# Lecture 18-2

Neural networks and CNN
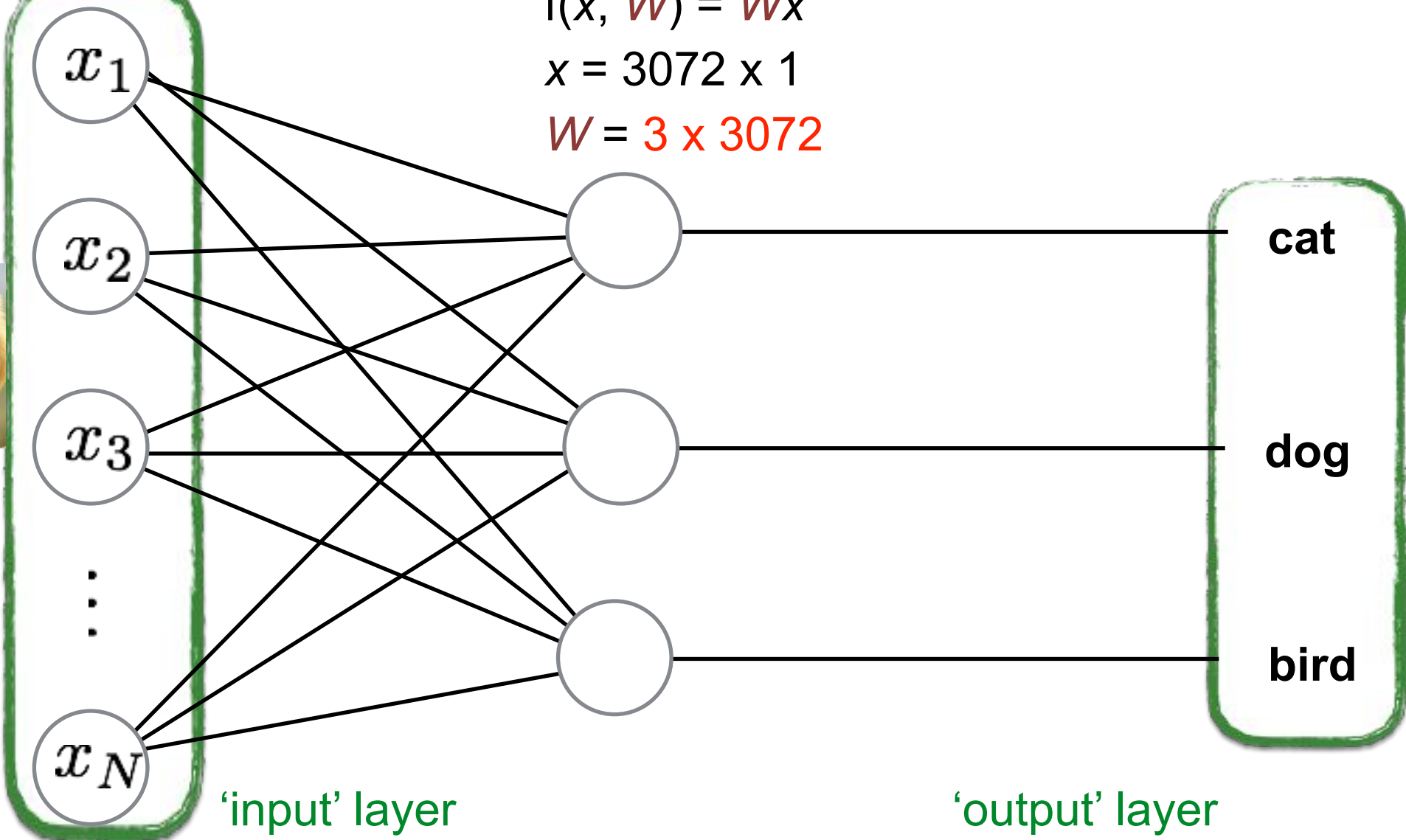
# Administrative

- A5 & A6 (bonus) are out

- Final Exam on 6/9 at 2:30 pm

- Makeup exam on 6/6

- Exam practice is out

# Recall: Linear Model



f($x$, $W$) = $Wx$

$x$ = 3072 x 1
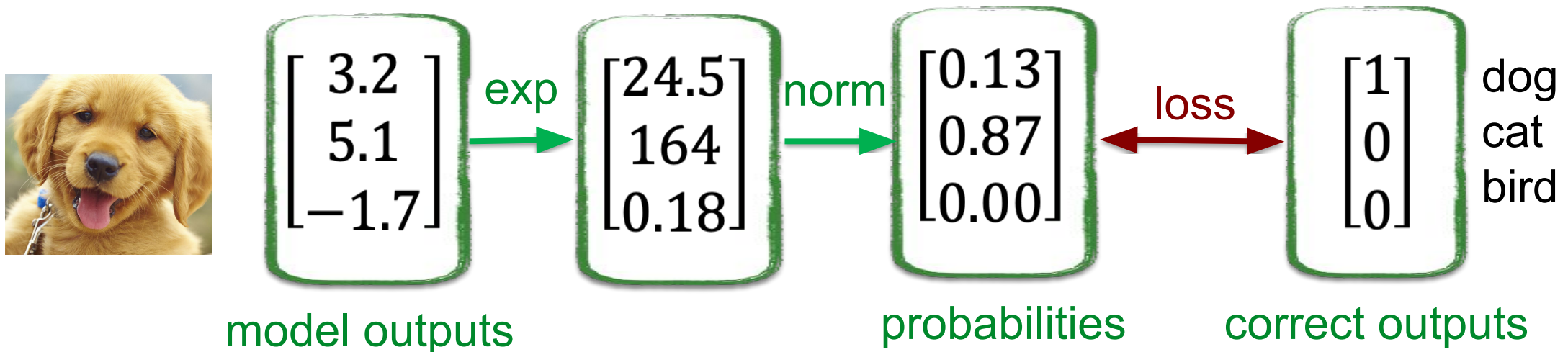
$W$ = 3 x 3072

'input' layer
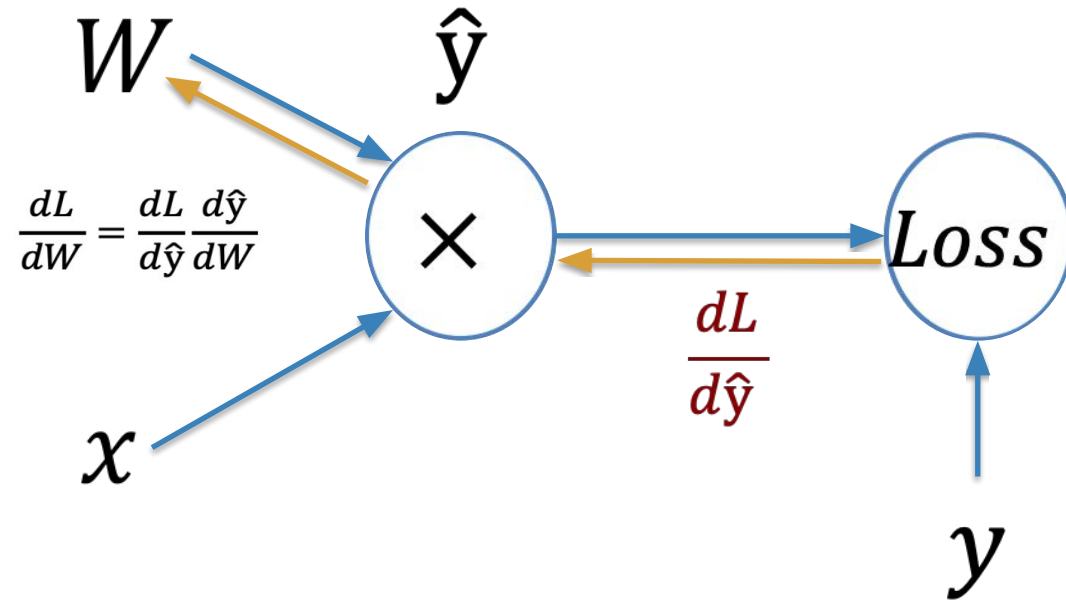
cat

dog

bird

'output' layer

# Recall: Softmax Classifier

$$L_i = -\log Prob[f(x_i, W) == y_i]$$

We need a mechanism to convert or normalize the output into probability range [0, 1]

Recall:  SOFTMAX: $Prob[f(x_i, W) == k] = \dfrac{e^{\hat{y}_k}}{\sum_j e^{\hat{y}_j}}$



$\begin{bmatrix} 3.2 \\ 5.1 \\ -1.7 \end{bmatrix}$ exp $\begin{bmatrix} 24.5 \\ 164 \\ 0.18 \end{bmatrix}$ norm $\begin{bmatrix} 0.13 \\ 0.87 \\ 0.00 \end{bmatrix}$ loss $\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$ dog cat bird

model outputs                    probabilities          correct outputs

# Recall: Gradient Descent through Backprop

$$\hat{y} = Wx$$
$$L = Loss(\hat{y}, y)$$

$$\frac{dL}{dW} = \frac{dL}{d\hat{y}} \frac{d\hat{y}}{dW}$$

$$\frac{dL}{dW} = \frac{dL}{d\hat{y}} \frac{d\hat{y}}{dW}$$

$$\frac{dL}{d\hat{y}}$$

Key Insight:
- visualize the computation as a graph flow
- Compute the forward pass to calculate the loss.
- Compute all gradients for each pair of nodes backwards
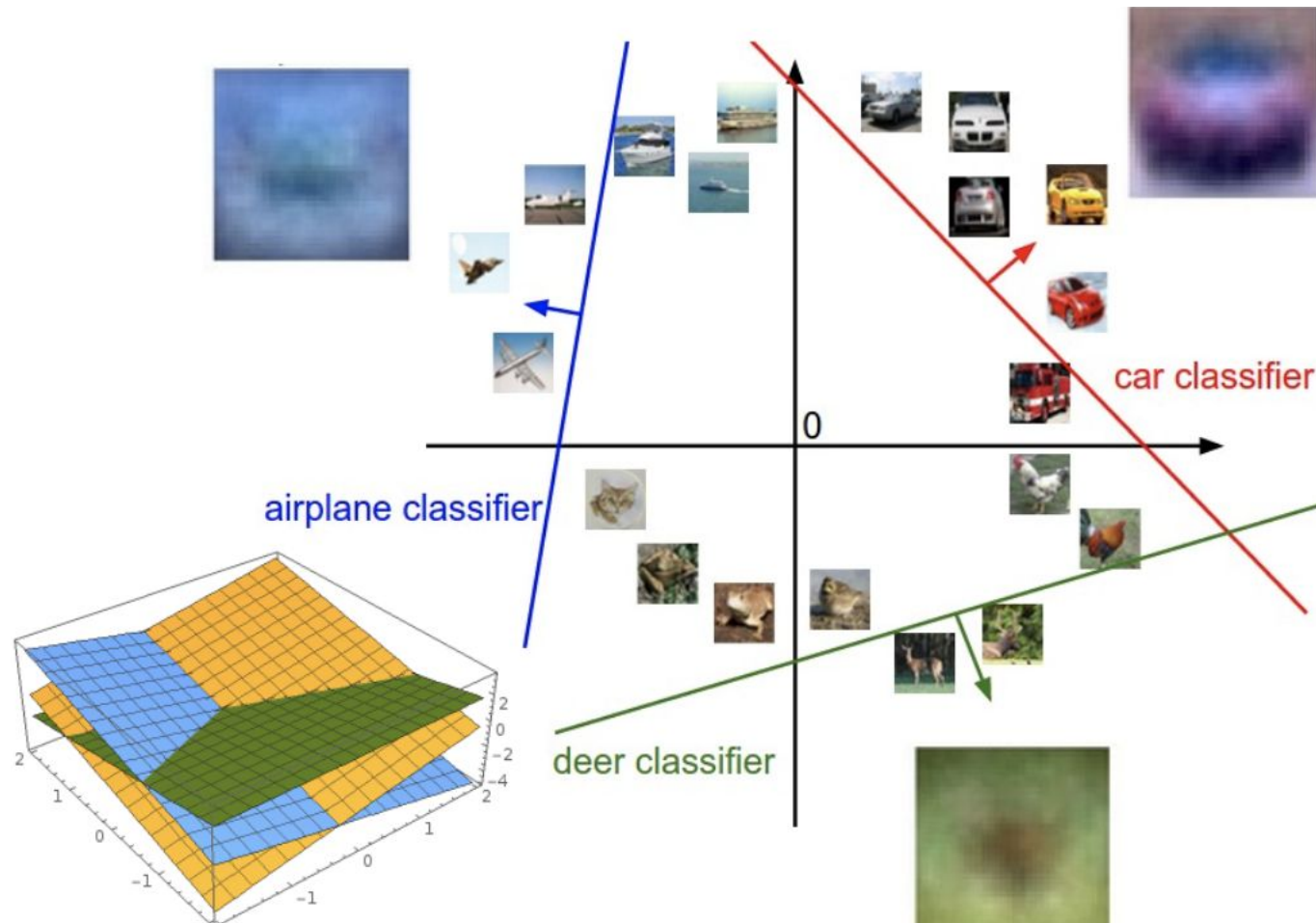
# Recall: we can featurize images into a vector

Image
Vector

Raw pixels
Raw pixels + (x,y)
PCA
LDA
BoW
BoW + spatial pyramids

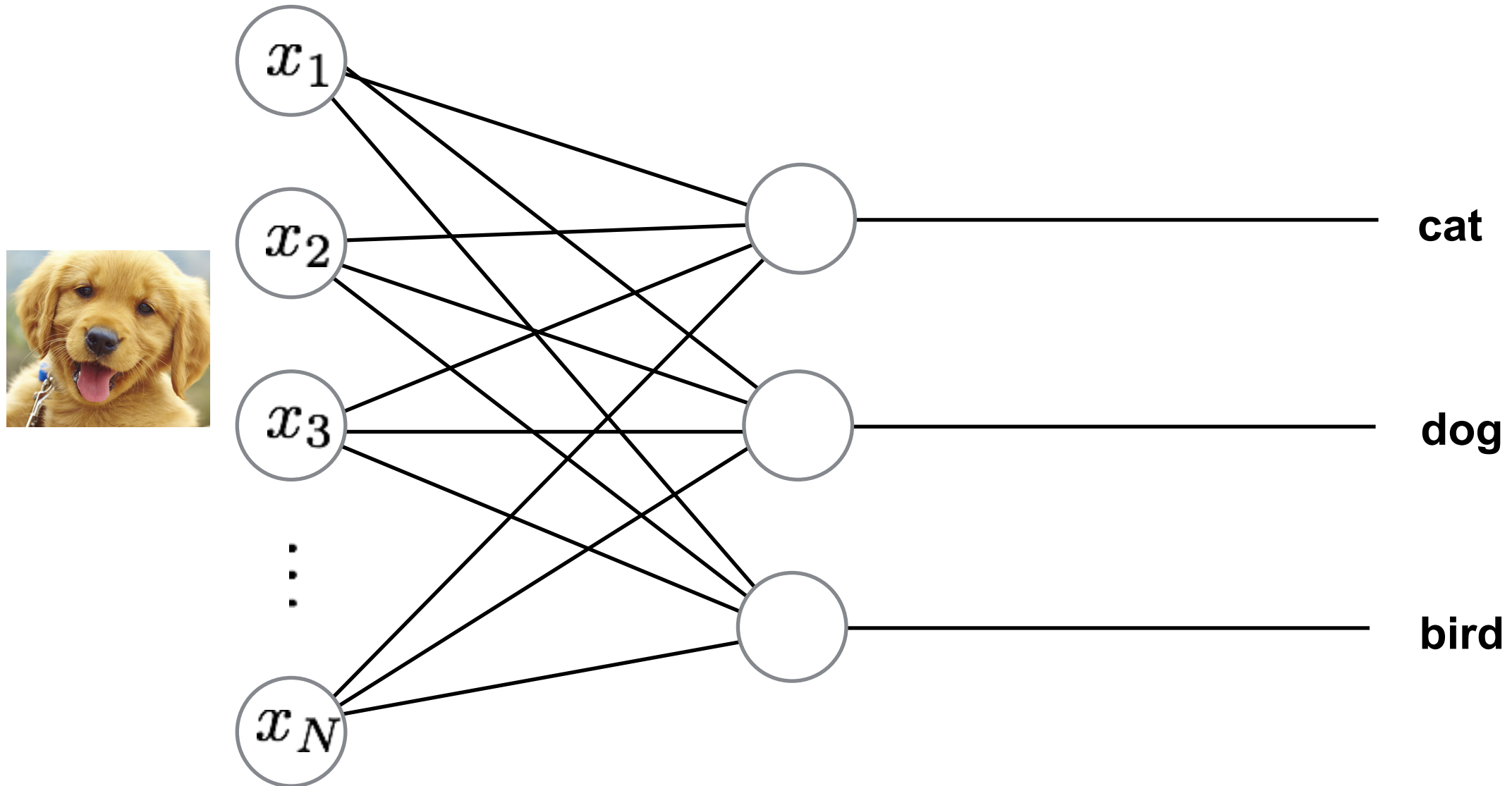# Interpreting the linear weights geometrically

Plot created using Wolfram Cloud
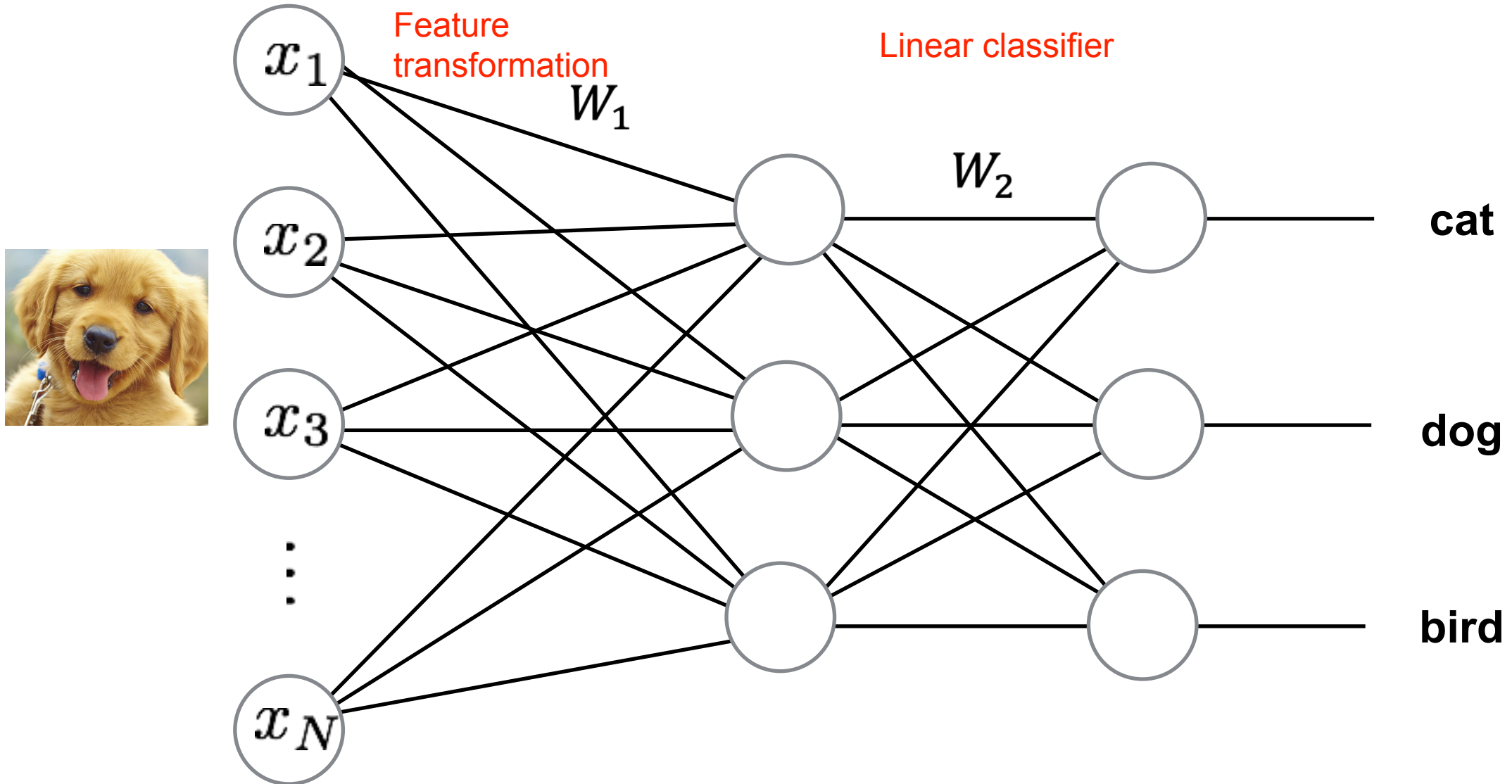
# Features sometimes might not be linearly separable

$$f(x, y) = (r(x, y), \theta(x, y))$$

but some mappings are!

# Remember our linear classifier

# Let's change the features by adding another layer



Feature transformation

Linear classifier
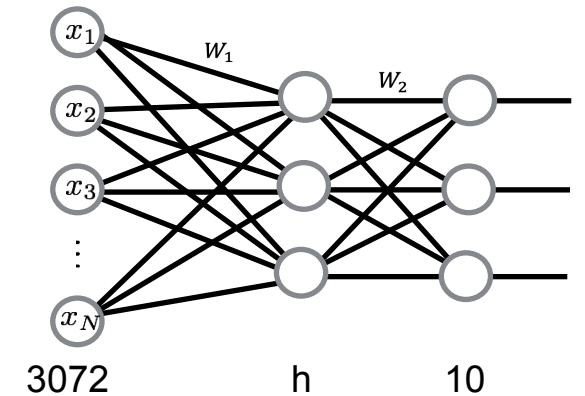
$W_1$

$W_2$

$x_1$
$x_2$
$x_3$
$x_N$

cat

dog

bird

# 2-layer network: mathematical formula

- Linear classifier: $y = Wx$

- 2-layer network: $y = W_2 \cdot \text{binarize}(W_1 x)$, where $\text{binarize}(W_1 x) = \begin{cases} 1 & \text{if } W_1 x > 0 \\ 0 & \text{otherwise} \end{cases}$

- 3-layer network: $y = W_3 \cdot \text{binarize}(W_2 \cdot \text{binarize}(W_1 x))$

The number of layers is a new hyperparameter!

# 2-layer network: mathematical formula

- Linear classifier: $y = Wx$

- 2-layer network: $y = W_2 \cdot \text{binarize}(W_1 x)$, where $\text{binarize}(W_1 x) = \begin{cases} 1 & \text{if } W_1 x > 0 \\ 0 & \text{otherwise} \end{cases}$

We know the size of $x$ = 1 x 3072 and y = 10 x 1, so what are $\boldsymbol{W_1}$ and $\boldsymbol{W_2}$

$$W_1 = h \times 3072 \qquad W_2 = 10 \times h$$

h is a new hyperparameter!



3072        h       10

# 2-layer network: mathematical formula

- Linear classifier: $y = Wx$

- 2-layer network: $y = W_2 \cdot \text{binarize}(W_1 x)$, where $\text{binarize}(W_1 x) = \begin{cases} 1 & \text{if } W_1 x > 0 \\ 0 & \text{otherwise} \end{cases}$

Why is the binarize necessary? Let's see what happen when we remove it:

$$y = W_2 W_1 x = Wx$$

Where: $W = W_2 W_1$

Activation is necessary to go from linear to non-linear models

# 2-layer network: mathematical formula

- Linear classifier: $y = Wx$

- 2-layer network: $y = W_2 \, \text{sigmoid}(W_1 x)$

Why is the binarize necessary?

- Neural science inspiration
- Non-differentiable

Let's approximate it with sigmoid

$$f(x) = \frac{1}{1 + e^{-x}}$$

# 2-layer network: mathematical formula

- Linear classifier:  $y = Wx$

- 2-layer network:  $y = W_2 \, \mathrm{ReLU}(W_1 x)$

Why is the <span style="color:red">sigmoid</span> necessary?

- Vanishing gradient

Let's <span style="color:red">replace</span> it with <span style="color:red">ReLU</span>

**ReLU**
$\max(0, x)$

# ReLU v.s. Sigmoid

ReLU is not from nowhere

- Connection between ReLU & Sigmoid

# 2-layer Neural Network



$W_1$

$W_2$

cat

dog

bird

ReLU, Sigmoid … etc

# Take-Home Exercise

Backprop on 2-layer neural network

$$\frac{dL}{dW} = \frac{dL}{d\hat{y}}\frac{d\hat{y}}{dW}$$

$W$ $\hat{y}$

$\times$ $Loss$

$\frac{dL}{d\hat{y}}$

$x$

$y$

1-layer case

2-layer case ???

# Fully Connected Layer

32x32x3 image -> stretch to 3072 x 1

**input**

1

3072

$$Wx$$

h x 3072
weights



$x_1$

$x_2$ $W_1$

$x_3$ $W_2$

$x_N$

**activation
(ReLU)**

1

h

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

# FC Layer Issues

## What if we are processing higher resolution image?

**input**

1    256x256x3

$Wx$

h x (256x256x3) weights

**activation (ReLU)**

1    h

Q: If h=128x128x3, how many parameters of a single layer?

$x_1$   $W_1$   $W_2$

$x_2$

$x_3$

$\vdots$

$x_N$

256x256x3    h

# FC Layer Issues

What if we are processing higher resolution image?

**input**

1

256x256x3

$Wx$

h x (256x256x3)
weights

**activation (ReLU)**

1

h

Q: If h=128x128x3, how many parameters of a single layer?

A: (256x256x3)x(128x128x3) ≈ 9.6 B

Too large to handle

$x_1$
$x_2$
$x_3$
$\vdots$
$x_N$

$W_1$

$W_2$

256x256x3          h

# Convolution Layer – A Special FC Layer

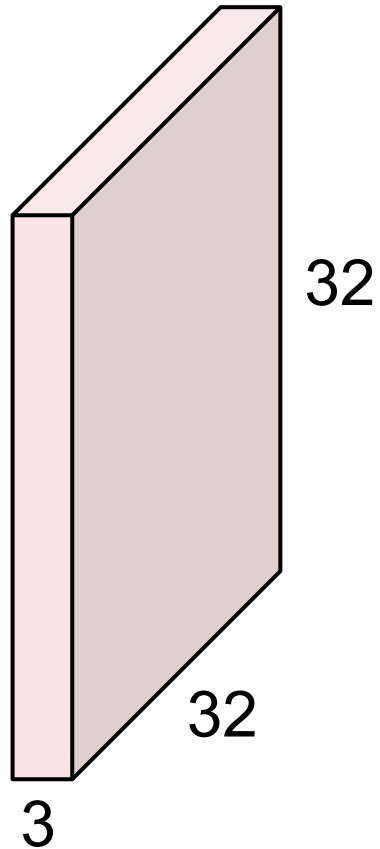32x32x3 image -> preserve spatial structure

32 height

32 width

3 depth

~~FC Layer: every output looks at the whole image~~

Main idea: every output only looks at small patches with small & shared number of parameters
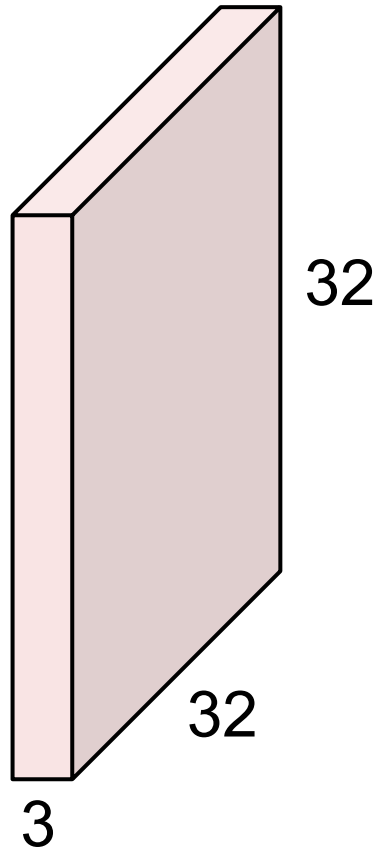
# Convolution Layer

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

32x32x**3** image

**Filters always extend the full depth of the input volume**

5x5x**3** filter

32

32

3

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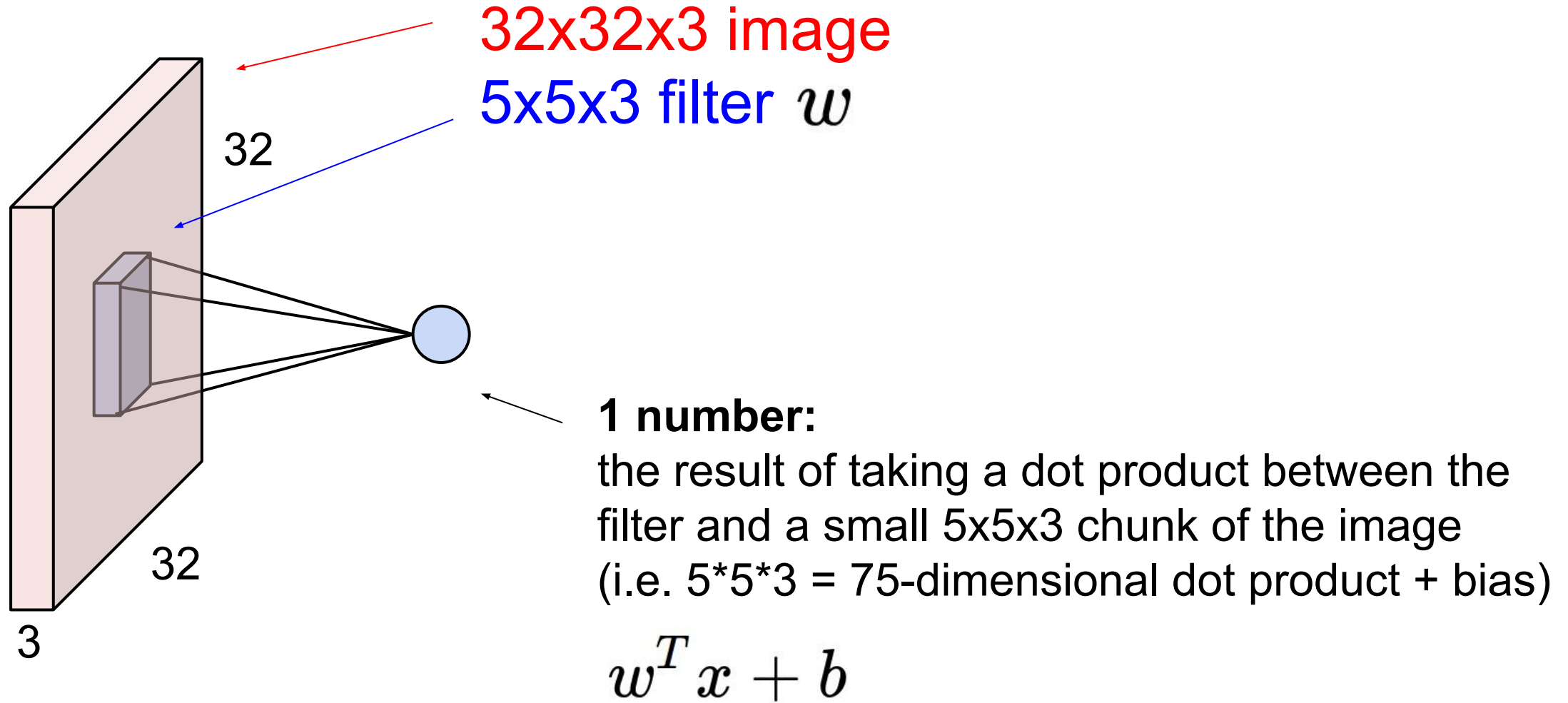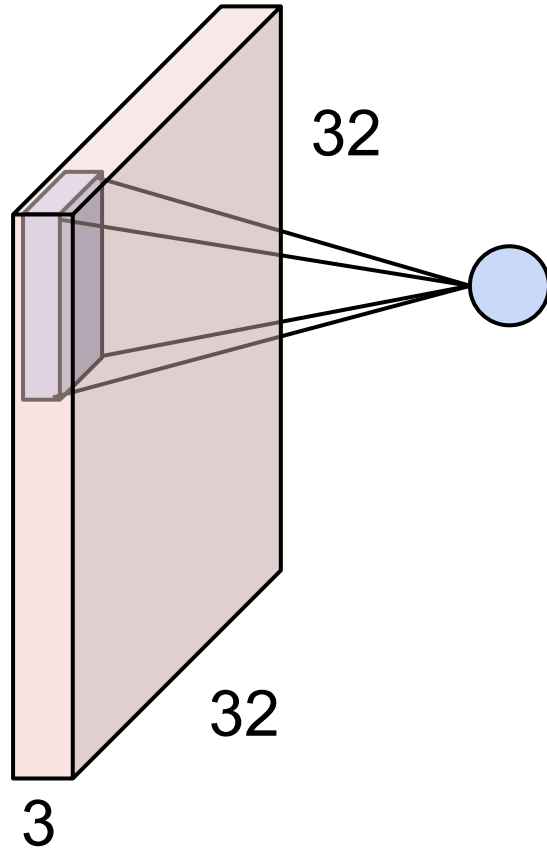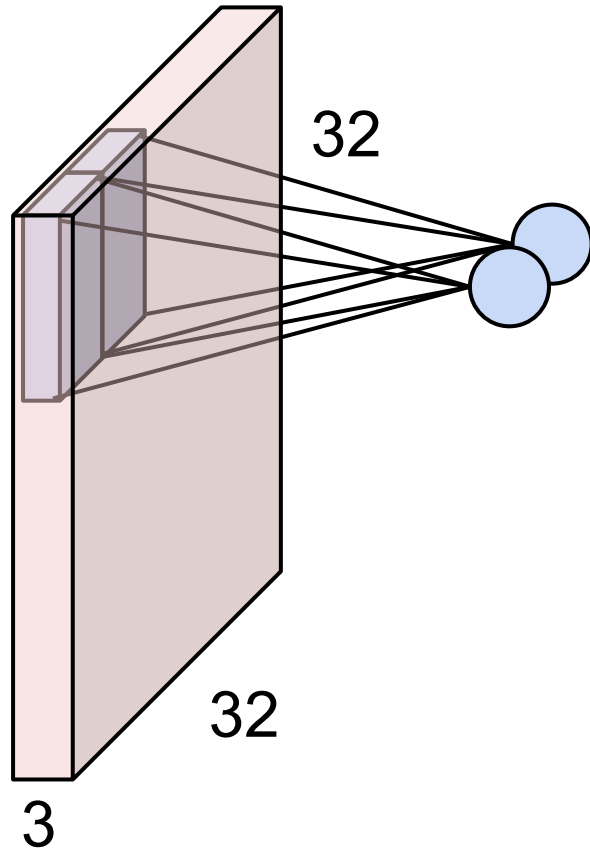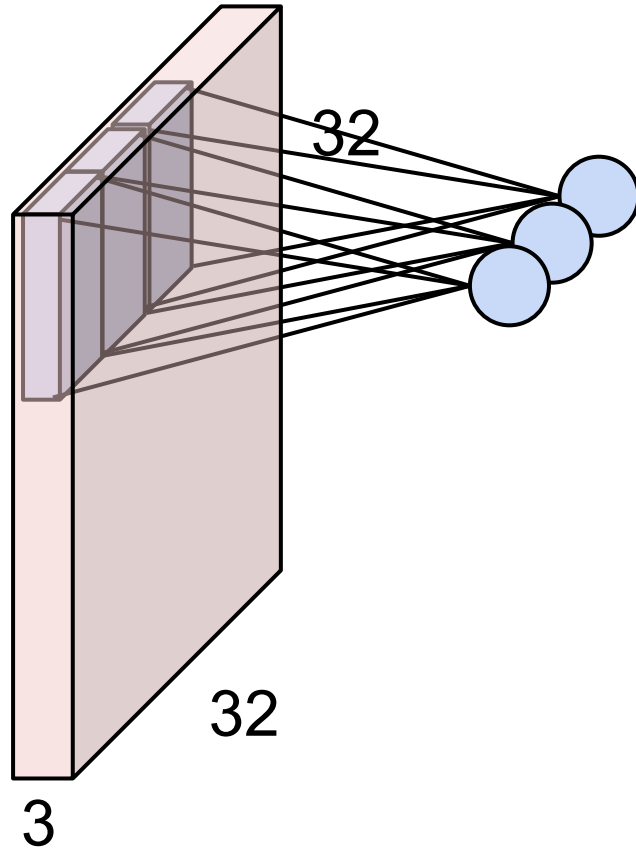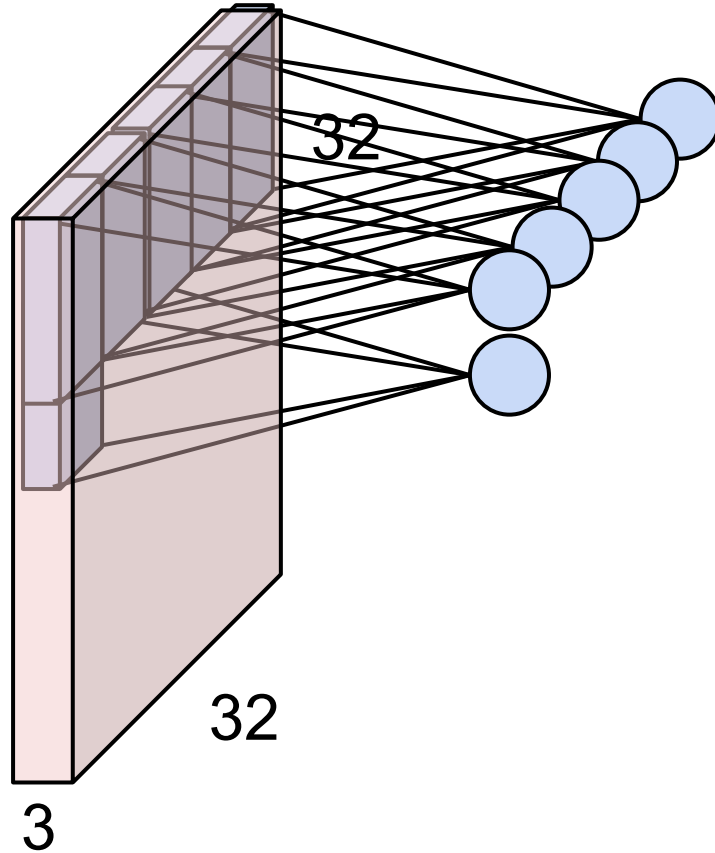
# Convolution Layer



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

Remark: A special case of a (32x32x3) x (28x28x1) linear layer!!

# Convolution Layer

consider a second, green filter

**32x32x3 image**
**5x5x3 filter**

32

32

3

**activation maps**

28

28

1

convolve (slide) over all
spatial locations

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

**activation maps**

32

32

3

Convolution Layer

28

28

6
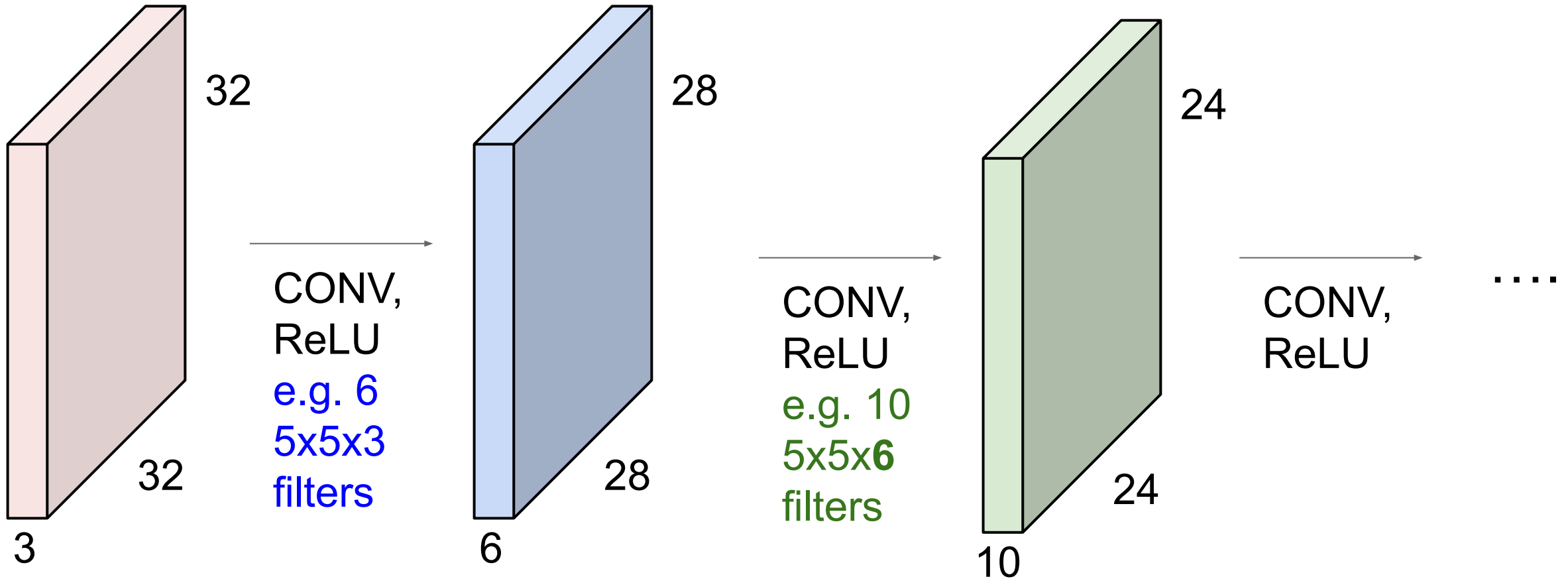
We stack these up to get a "new image" of size 28x28x6!

# ConvNet: Sequence of Convolution Layers, interspersed with activation functions



32
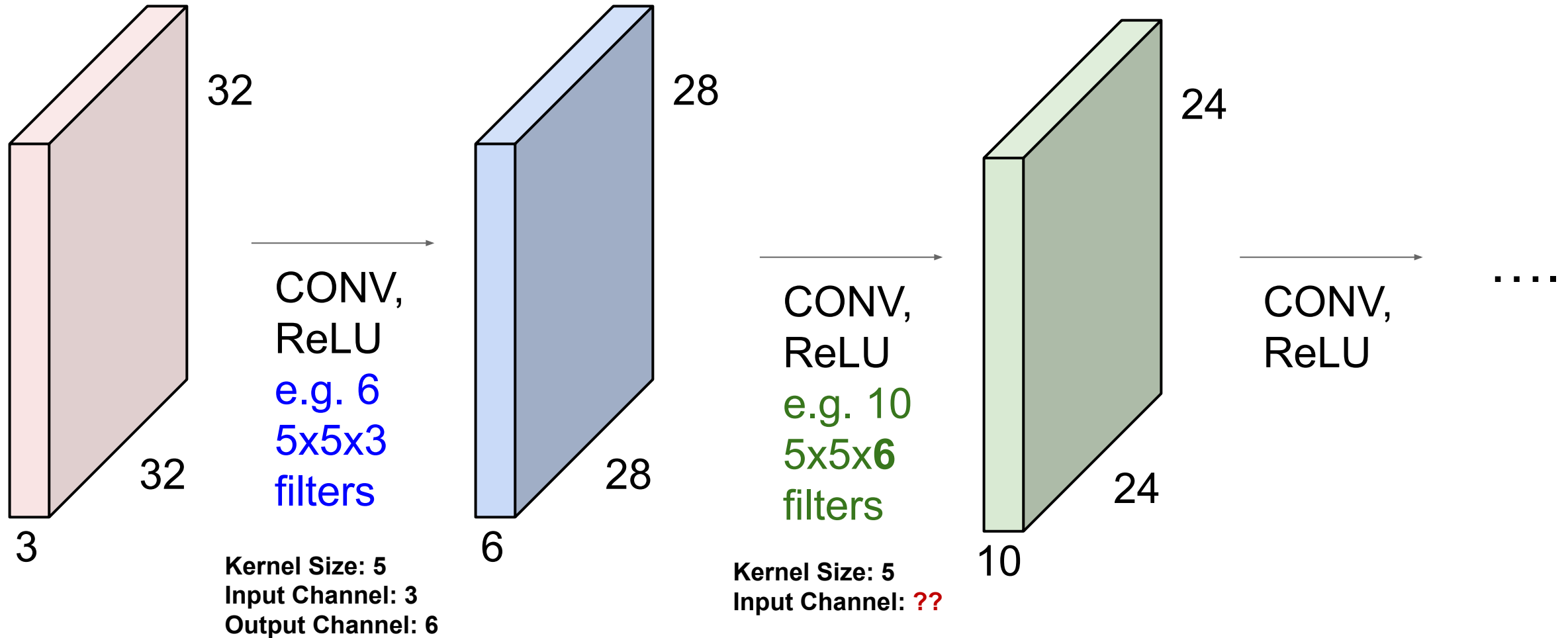
32

3

CONV,
ReLU
e.g. 6
5x5x3
filters

28

28

6

# ConvNet: Sequence of Convolution Layers, interspersed with activation functions



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

....

# ConvNet: Sequence of Convolution Layers, interspersed with activation functions



CONV, ReLU

e.g. 6 5x5x3 filters

**Kernel Size: 5**
**Input Channel: 3**
**Output Channel: 6**

CONV, ReLU

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

**Kernel Size: ??**

CONV, ReLU

....

32
32
3

28
28
6

24
24
10

# ConvNet: Sequence of Convolution Layers, interspersed with activation functions



32

32

3

CONV,
ReLU
e.g. 6
5x5x3
filters

**Kernel Size: 5**
**Input Channel: 3**
**Output Channel: 6**

28

28

6

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

**Kernel Size: 5**
**Input Channel: ??**

24

24

10

CONV,
ReLU

. . . .

# ConvNet: Sequence of Convolution Layers, interspersed with activation functions



32
32
3

CONV,
ReLU
e.g. 6
5x5x3
filters

**Kernel Size: 5**
**Input Channel: 3**
**Output Channel: 6**

28
28
6

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

**Kernel Size: 5**
**Input Channel: 6**
**Output Channel: ??**

24
24
10

CONV,
ReLU

....

# ConvNet: Sequence of Convolution Layers, interspersed with activation functions



32

32

3

CONV,
ReLU
e.g. 6
5x5x3
filters

**Kernel Size: 5**
**Input Channel: 3**
**Output Channel: 6**

28

28

6

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

**Kernel Size: 5**
**Input Channel: 6**
**Output Channel: 10**

24

24

10

CONV,
ReLU

....

A closer look at spatial dimensions:

32x32x3 image

5x5x3 filter

**Activation map (Output)**

32

32

3

convolve (slide) over all spatial locations

**28**

**28**

1

# A closer look at spatial dimensions:

(Step 1)

7

7
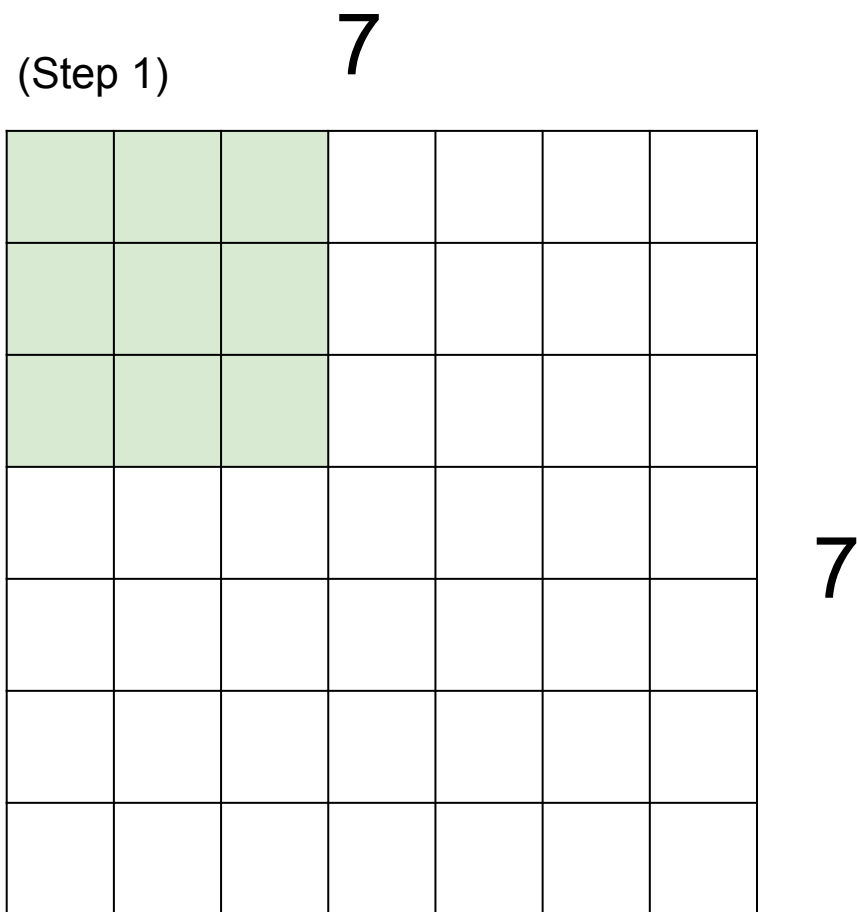
7x7 input (spatially)
assume 3x3 filter

# A closer look at spatial dimensions:

7

7

7x7 input (spatially)
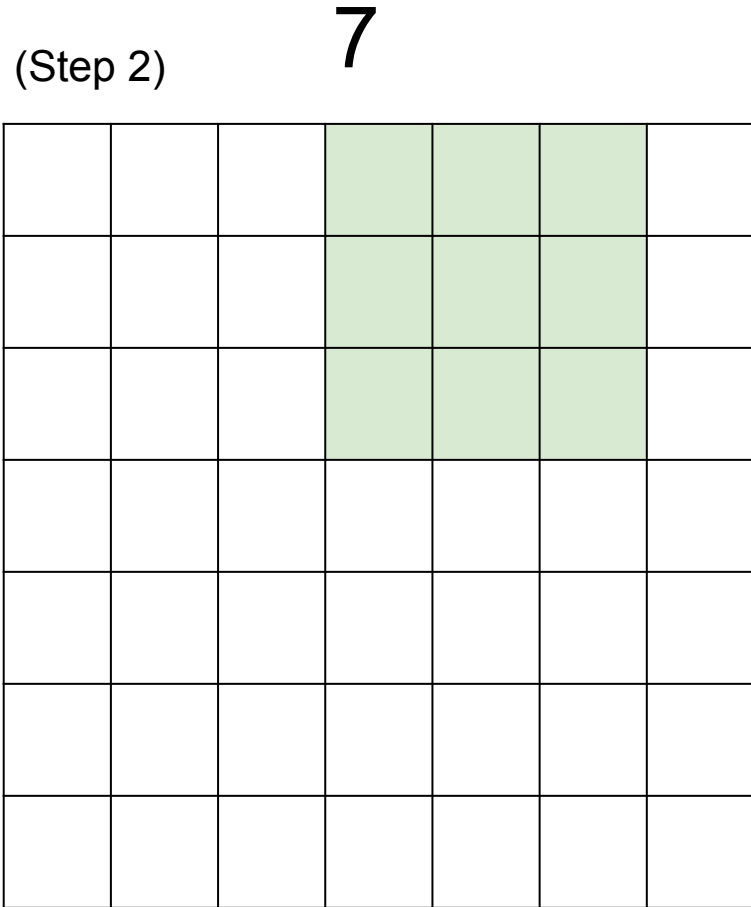assume 3x3 filter

A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter

(Step 3)

7

7

# A closer look at spatial dimensions:

(Step 4)

7

7

7x7 input (spatially)
assume 3x3 filter

A closer look at spatial dimensions:

(Step 5)

7

7

7x7 input (spatially)
assume 3x3 filter

**=> 5x5 output**

(7-3+1) x (7-3+1) = 5 x 5

A closer look at spatial dimensions:

(Step 1)

7

7

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

A closer look at spatial dimensions:

(Step 2)

7

7

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

A closer look at spatial dimensions:

(Step 3)

7

7



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

A closer look at spatial dimensions:

(Step 1)

7

7

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

A closer look at spatial dimensions:

(Step 2)

7

7



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

Next?

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

e.g. N = 7, K = 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 - K) / stride + 1

# In practice: Common to zero pad the border



e.g. input 7x7
**3x3** filter, applied with **stride 1**
**pad with 1 pixel** border => what is the output?

**7x7 output!**

(recall:)
$(N + 2P - K) / stride + 1$

# 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 KxK, and zero-padding with (K-1)/2. (will preserve size spatially)
e.g. K = 3 => zero pad with 1
       K = 5 => zero pad with 2
       K = 7 => zero pad with 3

**Why zero padding?**
E.g. 32x32 input convolved repeatedly with 5x5 filters shrinks volumes spatially!
(32 -> 28 -> 24 ...)

- can't stack deeply

- use stride to reduce size whenever we really wants

Examples time:

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

Let's assume output size is HxWxD.
What is D?

Examples time:

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

Let's assume output size is HxWxD.
What is D? **10**

Examples time:

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

Let's assume output size is HxWxC.
What is C? **10**
What is H or W?

Examples time:

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

Let's assume output size is HxWxC.
What is C? 10
What is H or W? (32+2*2-5)/1+1 = 32

Examples time:

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

Let's assume output size is HxWxD.
What is D? 10
What is H or W? (32+2*2-5)/1+1 = 32
So the total output size is: **32x32x10**

Examples time:

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

Number of parameters in this layer?

Exercise:

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**

# Convolution layer: summary

Let's assume input is $W_1$ x $H_1$ x $C_1$

Conv layer needs 4 hyperparameters:

- Number of filters $C_2$   (output channels)

- The filter size **K**

- The stride **S**

- The zero padding **P**

This will produce an output of $W_2$ x $H_2$ x $C_2$ where:

- $W_2 = (W_1 - K + 2P)/S + 1$

- $H_2 = (H_1 - K + 2P)/S + 1$
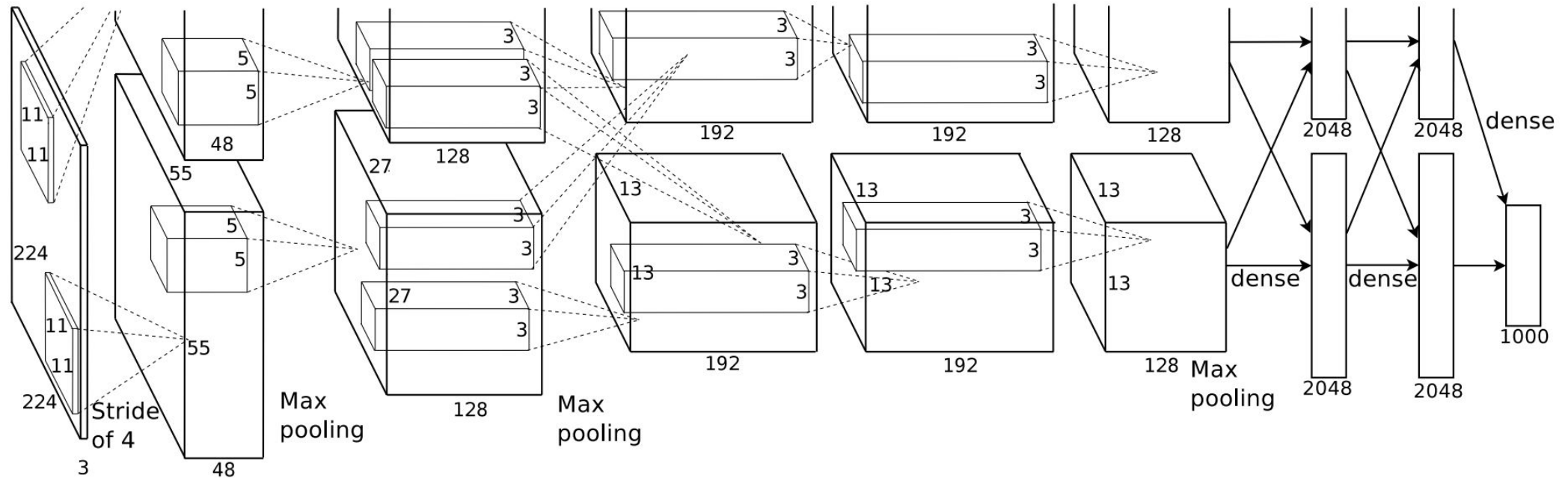
Number of parameters: $K^2 C_1 C_2$ and $C_2$ biases

# Receptive field



32

32

3

28

28

An activation map is a 28x28 sheet of neuron outputs:
1. Each is connected to a small region in the input
2. All of them share parameters

"5x5 filter" -> "5x5 receptive field for each neuron"

# 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

# The brain/neuron view of CONV Layer

Stack of Conv layers v.s. hierarchical organization in the visual processing system.

# What CNN learns?

# What CNN learns?

*[Zeiler and Fergus 2013]*

Visualization of VGG-16 by Lane McIntosh. VGG-16 architecture from [Simonyan and Zisserman 2014].



VGG-16 Conv1_1

VGG-16 Conv3_2

VGG-16 Conv5_3

one filter =>
one activation map

Activations:

example 5x5 filters
(32 total)

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.

# Recall:

Feature
Vector

Featurize

# Recall:

Feature
Vector

Featurize



e.g. 24 edge & blog filters
(Human priors)

# With CNN:

Feature
Vector

Featurize



CNN: Learning (more diverse) filters for you!!

# CNN (or Neural Network)

Representation (feature) learning + linear classifier

Feature Vector

Featurize



CNN: Learning (more diverse) filters for you!!

$x_1$

$x_2$

$x_3$

$x_N$

cat

dog

bird

Linear classifier (FC layer)

# 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

# Pooling layer: summary

Let's assume input is $W_1$ x $H_1$ x C
Pooling layer needs 2 hyperparameters:
- The kernel size **K**
- The stride **S**

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

Number of parameters: 0

# AlexNet (2012) – The start of modern deep learning

# AlexNet (2012) – The start of modern deep learning



Common Practice:
- Reduce the spatial dimension while increasing channels (why?)
- Output size is no larger than the input size (why?)

# AlexNet (2012) – The start of modern deep learning



224x224x3  v.s.   55x55x48

Common Practice:
- Reduce the spatial dimension while increasing channels
- Output size is no larger than the input size

# AlexNet (2012) – The start of modern deep learning



27x27x128  v.s.  13x13x192
Reduction ratio: (1/2)x(1/2)*1.5=0.375

Common Practice:
- Reduce the spatial dimension while increasing channels
- Output size is no larger than the input size

# Example: DPM is CNN

# Next time

Last Lecture: CV Frontier