



Lecture 11: Recognition

Instructor: Roni Sengupta

ULA: Andrea Dunn, William Li,
Liujie Zheng



Course Website:
Scan Me!

What is “Recognition”?



Next few slides adapted from Li, Fergus, & Torralba’s excellent [short course](#) on category and object recognition

What is “Recognition”?

- Verification: is that a lamp?



What is “Recognition”?

- Verification: is that a lamp?
- Detection: where are the people?



What is “Recognition”?

- Verification: is that a lamp?
- Detection: where are the people?
- Identification: is that Potala Palace?



What is “Recognition”?

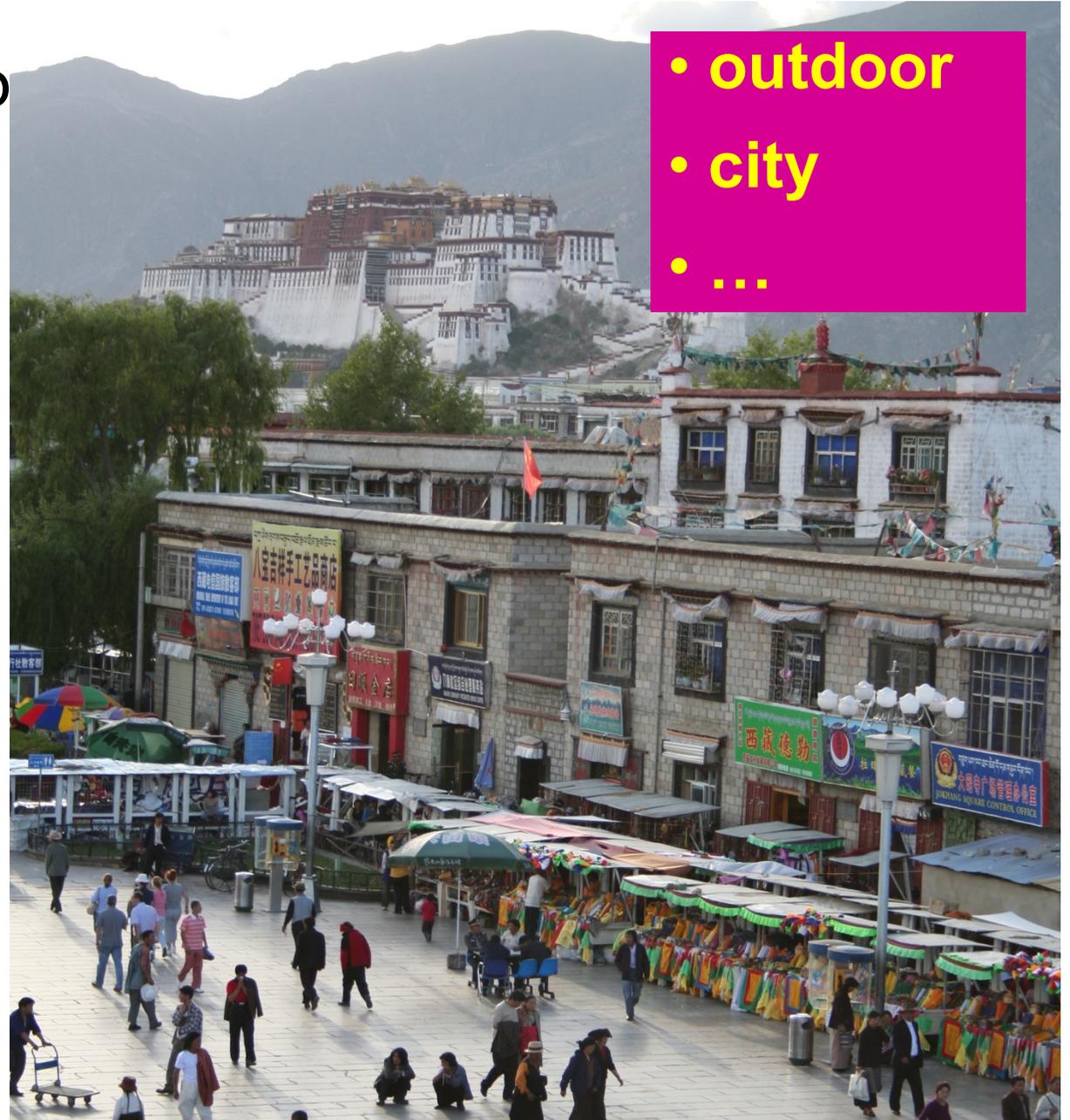
- Verification: is that a lamp?
- Detection: where are the people?
- Identification: is that Potala Palace?
- Object categorization



What is “Recognition”?

- Verification: is that a lamp?
- Detection: where are the people?
- Identification: is that Potala Palace?
- Object categorization
- Scene and context categorization

- outdoor
- city
- ...



What is “Recognition”?

- Verification: is that a lamp?
- Detection: where are the people?
- Identification: is that Potala Palace?
- Object categorization
- Scene and context categorization
- **Activity / Event Recognition**

what are these people doing?

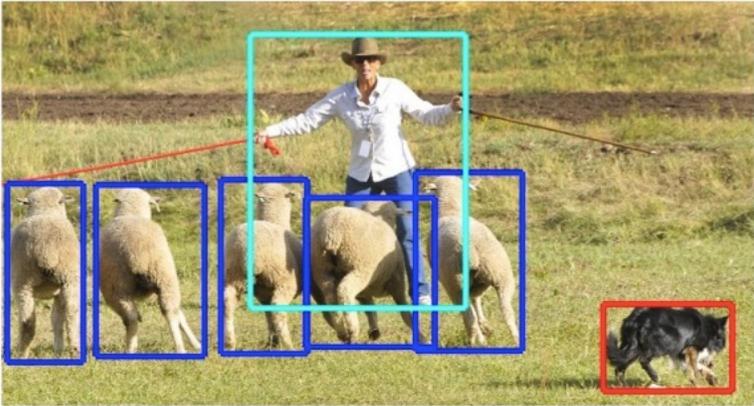


Recognition: What type of output?

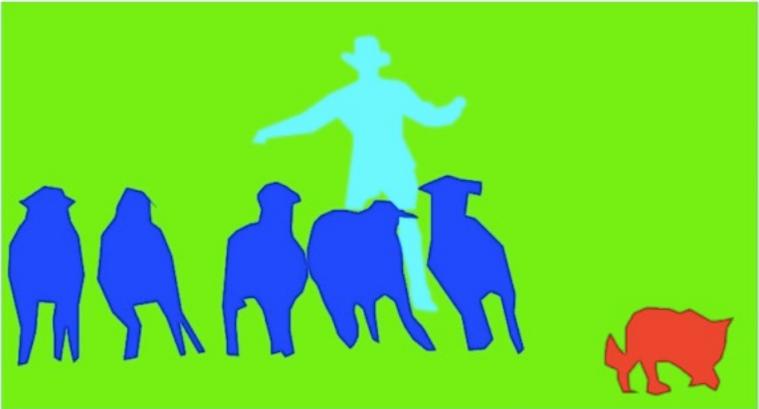
Image classification



Object detection



Semantic segmentation



Instance segmentation



Image Classification: a core task in computer vision

- Assume given set of discrete labels, e.g.
{cat, dog, cow, apple, tomato, truck, ... }

$f(\text{apple image}) = \text{“apple”}$

$f(\text{tomato image}) = \text{“tomato”}$

$f(\text{cow image}) = \text{“cow”}$

What Matters in Recognition?

- Data
 - More is always better (as long as it is good data)
 - Annotation is the hard part
- Representation
 - Low level: SIFT, HoG, GIST, edges
 - Mid level: Bag of words, sliding window, deformable model
 - High level: Contextual dependence
 - Deep learned features
- Learning Techniques
 - E.g. choice of classifier or inference method

What Matters in Recognition?

- Data
 - More is always better (as long as it is good data)
 - Annotation is the hard part
- Representation
 - Low level: SIFT, HoG, GIST, edges
 - Mid level: Bag of words, sliding window, deformable model
 - High level: Contextual dependence
 - Deep learned features
- Learning Techniques
 - E.g. choice of classifier or inference method

Data Sets

- PASCAL VOC
 - *Not* Crowdsourced, bounding boxes, 20 categories
- **ImageNet**
 - **Huge, Crowdsourced, Hierarchical, *Iconic* objects**
- SUN Scene Database, Places
 - *Not* Crowdsourced, 397 (or 720) scene categories
- LabelMe (Overlaps with SUN)
 - Sort of Crowdsourced, Segmentations, Open ended
- SUN *Attribute* database (Overlaps with SUN)
 - Crowdsourced, 102 attributes for every scene
- OpenSurfaces
 - Crowdsourced, materials
- **Microsoft COCO**
 - **Crowdsourced, large-scale objects**

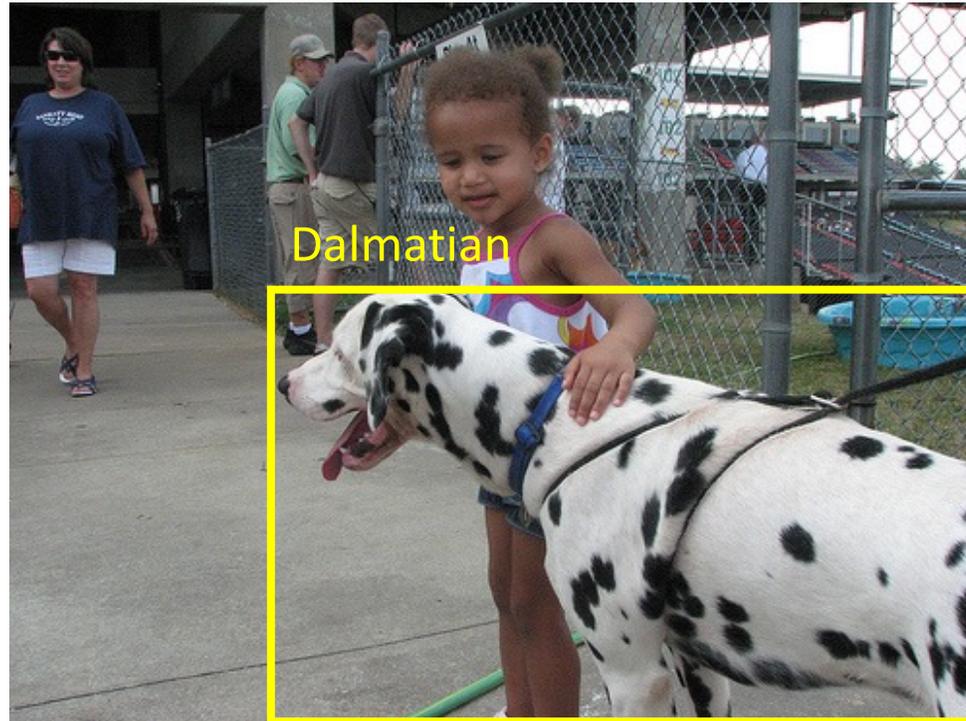
... and many more <https://paperswithcode.com/datasets?task=image-classification>

Large Scale Visual Recognition Challenge (ILSVRC)

2010-2017

IM  GENET

~~20 object classes~~ ————— ~~22,591 images~~
1000 object classes **1,431,167 images**



<http://image-net.org/challenges/LSVRC/{2010,2011,2012}>

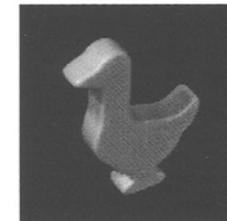
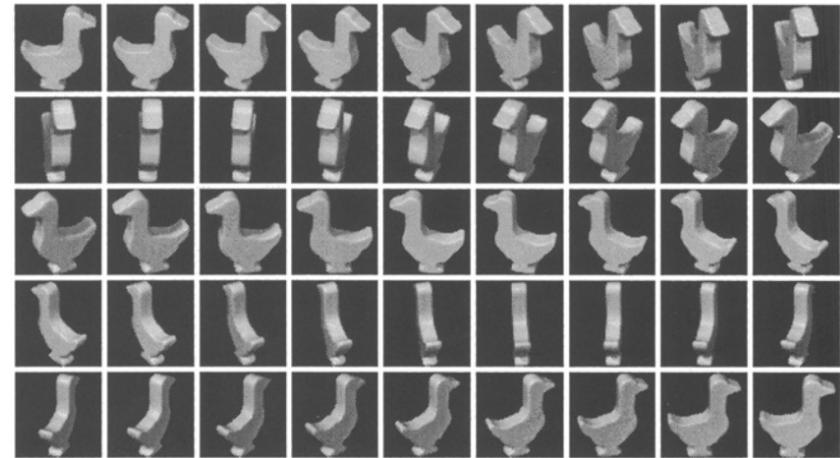
What Matters in Recognition?

- Data
 - More is always better (as long as it is good data)
 - Annotation is the hard part
- Representation
 - **Low level: SIFT, HoG, GIST, edges**
 - **Mid level: Bag of words, sliding window, deformable model**
 - **High level: Contextual dependence**
 - Deep learned features
- Learning Techniques
 - E.g. choice of classifier or inference method

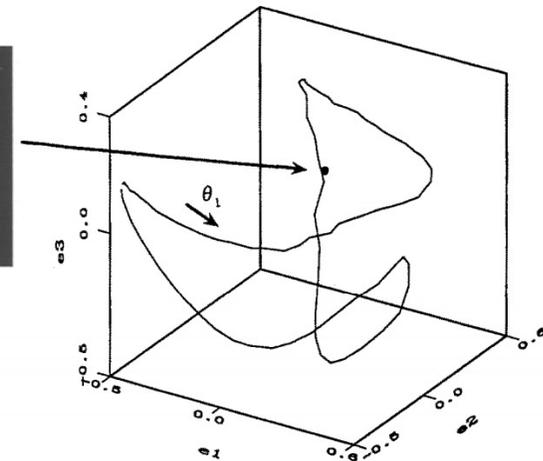
History of recognition: Appearance-based models



M. Turk and A. Pentland, [Face recognition using eigenfaces](#), CVPR 1991



(a)

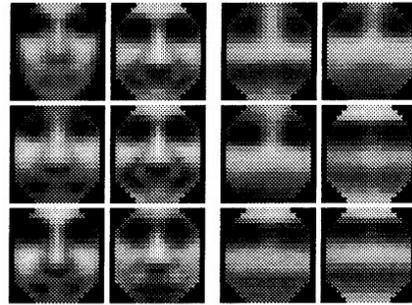


(b)

H. Murase and S. Nayar, [Visual learning and recognition of 3-d objects from appearance](#), IJCV 1995

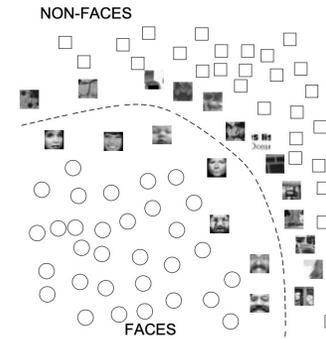
History of recognition: Features and classifiers

Appearance manifolds + neural network



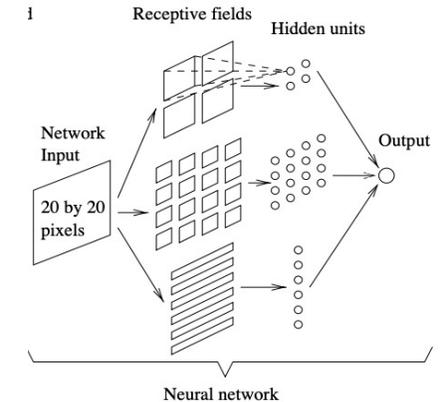
Face Centroids Non-Face Centroids
Sung & Poggio (1994)

Support vector machines



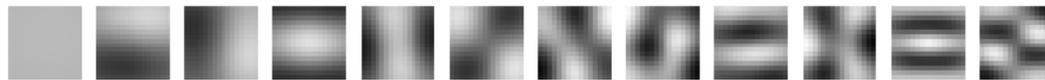
Osuna, Freund, Girosi (1997)

Neural network



Rowley, Baluja, Kanade (1998)

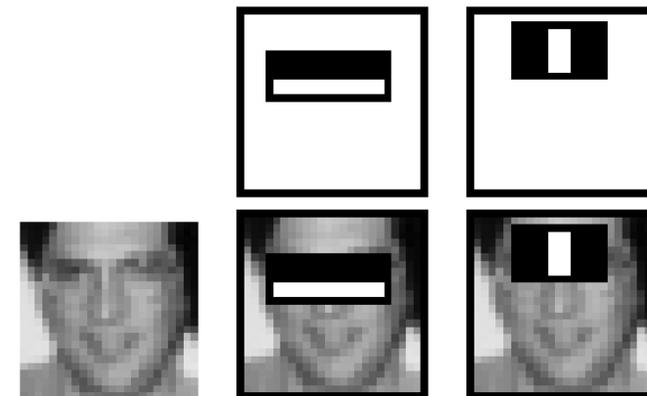
Statistics of feature responses, probabilistic classifier



$$\prod_{j=1}^{n_{magn}} \prod_{i=1}^{n_{subs}} \frac{P(q1_j^i | \text{object}) P(pos_j^i | q2_i, \text{object})}{\frac{P(q1_j^i | \overline{\text{object}})}{n_{subs}}} > \lambda = \frac{P(\overline{\text{object}})}{P(\text{object})}$$

Schneiderman & Kanade (1998)

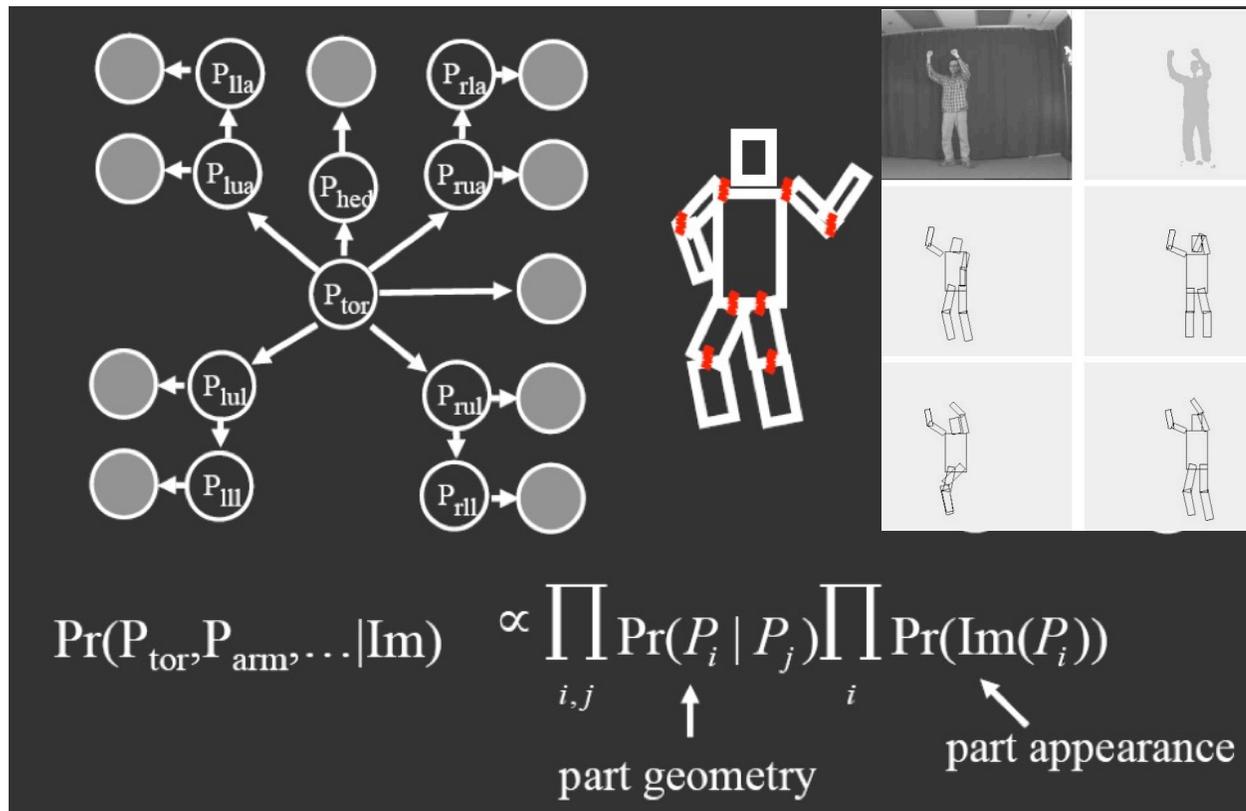
Rectangle features, boosting



Viola & Jones (2001)

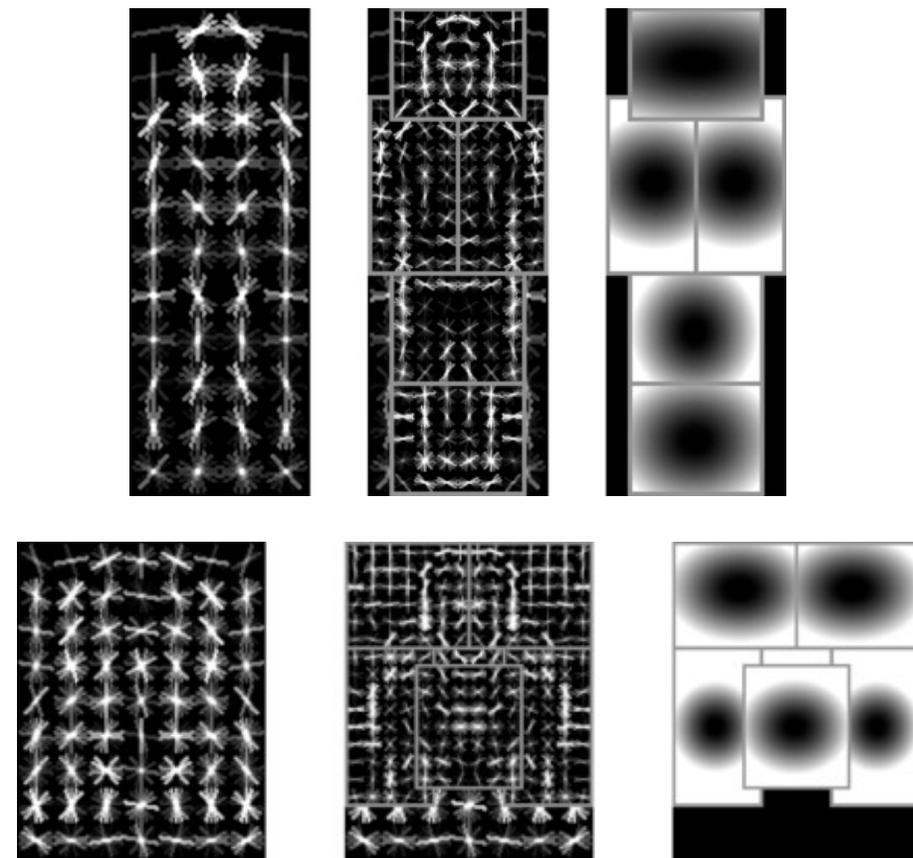
History of recognition: Deformable templates

Pictorial structures revisited



Felzenszwalb & Huttenlocher (2000)

Discriminatively trained deformable part-based models

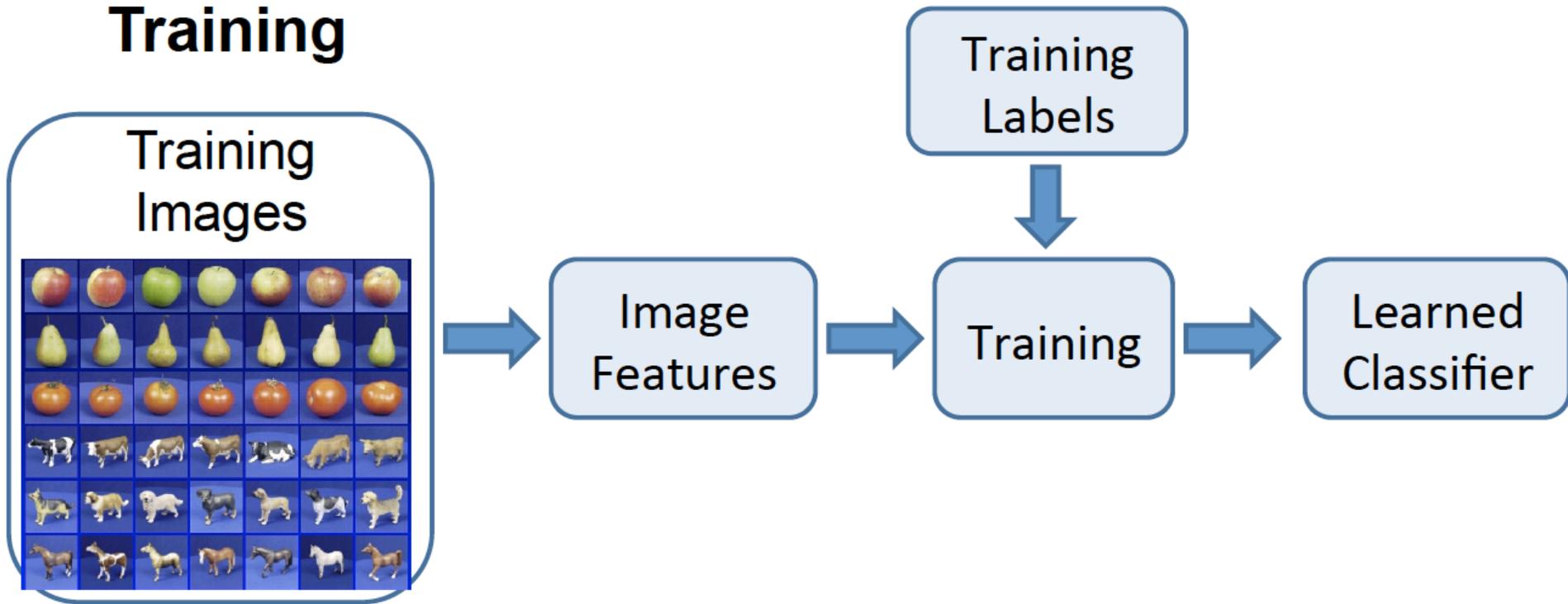


Felzenszwalb et al. (2008)

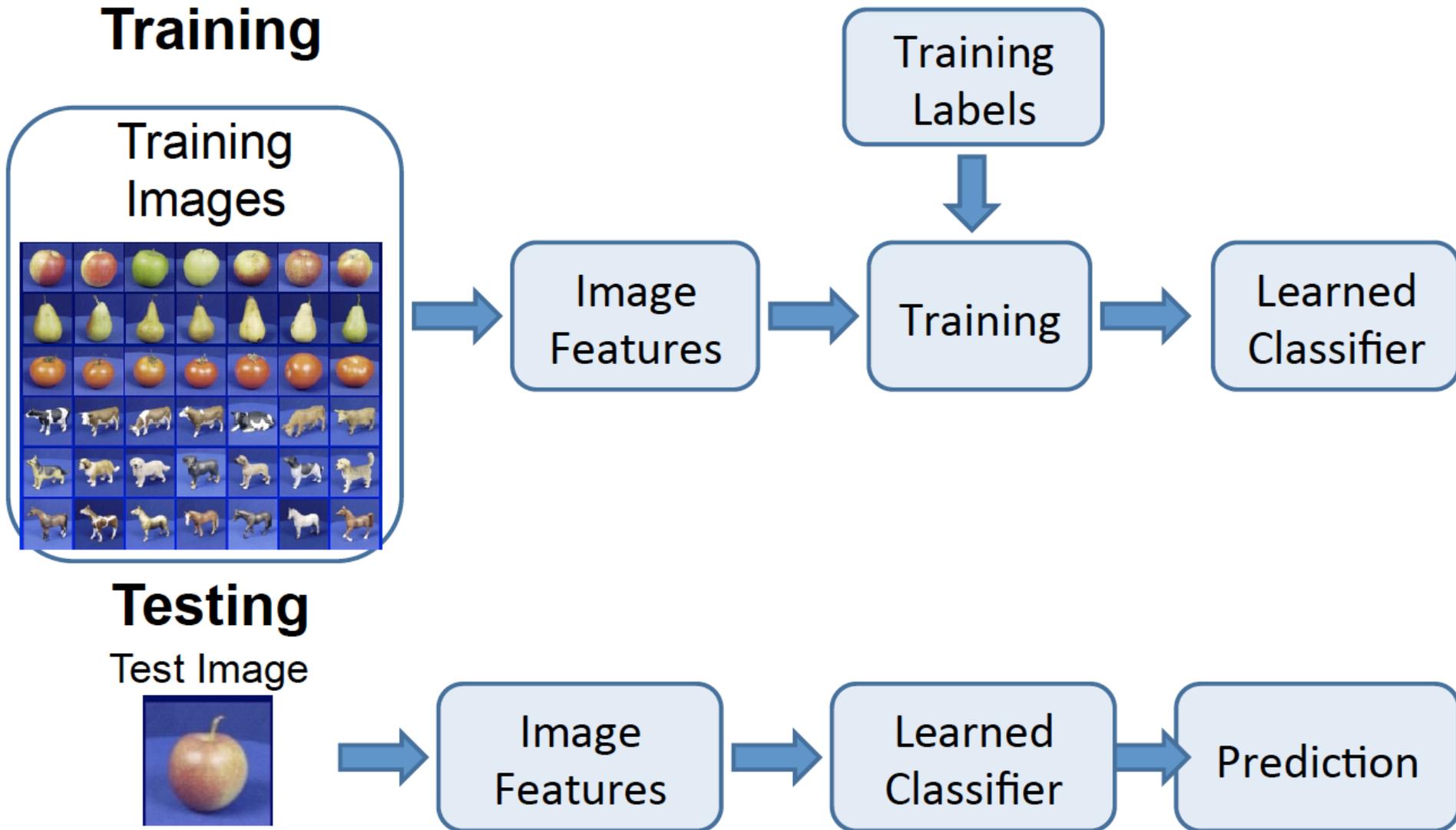
What Matters in Recognition?

- Data
 - More is always better (as long as it is good data)
 - Annotation is the hard part
- Representation
 - Low level: SIFT, HoG, GIST, edges
 - Mid level: Bag of words, sliding window, deformable model
 - High level: Contextual dependence
 - Deep learned features
- Learning Techniques
 - E.g. choice of classifier or inference method

Training & Testing a Classifier



Training & Testing a Classifier



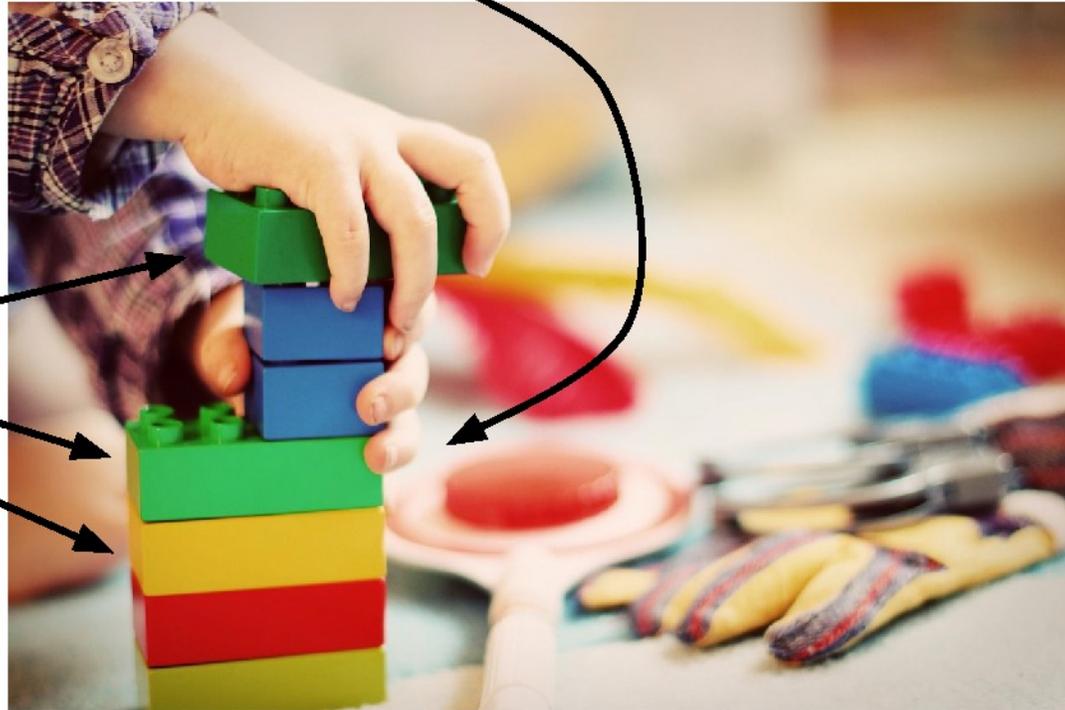
Classifiers

- Nearest Neighbor
- kNN (“k-Nearest Neighbors”)
- Linear Classifier
- Neural Network
- Deep Neural Network
- ...

Linear Classifiers

Neural Network

Linear
classifiers



[This image is CC0.1.0 public domain](#)

Score functions



class scores

Parametric Approach



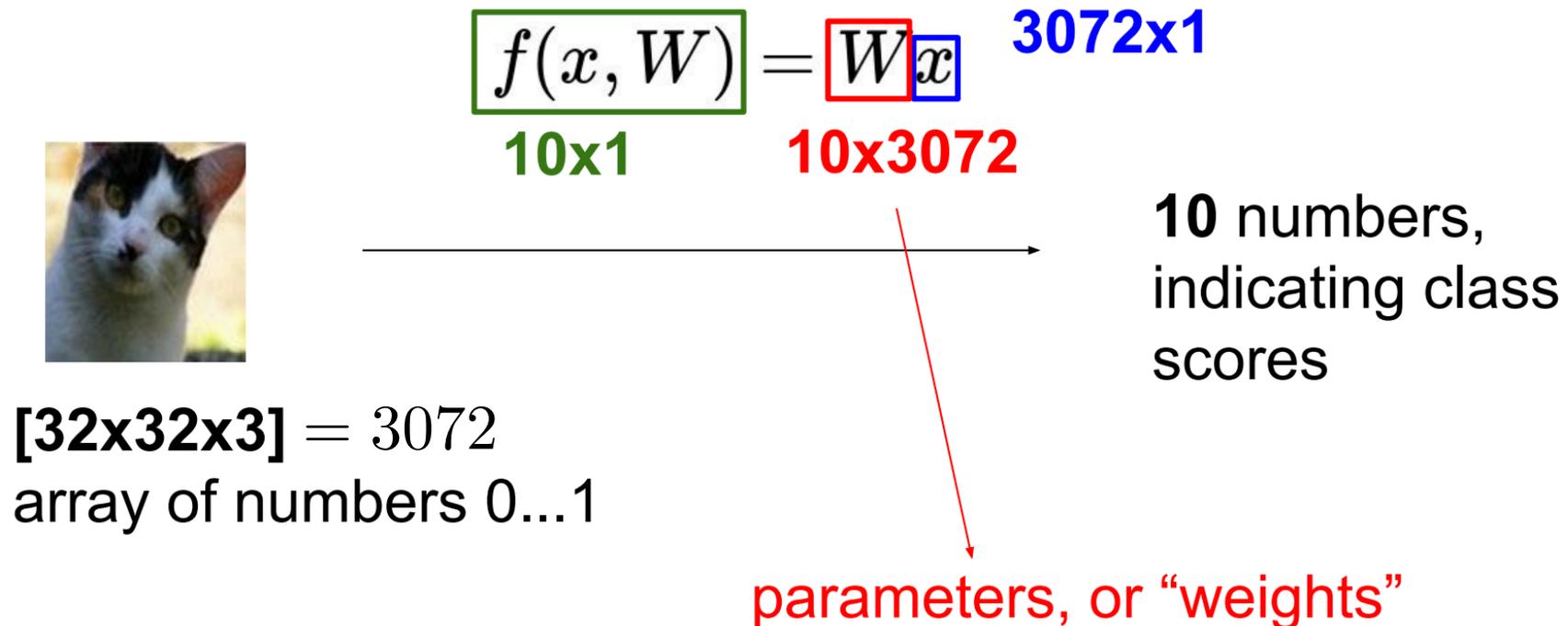
[32x32x3] = 3072
array of numbers 0...1
(3072 numbers total)

image parameters

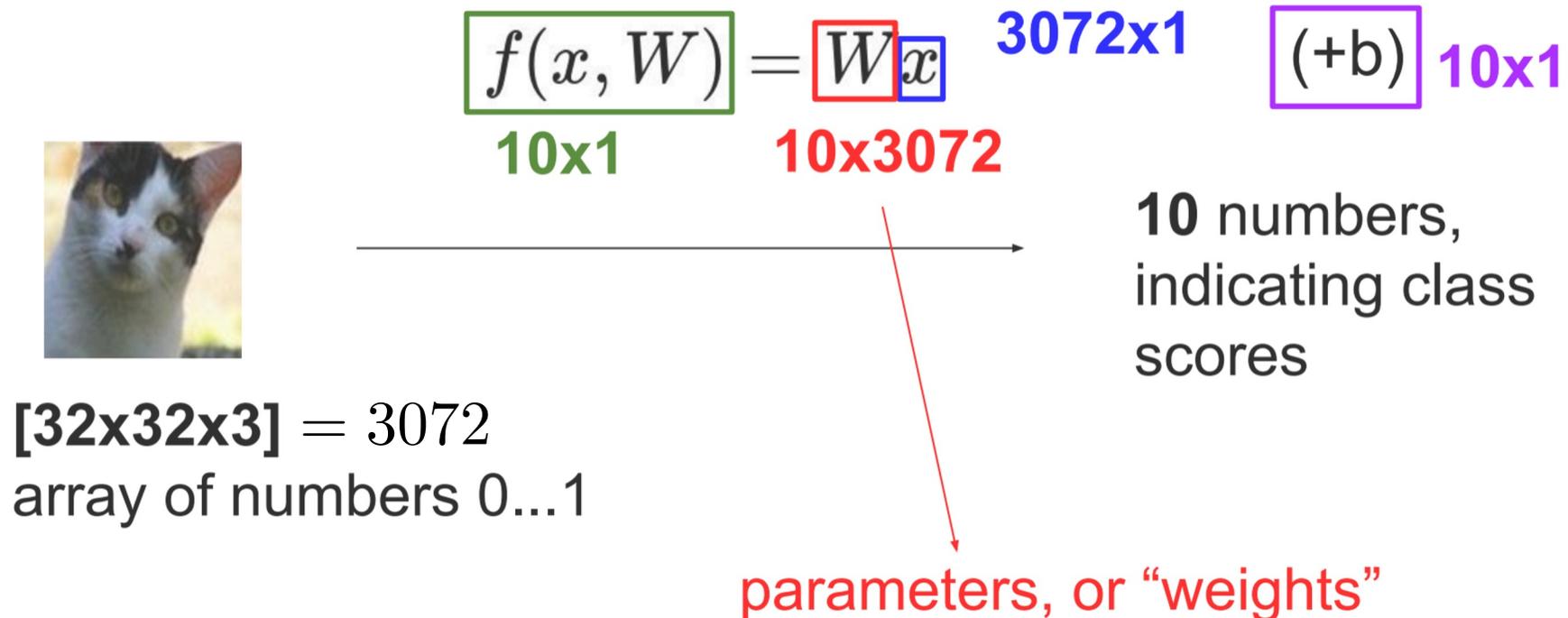
$f(\mathbf{x}, \mathbf{W})$

10 numbers,
indicating class
scores

Parametric Approach: Linear Classifier



Parametric Approach: Linear Classifier



Linear Classifier

define a **score function**

$$f(x_i, W, b) = Wx_i + b$$

data (image)

class scores

“weights”

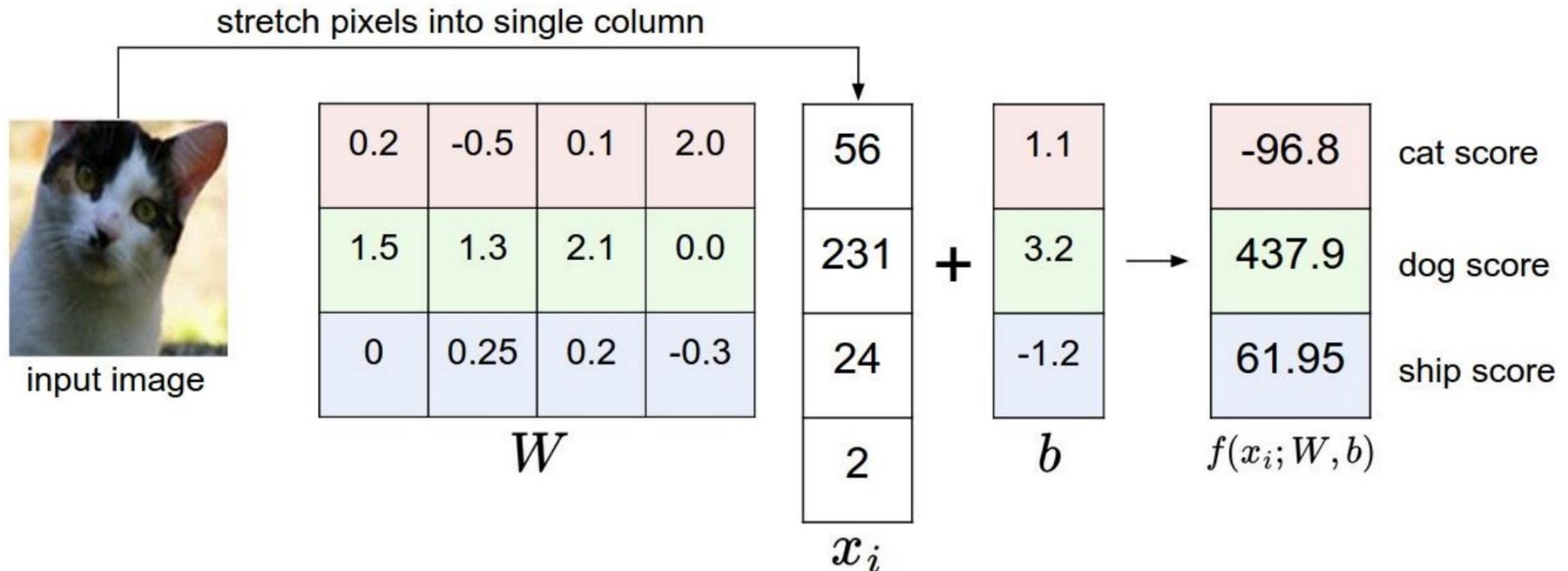
“bias vector”

“parameters”

The diagram shows the equation $f(x_i, W, b) = Wx_i + b$ centered on the slide. Several annotations with arrows point to parts of the equation: 'data (image)' points to x_i ; 'class scores' points to $f(x_i, W, b)$; '“weights”' points to W ; '“bias vector”' points to b ; and '“parameters”' points to both W and b .

Interpretation: Algebraic

Example with an image with 4 pixels, and 3 classes (cat/dog/ship)

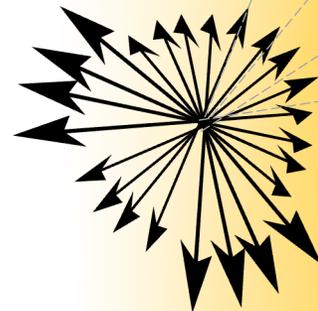


Interpretation: Geometric

- Parameters define a hyperplane for each class:

$$f(x_i, W, b) = Wx_i + b$$

- We can think of each class score as defining a distribution that is proportional to distance from the corresponding hyperplane



The Space of
All Images

Simpler example: binary classification

- Two classes (e.g., “cat” and “not cat”)
 - AKA “positive” and “negative” classes



cat

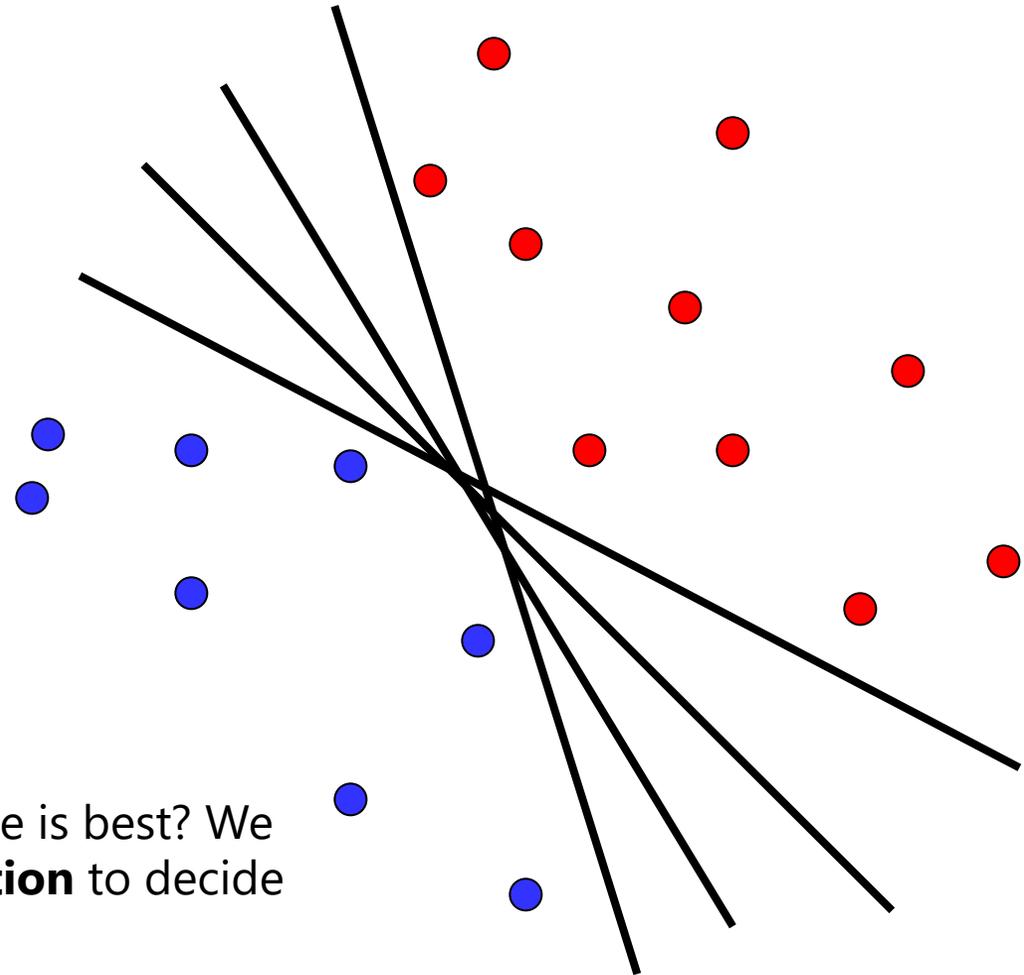
not cat

Simpler example: binary classification

- Find linear function (*hyperplane*) to separate positive and negative examples

$$\mathbf{x}_i \text{ positive : } \mathbf{x}_i \cdot \mathbf{w} + b \geq 0$$

$$\mathbf{x}_i \text{ negative : } \mathbf{x}_i \cdot \mathbf{w} + b < 0$$



Which hyperplane is best? We need a **loss function** to decide

Linear classification



airplane	-3.45	-0.51	3.42
automobile	-8.87	6.04	4.64
bird	0.09	5.31	2.65
cat	2.9	-4.22	5.1
deer	4.48	-4.19	2.64
dog	8.02	3.58	5.55
frog	3.78	4.49	-4.34
horse	1.06	-4.37	-1.5
ship	-0.36	-2.09	-4.79
truck	-0.72	-2.93	6.14

Cat image by Nikita is licensed under CC-BY 2.0; Car image is CC0 1.0 public domain; Frog image is in the public domain

Output scores

TODO:

1. Define a **loss function** that quantifies our unhappiness with the scores across the training data.
2. Come up with a way of efficiently finding the parameters that minimize the loss function.
(optimization)

Loss functions

Suppose: 3 training examples, 3 classes.
With some W the scores $f(x, W) = Wx$ are:



cat	3.2	1.3	2.2
car	5.1	4.9	2.5
frog	-1.7	2.0	-3.1

A **loss function** tells how good our current classifier is

Given a dataset of examples

$$\{(x_i, y_i)\}_{i=1}^N$$

Where x_i is image and
 y_i is (integer) label

Loss over the dataset is a sum of loss over examples:

$$L = \frac{1}{N} \sum_i L_i(f(x_i, W), y_i)$$

Softmax classifier

- Interpret Scores as unnormalized log probabilities of classes

$$f(x_i, W) = Wx_i \quad (\text{score function})$$

$$\frac{e^{f_{y_i}}}{\sum_j e^{f_j}}$$

softmax function

Squashes values into *probabilities* ranging from 0 to 1

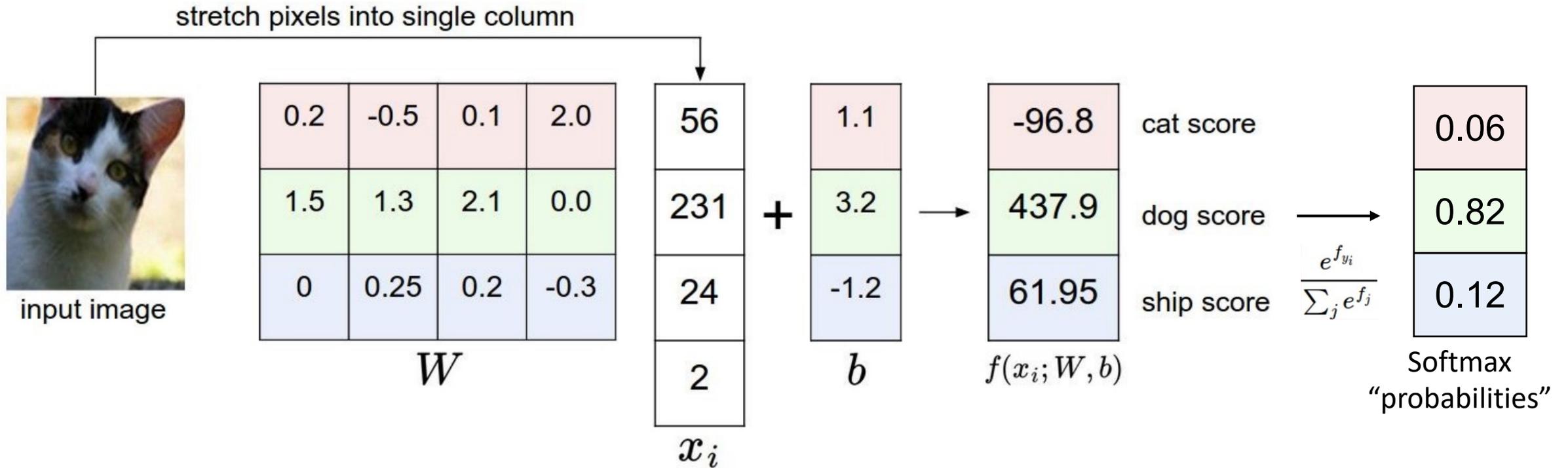
$$P(y_i | x_i; W)$$

Example with three classes:

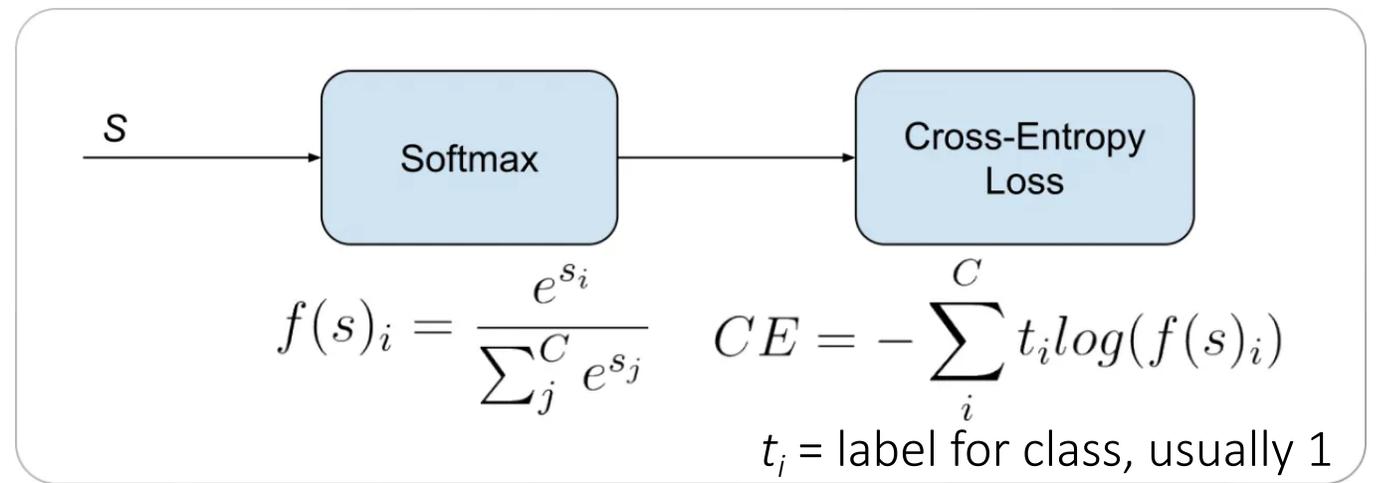
$$[1, -2, 0] \rightarrow [e^1, e^{-2}, e^0] = [2.71, 0.14, 1] \rightarrow [0.7, 0.04, 0.26]$$

Softmax classifier

Example with an image with 4 pixels, and 3 classes (cat/dog/ship)



Cross-entropy loss



$$f(x_i, W) = Wx_i \quad (\text{score function})$$

Summary

- Have score function and loss function
 - Currently, score function is based on linear classifier
 - Next, will generalize to convolutional neural networks
- Find W and b to minimize loss

$$L = \frac{1}{N} \sum_i -\log \left(\frac{e^{f_{y_i}}}{\sum_j e^{f_j}} \right) + \lambda \sum_k \sum_l W_{k,l}^2$$

Average of cross-entropy loss
over all training examples

Regularization term

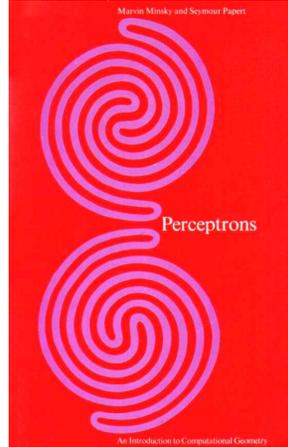
(Deep) Neural Networks

History of recognition: Neural networks

Perceptrons

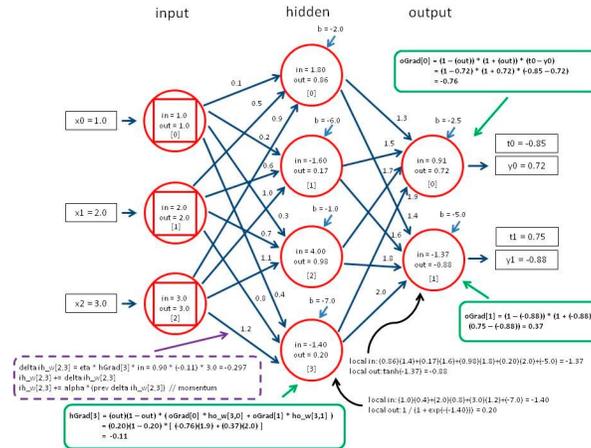


Rosenblatt (1958)



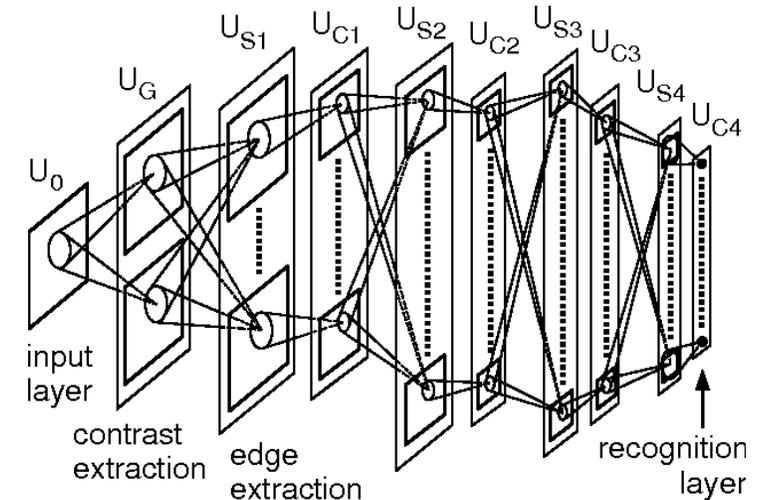
Minsky & Papert (1969)

Back-propagation



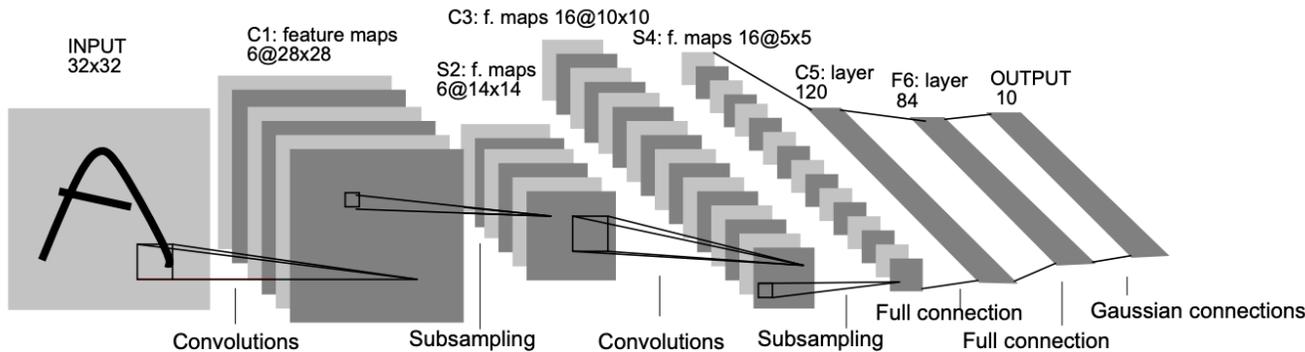
Rumelhart, Hinton & Williams (1986)

Neocognitron



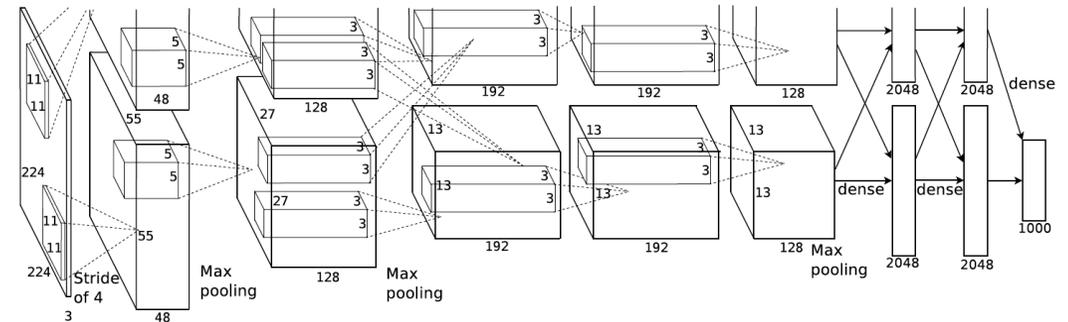
Fukushima (1980)

LeNet-5



LeCun et al. (1998)

AlexNet



Krizhevsky et al. (2012)

Neural networks

(**Before**) Linear score function: $f = Wx$

Neural networks

(**Before**) Linear score function:

$$f = Wx$$

(**Now**) 2-layer Neural Network

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

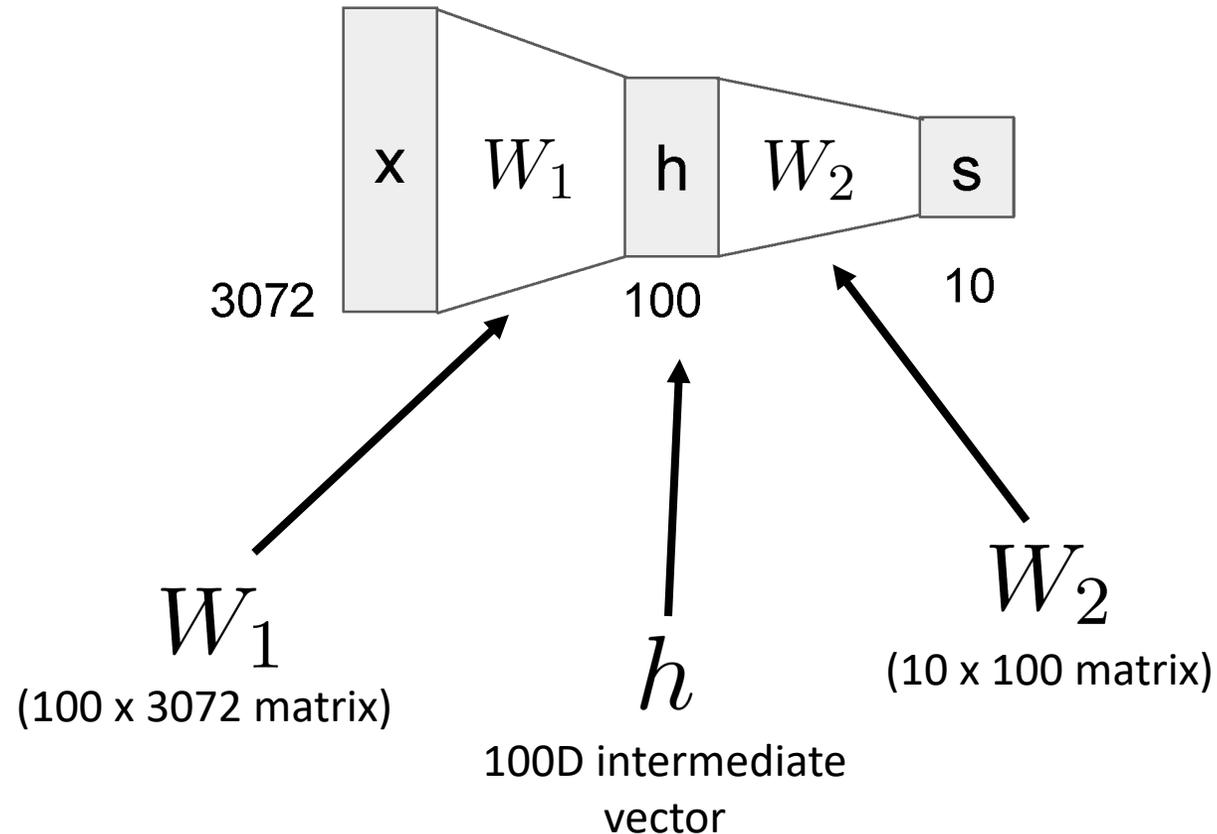


Non-linear Activation Function
(many other choices exist)

Neural networks

(**Before**) Linear score function: $f = Wx$

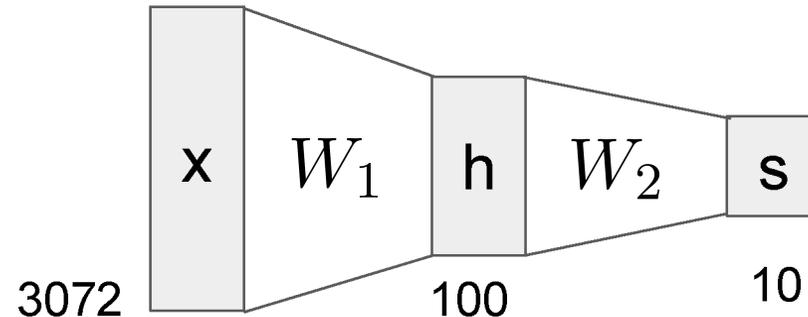
(**Now**) 2-layer Neural Network $f = W_2 \max(0, W_1 x)$



Neural networks

(**Before**) Linear score function: $f = Wx$

(**Now**) 2-layer Neural Network $f = W_2 \max(0, W_1 x)$



- Total number of weights to learn:

$$3,072 \times 100 + 100 \times 10 = 308,200$$

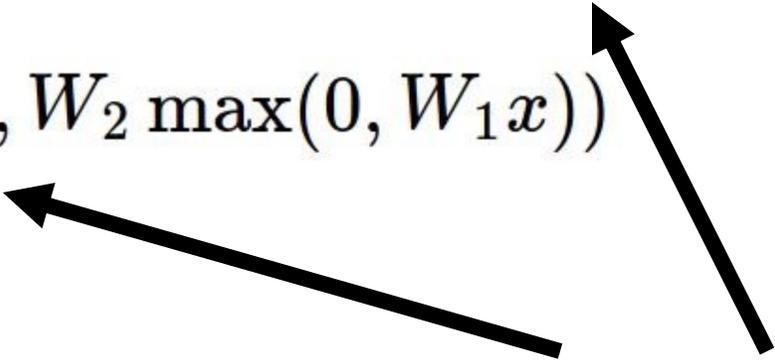
Neural networks

(**Before**) Linear score function: $f = Wx$

(**Now**) 2-layer Neural Network
or 3-layer Neural Network $f = W_2 \max(0, W_1 x)$

$$f = W_3 \max(0, W_2 \max(0, W_1 x))$$

also called “Multi-Layer
Perceptrons” (MLPs)

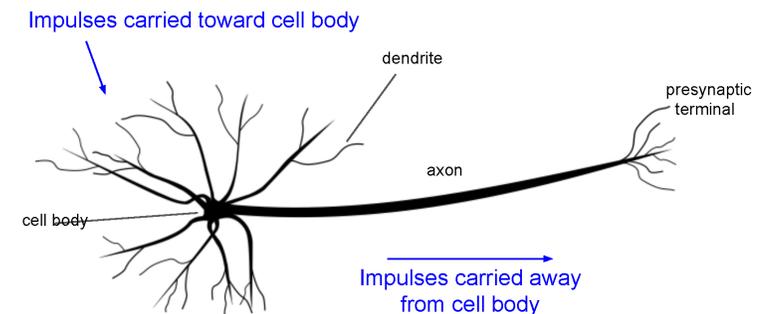


Neural networks

- Very coarse generalization of neural networks:
 - Linear functions chained together and separated by non-linearities (*activation functions*), e.g. “max”

$$f = W_3 \max(0, W_2 \max(0, W_1 x))$$

- Why separate linear functions with non-linear functions?
- *Very roughly* inspired by real neurons

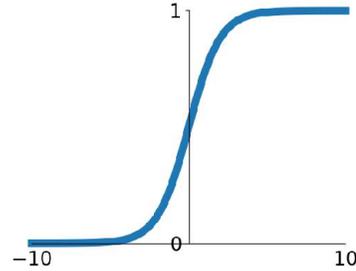


This image by Felipe Perucho is licensed under [CC-BY 3.0](https://creativecommons.org/licenses/by/3.0/)

Activation functions

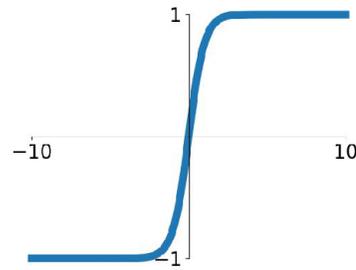
Sigmoid

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



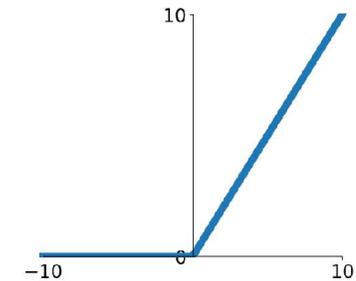
tanh

$$\tanh(x)$$



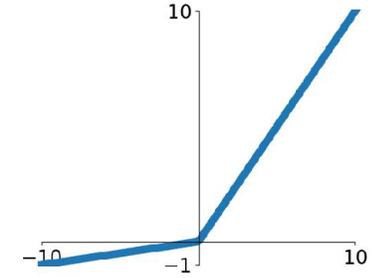
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$

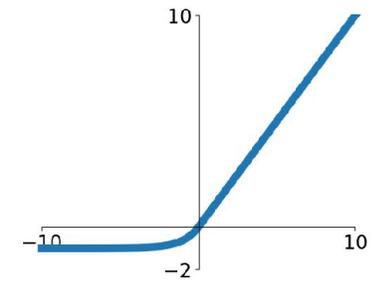


Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

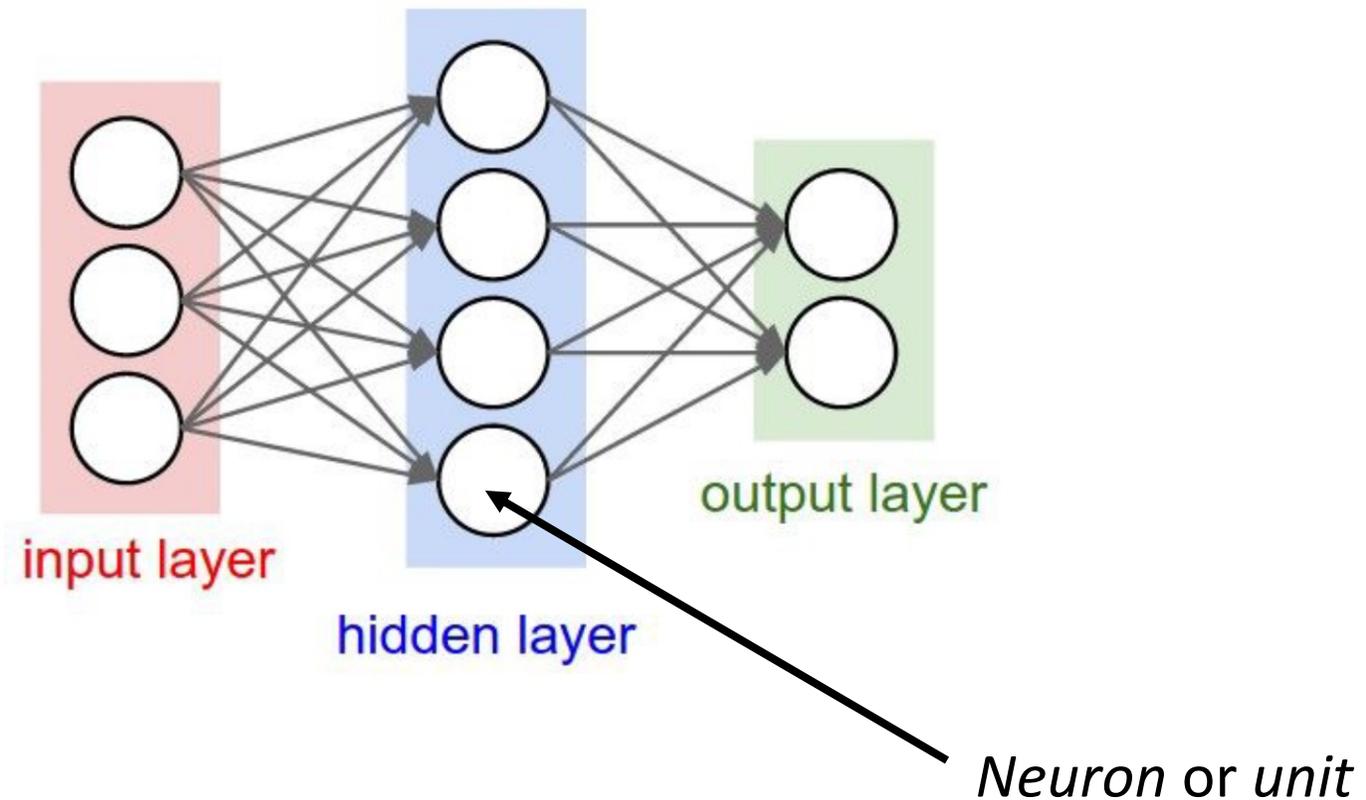
ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$

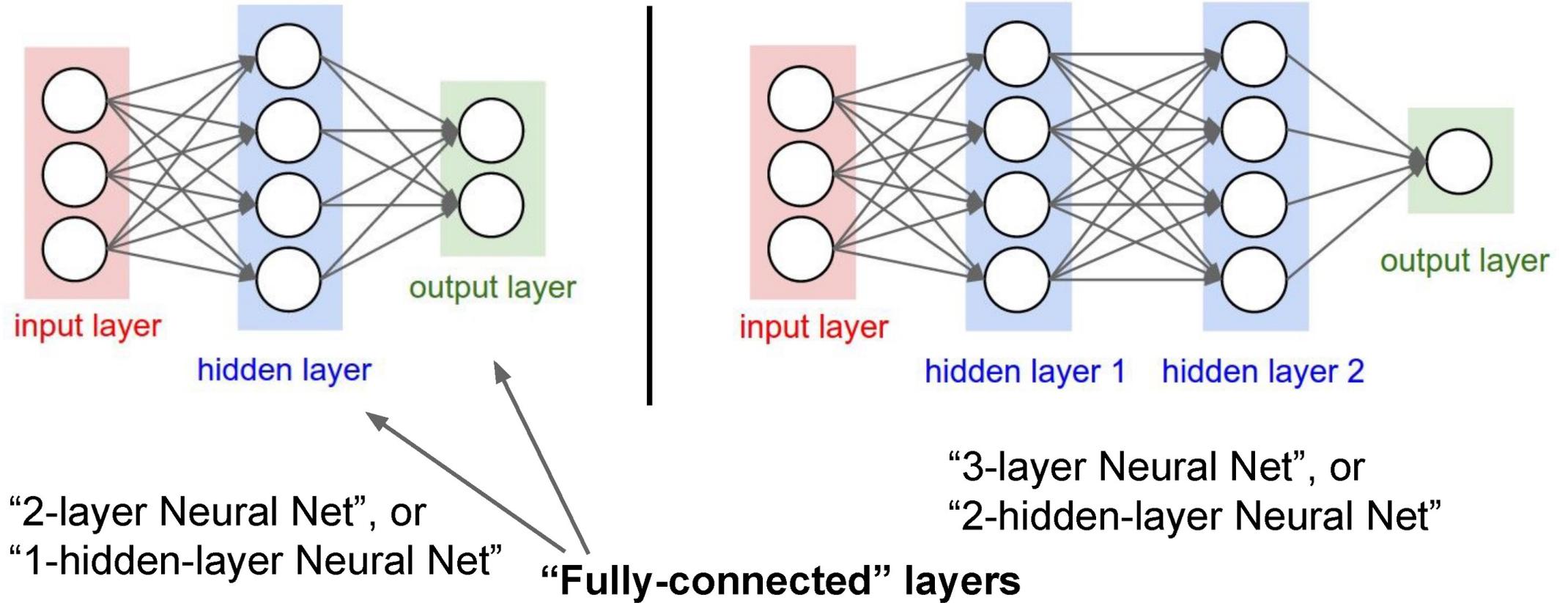


Neural network architecture

- Computation graph for a 2-layer neural network



Neural networks: Architectures

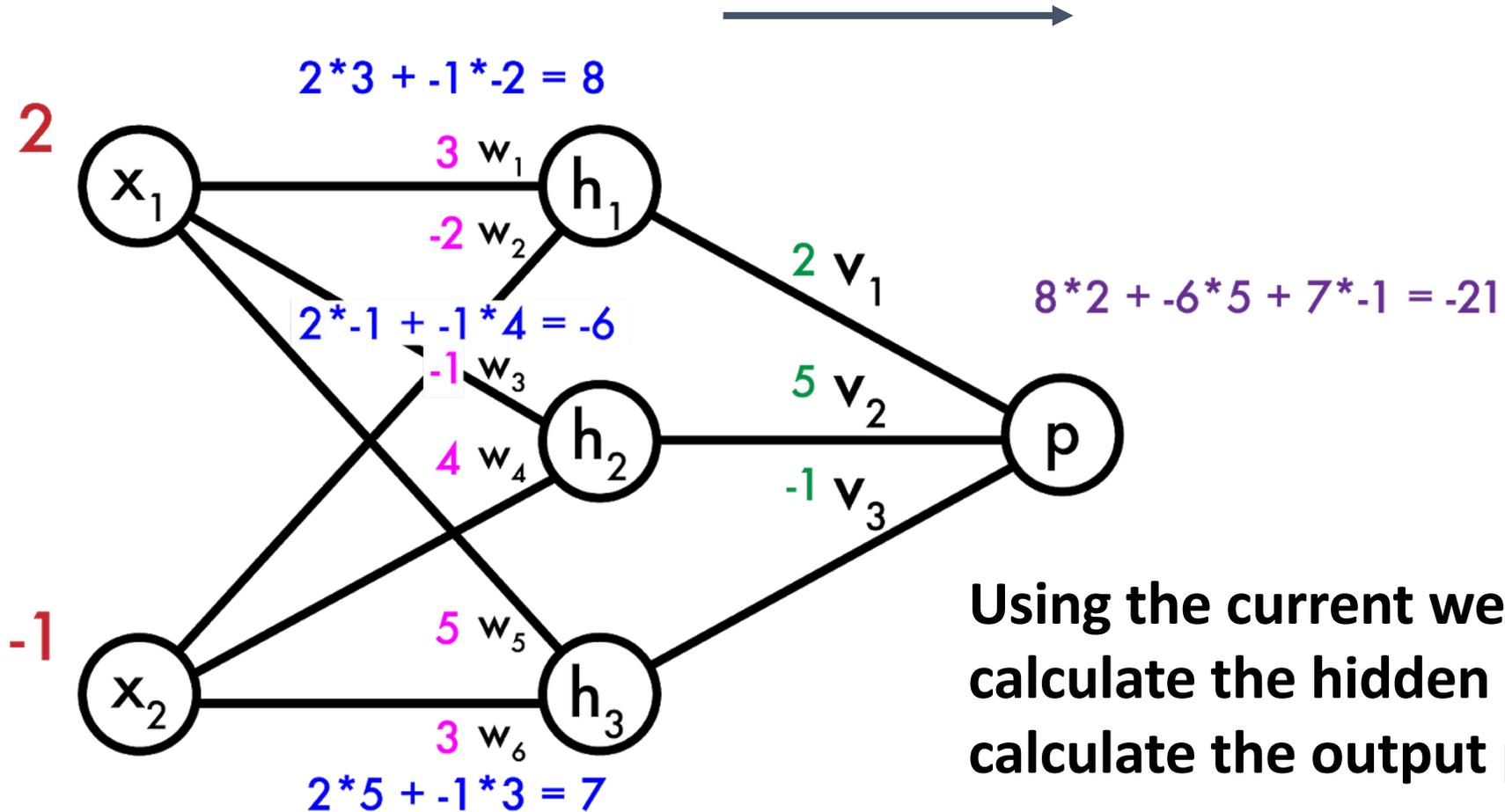


Deep networks typically have many layers and potentially millions of parameters

Optimizing parameters with gradient descent

- How do we find the best \mathbf{W} and \mathbf{b} parameters?
- In general: *gradient descent*
 1. Start with a guess of a good \mathbf{W} and \mathbf{b} (or randomly initialize them)
 2. Compute the loss function for this initial guess and the *gradient* of the loss function
 3. Step some distance in the negative gradient direction (direction of steepest descent)
 4. Repeat steps 2 & 3
- Note: efficiently performing step 2 for deep networks is called *backpropagation*

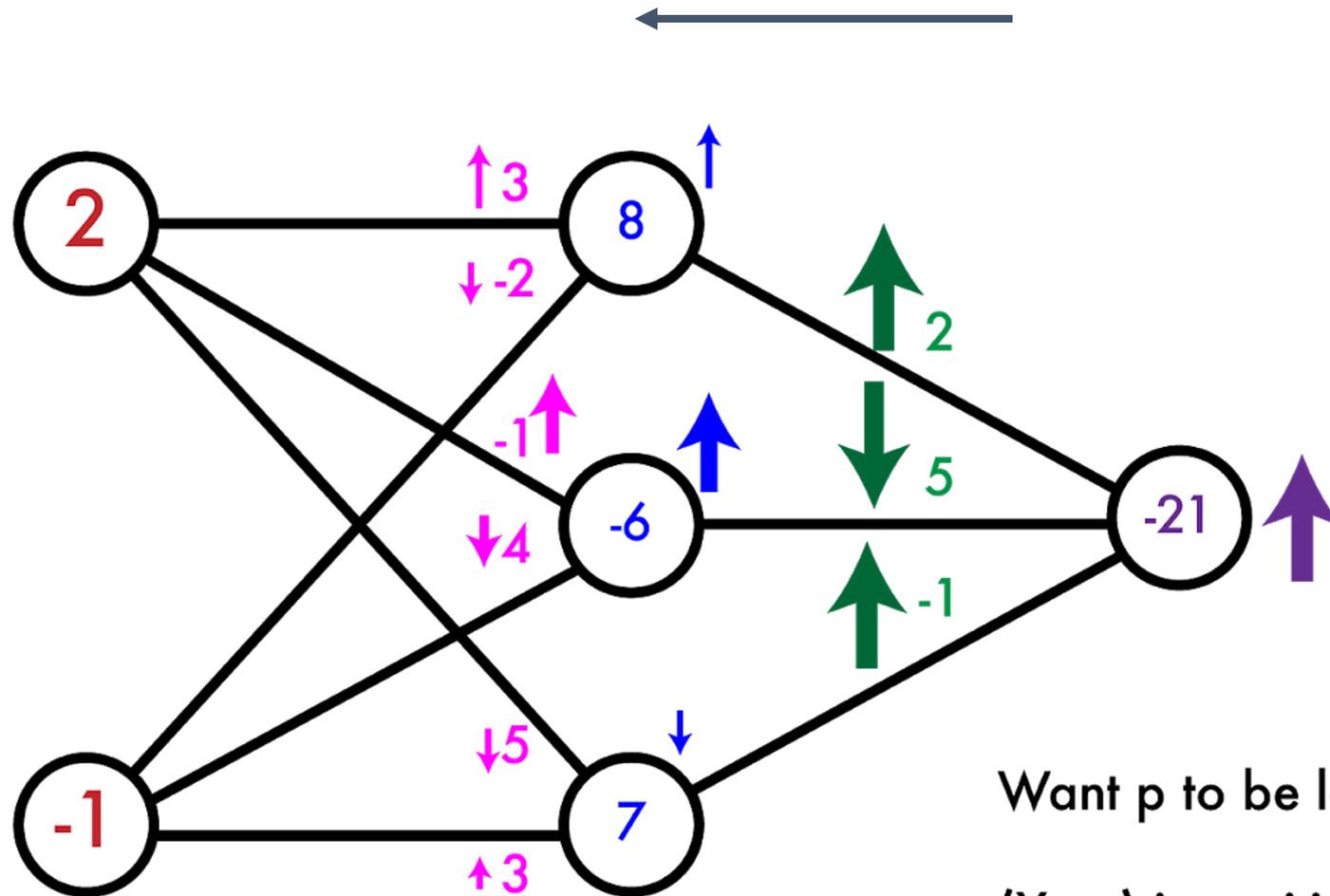
The Learning Cycle: Forward Propagation



Using the current weights of the model calculate the hidden layer neurons, then calculate the output p.



The Learning Cycle: Backward Propagation

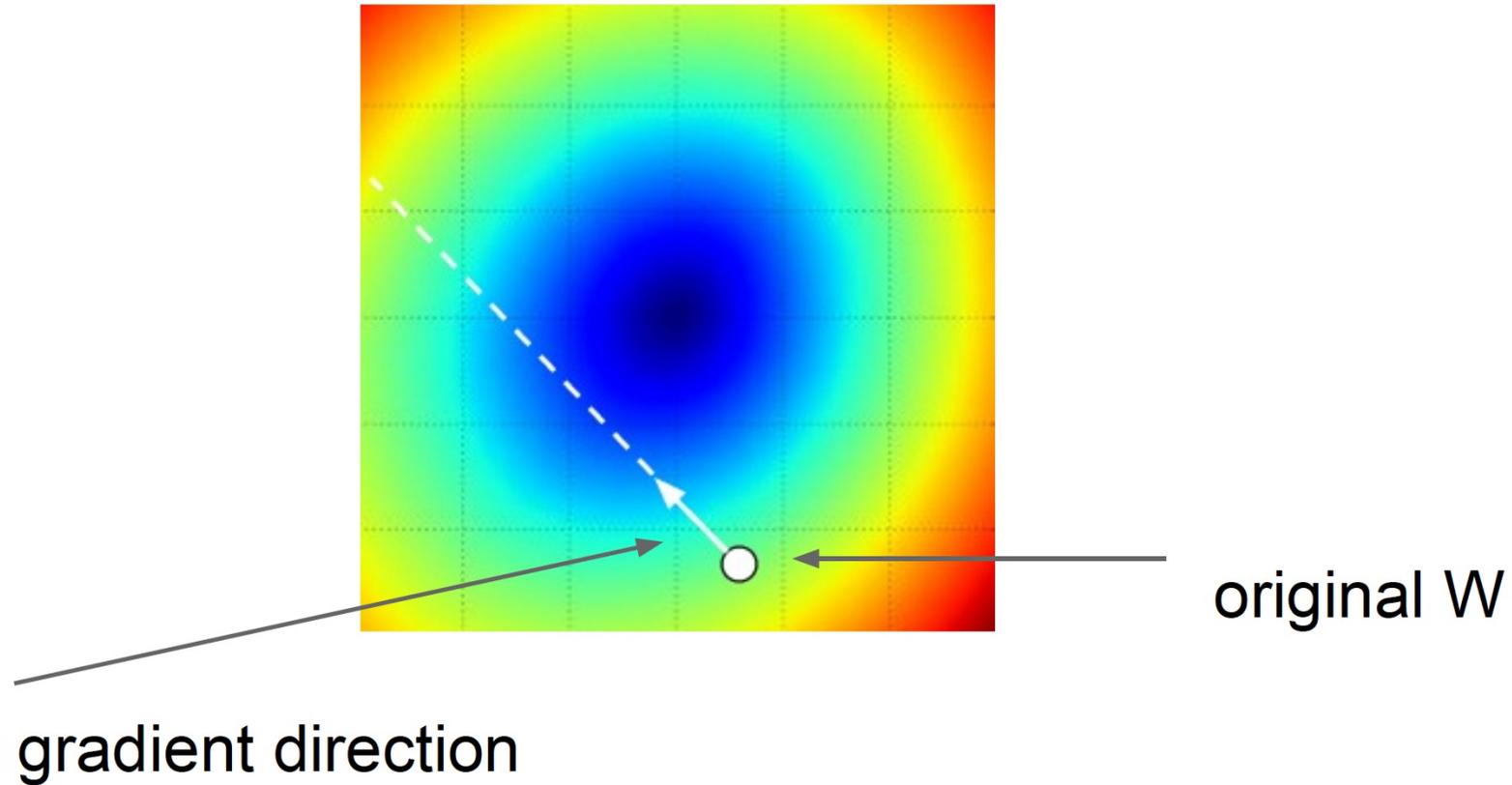


Want p to be larger...

$(Y - p)$ is positive

How do we change our weights?



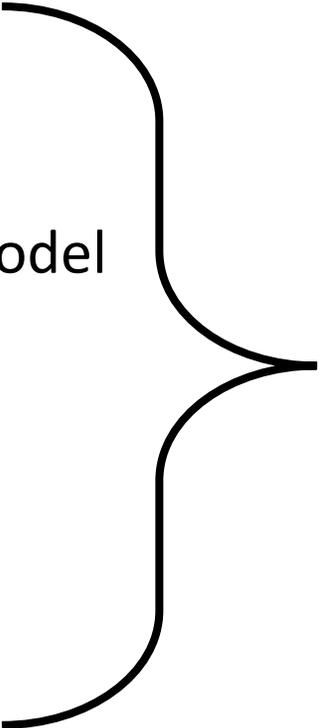


Gradient descent: walk in the direction opposite gradient

- **Q:** How far?
- **A:** Step size: *learning rate*
- Too big: will miss the minimum
- Too small: slow convergence

What Matters in Recognition?

- Data
 - More is always better (as long as it is good data)
 - Annotation is the hard part
- Representation
 - Low level: SIFT, HoG, GIST, edges
 - Mid level: Bag of words, sliding window, deformable model
 - High level: Contextual dependence
 - Deep learned features
- Learning Techniques
 - E.g. choice of classifier or inference method



Combine!

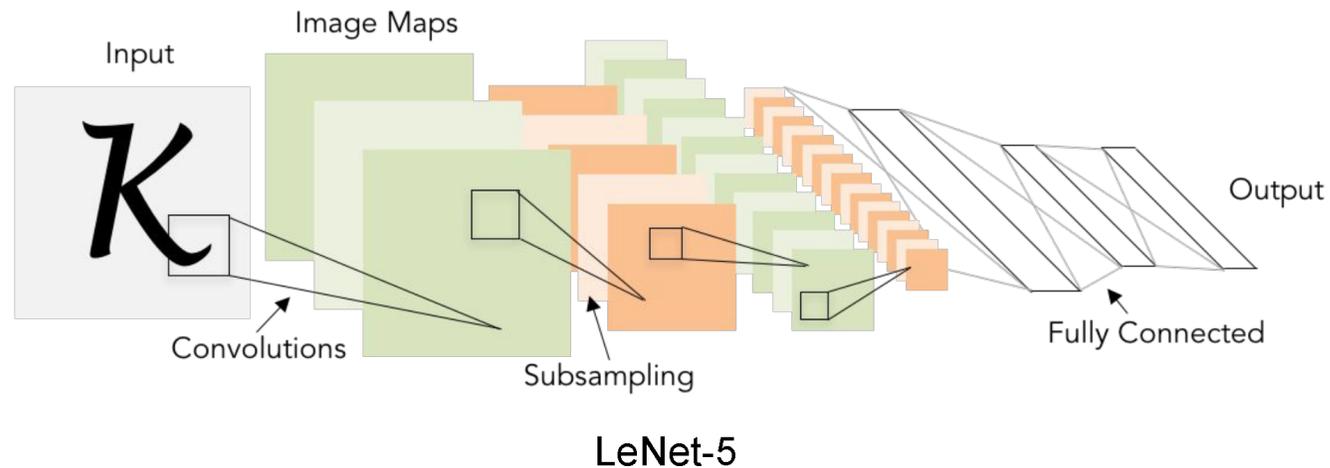
Convolutional Neural Networks (CNNs)

Convolutional neural networks

A bit of history:

Gradient-based learning applied to document recognition

[LeCun, Bottou, Bengio, Haffner 1998]



Fast-forward to today: ConvNets are everywhere

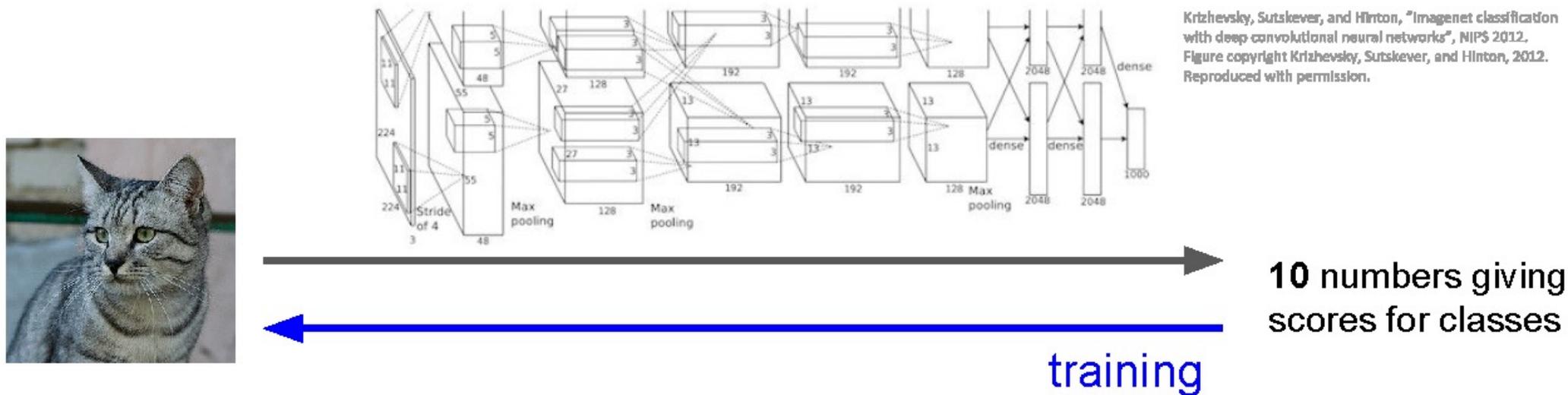
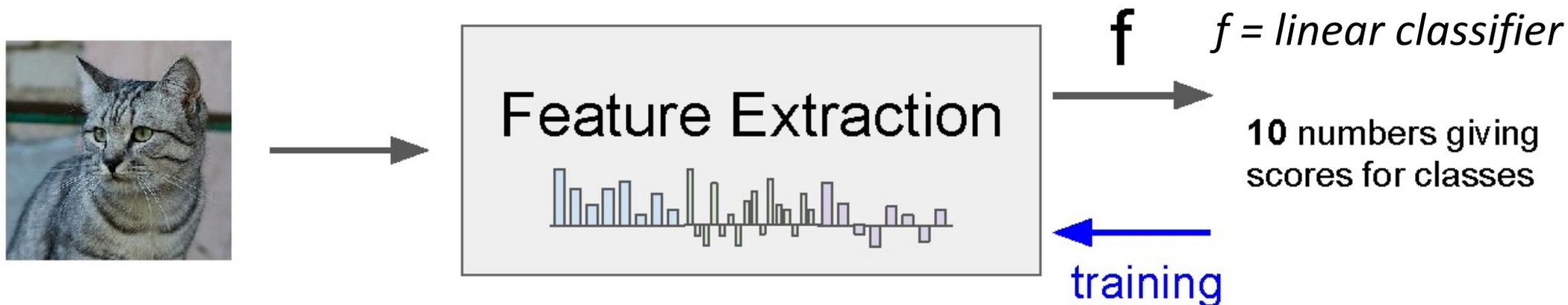
Classification

Retrieval



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

Image features vs ConvNets



Last layer of most CNNs is a linear classifier

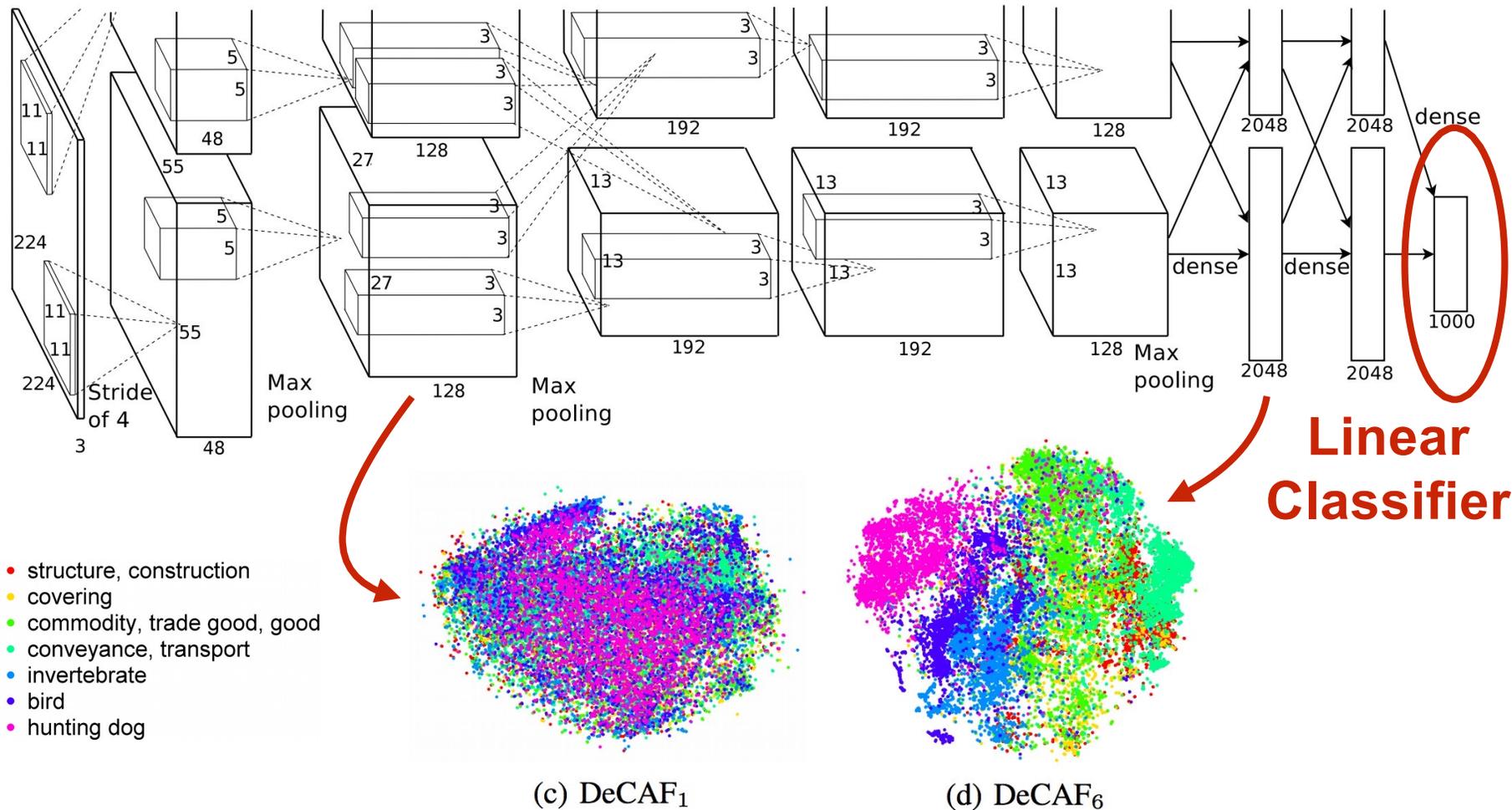


*Input
Pixels*

*Perform everything with a big neural
network, trained end-to-end*

Key: perform enough processing so that by the time you get to the end of the network, the classes are linearly separable

Visualizing AlexNet in 2D with t-SNE



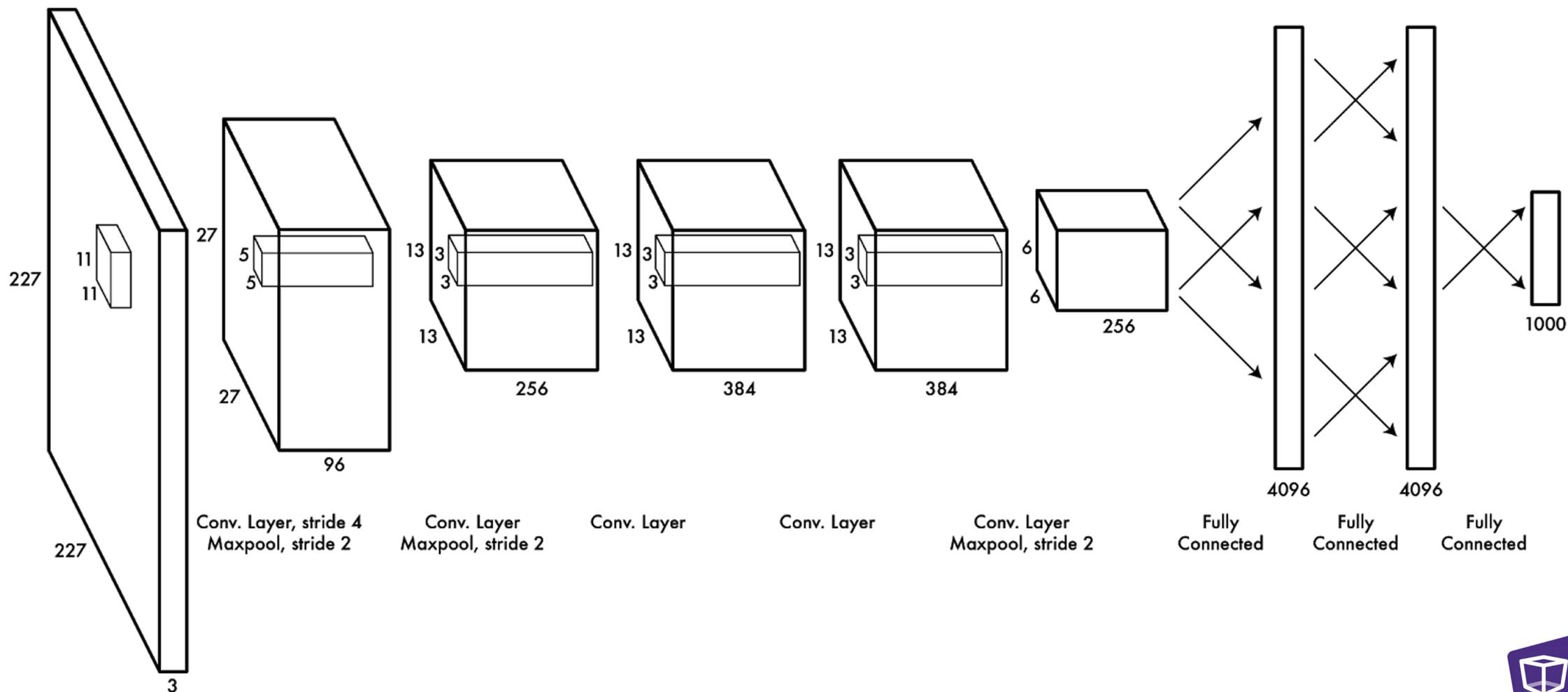
(2D visualization using t-SNE)

[Donahue, "DeCAF: DeCAF: A Deep Convolutional ...", arXiv 2013]

Convolutional neural networks

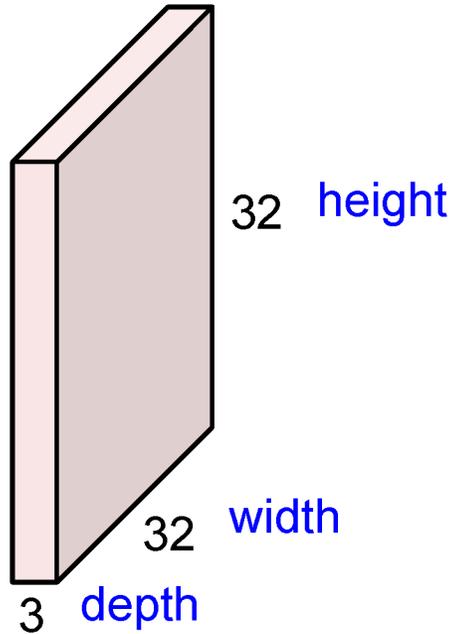
- Layer types:
 - *Convolutional layer*
 - Pooling layer
 - Fully-connected layer

AlexNet: An Early Example



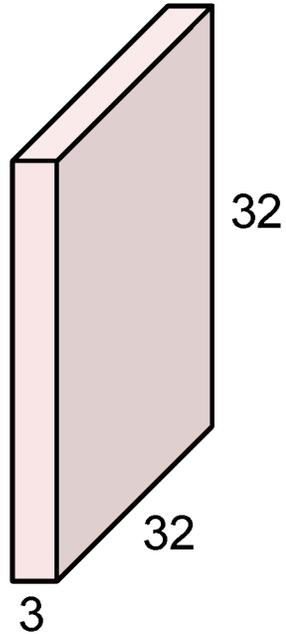
Convolution Layer

32x32x3 image -> preserve spatial structure

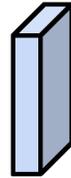


Convolution Layer

32x32x3 image



5x5x3 filter



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

Convolution (Recap)

- Same as cross-correlation, except that the kernel is “flipped” (horizontally and vertically)

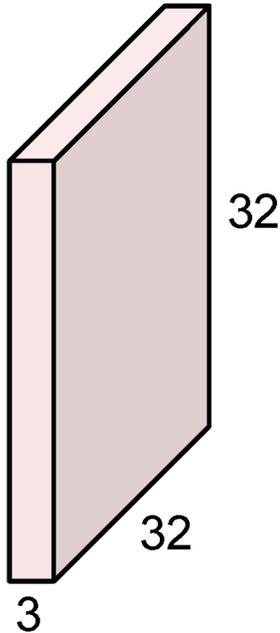
$$G[i, j] = \sum_{u=-k}^k \sum_{v=-k}^k H[u, v] F[i - u, j - v]$$

$$G[i, j] = \sum_{u=-k}^k \sum_{v=-k}^k H[u, v] F[i + u, j + v] \quad \text{Cross-correlation}$$

- Convolution is **commutative** and **associative**

Convolution Layer

32x32x3 image



Filters always extend the full depth of the input volume

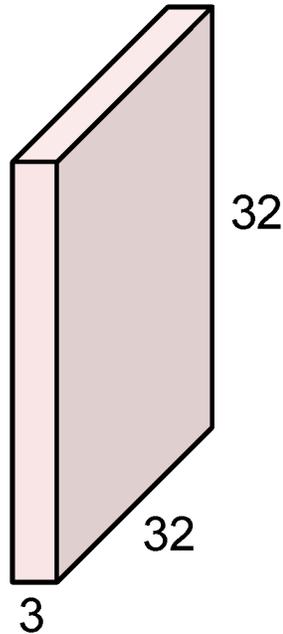
5x5x3 filter



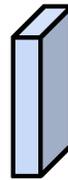
Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”

Convolution Layer

32x32x3 image



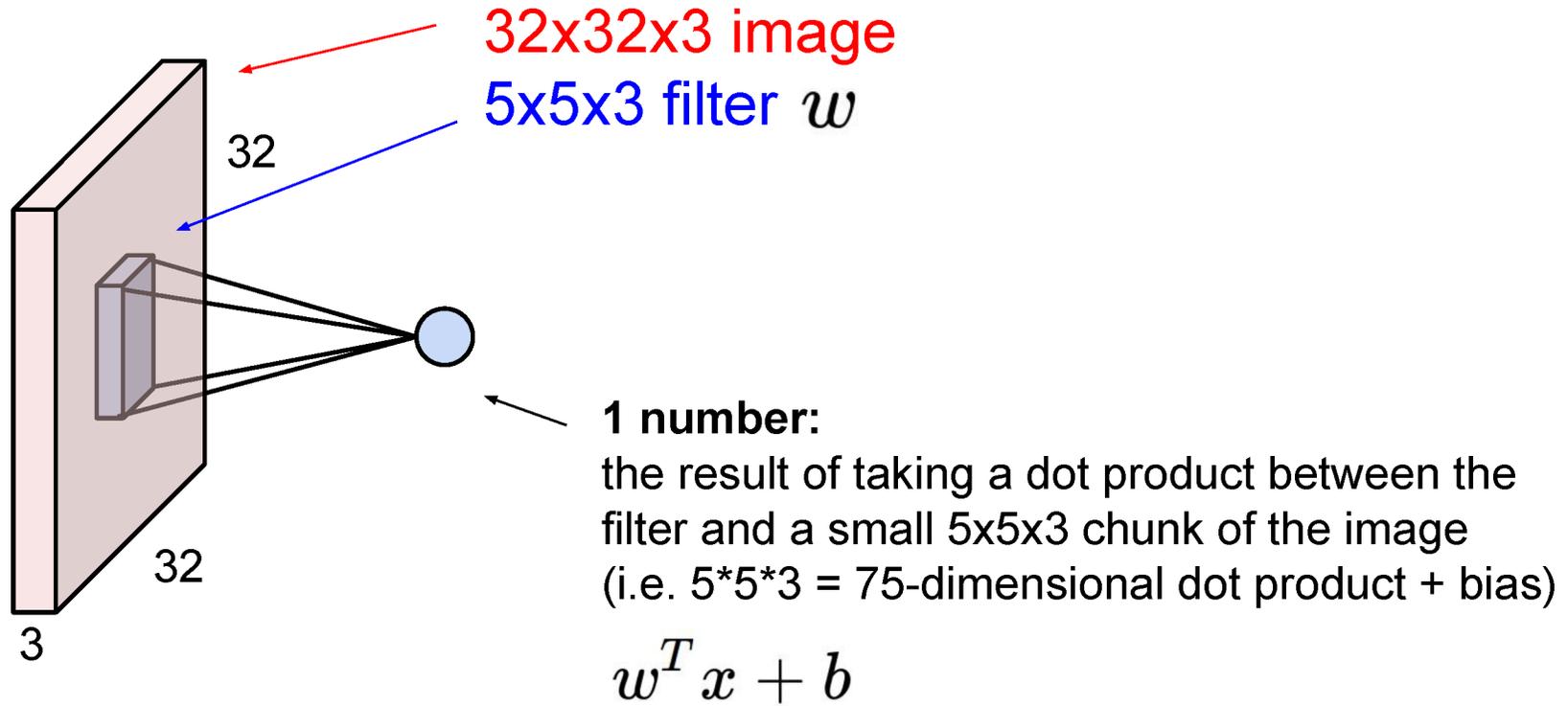
5x5x3 filter



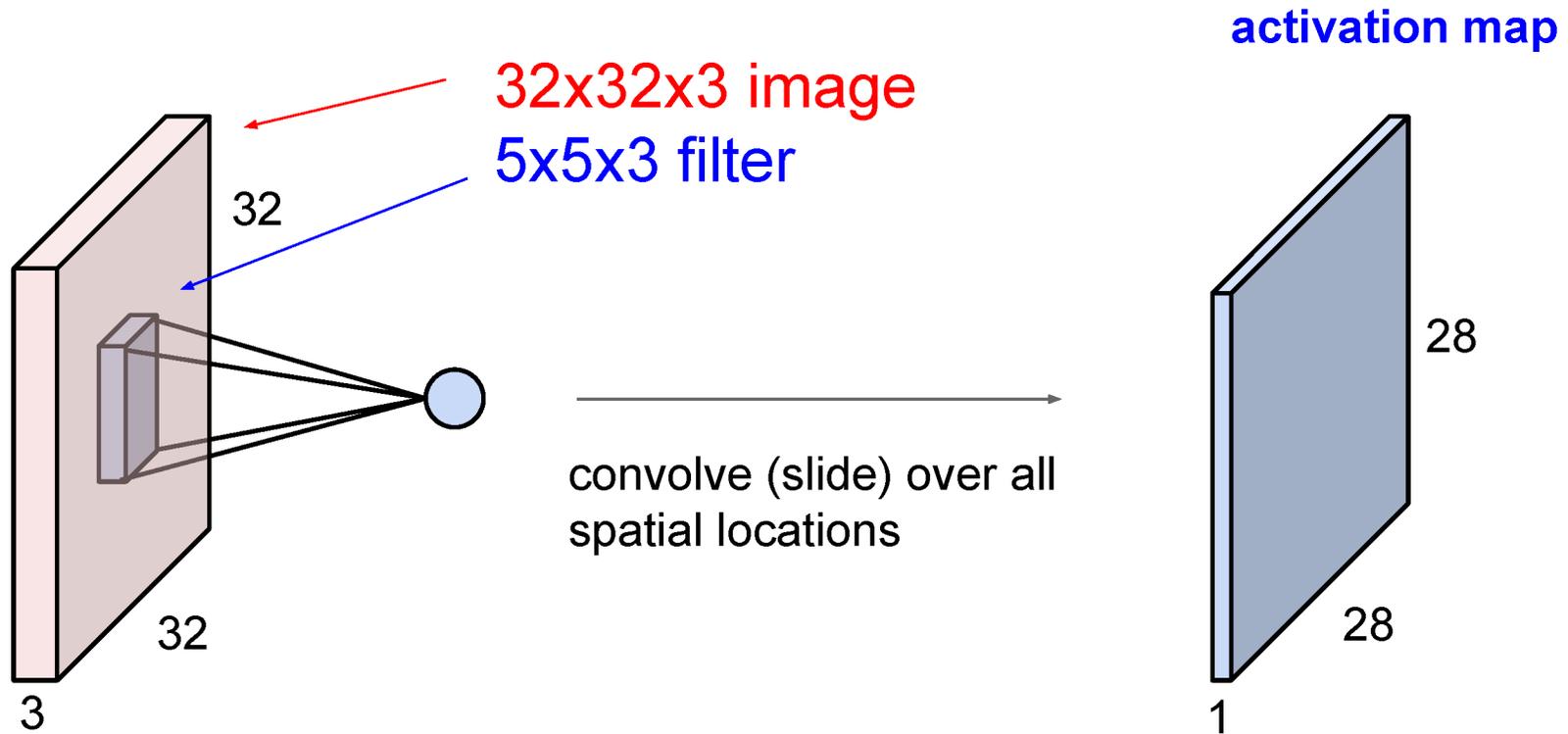
Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”

Number of weights: $5 \times 5 \times 3 + 1 = \mathbf{76}$
(vs. 3072 for a fully-connected layer)
(+1 for bias)

Convolution Layer

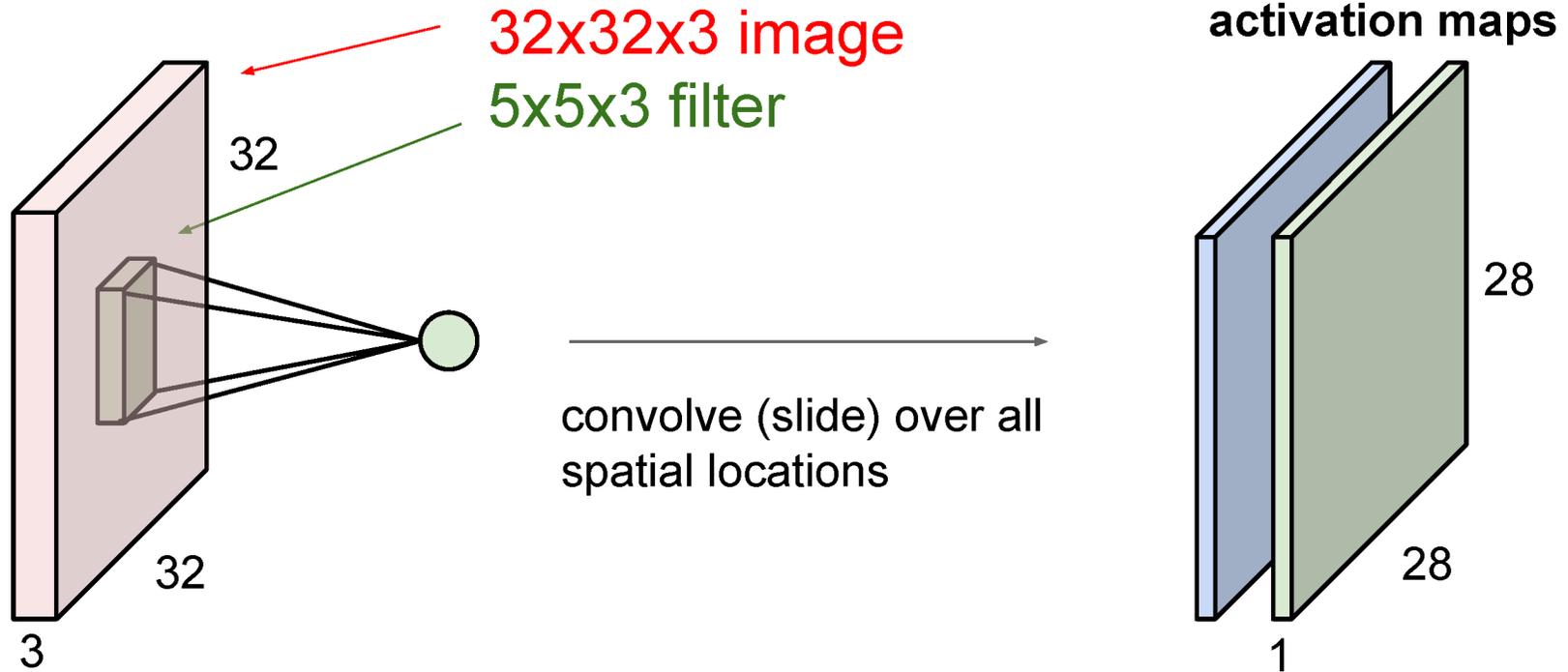


Convolution Layer

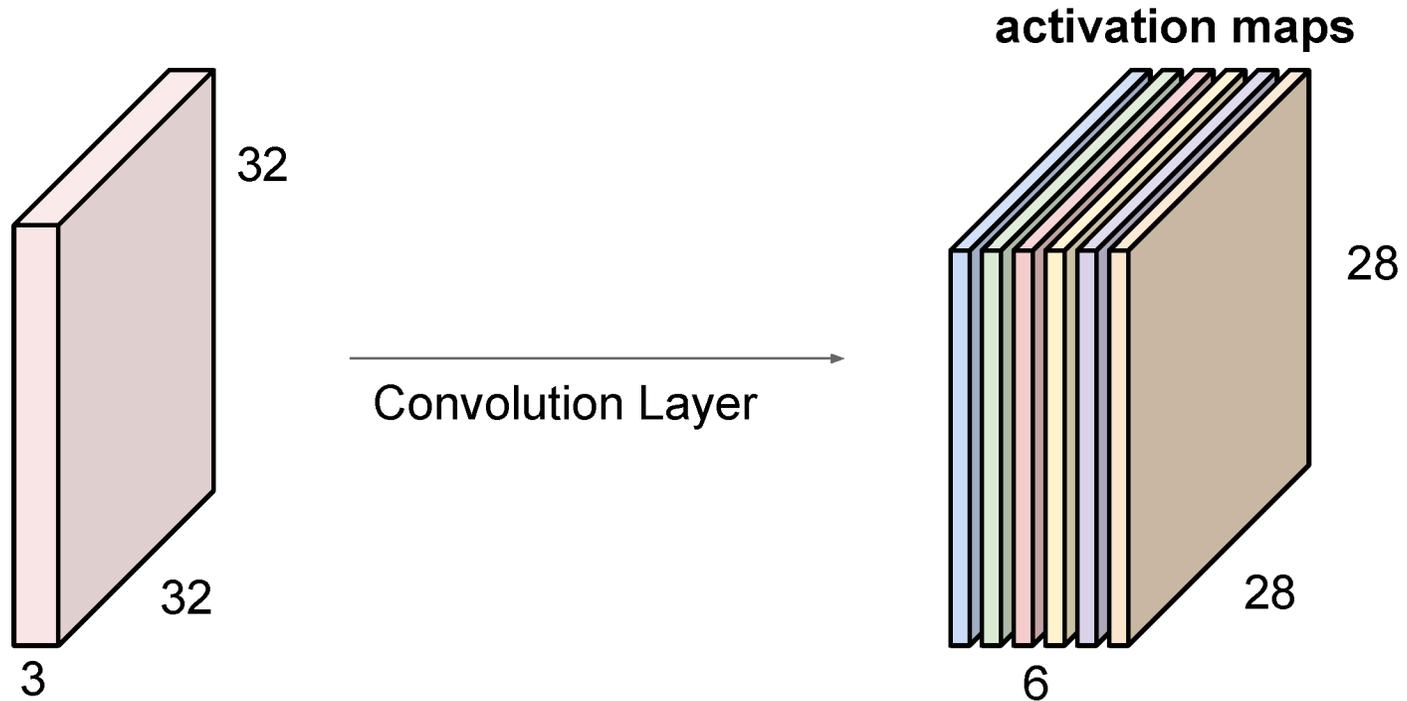


Convolution Layer

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 28x28x6!

(total number of parameters: $6 \times (75 + 1) = 456$)

Padding & Stride in CNN

```
CLASS torch.nn.Conv2d(in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1, groups=1, bias=True, padding_mode='zeros', device=None, dtype=None) \[SOURCE\]
```

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

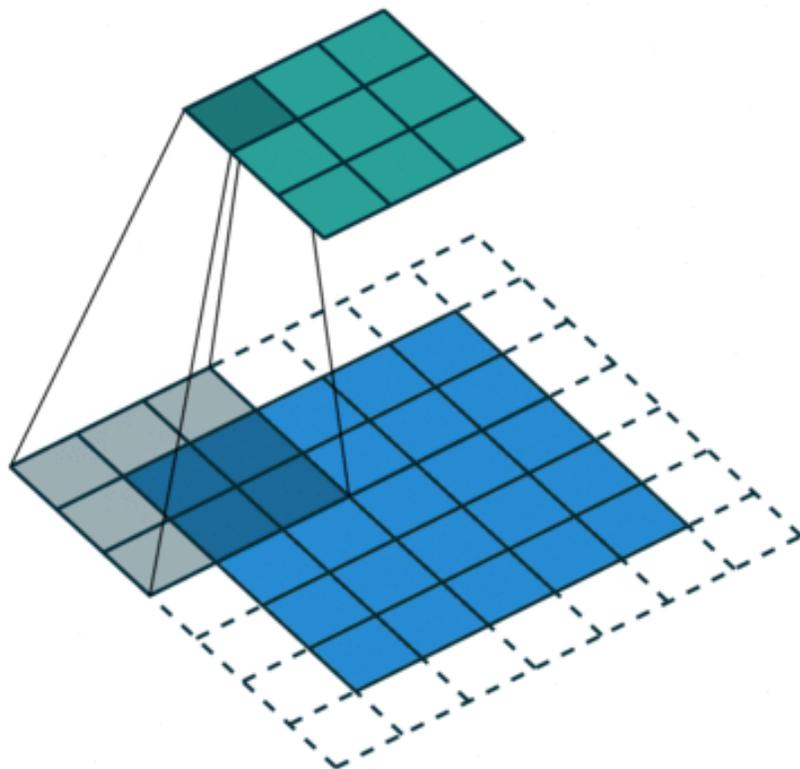
In the simplest case, the output value of the layer with input size (N, C_{in}, H, W) and output $(N, C_{out}, H_{out}, W_{out})$ can be precisely described as:

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

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.

$$H_{out} = \left\lfloor \frac{H_{in} + 2 \times \text{padding}[0] - \text{dilation}[0] \times (\text{kernel_size}[0] - 1) - 1}{\text{stride}[0]} + 1 \right\rfloor$$

$$W_{out} = \left\lfloor \frac{W_{in} + 2 \times \text{padding}[1] - \text{dilation}[1] \times (\text{kernel_size}[1] - 1) - 1}{\text{stride}[1]} + 1 \right\rfloor$$

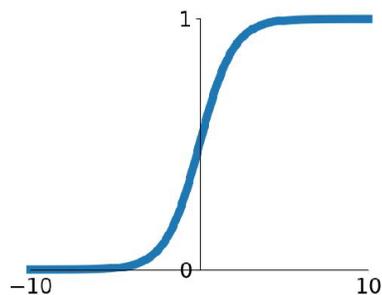


padding=1, stride=2

Activation functions

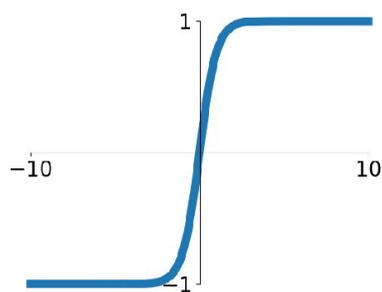
Sigmoid

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



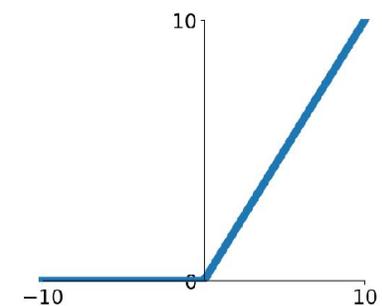
tanh

$$\tanh(x)$$



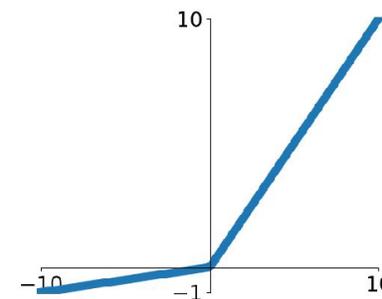
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$

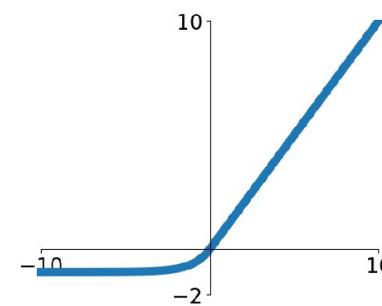


Maxout

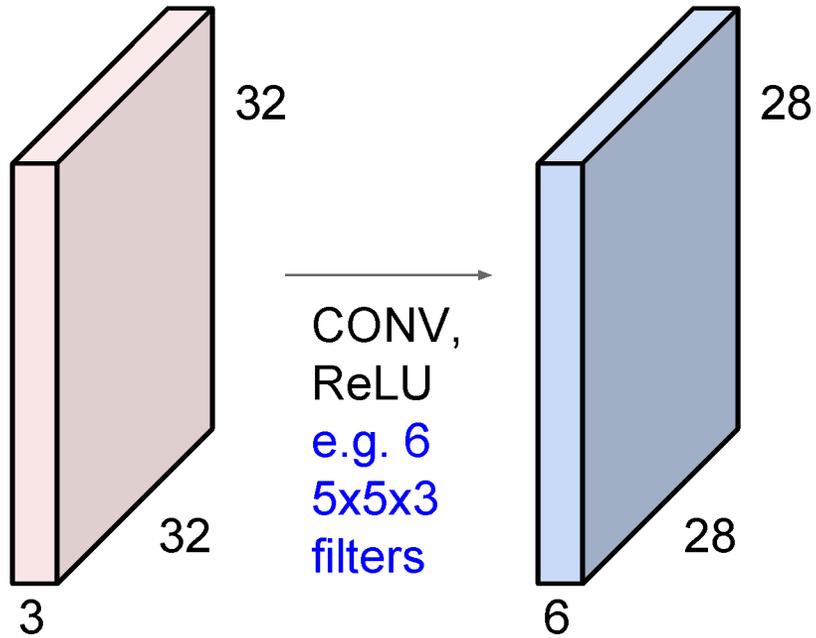
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ELU

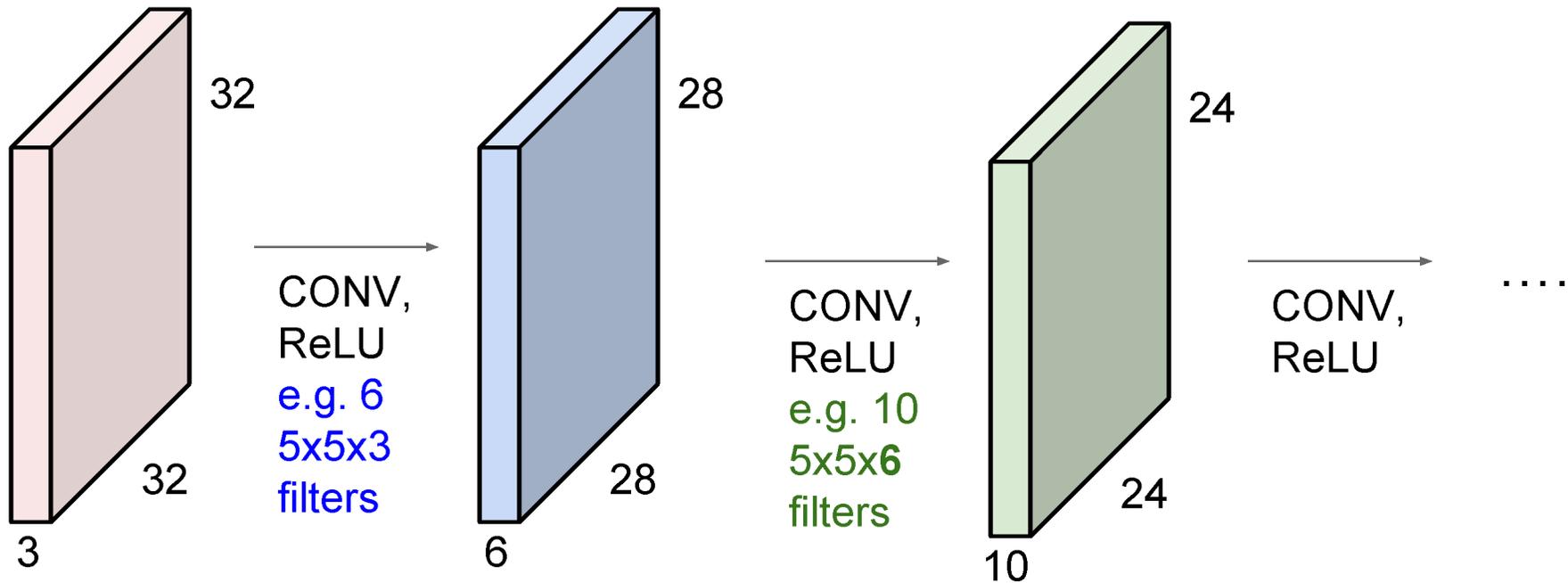
$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



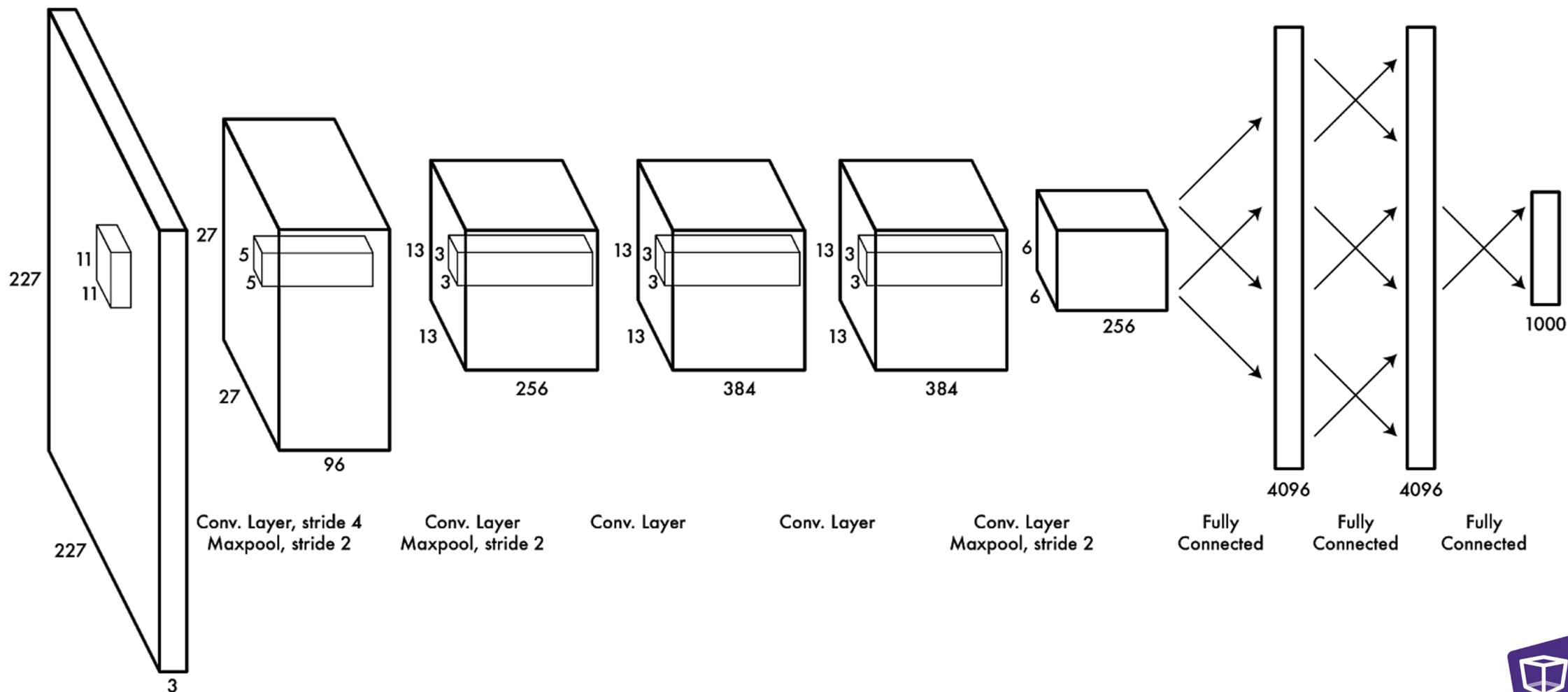
Preview: ConvNet is a sequence of Convolution Layers, interspersed with activation functions



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



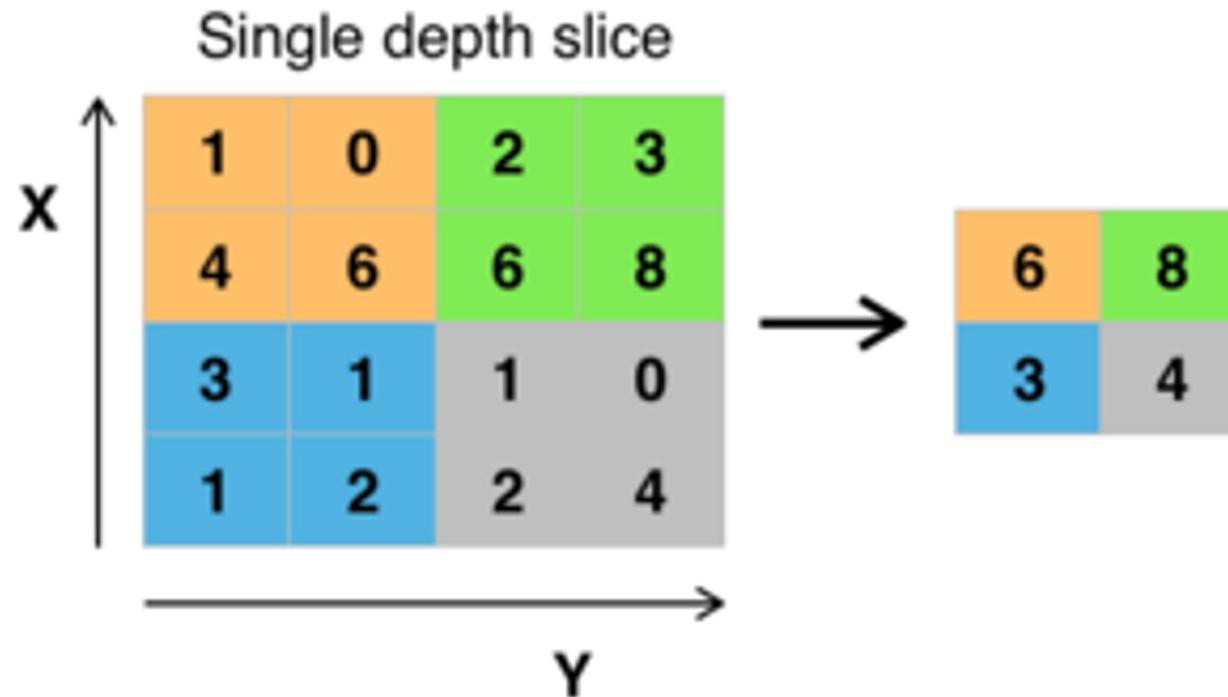
AlexNet: An Early Example



How Else to Shrink the Model Size?

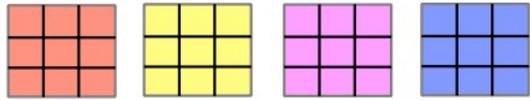
Pooling Layer:

- Max Pooling
- Other pooling options like average pooling are also used



Convolutional Networks

Learnable 3x3 Convolutional Kernels



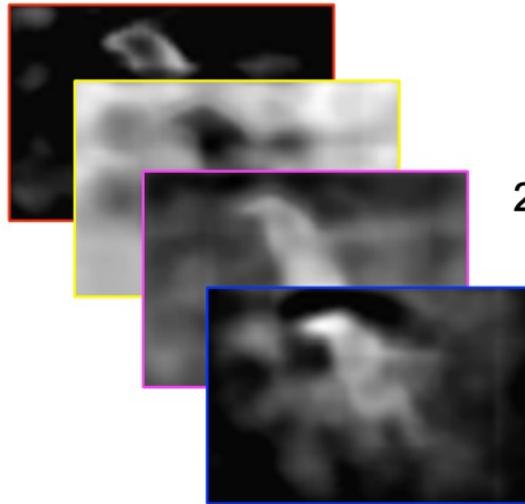
Input Image (Grayscale)



80 x 120 x 1

2D Conv.

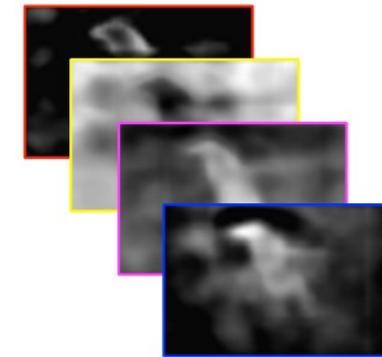
Conv. Feature Maps



80 x 120 x 4

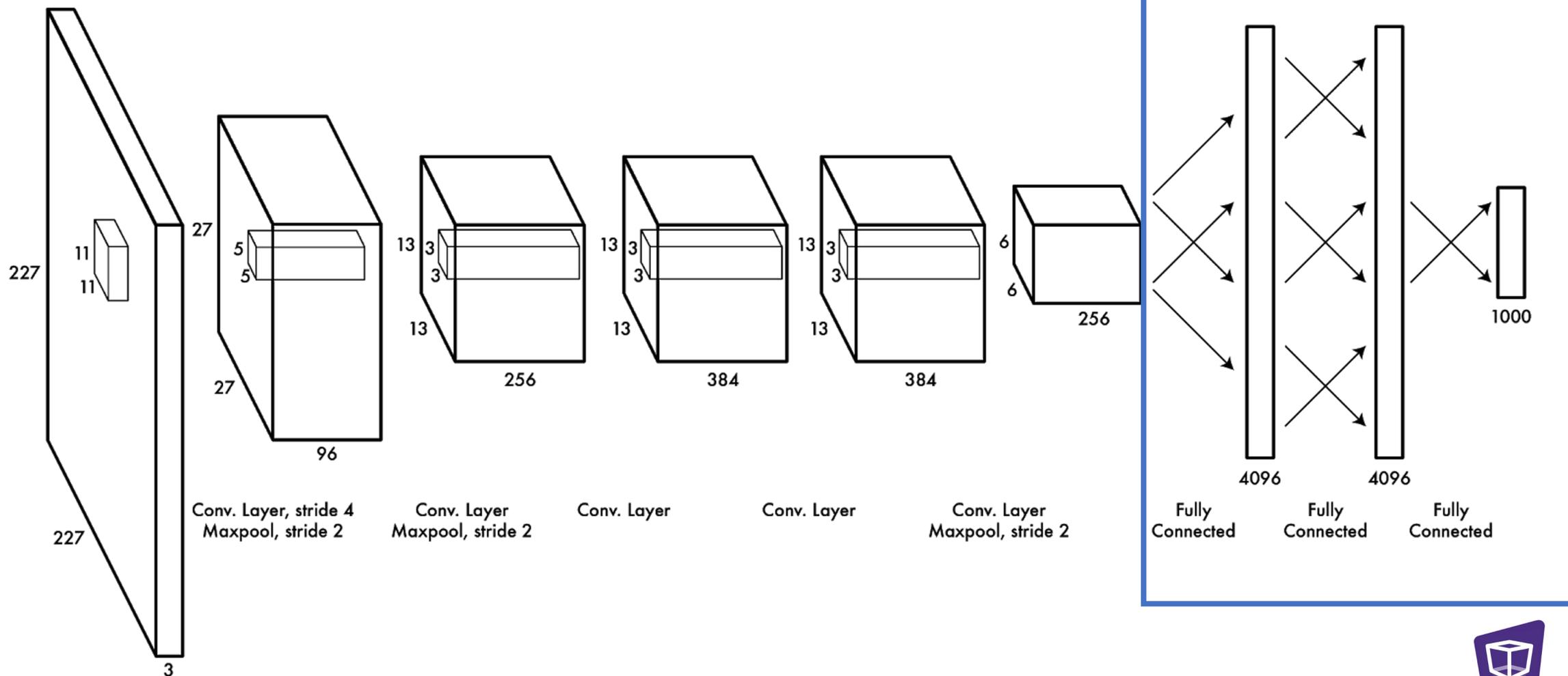
2D Max Pooling

Pooled Feature Maps



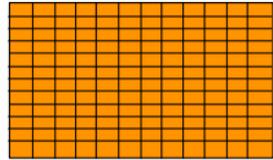
40 x 60 x 4

AlexNet: An Early Example



Convolutional Networks

Learnable FC Layer Weight Matrix

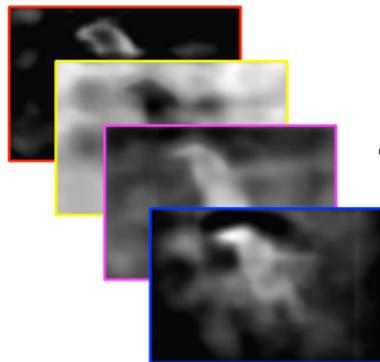


$$W \in \mathbb{R}^{d \times C}$$

d - feature dimensionality (4800 in this example)

C - number of classes

Pooled Feature Maps



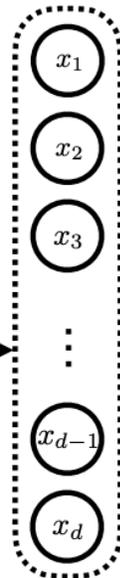
40 x 60 x 4

2D Conv.



40 x 60 x 2

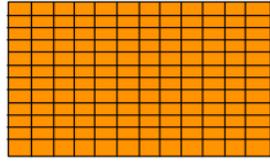
Flatten



1 x 4800

Convolutional Networks

Learnable FC Layer Weight Matrix

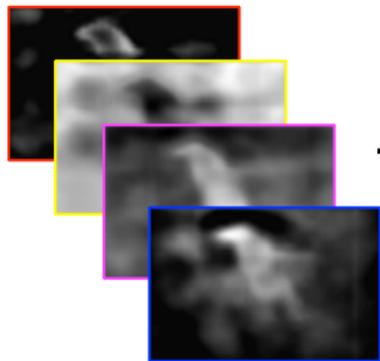


$$W \in \mathbb{R}^{d \times C}$$

d - feature dimensionality (4800 in this example)

C - number of classes

Pooled Feature Maps



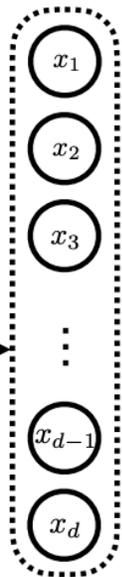
40 x 60 x 4

2D Conv.



40 x 60 x 2

Flatten



1 x 4800

FC Layer + Softmax



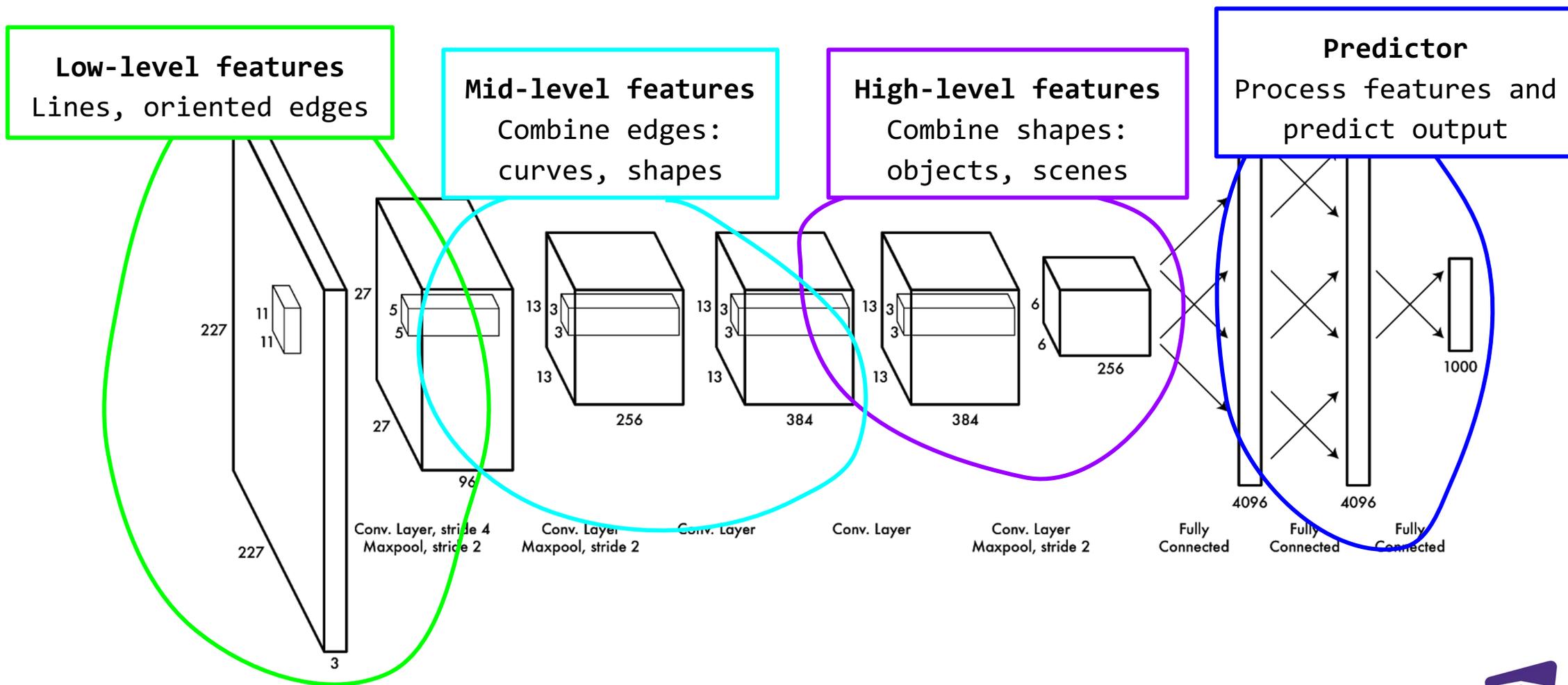
1 x C

Cat: 0.01

Dog: 0.03

Penguin: 0.91

Where Models Learn Features of an Image



Preview

[Zeiler and Fergus 2013]

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

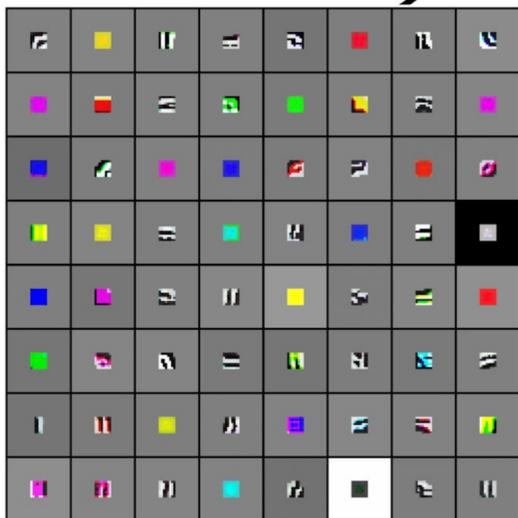


Low-level features

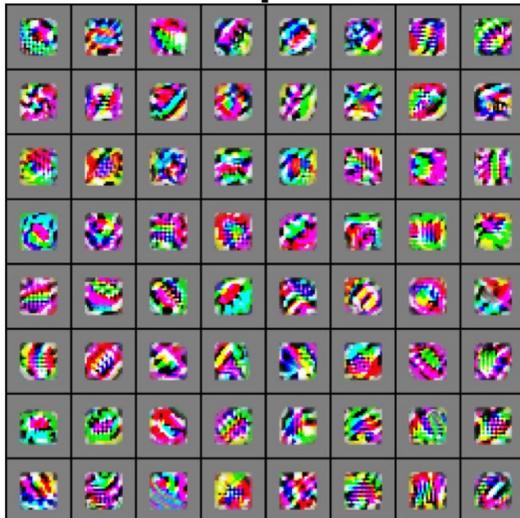
Mid-level features

High-level features

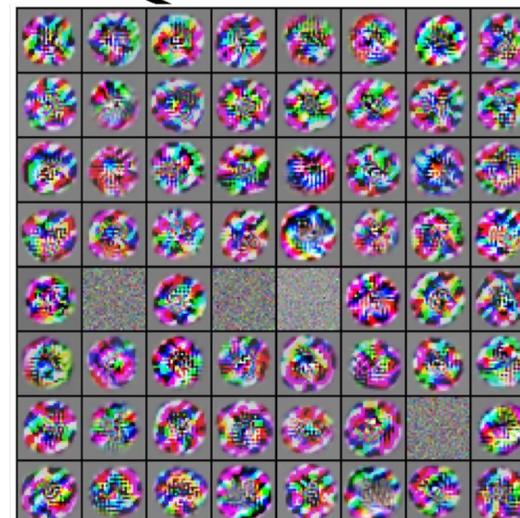
Linearly separable classifier



VGG-16 Conv1_1



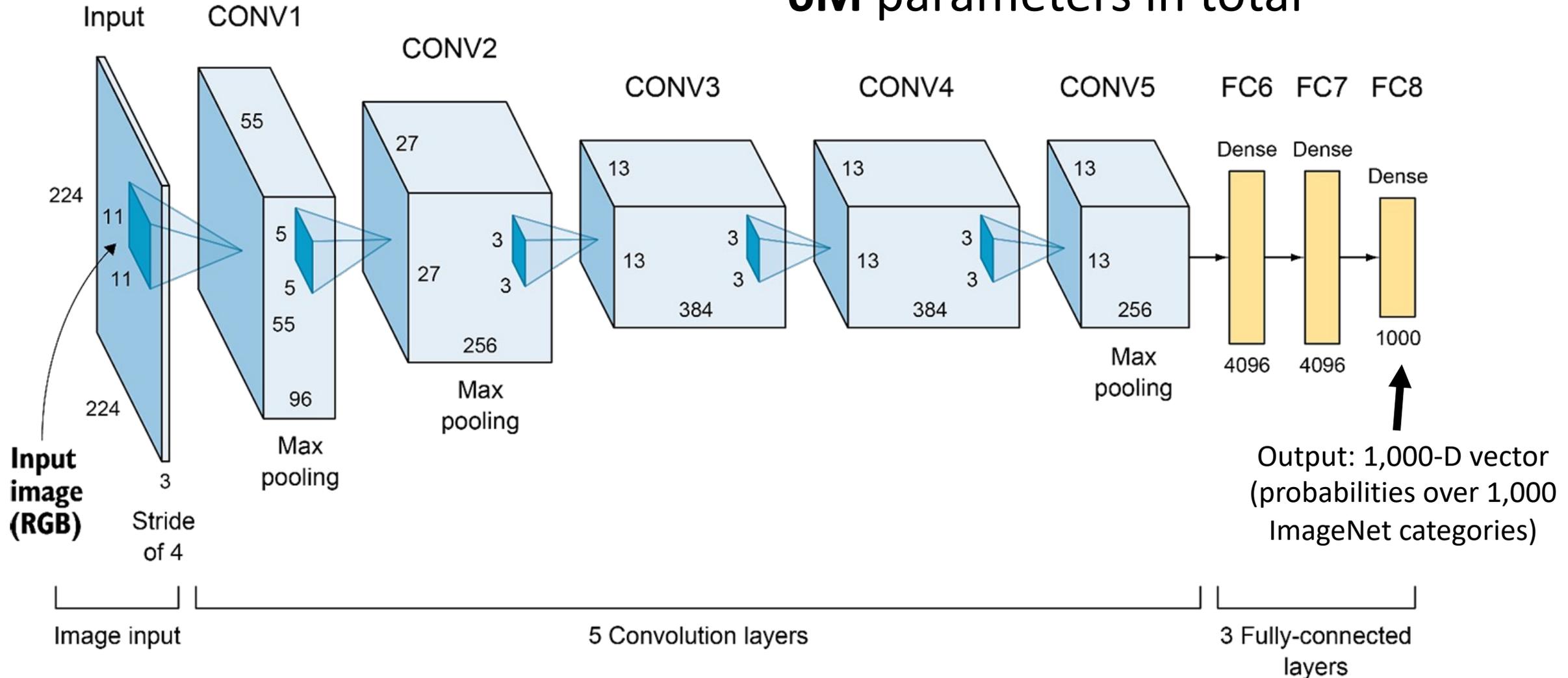
VGG-16 Conv3_2



VGG-16 Conv5_3

AlexNet (2012)

6M parameters in total



Big picture

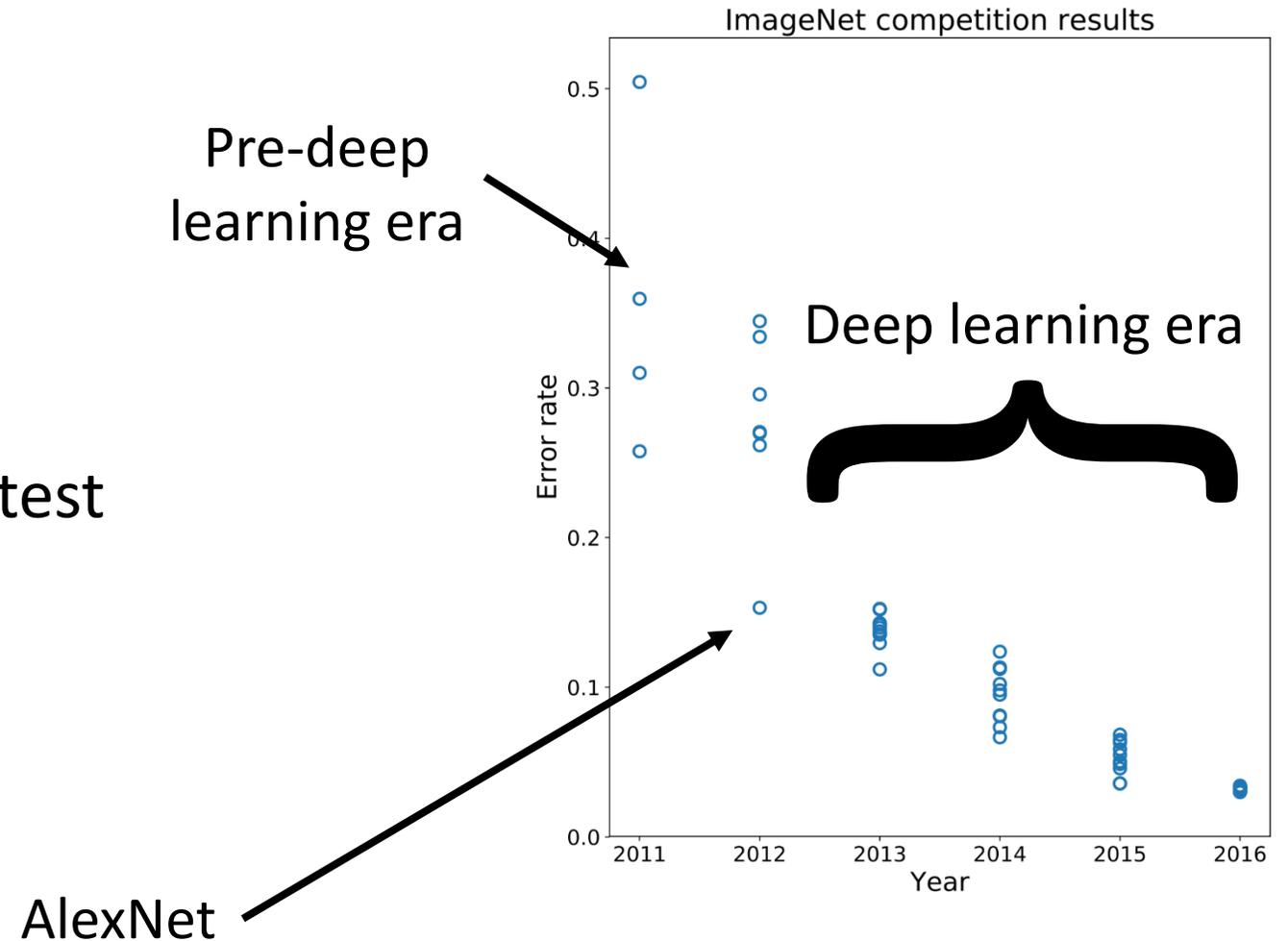
- A convolutional neural network can be thought of as a function from images to class scores
 - With millions of adjustable weights...
 - ... leading to a very non-linear mapping from images to features / class scores.
 - We will set these weights based on classification accuracy on training data...
 - ... and hopefully our network will generalize to new images at test time

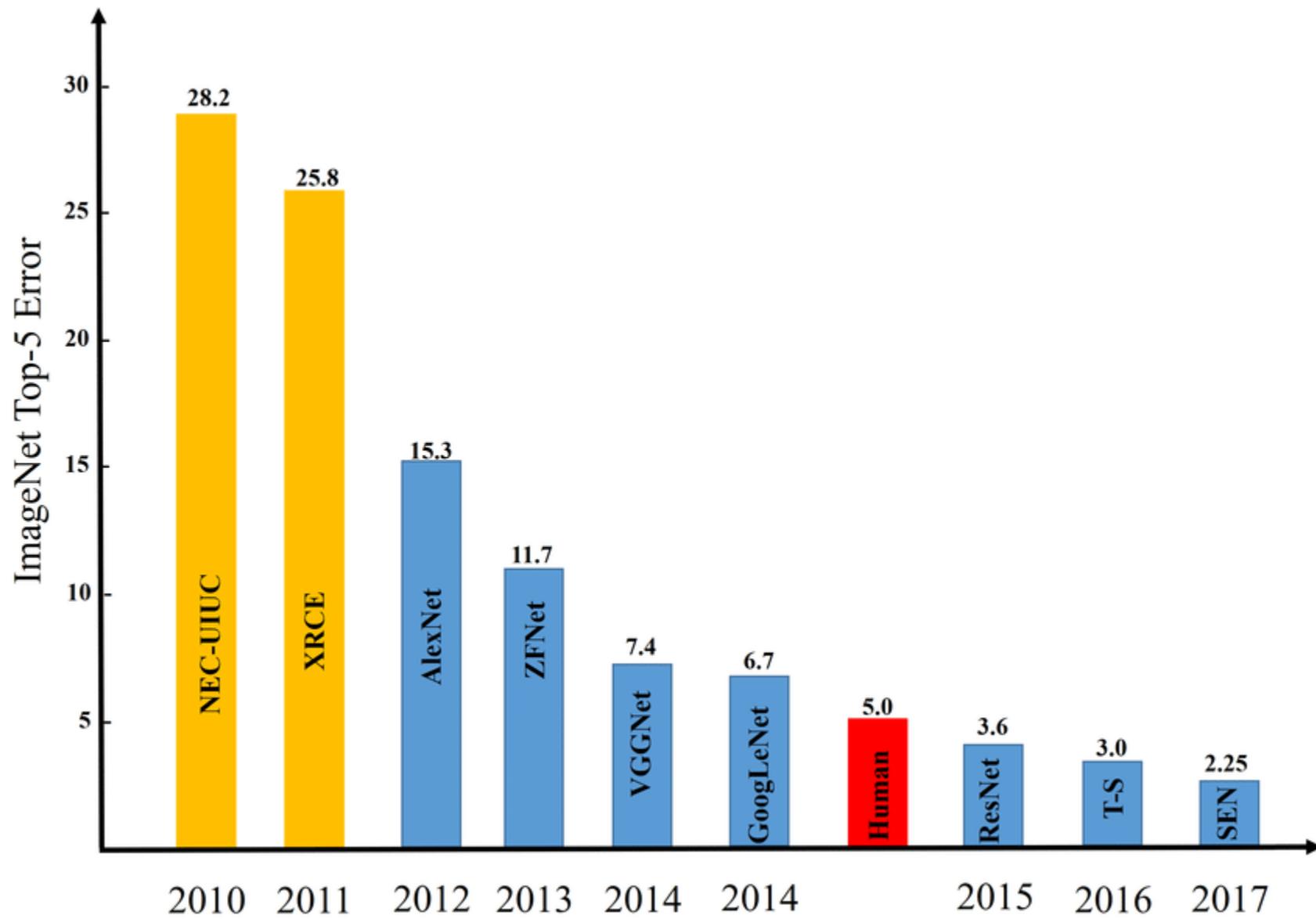
Data is key—enter ImageNet

- ImageNet (and the ImageNet Large-Scale Visual Recognition Challenge, aka **ILSVRC**) has been key to training deep learning methods
 - J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li and L. Fei-Fei, ImageNet: A Large-Scale Hierarchical Image Database. *CVPR*, 2009.
- **ILSVRC**: 1,000 object categories, each with ~700-1300 training images. Test set has 100 images per categories (100,000 total).
- Standard ILSVRC error metric: top-5 error
 - if the correct answer for a given test image is in the top 5 categories, your answer is judged to be correct

Performance improvements on ILSVRC

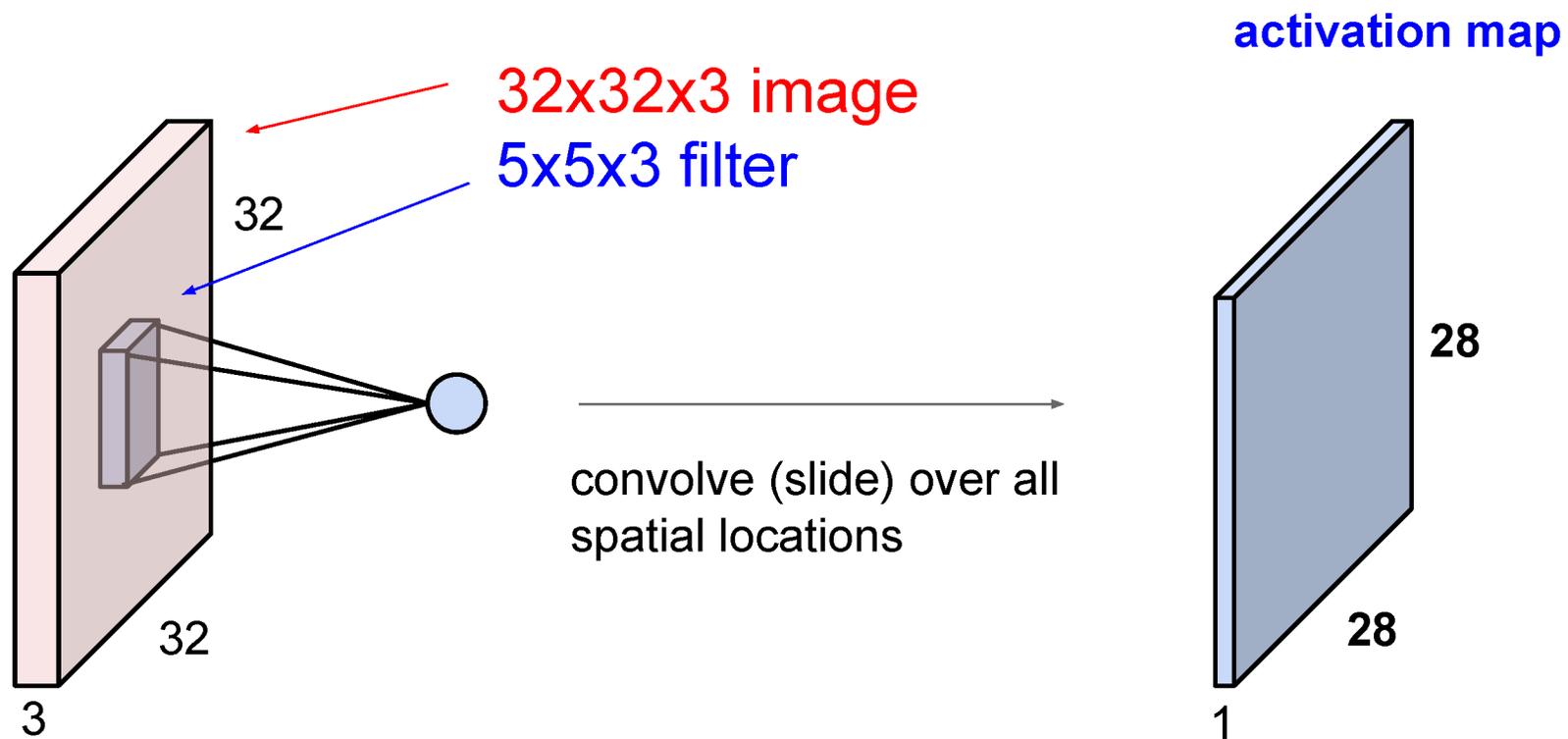
- ImageNet Large-Scale Visual Recognition Challenge
- Held from 2011-2017
- 1000 categories, 1000 training images per category
- Test performance on held-out test set of images





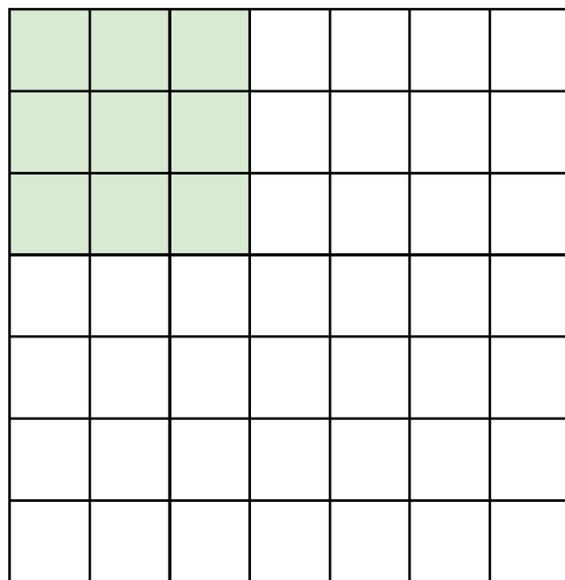
Closer look at Convolution (Extra Slides)

A closer look at spatial dimensions:



A closer look at spatial dimensions:

7

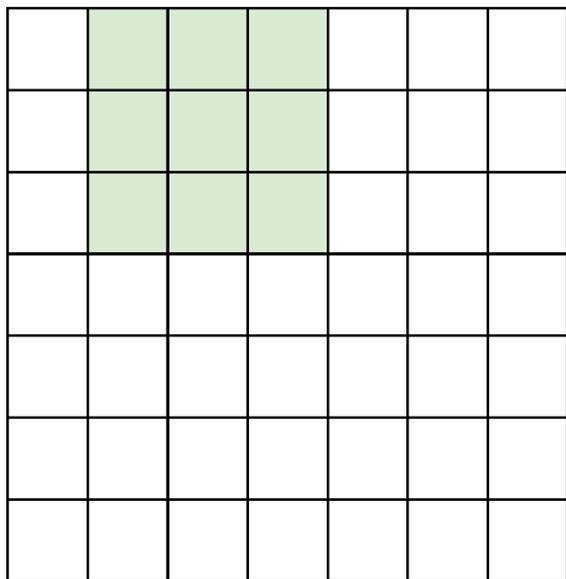


7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

7

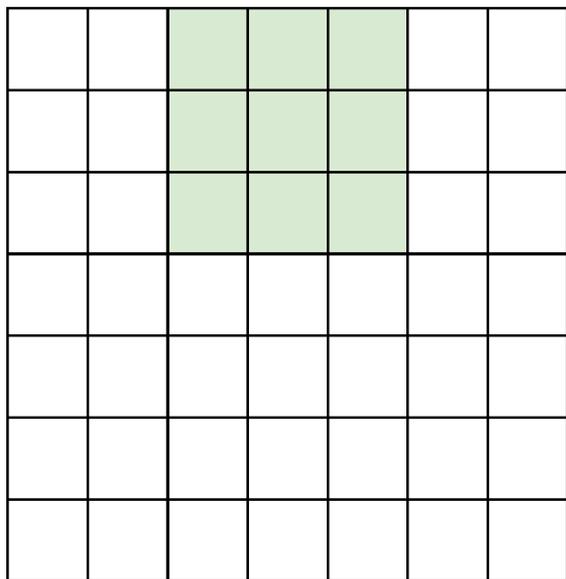


7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

7

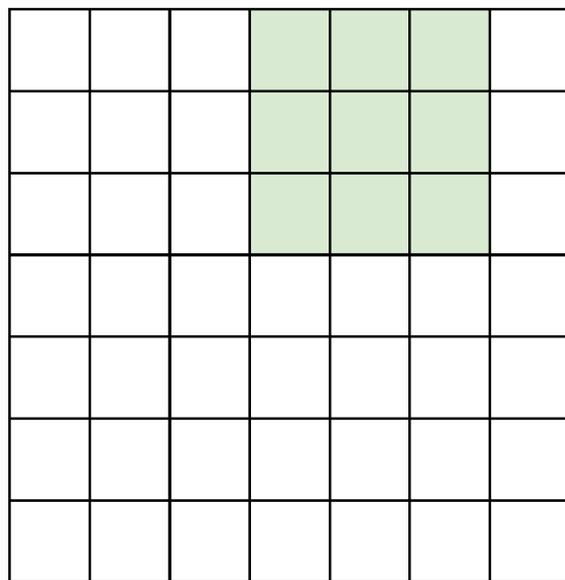


7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

7

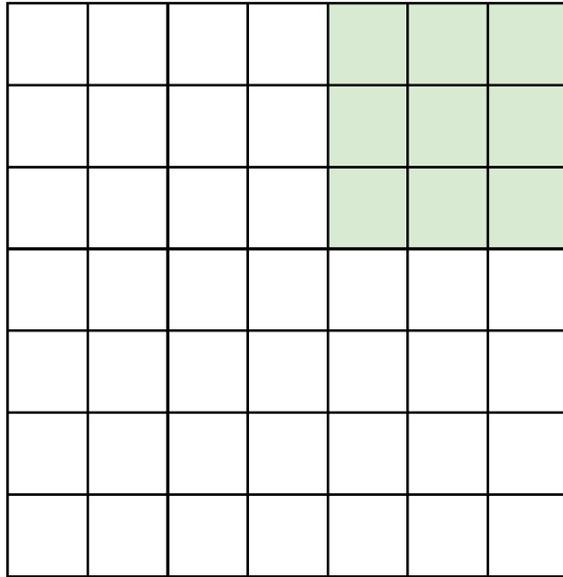


7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

7

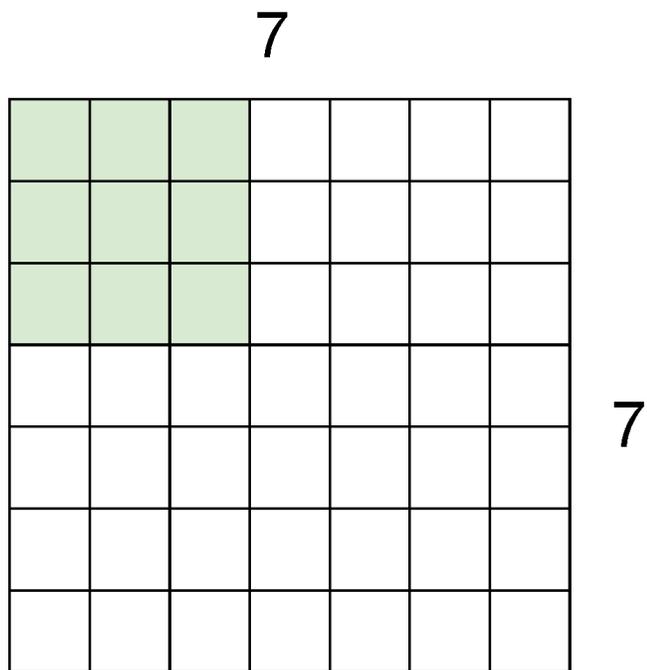


7

7x7 input (spatially)
assume 3x3 filter

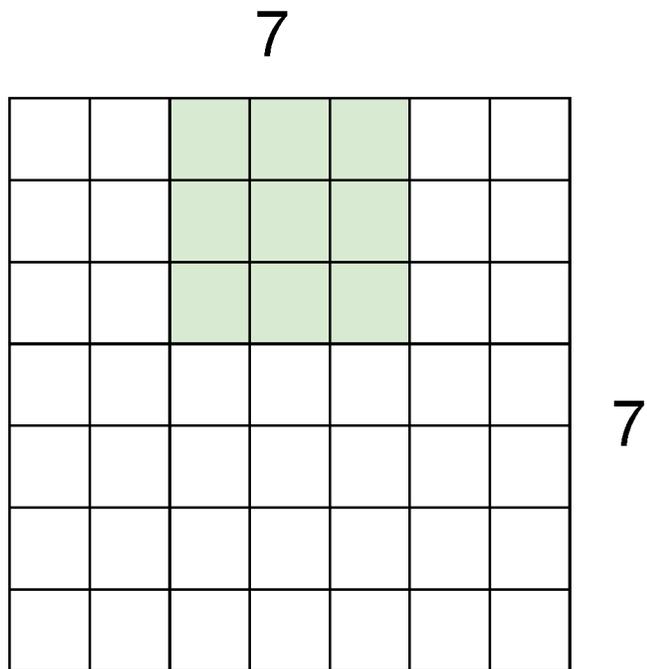
=> 5x5 output

A closer look at spatial dimensions:



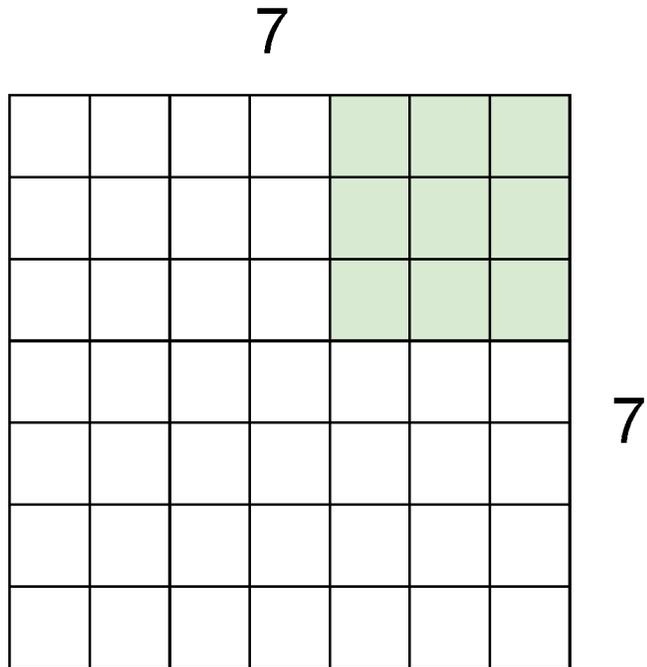
7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**

A closer look at spatial dimensions:



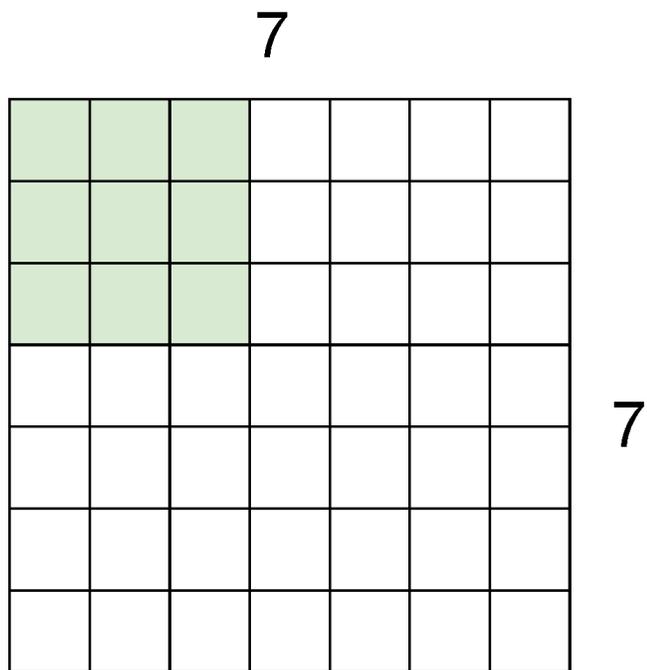
7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**

A closer look at spatial dimensions:



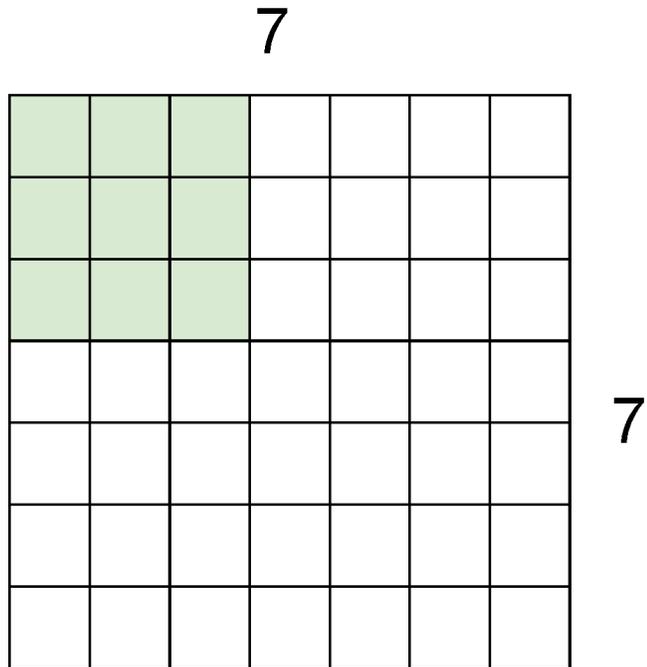
7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**
=> 3x3 output!

A closer look at spatial dimensions:



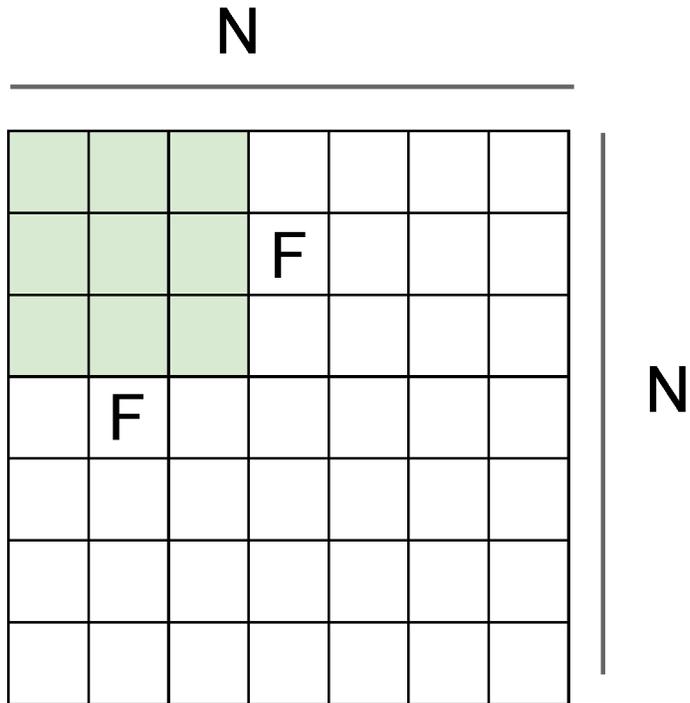
7x7 input (spatially)
assume 3x3 filter
applied **with stride 3?**

A closer look at spatial dimensions:



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

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

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

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

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

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:)

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

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!

Output filter size: $(N + 2 * \text{pad} - F) / S + 1$

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 $F \times F$, and zero-padding with $(F-1)/2$. (will preserve size spatially)

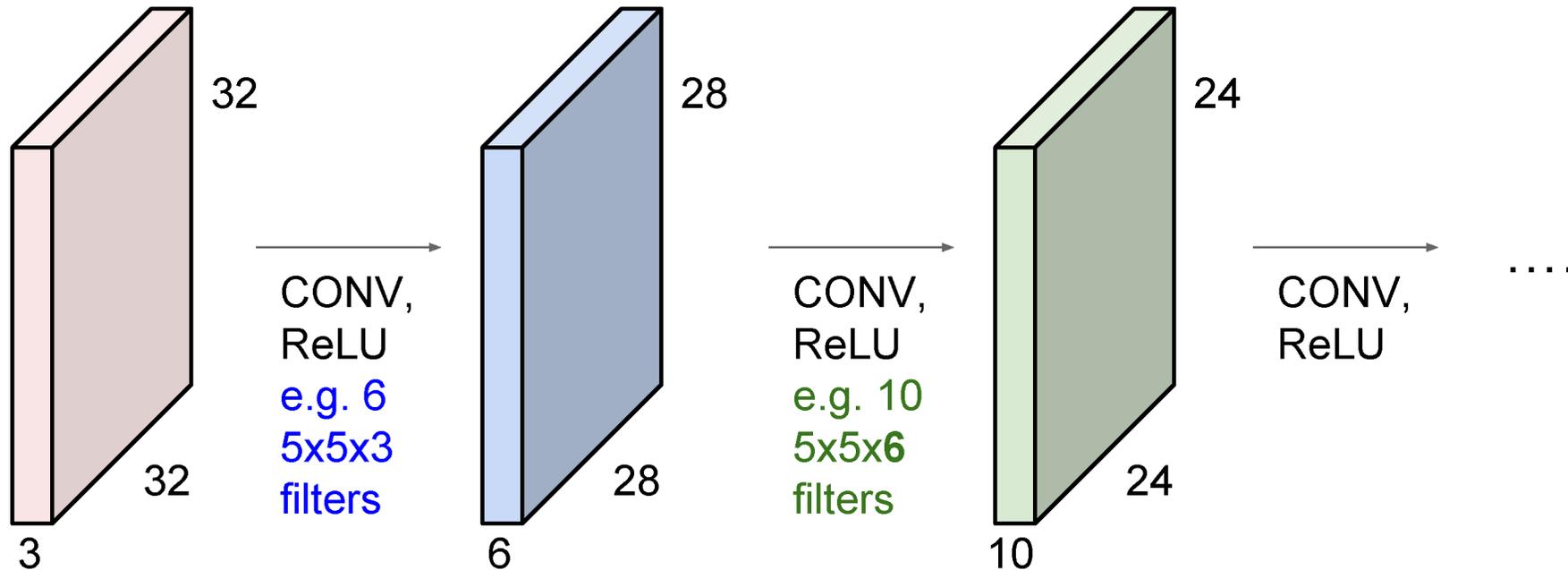
e.g. $F = 3 \Rightarrow$ zero pad with 1

$F = 5 \Rightarrow$ 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.

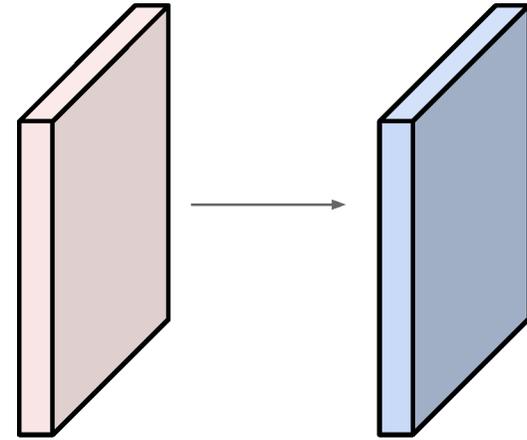


Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2

Output volume size: ?



Examples time:

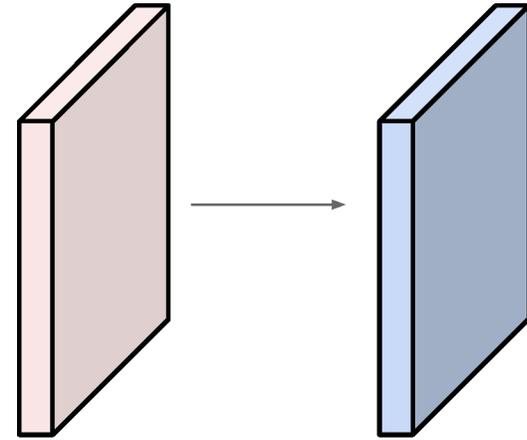
Input volume: **32x32x3**

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

Output volume size:

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

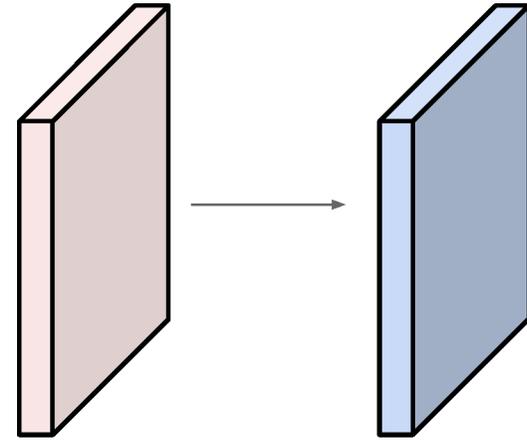
32x32x10



Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2

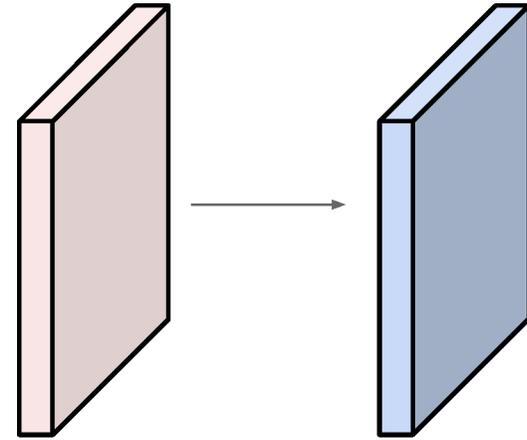


Number of parameters in this layer?

Examples time:

Input volume: **32x32x3**

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

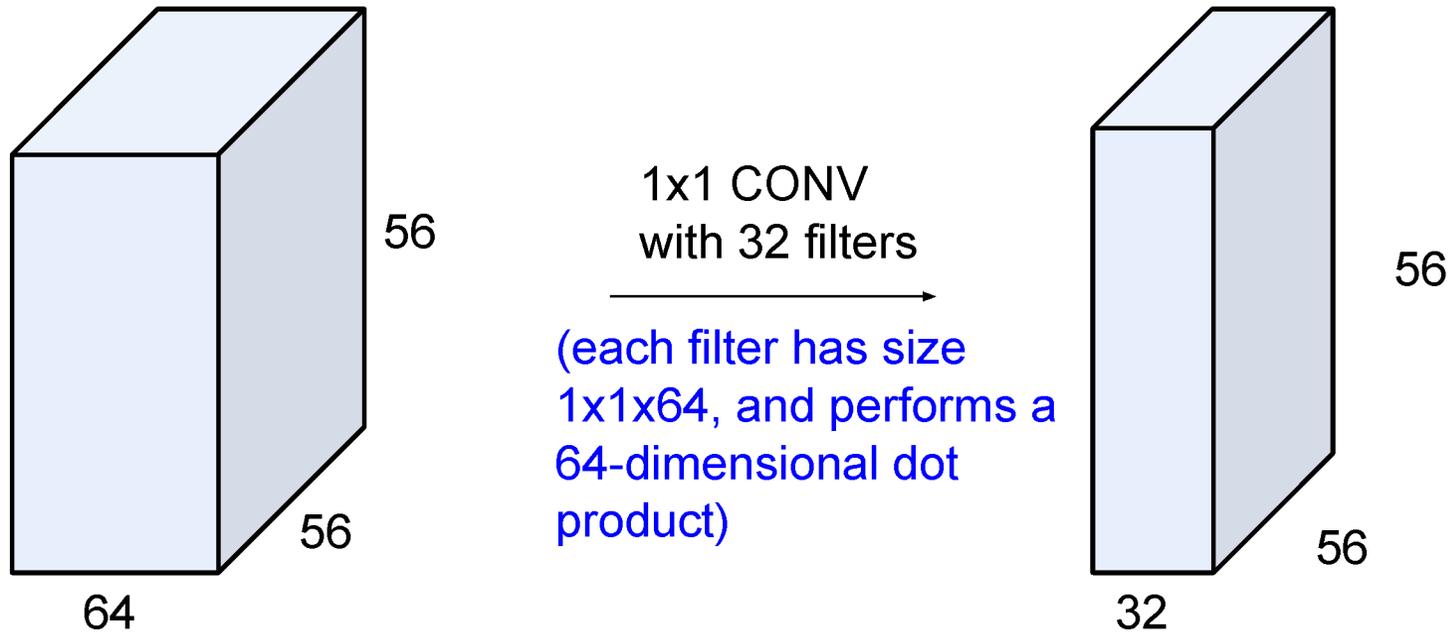


Number of parameters in this layer?

each filter has $5*5*3 + 1 = 76$ params (+1 for bias)

$\Rightarrow 76*10 = 760$

(btw, 1x1 convolution layers make perfect sense)



Convolutional layer—properties

- Small number of parameters to learn compared to a fully connected layer
- Preserves spatial structure—output of a convolutional layer is shaped like an image
- **Translation equivariant:** passing a translated image through a convolutional layer is (almost) equivalent to translating the convolution output (but be careful of image boundaries)

Self-study

[ConvNetJS demo: training on CIFAR-10]

[ConvNetJS](#) CIFAR-10 demo

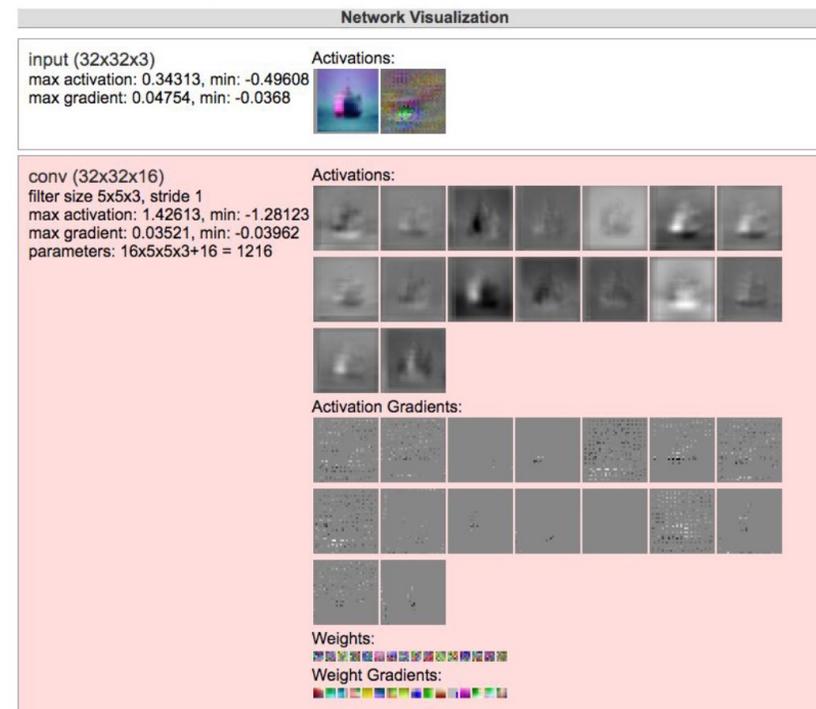
Description

This demo trains a Convolutional Neural Network on the [CIFAR-10 dataset](#) in your browser, with nothing but Javascript. The state of the art on this dataset is about 90% accuracy and human performance is at about 94% (not perfect as the dataset can be a bit ambiguous). I used [this python script](#) to parse the [original files](#) (python version) into batches of images that can be easily loaded into page DOM with img tags.

This dataset is more difficult and it takes longer to train a network. Data augmentation includes random flipping and random image shifts by up to 2px horizontally and vertically.

By default, in this demo we're using Adadelta which is one of per-parameter adaptive step size methods, so we don't have to worry about changing learning rates or momentum over time. However, I still included the text fields for changing these if you'd like to play around with SGD+Momentum trainer.

Report questions/bugs/suggestions to [@karpathy](#).



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

Train a neural network for classification
on CIFAR10 dataset in google colab

[Google Colab Page](#)

Where to Look for More Information

- Explore existing computer vision and machine learning frameworks
 - <https://pytorch.org/>
 - <https://www.tensorflow.org/>
 - <https://keras.io/>
 - <https://opencv.org/>
- Watch more in-depth lecture series
 - [The Ancient Secrets of Computer Vision - Joseph Redmon](#)
 - [Deep Learning Specialization - Andrew Ng](#)
- Checkout other online courses and guides
 - <https://ai.google/education/>
 - <https://www.udacity.com/course/deep-learning-pytorch--ud188>



Slide Credits

- [CS5670, Introduction to Computer Vision](#), Cornell Tech, by Noah Snavely.
- CS 543 [Computer Vision](#), by Stevlana Lazebnik, UIUC.
- EECS 442 [Computer Vision](#), by Justin Johnson & David Fouhey, U Michigan.