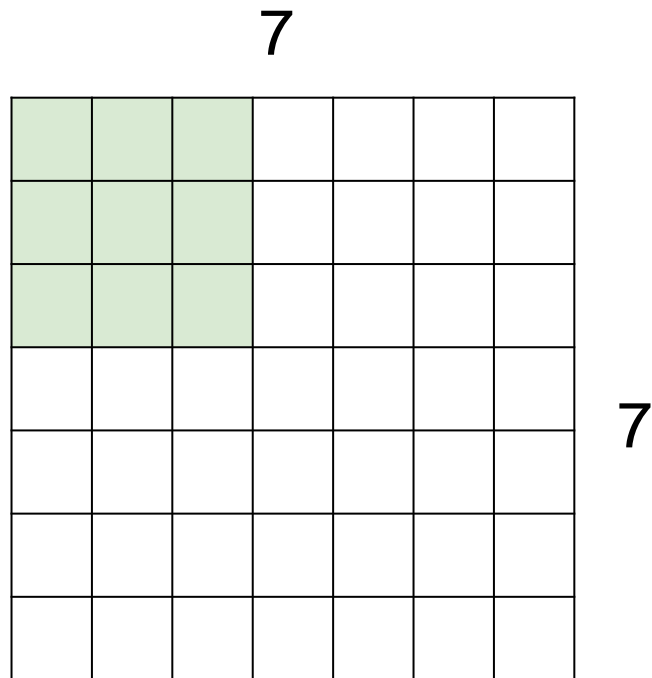
Problem: For large images we need many layers
for each output to "see" the whole image image

Solution: Downsample inside the network

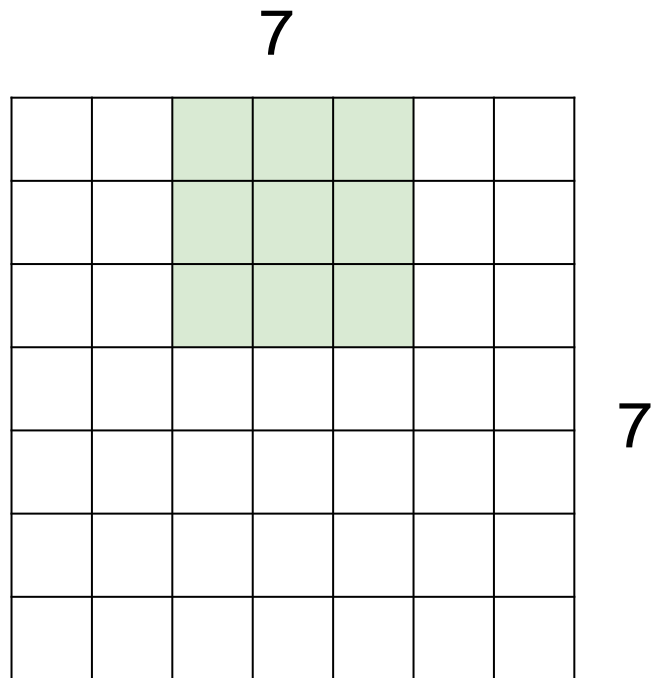# Solution: Strided Convolution

7



7

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

# Solution: Strided Convolution

7



7

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

**=> 3x3 output!**

# Convolution layer: summary

Let's assume input is $W_1$ x $H_1$ x C

Conv layer needs 4 hyperparameters:
- Number of filters **K**
- The filter size **F**
- The stride **S**
- The zero padding **P**

This will produce an output of $W_2$ x $H_2$ x K
where:
- $W_2 = (W_1 - F + 2P)/S + 1$
- $H_2 = (H_1 - F + 2P)/S + 1$

Number of parameters: $F^2CK$ and K biases

# Convolution layer: summary

Let's assume input is $W_1$ x $H_1$ x C
Conv layer needs 4 hyperparameters:

- Number of filters **K**
- The filter size **F**
- The stride **S**
- The zero padding **P**

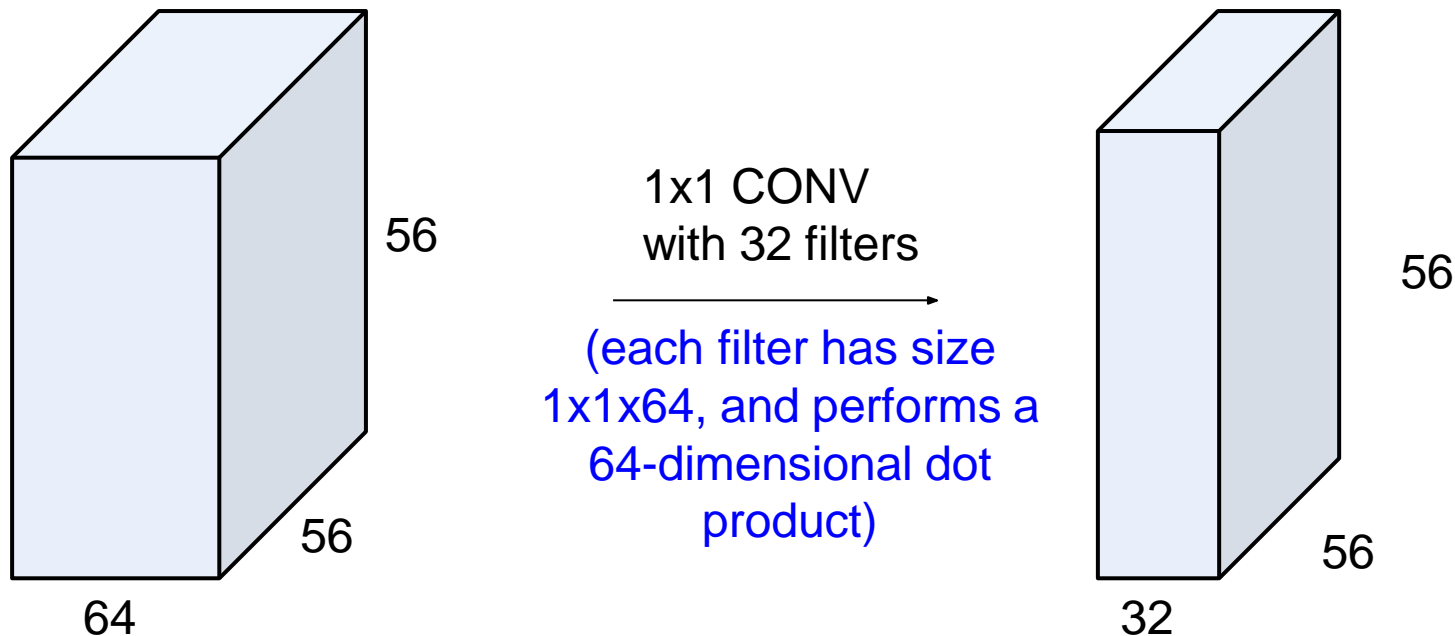This will produce an output of $W_2$ x $H_2$ x K
where:

- $W_2 = (W_1 - F + 2P)/S + 1$
- $H_2 = (H_1 - F + 2P)/S + 1$

Number of parameters: $F^2CK$ and K biases

K = (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 = ? (whatever fits)
- F = 1, S = 1, P = 0

# (btw, 1x1 convolution layers make perfect sense)



56

1x1 CONV
with 32 filters

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

56

56

64

56

32

# (btw, 1x1 convolution layers make perfect sense)



1x1 CONV
with 32 filters

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

56

56

64

56

56

32

# Example: CONV layer in PyTorch

Conv2d

Applies a 2D convolution over an input signal composed of several input planes.

In the simplest case, the output value of the layer with input size $(N, C_{in}, H, W)$ and output $(N, C_{out}, H_{out}, W_{out})$ can be precisely described as:

$$out(N_i, C_{out_j}) = bias(C_{out_j}) + \sum_{k=0}^{C_{in}-1} weight(C_{out_j}, k) \star input(N_i, k)$$

where $\star$ is the valid 2D cross-correlation operator, $N$ is a batch size, $C$ denotes a number of channels, $H$ is a height of input planes in pixels, and $W$ is width in pixels.

- stride controls the stride for the cross-correlation, a single number or a tuple.
- padding controls the amount of implicit zero-paddings on both sides for padding number of points for each dimension.
- dilation controls the spacing between the kernel points; also known as the à trous algorithm. It is harder to describe, but this link has a nice visualization of what dilation does.
- groups controls the connections between inputs and outputs. in_channels and out_channels must both be divisible by groups. For example,
    - At groups=1, all inputs are convolved to all outputs.
    - At groups=2, the operation becomes equivalent to having two conv layers side by side, each seeing half the input channels, and producing half the output channels, and both subsequently concatenated.
    - At groups= in_channels, each input channel is convolved with its own set of filters, of size: $\left\lfloor \frac{C_{out}}{C_{in}} \right\rfloor$.

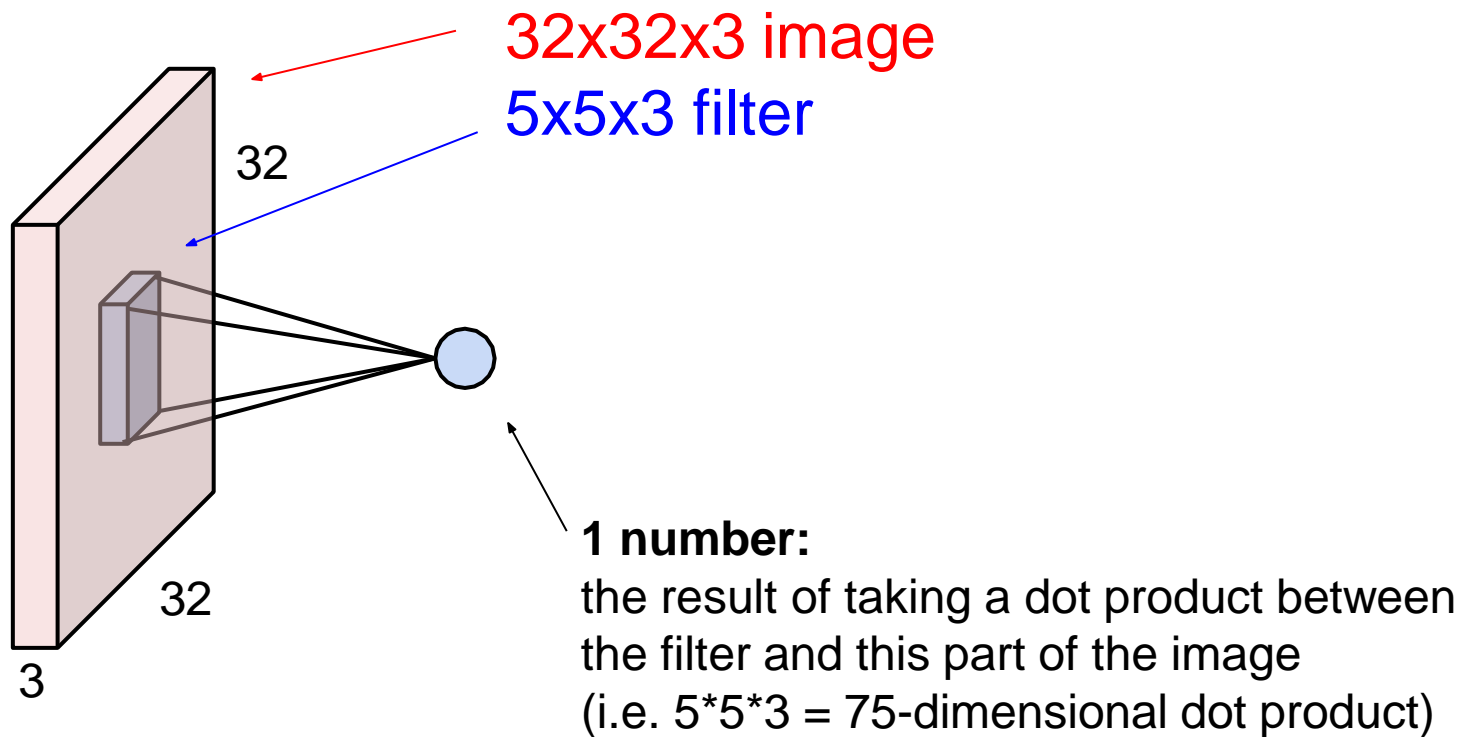The parameters kernel_size, stride, padding, dilation can either be:

- a single int – in which case the same value is used for the height and width dimension
- a tuple of two ints – in which case, the first int is used for the height dimension, and the second int for the width dimension
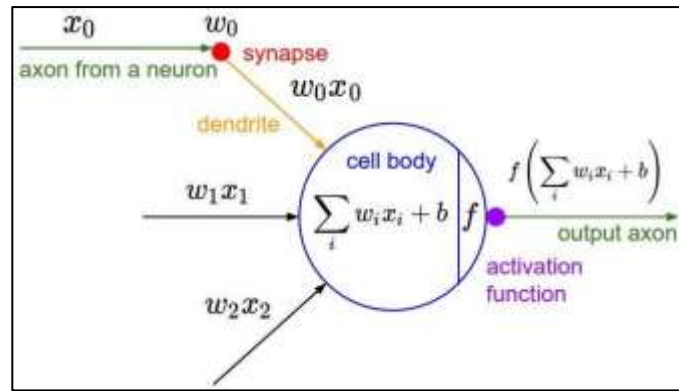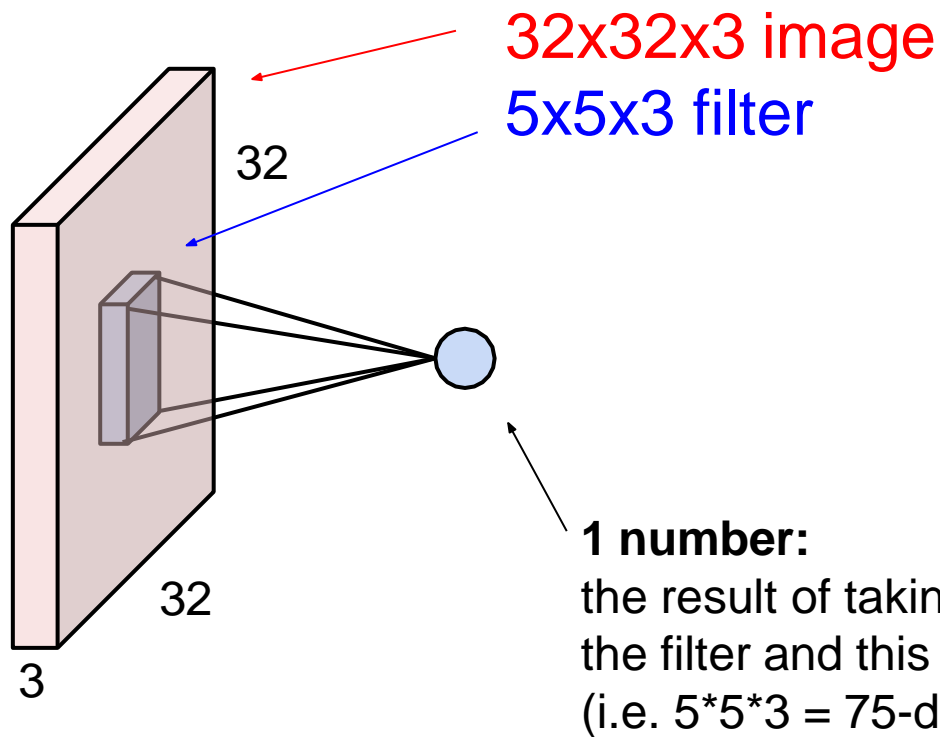
Conv layer needs 4 hyperparameters:
- Number of filters **K**
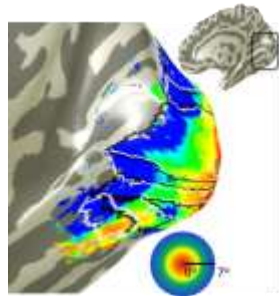- The filter size **F**
- The stride **S**
- The zero padding **P**
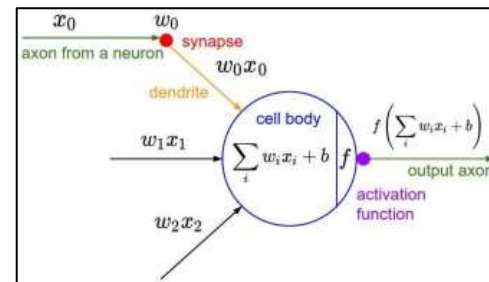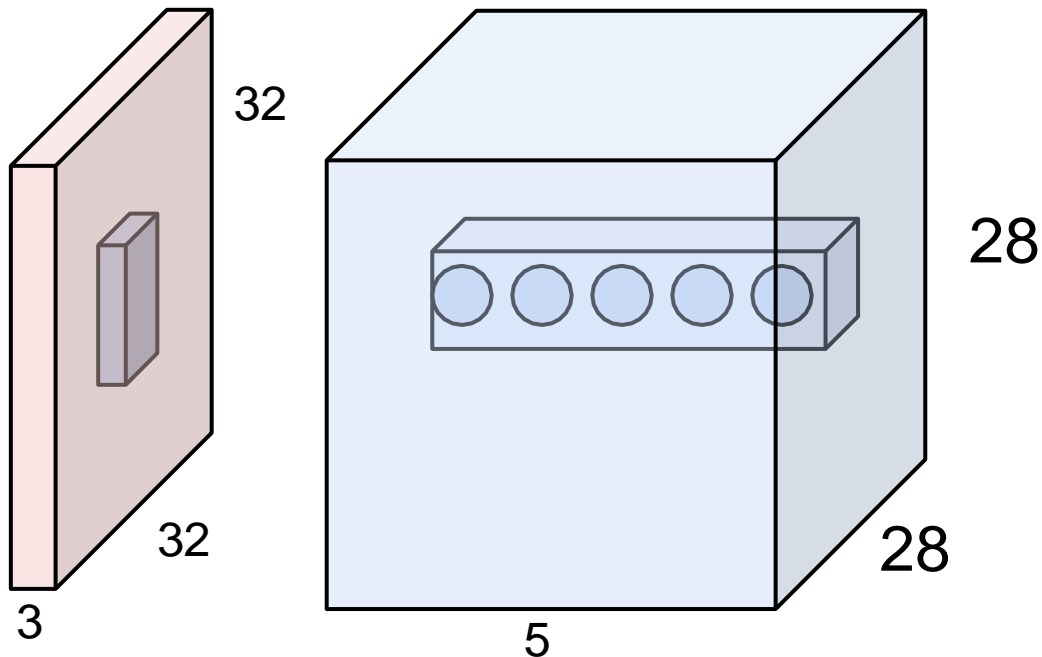
# The brain/neuron view of CONV Layer



32x32x3 image

5x5x3 filter

**1 number:**
the result of taking a dot product between
the filter and this part of the image
(i.e. 5*5*3 = 75-dimensional dot product)

# The brain/neuron view of CONV Layer

32x32x3 image

5x5x3 filter

32



It's just a neuron with local connectivity...

**1 number:**
the result of taking a dot product between the filter and this part of the image
(i.e. 5*5*3 = 75-dimensional dot product)
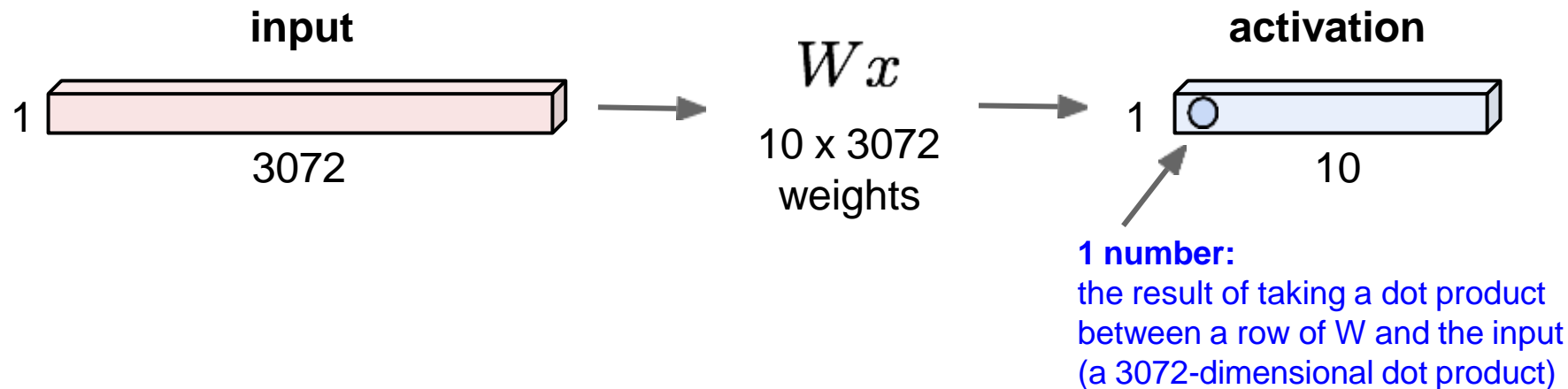
32

3

# The brain/neuron view of CONV Layer



E.g. with 5 filters,
CONV layer consists of
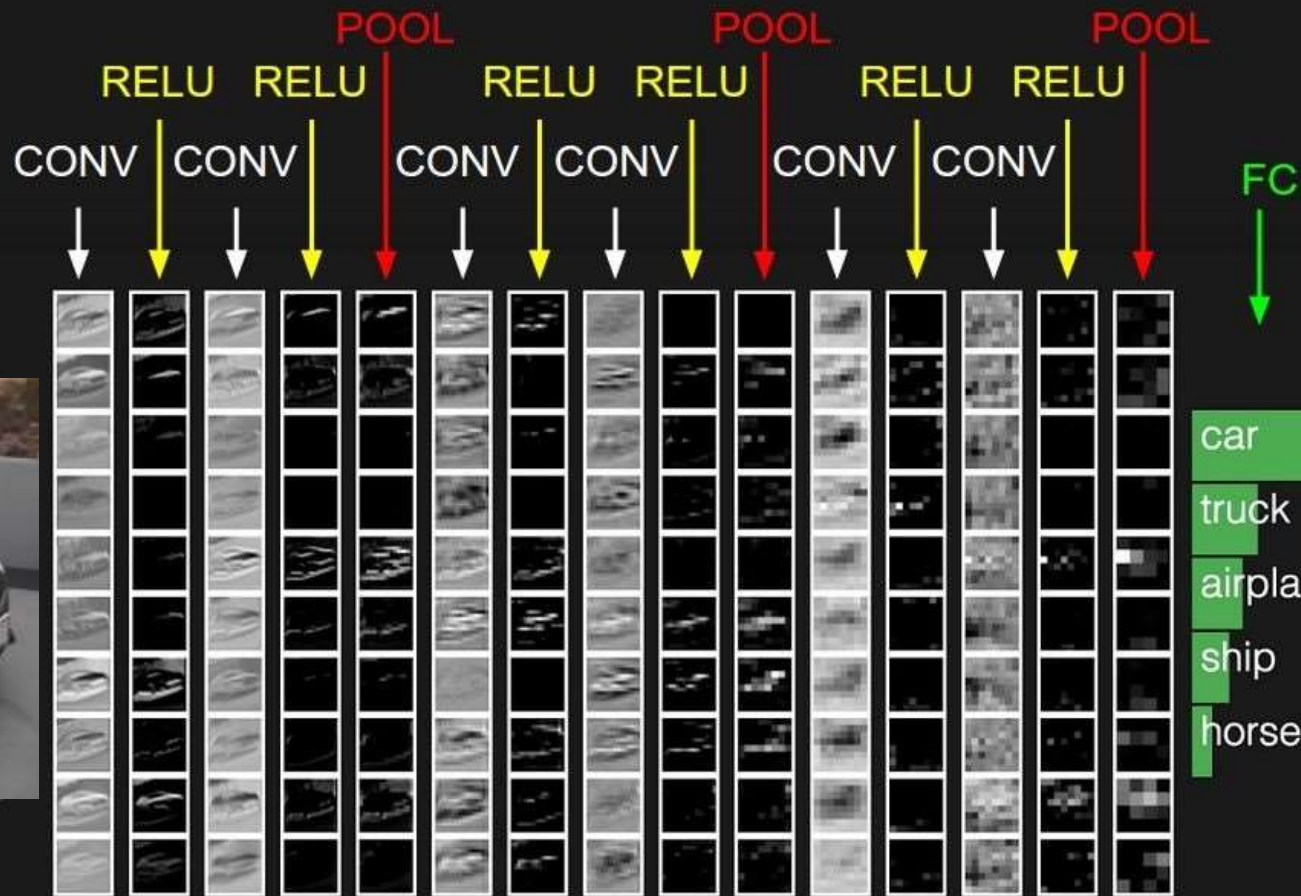neurons arranged in a 3D grid
(28x28x5)

There will be 5 different
neurons all looking at the same
region in the input volume

# Reminder: Fully Connected Layer
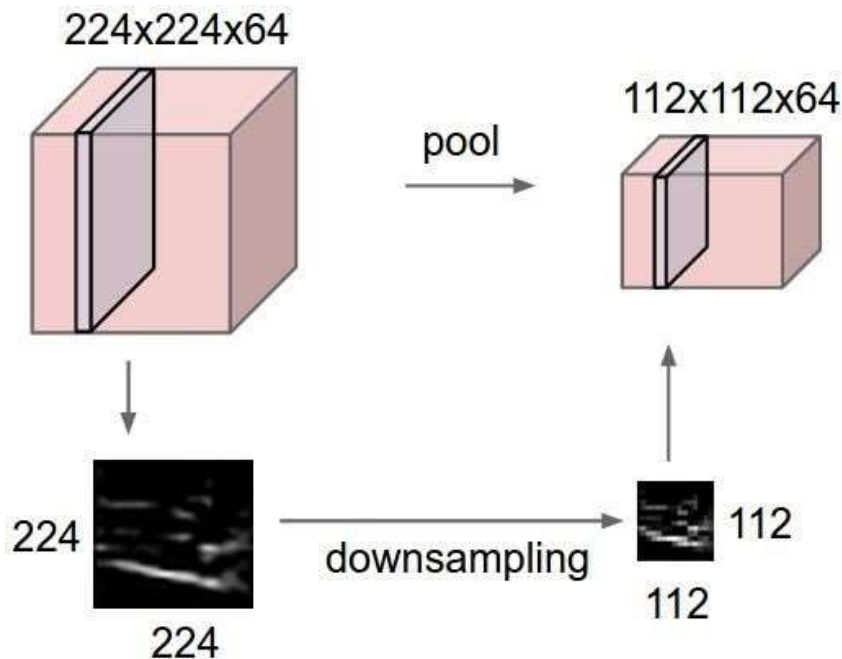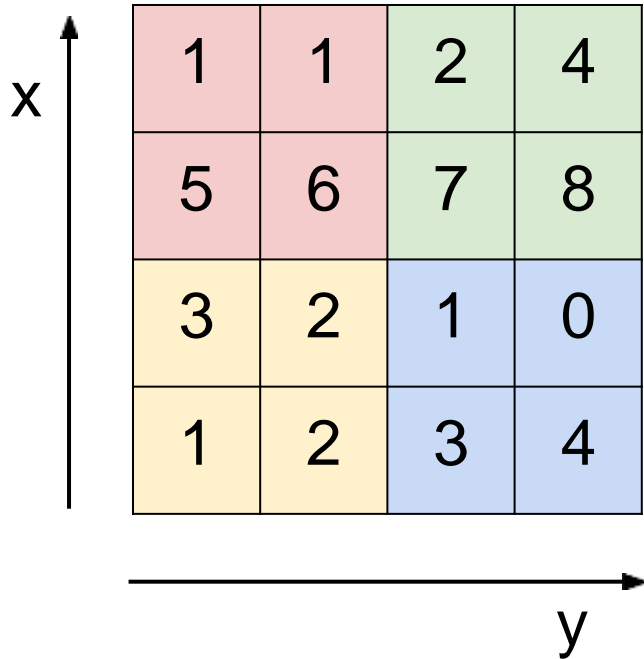
32x32x3 image -> stretch to 3072 x 1

**input**

1

3072

$Wx$

10 x 3072 weights

**activation**

1

10

**1 number:**
the result of taking a dot product between a row of W and the input (a 3072-dimensional dot product)

POOL     POOL     POOL

RELU RELU   RELU RELU   RELU RELU

CONV CONV   CONV CONV   CONV CONV    FC

car
truck
airplane
ship
horse

# Pooling layer

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

# MAX POOLING
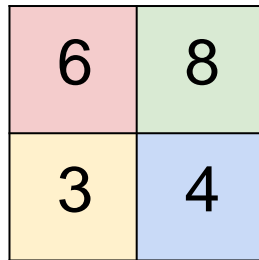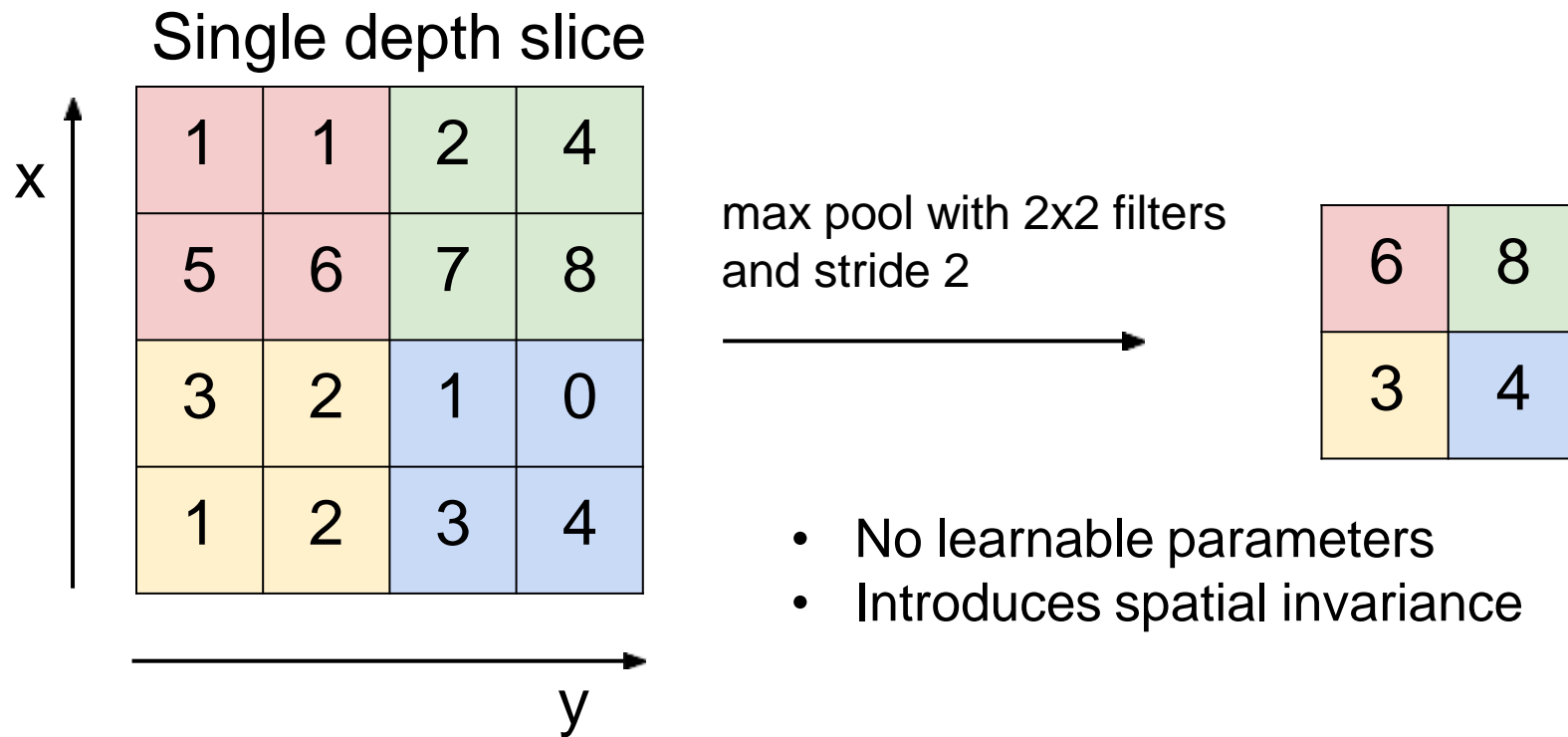
Single depth slice



max pool with 2x2 filters
and stride 2

# MAX POOLING

Single depth slice

x

| 1 | 1 | 2 | 4 |
|---|---|---|---|
| 5 | 6 | 7 | 8 |
| 3 | 2 | 1 | 0 |
| 1 | 2 | 3 | 4 |

y

max pool with 2x2 filters
and stride 2

⟶

| 6 | 8 |
|---|---|
| 3 | 4 |

- No learnable parameters
- Introduces spatial invariance

# Pooling layer: summary

Let's assume input is $W_1$ x $H_1$ x C
Conv layer needs 2 hyperparameters:
- The spatial extent **F**
- The stride **S**

This will produce an output of $W_2$ x $H_2$ x C where:
- $W_2 = (W_1 - F)/S + 1$
- $H_2 = (H_1 - F)/S + 1$

Number of parameters: 0

# Fully Connected Layer (FC layer)

- Contains neurons that connect to the entire input volume, as in ordinary Neural Networks
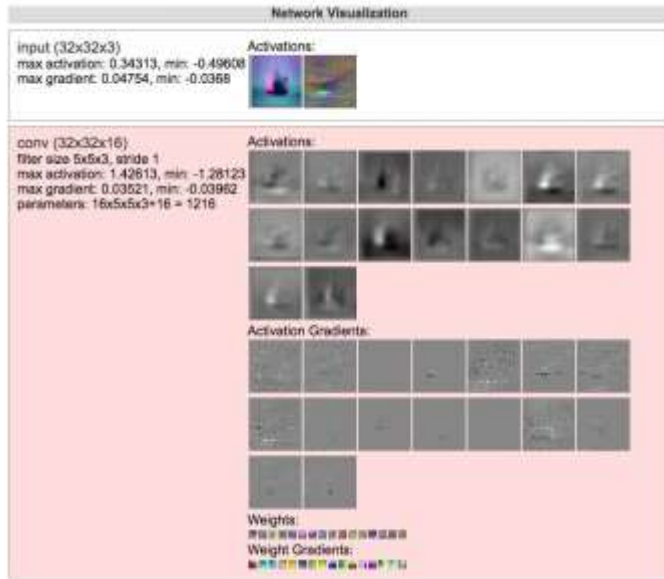
# [ConvNetJS demo: training on CIFAR-10]



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

# Summary

- ConvNets stack CONV,POOL,FC layers
- Trend towards smaller filters and deeper architectures
- Trend towards getting rid of POOL/FC layers (just CONV)
- Historically architectures looked like
  **[(CONV-RELU)*N-POOL?]*M-(FC-RELU)*K,SOFTMAX**
  where N is usually up to ~5, M is large, 0 <= K <= 2.
- But recent advances such as ResNet/GoogLeNet have challenged this paradigm

# Next time: CNN Architectures