



CS231n: Deep Learning for Computer Vision

Lecture 1: Introduction

Computer Vision is everywhere!



Left to right:
[Image by Roger H Goun](#) is licensed under CC BY 2.0
[Image is CCO 1.0 public domain](#)
[Image is CCO 1.0 public domain](#)
[Image is CCO 1.0 public domain](#)



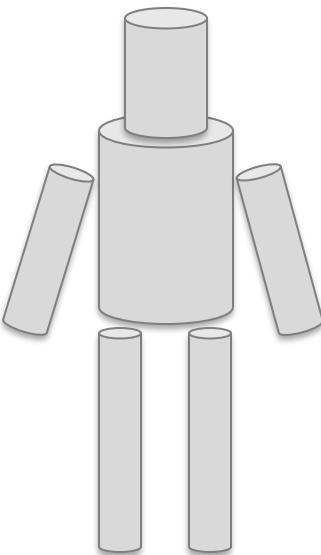
Left to right:
[Image is free to use](#)
[Image is CCO 1.0 public domain](#)
[Image by NASA](#) is licensed under CC BY 2.0
[Image is CCO 1.0 public domain](#)



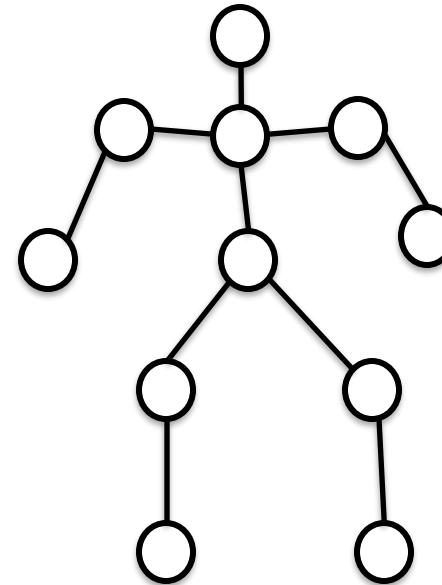
Bottom row, left to right:
[Image is CCO 1.0 public domain](#)
[Image by Derek Keats](#) is licensed under CC BY 2.0; changes made
[Image is public domain](#)
[Image is licensed under CC-BY 2.0; changes made](#)

Where did we come from?

Recognition via Parts (1970s)



Generalized Cylinders,
Brooks and Binford,
1979



Pictorial Structures,
Fischler and Elshlager, 1973



Recognition via Edge Detection (1980s)



1959
Hubel & Wiesel

1963
Roberts

1970s
David Marr

1979
Gen. Cylinders

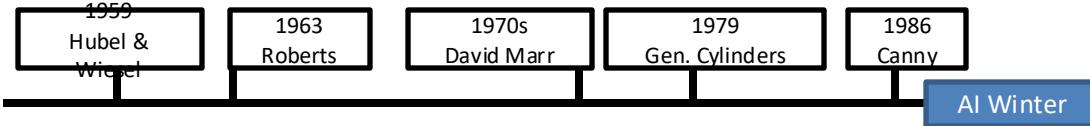
1986
Canny

John Canny, 1986
David Lowe, 1987

Image is CC0 1.0 public domain

Arriving at an “AI winter”

- Enthusiasm (and funding!) for AI research dwindled
- “Expert Systems” failed to deliver on their promises
- But subfields of AI continues to grow
 - Computer vision, NLP, robotics, compbio, etc.



[Left Image](#) is CC BY 3.0

[Middle Image](#) is public domain

[Right Image](#) is CC-BY 2.0; changes made

Slide inspiration: Justin Johnson

Recognition via Grouping (1990s)



1959
Hubel & Wiesel

1963
Roberts

1970s
David Marr

1979
Gen. Cylinders

1986
Canny
1997
Norm. Cuts

AI Winter

Normalized Cuts, Shi and Malik, 1997

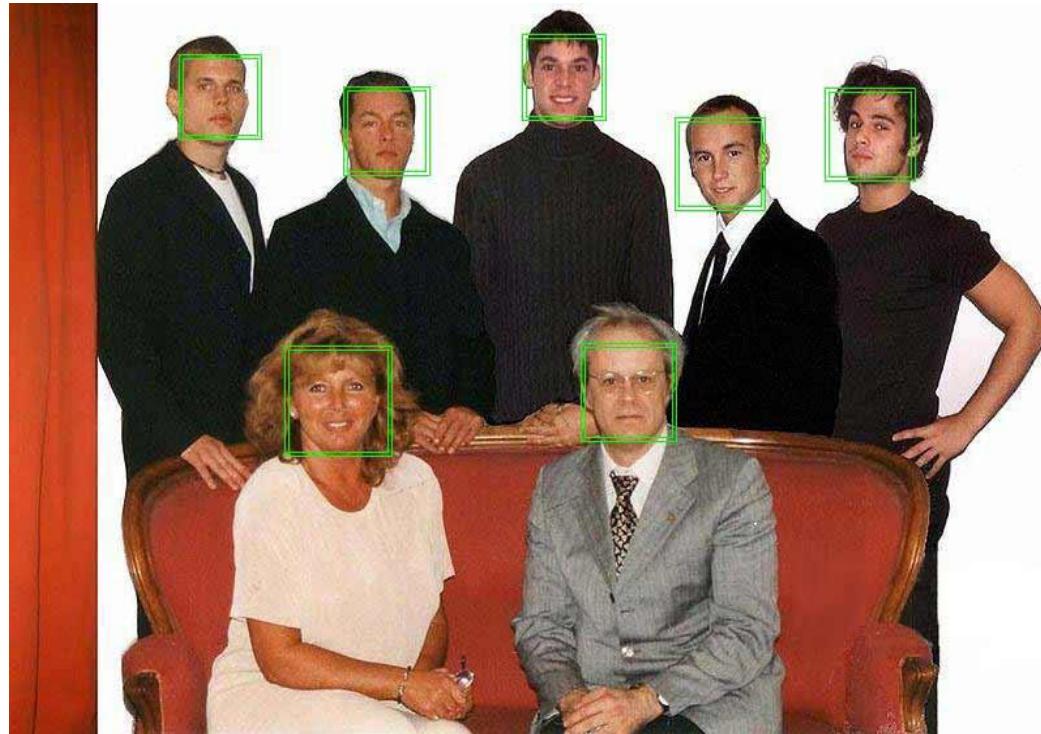
[Left Image](#) is CC BY 3.0
[Middle Image](#) is public domain

[Right Image](#) is CC-BY 2.0; changes made

Face Detection

Viola and Jones, 2001

One of the first successful applications of machine learning to vision



1959
Hubel & Wiesel

1963
Roberts

1970s
David Marr

1979
Gen. Cylinders

1986
Canny

1997
Norm. Cuts

1999
SIFT

2001
V&J

AI Winter

Slide inspiration: Justin Johnson

Caltech 101 images



1959
Hubel & Wiesel

1963
Roberts

1970s
David Marr

1979
Gen. Cylinders

1986
Canny

1997
Norm. Cuts

1999
SIFT

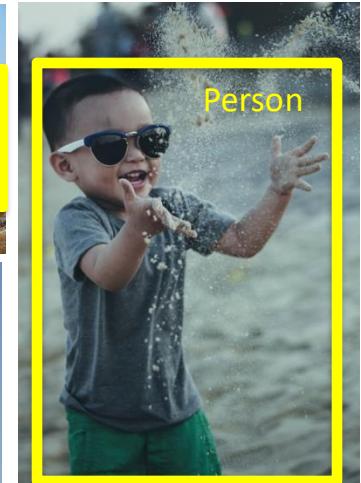
2001
V&J

2004, 2007
Caltech101;
PASCAL

AI Winter

PASCAL Visual Object Challenge

[Image is CC0 1.0 public domain](#)



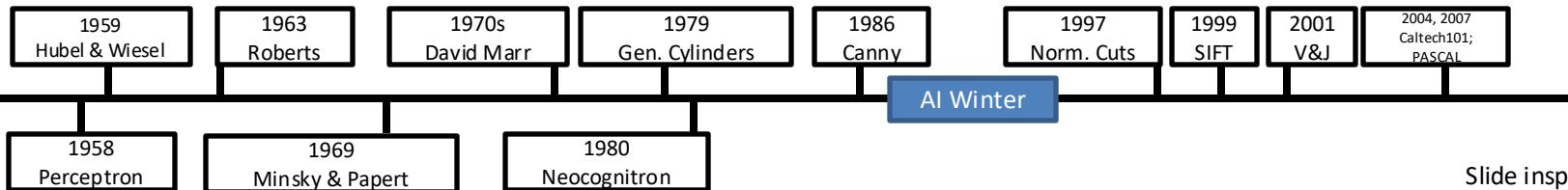
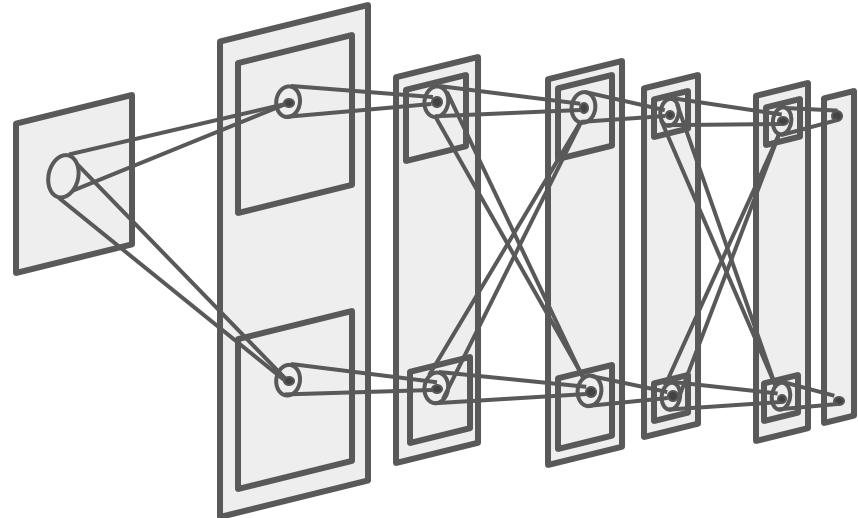
[Image is CC0 1.0 public domain](#)

Neocognitron: Fukushima, 1980

Computational model the visual system,
directly inspired by Hubel and Wiesel's
hierarchy of complex and simple cells

Interleaved simple cells (convolution)
and complex cells (pooling)

No practical training algorithm



Slide inspiration: Justin Johnson

Backprop: Rumelhart, Hinton, and Williams, 1986

Introduced backpropagation for computing gradients in neural networks

Successfully trained perceptrons with multiple layers

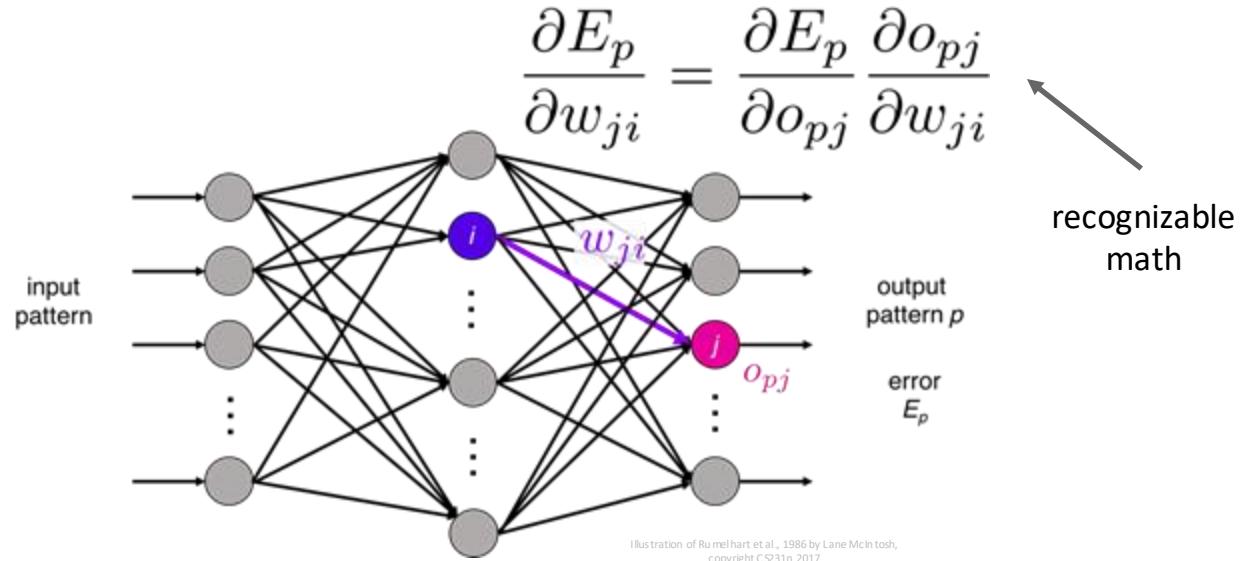
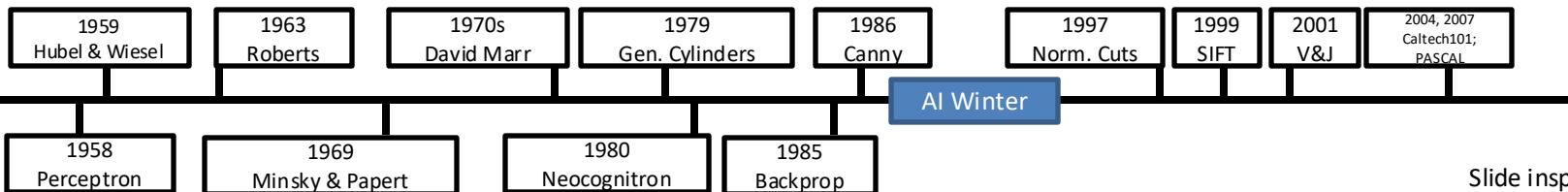
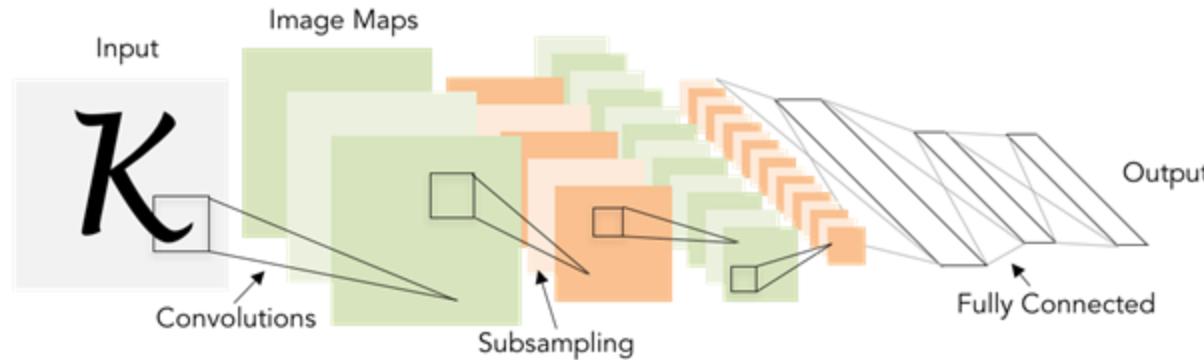


Illustration of Rumelhart et al., 1986 by Lane McIntosh,
copyright CS231n 2017



Convolutional Networks: LeCun et al, 1998

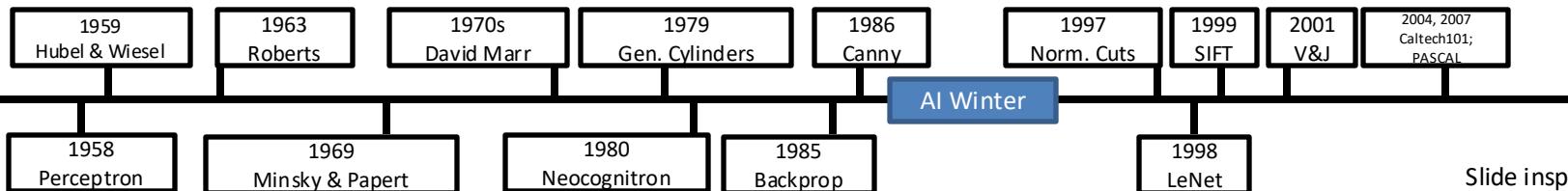


Applied backprop algorithm to a Neocognitron-like architecture

Learned to recognize handwritten digits

Was deployed in a commercial system by NEC, processed handwritten checks

Very similar to our modern convolutional networks!



Slide inspiration: Justin Johnson

2000s: “Deep Learning”

People tried to train neural networks that were deeper and deeper

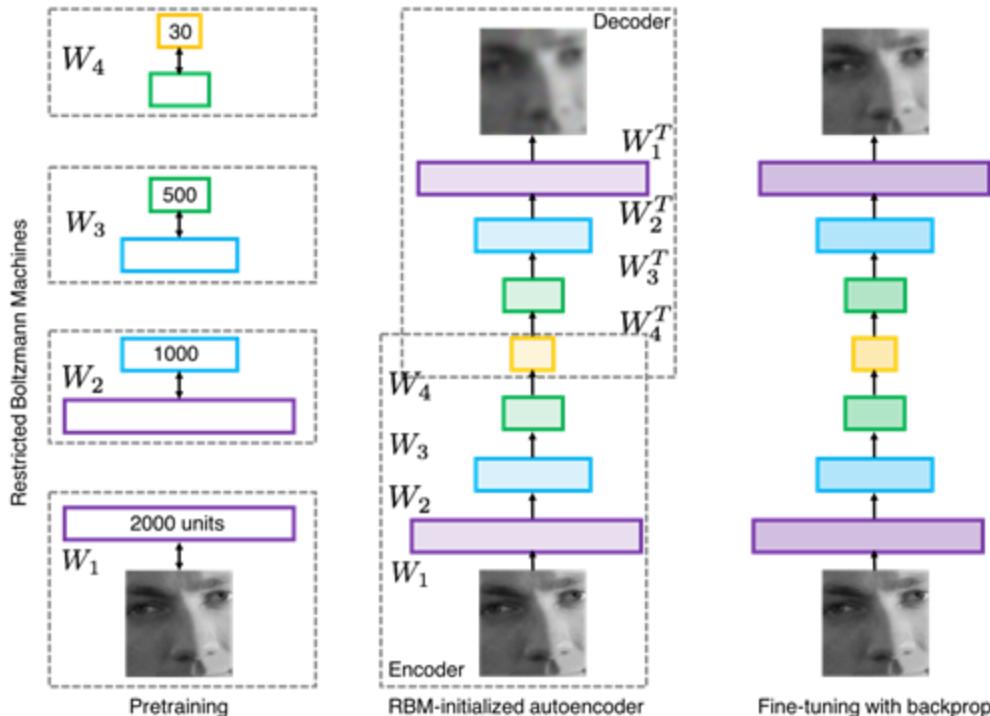
Not a mainstream research topic at this time

Hinton and Salakhutdinov, 2006

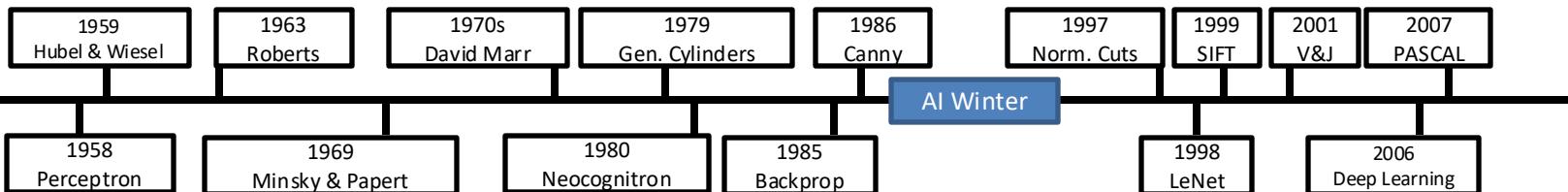
Bengio et al, 2007

Lee et al, 2009

Glorot and Bengio, 2010



Fine-tuning with backprop



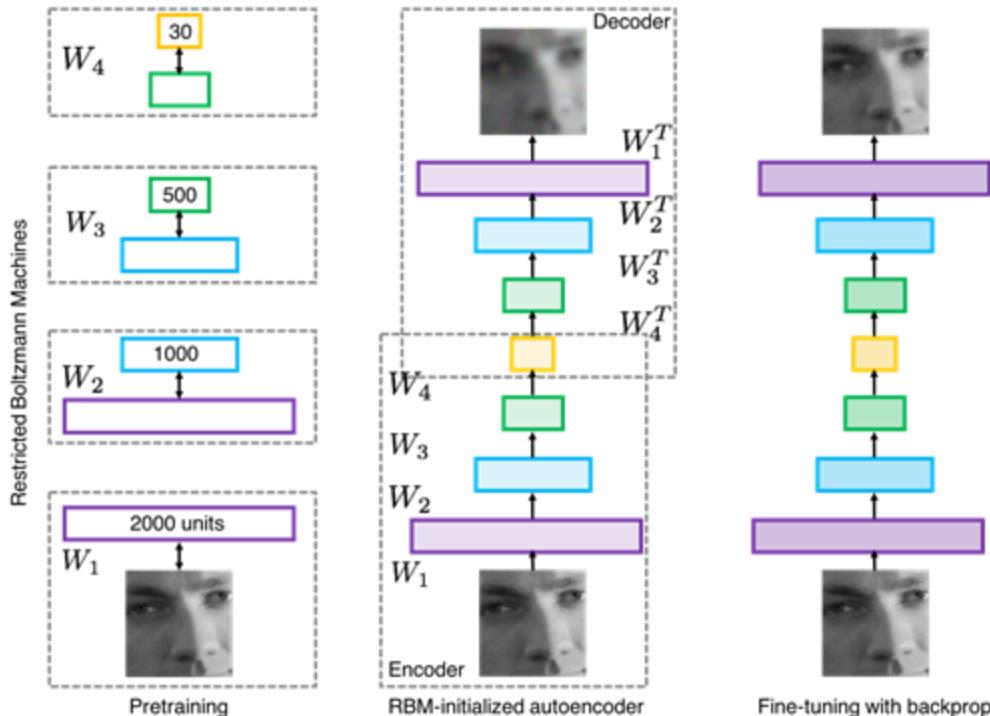
2000s: “Deep Learning”

People tried to train neural networks that were deeper and deeper

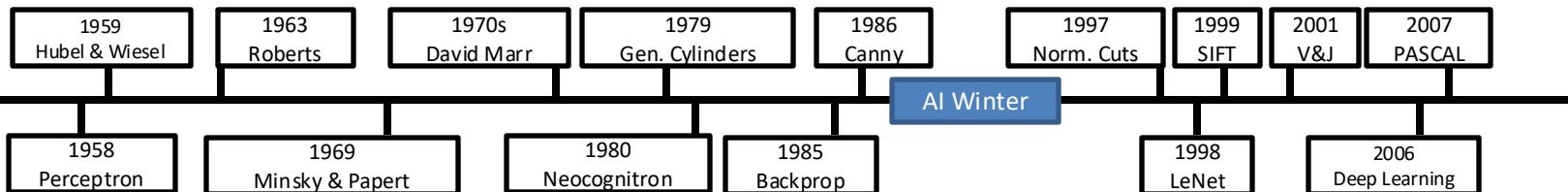
Not a mainstream research topic at this time

No good dataset to work on

Hinton and Salakhutdinov, 2006
Bengio et al, 2007
Lee et al, 2009
Glorot and Bengio, 2010



Slide inspiration: Justin Johnson



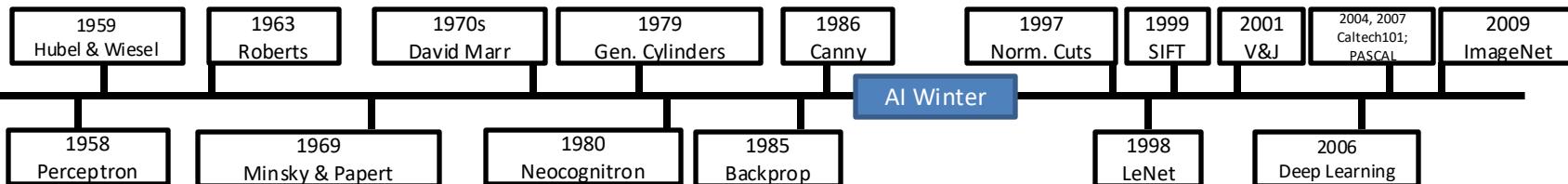
IMAGENET Large Scale Visual Recognition Challenge

The Image Classification Challenge:
1,000 object classes
1,431,167 images

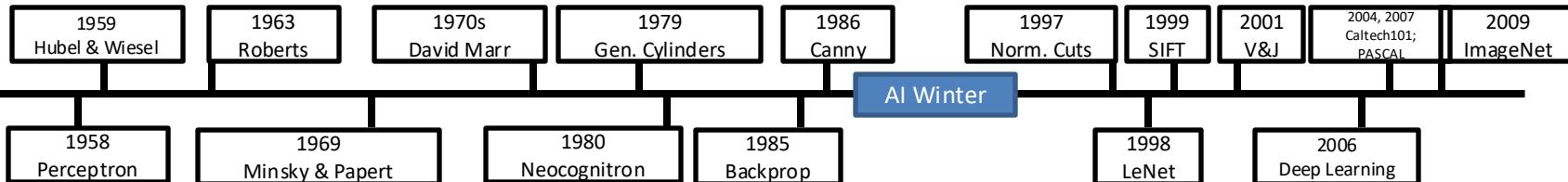
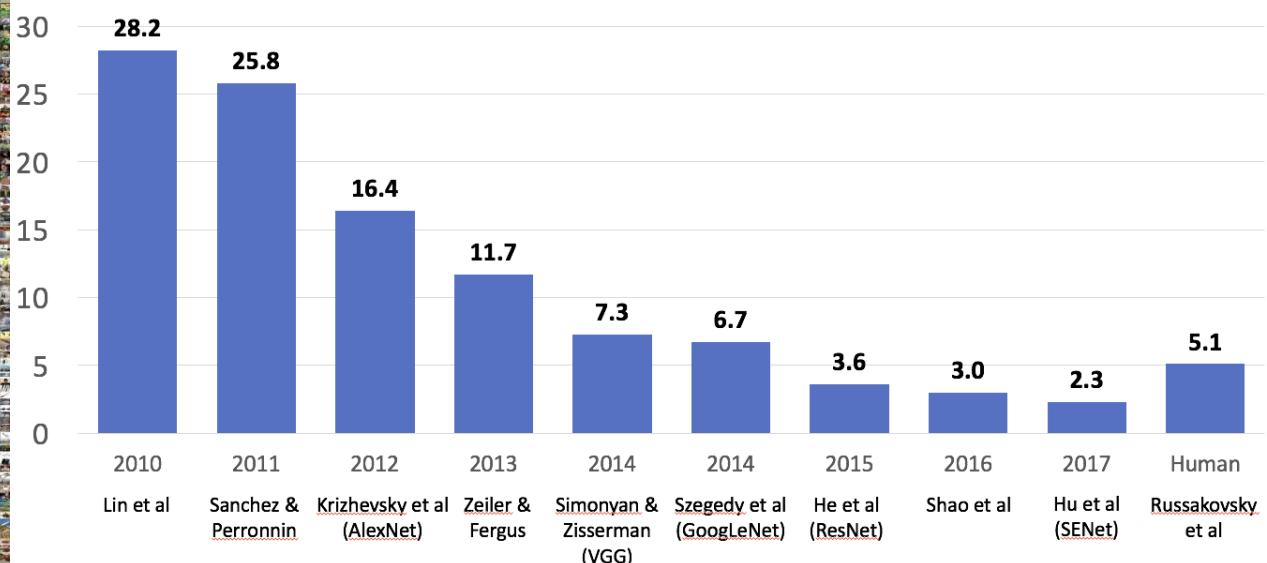


Output:
Scale
T-shirt
Steel drum
Drumstick
Mud turtle

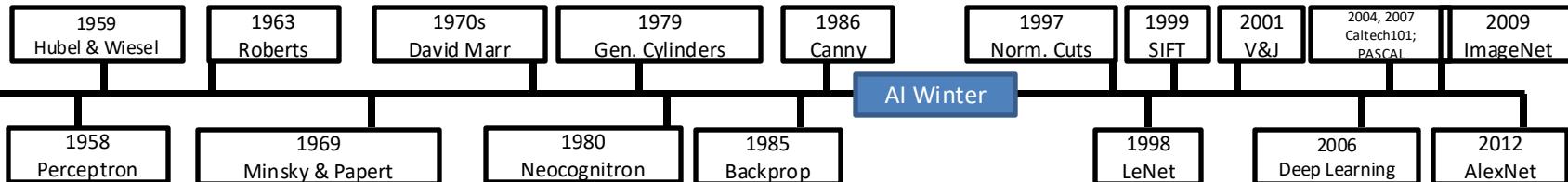
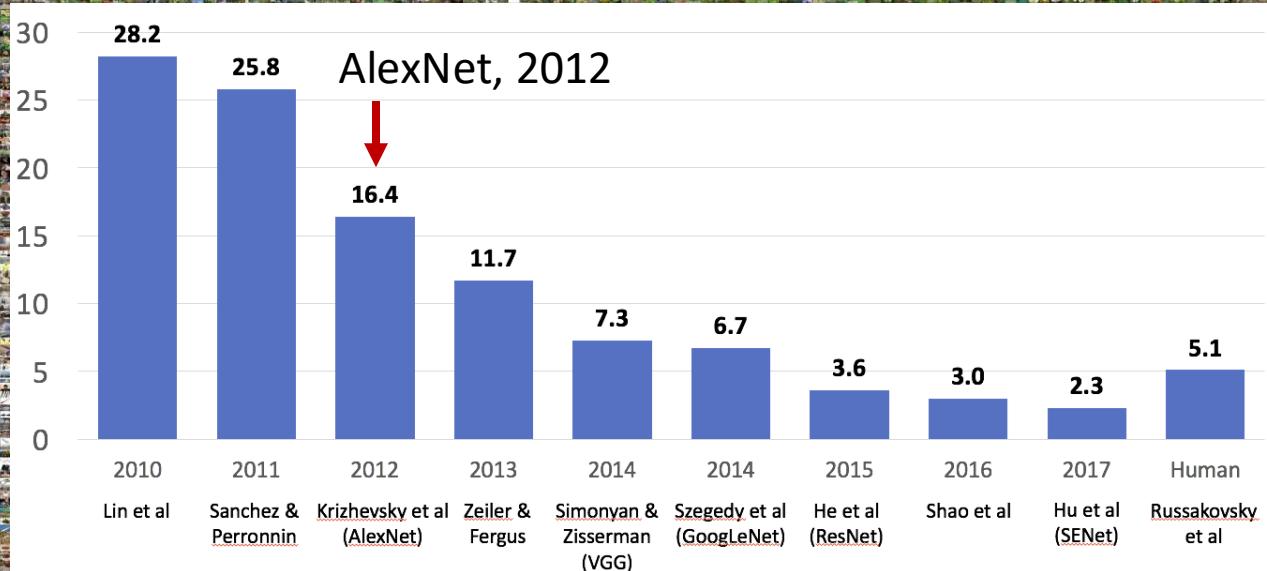
Deng et al, 2009
Russakovsky et al. IJCV 2015



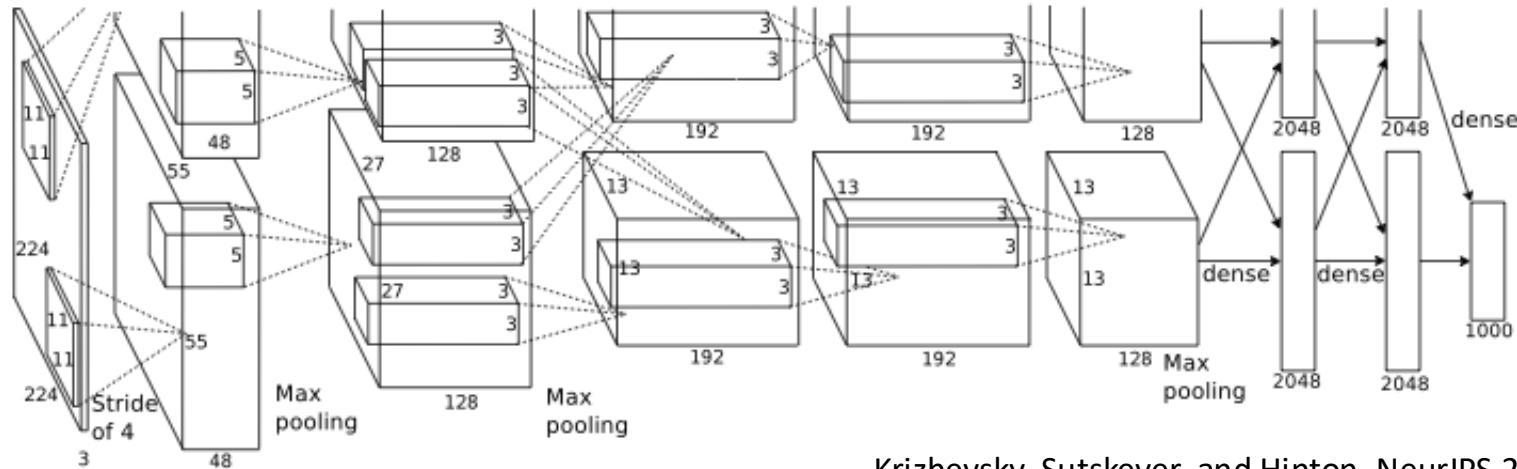
IMAGENET Large Scale Visual Recognition Challenge



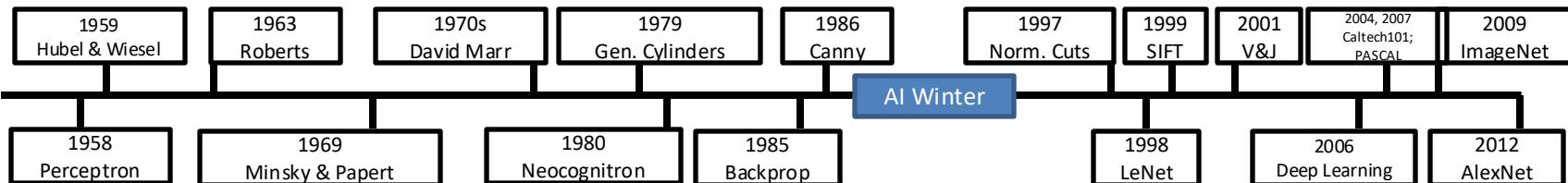
IMAGENET Large Scale Visual Recognition Challenge



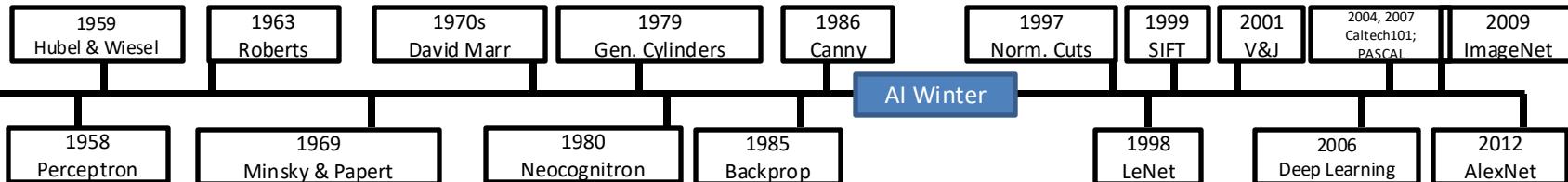
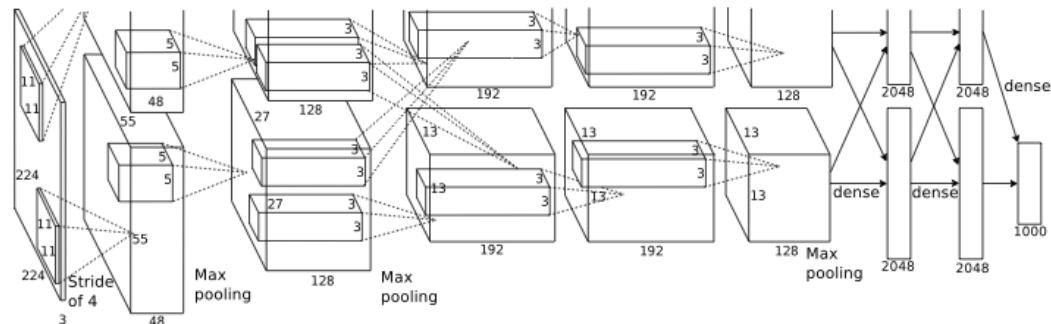
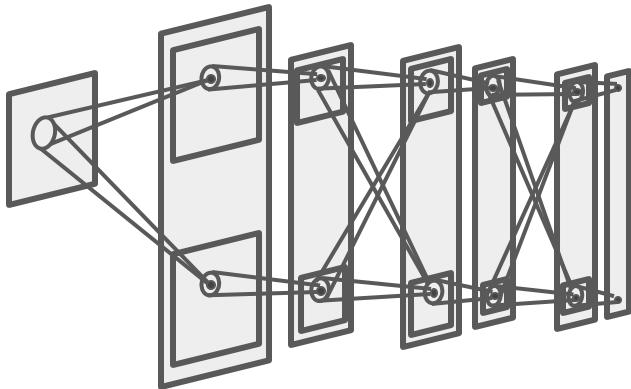
AlexNet: Deep Learning Goes Mainstream



Krizhevsky, Sutskever, and Hinton, NeurIPS 2012

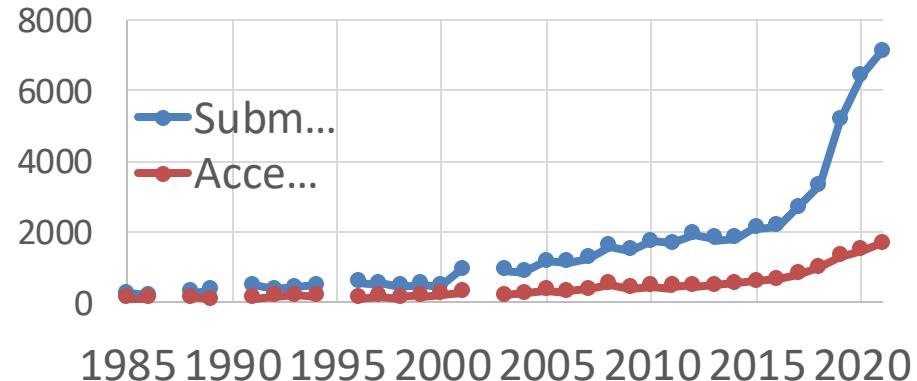


AlexNet vs. Neocognitron: 32 years apart

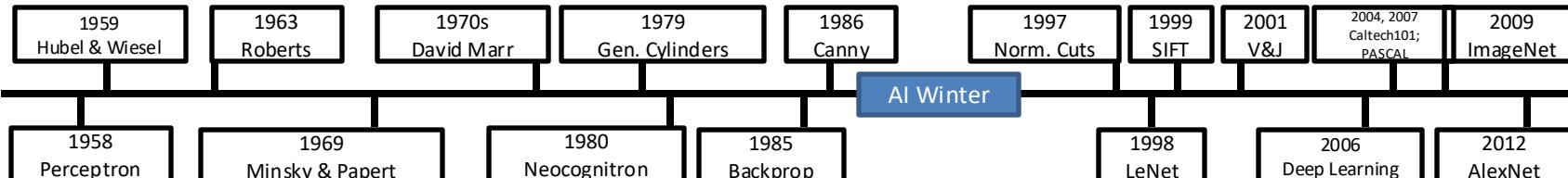


2012 to Present: Deep Learning Explosion

CVPR Papers



Publications at top Computer Vision conference



2012 to Present: Deep Learning is Everywhere

Image Classification



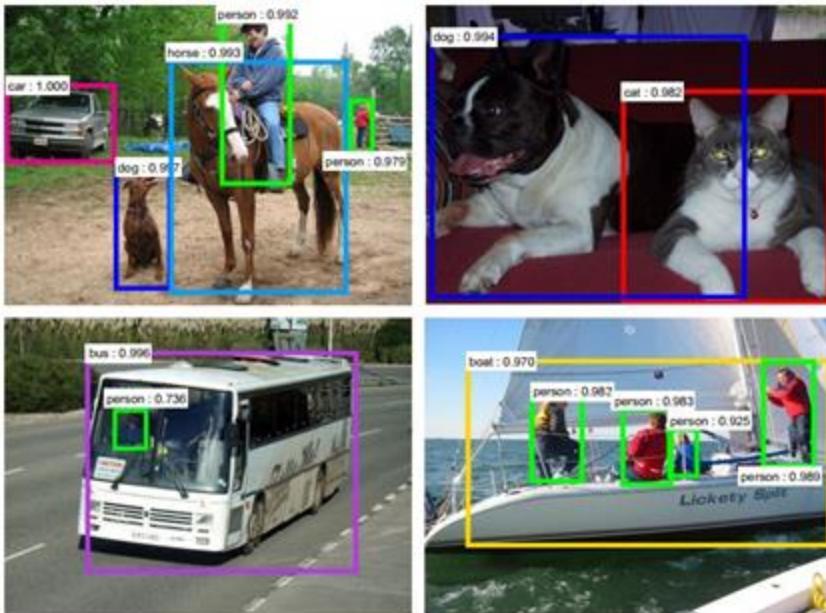
Image Retrieval



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

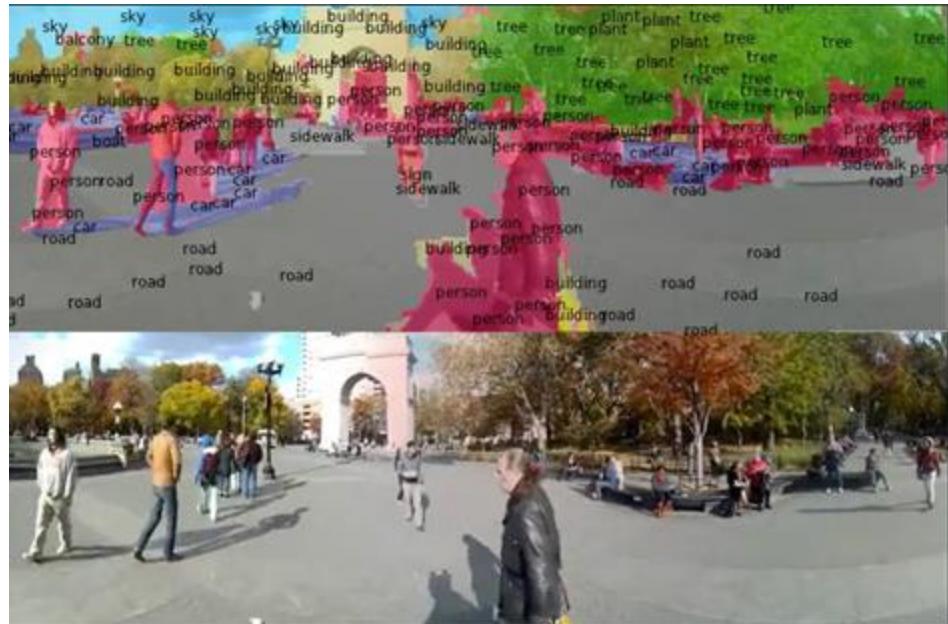
2012 to Present: Deep Learning is Everywhere

Object Detection



Ren, He, Girshick, and Sun, 2015

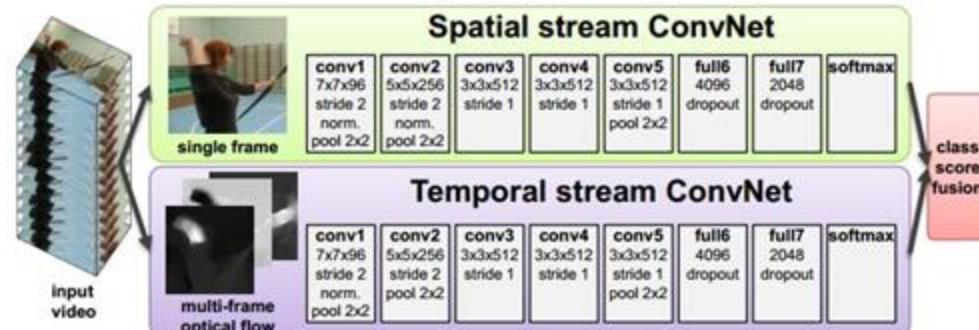
Image Segmentation



Fabaret et al, 2012

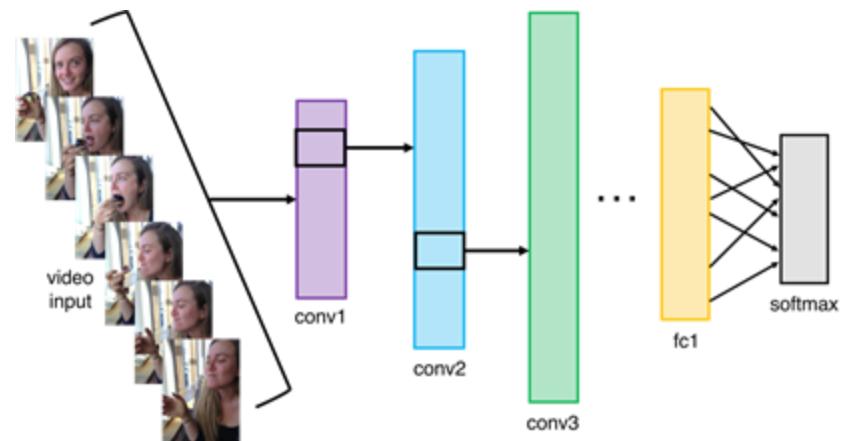
2012 to Present: Deep Learning is Everywhere

Video Classification



Simonyan et al, 2014

Activity Recognition

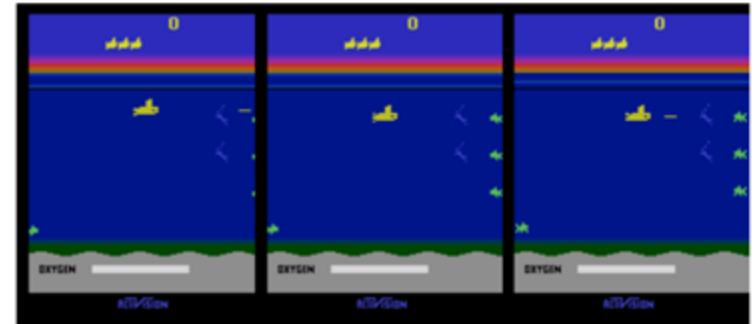
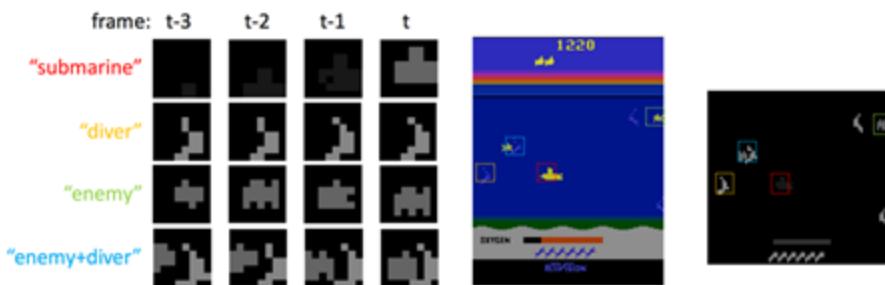


2012 to Present: Deep Learning is Everywhere

Pose Recognition (Toshev and Szegedy, 2014)

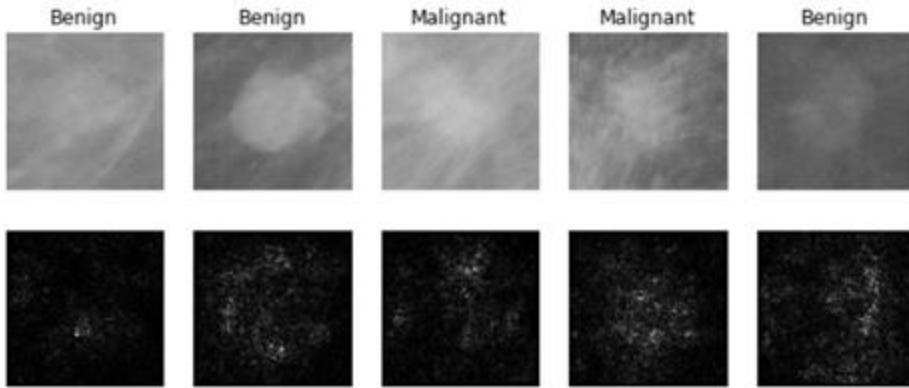


Playing Atari games (Guo et al, 2014)



2012 to Present: Deep Learning is Everywhere

Medical Imaging



Levy et al, 2016

Figure reproduced with permission

Whale recognition



Galaxy Classification



Dieleman et al, 2014

From left to right: [public domain by NASA](#), [usage permitted by ESA/Hubble](#), [public domain by NASA](#), and [public domain](#).

Kaggle Challenge

This image by Christin Khan is in the public domain and originally came from the U.S. NOAA.

2012 to Present: Deep Learning is Everywhere



Image Captioning

Vinyals et al, 2015

Karpathy and Fei-Fei,
2015

*A white teddy bear
sitting in the grass*

*A man in a baseball
uniform throwing a ball*

*A woman is holding
a cat in her hand*



*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/teddy-plush-bears-cute-teddy-bear-1643010/>
<https://pixabay.com/en/antique-cat-1623436/>
<https://pixabay.com/en/surf-wave-summer-sport-ltoral-11668716/>
<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-left-field-1045263/>

Captions generated by Justin Johnson using NeuralTalk2

2012 to Present: Deep Learning is Everywhere

TEXT PROMPT

an armchair in the shape of an avocado. an armchair imitating an avocado.

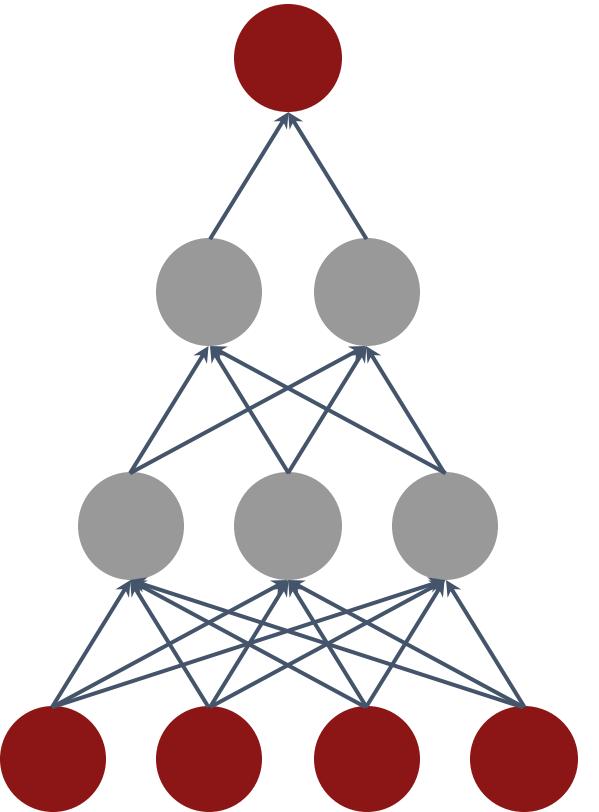
AI-GENERATED IMAGES



Slide inspiration: Justin Johnson



Computation
April 1, 2025



Algorithms



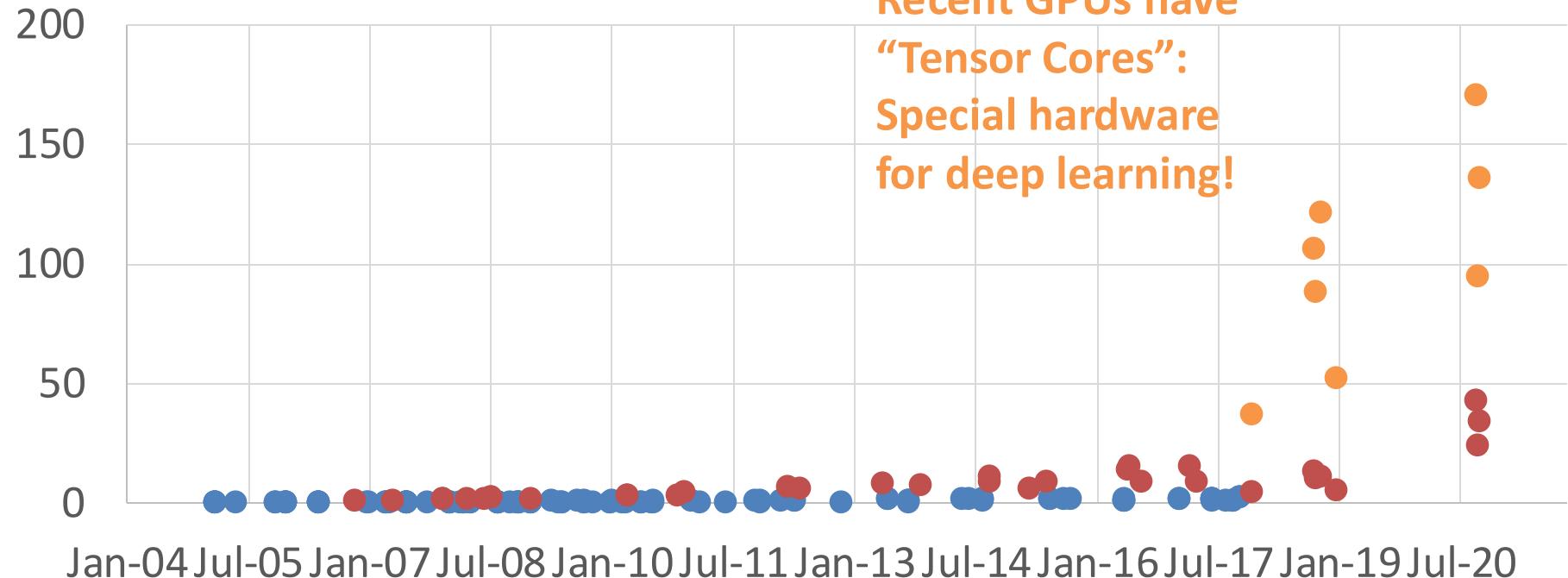
Data

GFLOP per Dollar

● CPU ● GPU (FP32) ● GPU (Tensor Core)

Recent GPUs have

“Tensor Cores”:
Special hardware
for deep learning!



Despite the successes, computer
vision still has a long way to go

Computer Vision Can Cause Harm

Harmful Stereotypes

Gender Classifier	Darker Male	Darker Female	Lighter Male	Lighter Female	Largest Gap
Microsoft	94.0%	79.2%	100%	98.3%	20.8%
FACE++	99.3%	65.5%	99.2%	94.0%	33.8%
IBM	88.0%	65.3%	99.7%	92.9%	34.4%

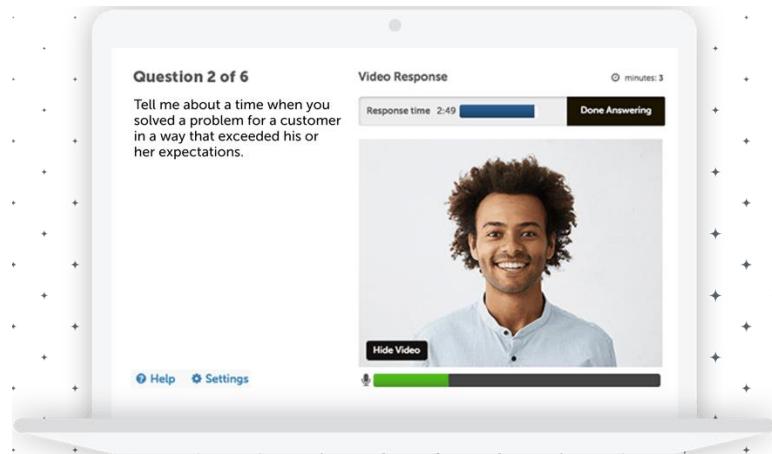


Affect people's lives

Technology

A face-scanning algorithm increasingly decides whether you deserve the job

HireVue claims it uses artificial intelligence to decide who's best for a job. Outside experts call it 'profoundly disturbing.'



Source: <https://www.washingtonpost.com/technology/2019/10/22/ai-hiring-face-scanning-algorithm-increasingly-decides-whether-you-deserve-job/>
<https://www.hirevue.com/platform/online-video-interviewing-software>

Example Credit: Timnit Gebru

Computer Vision Can Save Lives

How to take care of seniors
while keeping them safe?



Early Symptom Detection
of COVID-19



Monitor Patients with
Mild Symptoms



Manage Chronic Conditions

Versatile



Mobility



Infection



Sleep



Diet

Scalable



Low-cost



Burden-free



And there is a lot we don't know how to do



https://fedandfit.com/wp-content/uploads/2020/06/summer-activities-for-kids_optimized-scaled.jpeg



This image is
copyright-free [United
States govern
ment
work](#)

Lecture 2: Image Classification with Linear Classifiers

Image Classification

A Core Task in Computer Vision

Today:

- The image classification task
- Two basic data-driven approaches to image classification
 - K-nearest neighbor and linear classifier

Image Classification: A core task in Computer Vision



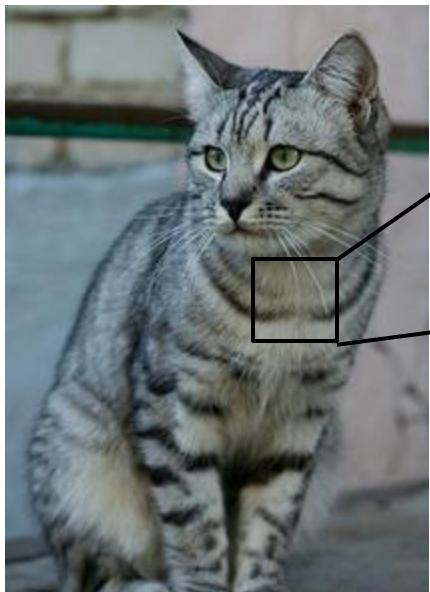
(assume given a set of possible labels)
{dog, cat, truck, plane, ...}



cat

This image by [Nikita](#) is
licensed under [CC-BY2.0](#)

The Problem: Semantic Gap



[105 112 188 111 104 99 186 99 96 103 112 119 104 97 93 87]
[91 98 182 106 184 79 98 103 99 105 123 136 110 105 94 85]
[76 85 98 105 128 105 87 96 95 99 115 112 186 103 99 85]
[99 81 81 93 120 131 127 106 95 98 102 99 96 93 101 94]
[106 91 61 64 69 91 88 85 101 107 109 98 75 84 96 95]
[114 108 85 55 55 69 64 54 64 87 112 129 98 74 84 91]
[133 137 147 183 65 81 80 65 52 54 74 84 102 93 85 82]
[128 137 144 148 109 95 86 76 62 65 63 63 60 73 86 101]
[125 133 148 137 119 121 117 94 65 79 88 65 54 64 72 98]
[127 125 131 147 133 127 126 131 111 96 89 75 61 64 72 84]
[115 114 109 123 150 148 131 118 113 109 100 92 74 65 72 78]
[89 93 98 97 108 147 131 118 113 114 113 109 106 95 77 88]
[63 77 86 81 77 79 102 123 117 115 117 125 125 130 115 87]
[62 65 82 89 78 71 88 101 124 126 119 101 107 114 131 119]
[63 65 75 88 89 71 62 81 120 138 135 105 81 98 118 118]
[87 65 71 87 106 95 69 45 76 130 126 107 92 94 105 112]
[118 97 82 86 117 123 116 66 41 51 95 93 89 95 102 107]
[164 146 112 88 82 120 124 104 76 48 45 66 88 101 102 109]
[157 170 157 128 93 86 114 134 112 97 69 55 70 82 99 94]
[130 128 134 161 139 100 109 118 121 134 114 87 65 53 69 86]
[128 112 96 117 150 144 128 115 104 107 102 93 87 81 72 79]
[123 107 96 86 83 112 153 149 122 109 104 75 80 107 112 99]
[122 121 102 88 82 86 94 117 145 148 153 102 58 78 92 107]
[122 164 148 103 71 56 78 83 93 103 119 139 102 61 69 84]

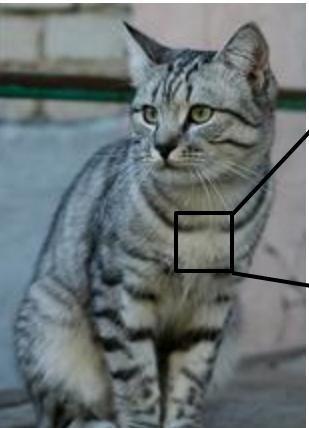
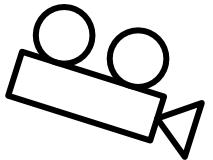
What the computer sees

An image is a tensor of integers between [0, 255]:

e.g. 800 x 600 x 3
(3 channels RGB)

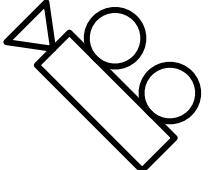
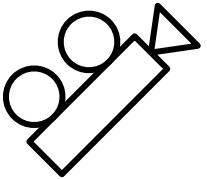
This image by [Nikita](#) is
licensed under [CC-BY2.0](#)

Challenges: Viewpoint variation



1149	152	189	111	184	99	186	99	96	183	112	119	184	97	93	87
1195	98	182	386	384	79	98	692	99	185	123	136	128	185	94	851
1176	85	98	385	328	189	87	94	95	99	115	132	186	183	99	853
1198	85	85	82	128	131	127	149	95	94	180	99	96	93	181	941
1180	85	85	85	85	85	85	85	85	85	180	85	85	85	85	853
1154	188	85	55	55	69	64	54	64	87	112	129	98	74	84	913
1133	137	147	383	65	81	88	65	52	54	74	84	182	93	85	821
1138	137	144	348	589	95	86	76	62	65	63	63	69	73	64	5813
1129	132	132	129	129	125	125	125	84	85	79	85	54	64	72	581
1132	132	132	132	132	132	132	132	132	132	132	132	132	132	132	132
1115	154	189	323	558	548	571	218	513	189	189	92	74	85	73	781
1189	89	99	97	589	547	531	518	513	114	113	189	186	95	77	681
1163	77	88	85	77	79	581	123	117	115	117	125	125	138	125	871
1182	82	82	82	82	82	82	82	82	82	82	82	82	82	82	82
1189	85	75	88	89	85	62	52	52	52	52	138	135	185	85	58
1187	85	75	87	186	95	69	49	76	138	136	187	92	94	185	3321
1118	97	82	86	157	323	515	66	41	52	95	93	89	89	89	582
1164	148	137	137	82	129	521	584	76	48	45	64	88	185	182	1891
1115	129	129	129	129	129	129	129	129	129	129	129	129	129	129	129
1138	128	134	181	339	549	189	518	521	134	114	87	65	53	89	861
1128	112	96	127	558	544	526	115	149	187	182	93	93	87	85	72
1123	187	96	86	83	512	553	549	522	149	184	75	88	187	132	993
1122	125	182	88	82	86	94	117	145	148	153	182	58	76	92	587
1122	164	148	148	71	56	78	83	93	181	119	182	61	69	841	

All pixels change when
the camera moves!



[This image](#) by [Nikita](#) is
licensed under [CC-BY2.0](#)

Challenges: Illumination



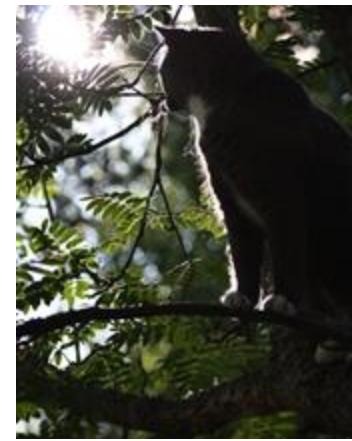
[This image](#) is CC0 1.0 public domain



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



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



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

Challenges: Background Clutter



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



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

Challenges: Occlusion



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

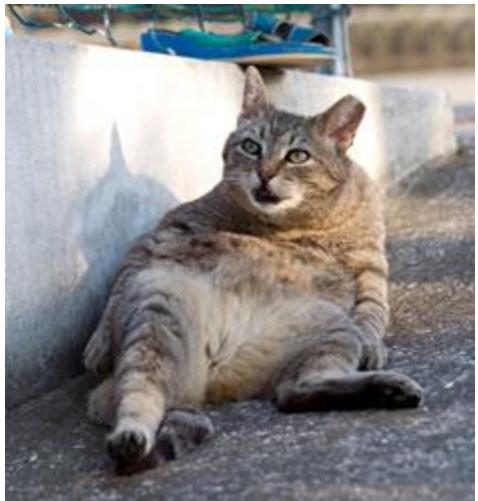


[This image](#) is [CC0 1.0](#) public domain



[This image](#) by [jonsson](#) is licensed
under [CC-BY2.0](#)

Challenges: Deformation



[This image](#) by [Umberto Salvagnin](#) is licensed under [CC-BY 2.0](#)



[This image](#) by [Umberto Salvagnin](#) is licensed under [CC-BY 2.0](#)



[This image](#) by [sare bear](#) is licensed under [CC-BY 2.0](#)



[This image](#) by [Tom Thai](#) is licensed under [CC-BY 2.0](#)

Challenges: Intraclass variation



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

Challenges: Context

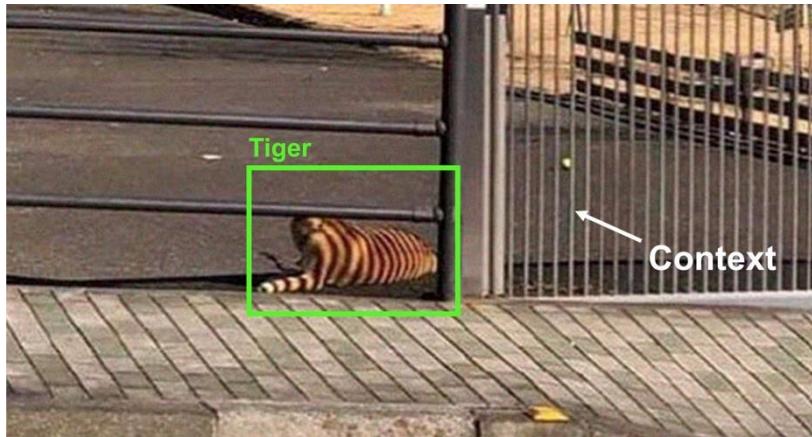


Image source: https://www.linkedin.com/posts/ralph-aboujaoude-diaz-40838313_technology-artificialintelligence-computervision-activity-6912446088364875776-h-lq?utm_source=linkedin_share&utm_medium=member_desktop_web

An image classifier

```
def classify_image(image):  
    # Some magic here?  
    return class_label
```

Unlike e.g. sorting a list of numbers,
no obvious way to hard-code the algorithm for
recognizing a cat, or other classes.

Linear Classifier

Parametric Approach: Linear Classifier

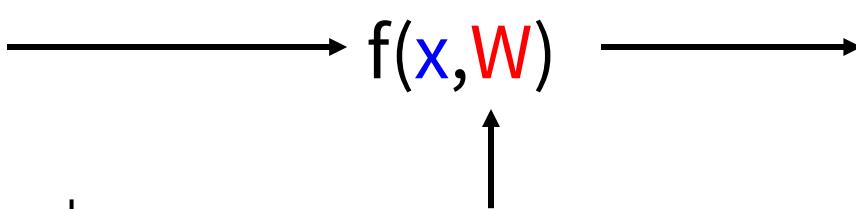


Image

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

10x1 3072x1
10x3072 10x1

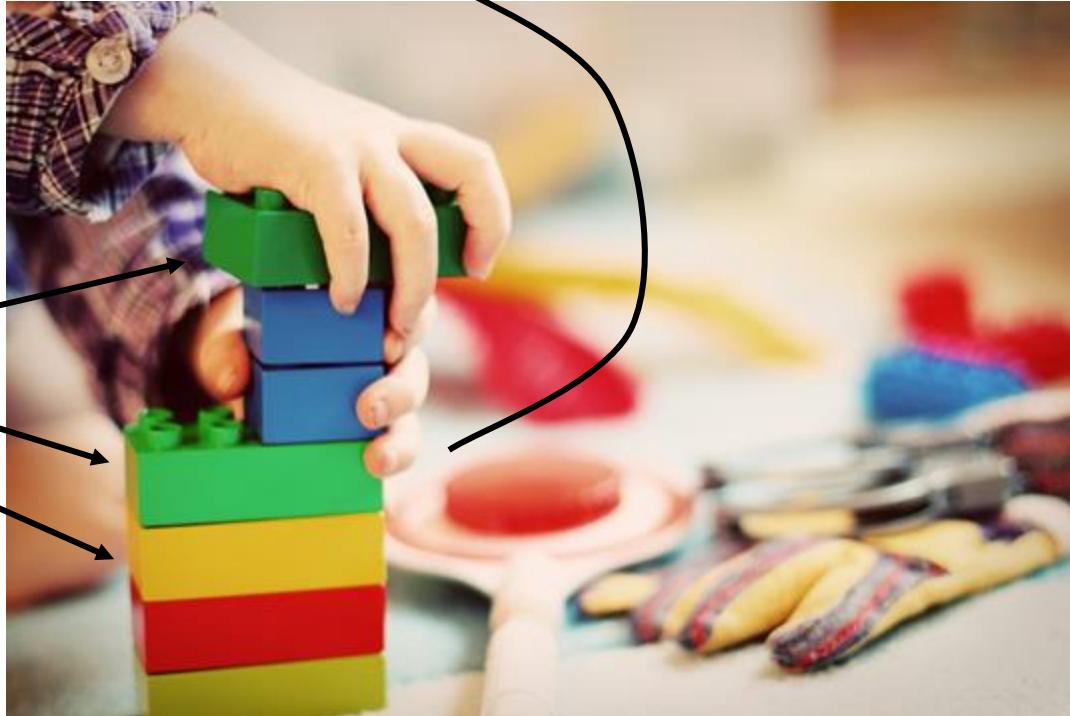
Array of 32x32x3 numbers
(3072 numbers total)



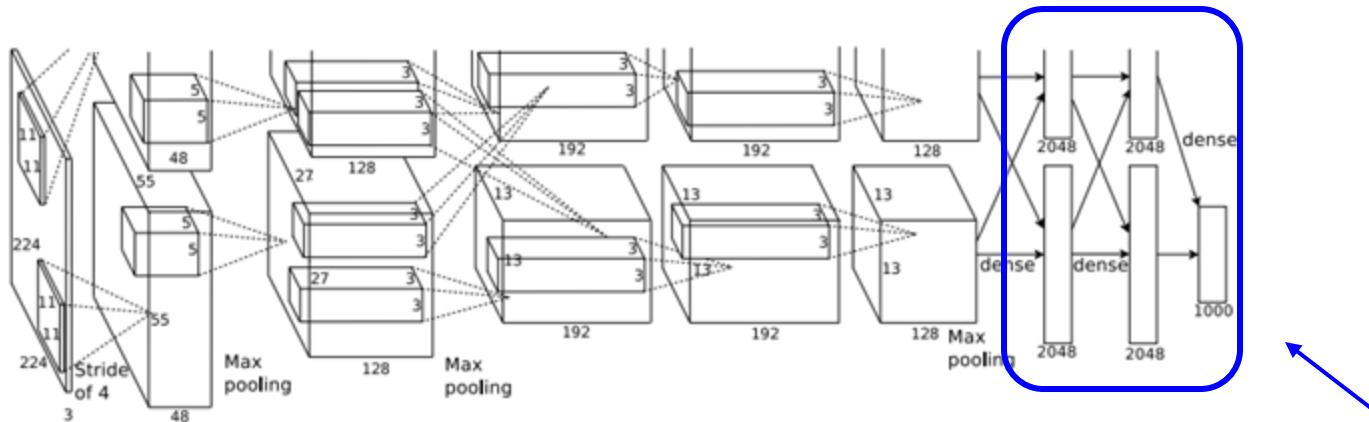
W
parameters
or weights

Neural Network

Linear
classifiers

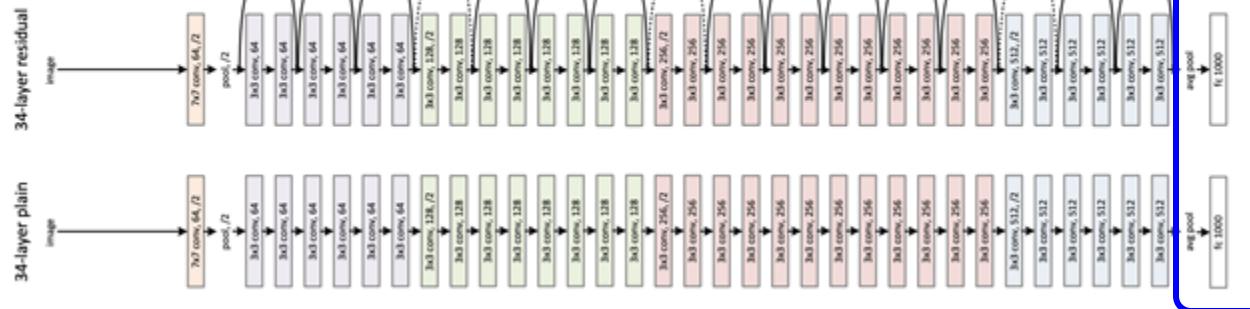


[This image](#) is [CC0 1.0](#) public domain



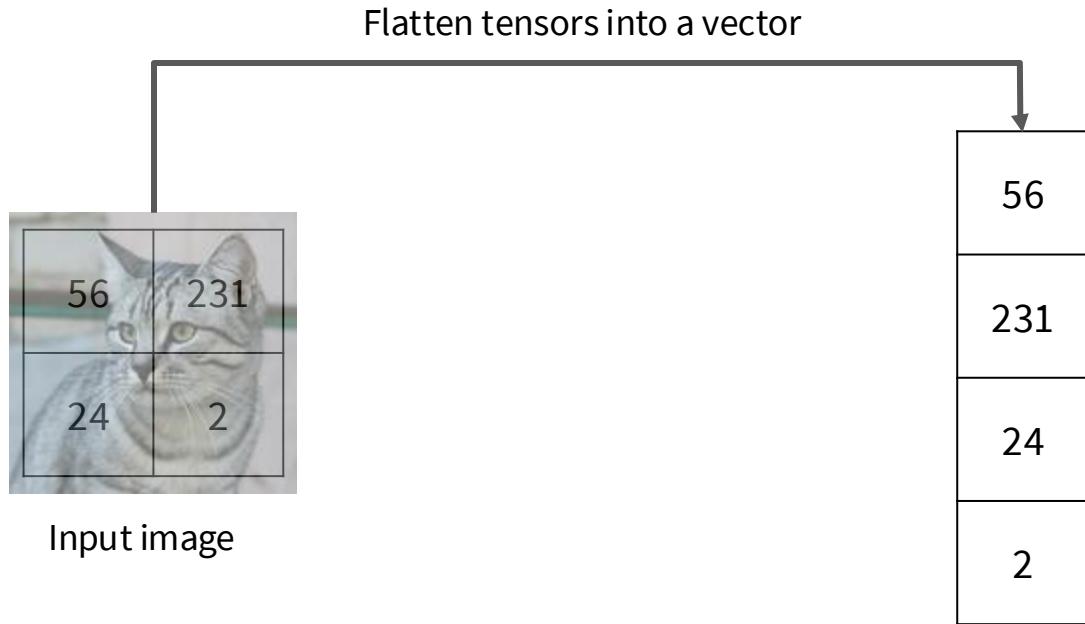
[Krizhevsky et al. 2012]

Linear layers

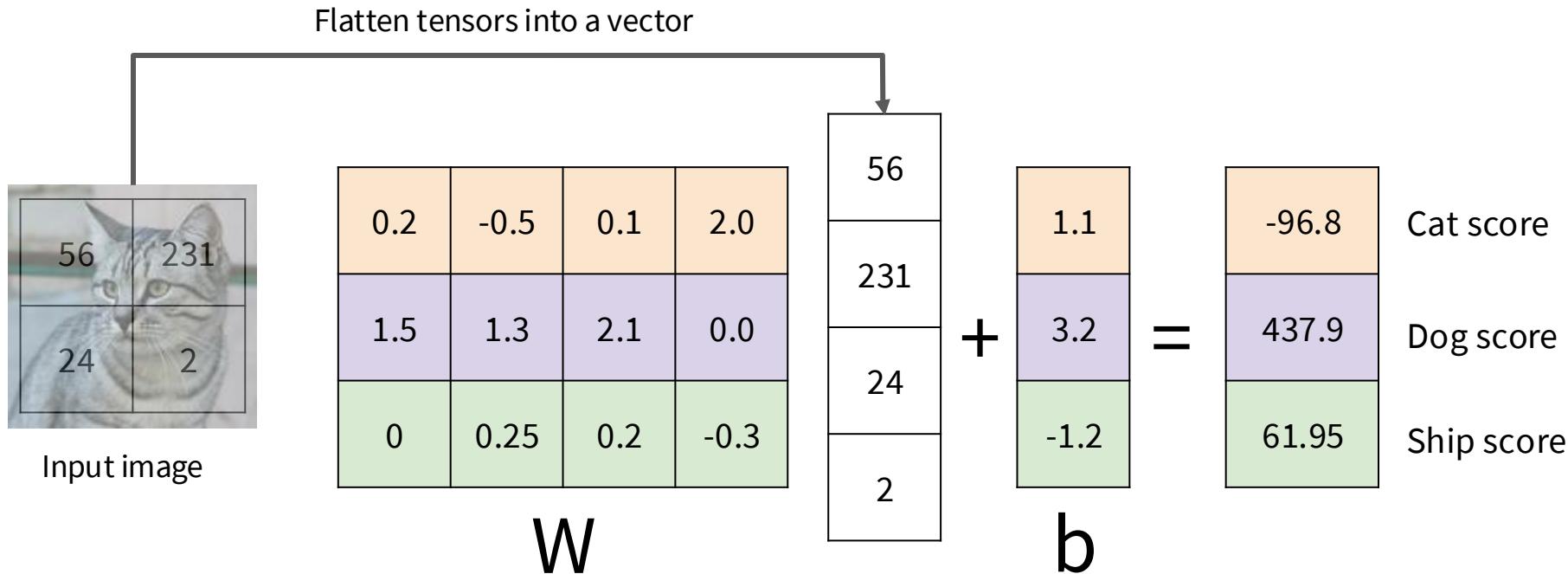


[He et al. 2015]

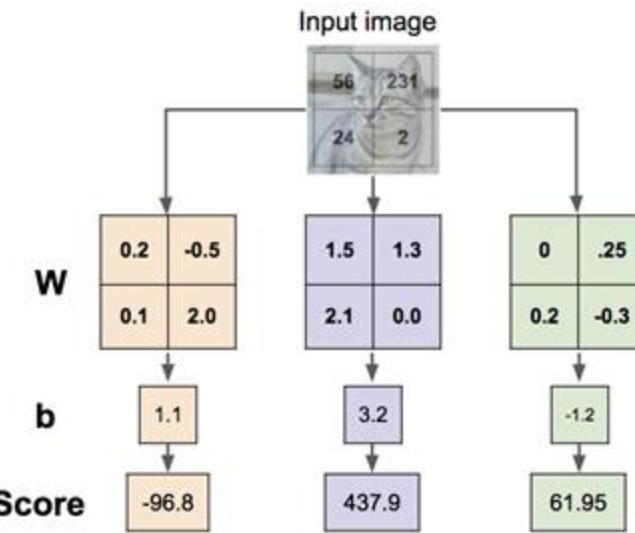
Example with an image with 4 pixels, and 3 classes (cat/dog/ship)



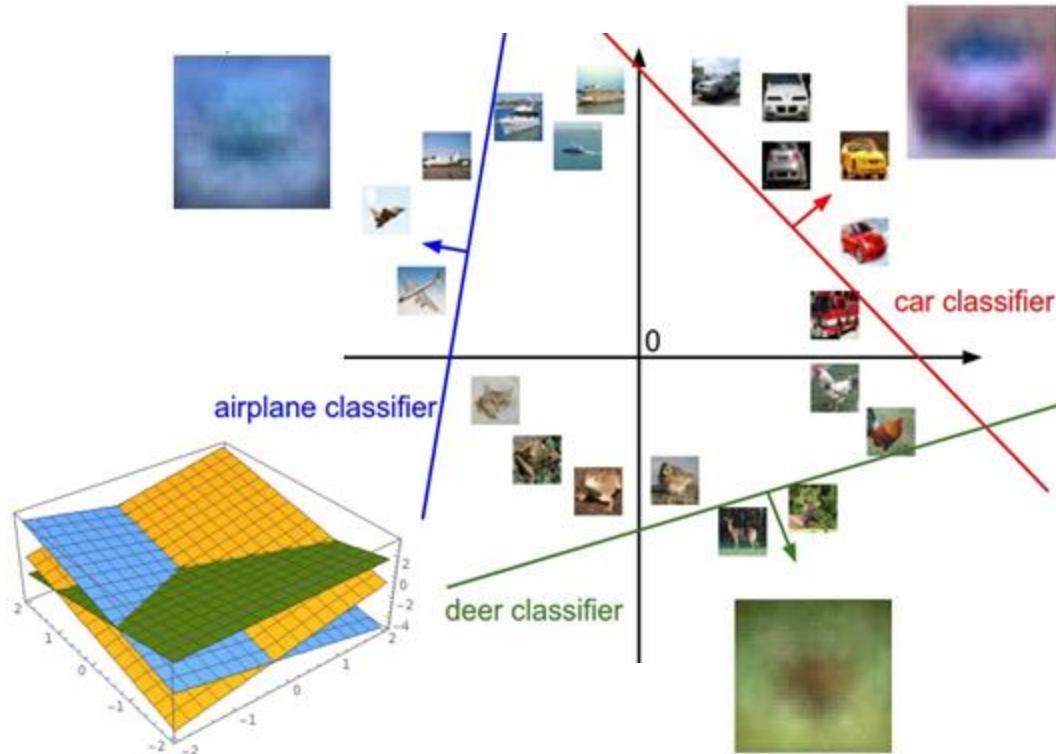
Example with an image with 4 pixels, and 3 classes (**cat/dog/ship**) Algebraic Viewpoint



Interpreting a Linear Classifier: Visual Viewpoint



Interpreting a Linear Classifier: Geometric Viewpoint



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



Array of 32x32x3 numbers
(3072 numbers total)

Plot created using [Wolfram Cloud](#)

[Cat image](#) by [Nikita](#) is licensed under [CC-BY 2.0](#)

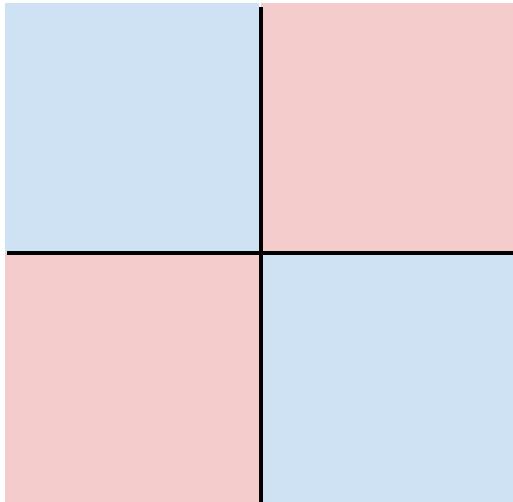
Hard cases for a linear classifier

Class 1:

First and third quadrants

Class 2:

Second and fourth quadrants

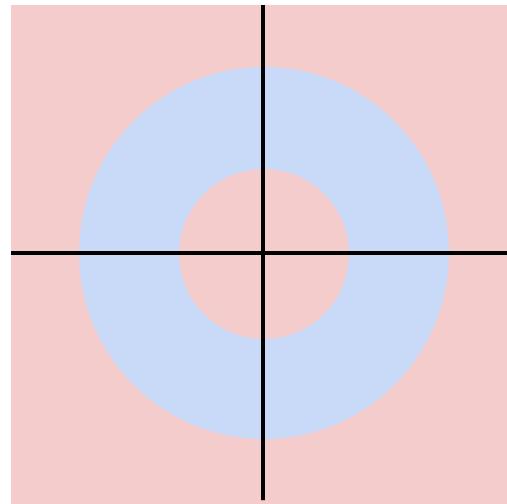


Class 1:

$1 \leq L_2 \text{ norm} \leq 2$

Class 2:

Everything else

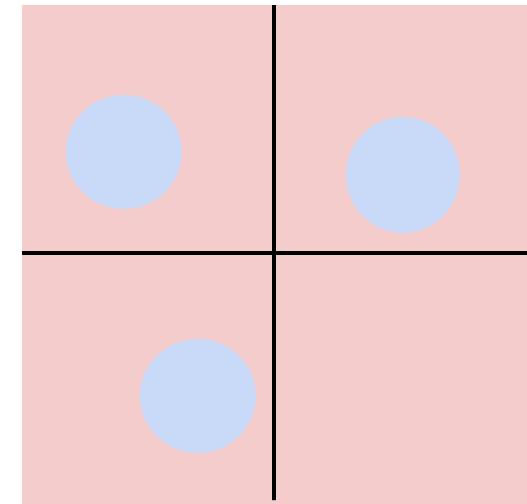


Class 1:

Three modes

Class 2:

Everything else



Softmax classifier

Softmax Classifier (Multinomial Logistic Regression)



Want to interpret raw classifier scores as probabilities

$$s = f(x_i; W)$$

$$P(Y = k|X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

Softmax Function

cat
car
frog

3.2
5.1
-1.7

Unnormalized log-probabilities / logits

exp

24.5
164.0
0.18

unnormalized probabilities

normalize

0.13
0.87
0.00

probabilities

$$L_i = -\log P(Y = y_i|X = x_i)$$

$$\rightarrow L_i = -\log(0.13) \\ = 2.04$$

Maximum Likelihood Estimation
Choose weights to maximize the likelihood of the observed data
(See CS 229 for details)

Softmax Classifier (Multinomial Logistic Regression)



Want to interpret raw classifier scores as probabilities

$$s = f(x_i; W)$$

$$P(Y = k|X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

Softmax Function

cat
car
frog

3.2
5.1
-1.7

Unnormalized log-probabilities / logits

Probabilities must be ≥ 0

24.5
164.0
0.18

unnormalized probabilities

exp

normalize

Probabilities must sum to 1

0.13
0.87
0.00

probabilities

$$L_i = -\log P(Y = y_i|X = x_i)$$

compare

1.00
0.00
0.00

Correct probs

Softmax Classifier (Multinomial Logistic Regression)



Want to interpret raw classifier scores as probabilities

$$s = f(x_i; W)$$

$$P(Y = k|X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

Softmax Function

cat
car
frog

3.2
5.1
-1.7

Unnormalized log-probabilities / logits

exp

24.5
164.0
0.18

unnormalized probabilities

normalize

0.13
0.87
0.00

probabilities

$$L_i = -\log P(Y = y_i|X = x_i)$$

compare ←
Kullback-Leibler divergence
 $D_{KL}(P||Q) = \sum_y P(y) \log \frac{P(y)}{Q(y)}$

1.00

0.00

0.00

Correct probs

Softmax Classifier (Multinomial Logistic Regression)



Want to interpret raw classifier scores as probabilities

$$s = f(x_i; W)$$

$$P(Y = k|X = x_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

Softmax Function

cat
car
frog

3.2
5.1
-1.7

Unnormalized log-probabilities / logits

exp

24.5
164.0
0.18

unnormalized probabilities

normalize

0.13
0.87
0.00

probabilities

$$L_i = -\log P(Y = y_i|X = x_i)$$

1.00
0.00
0.00

Cross Entropy

$$H(P, Q) = H(p) + D_{KL}(P||Q)$$

Correct probs

Softmax Classifier (Multinomial Logistic Regression)



Want to interpret raw classifier scores as probabilities

$$\mathbf{s} = f(\mathbf{x}_i; \mathbf{W})$$

$$P(Y = k | X = \mathbf{x}_i) = \frac{e^{s_k}}{\sum_j e^{s_j}}$$

Softmax
Function

Maximize probability of correct class

$$L_i = -\log P(Y = y_i | X = \mathbf{x}_i)$$

Putting it all together:

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

cat	3.2
car	5.1
frog	-1.7

Lecture 5: Image Classification with CNNs

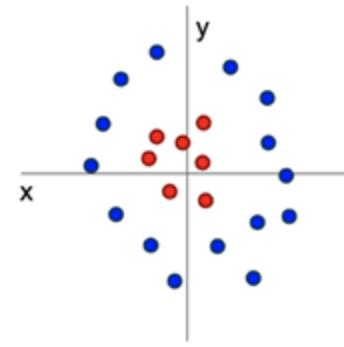
Problem: Linear Classifiers are not very powerful

Visual Viewpoint



Linear classifiers learn
one template per class

Geometric Viewpoint



Linear classifiers can
only draw linear
decision boundaries

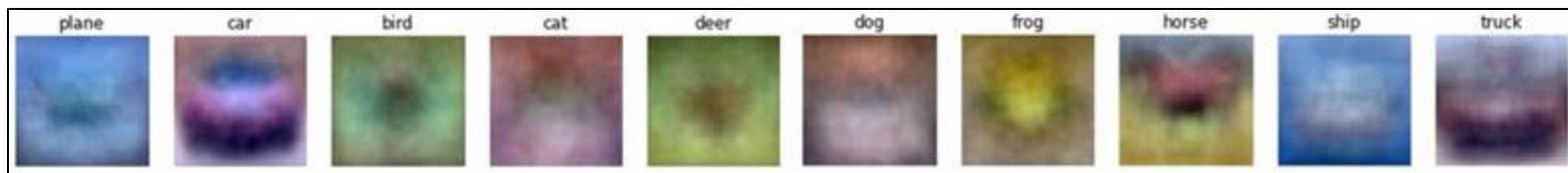
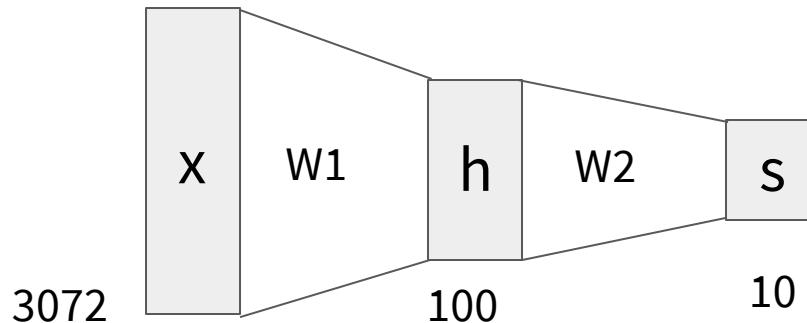
Last time: Neural Networks

Linear score function:

$$f = Wx$$

2-layer Neural Network

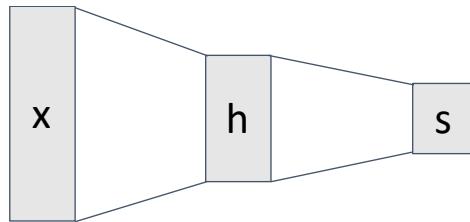
$$f = W_2 \max(0, W_1 x)$$



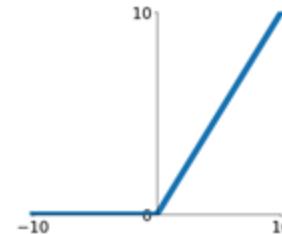
Today: Convolutional Networks

We have
already
seen these

Fully-Connected Layer



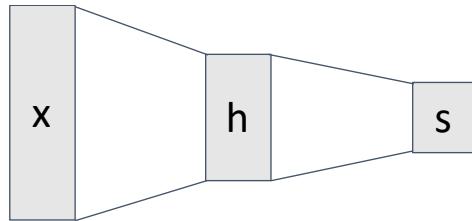
Activation Function



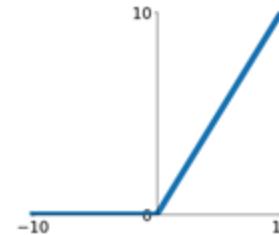
Today: Convolutional Networks

We have
already
seen these

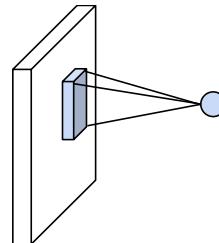
Fully-Connected Layer



Activation Function

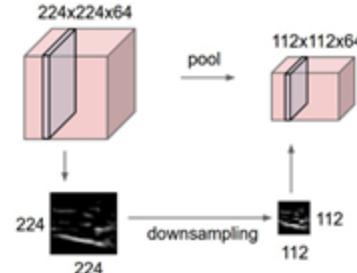


Convolution Layer



Today: Image-specific operators

Pooling Layer



Pixel space

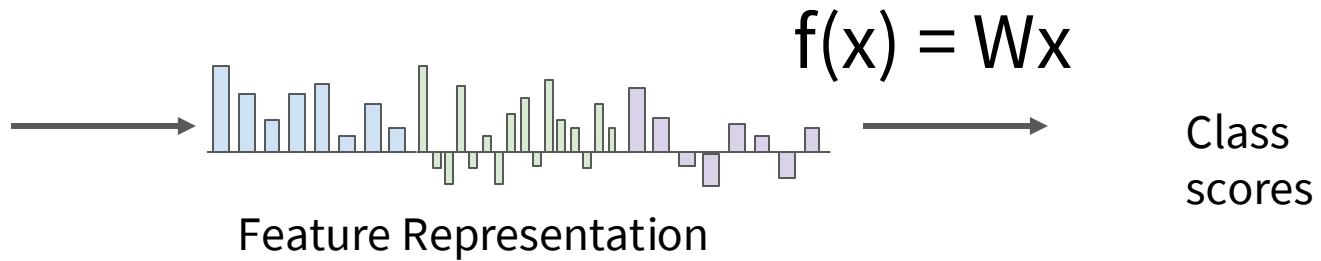


$$f(x) = Wx$$

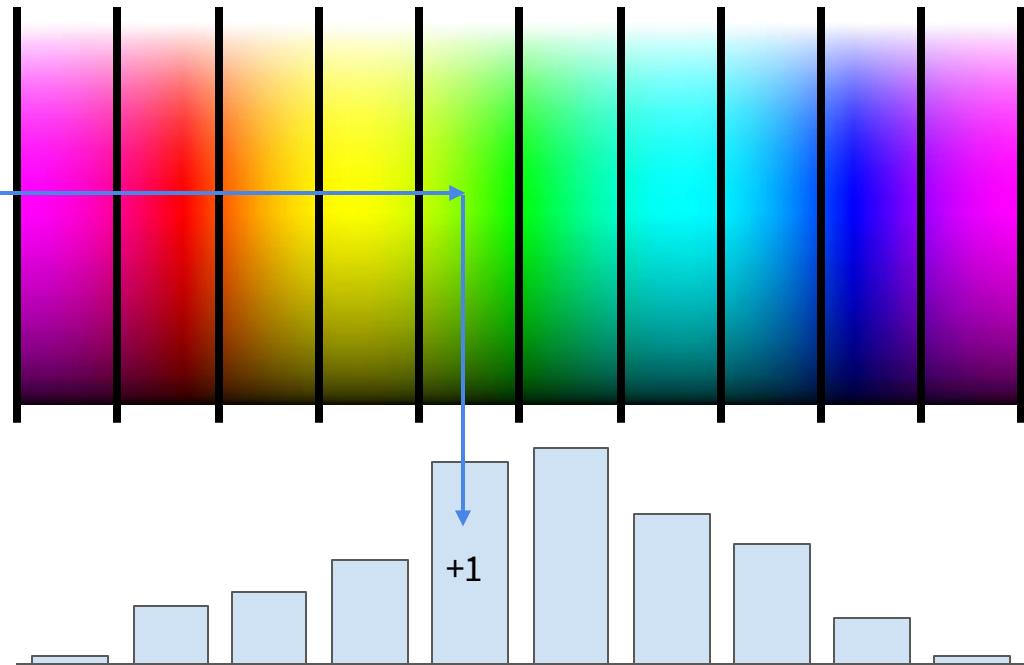
Class scores



Image features



Example: Color Histogram

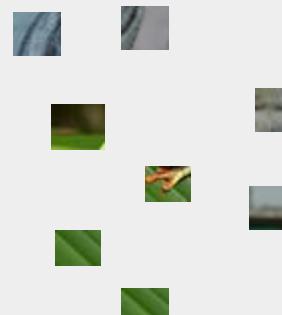


Example: Bag of Words

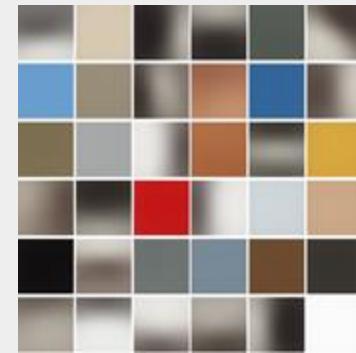
Step 1: Build codebook



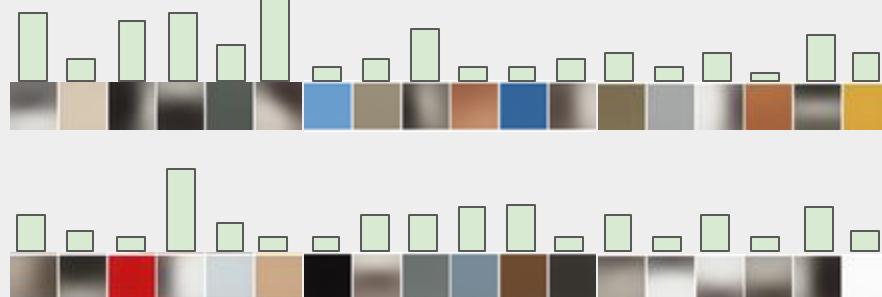
Extract random patches



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

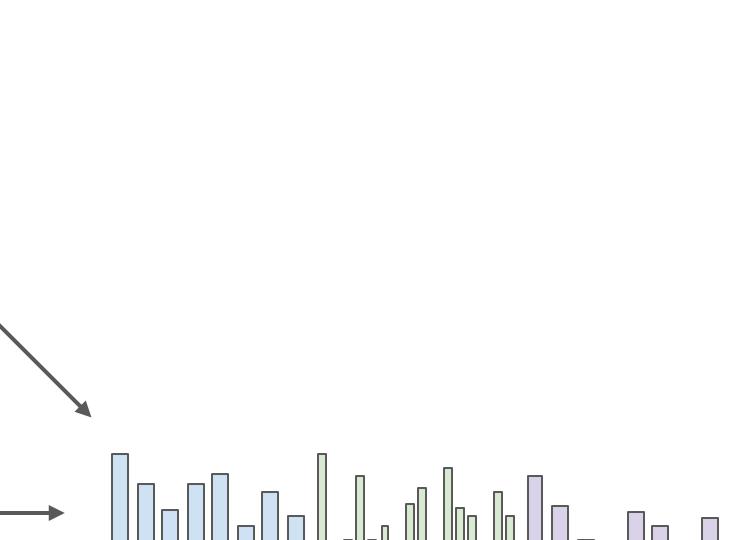
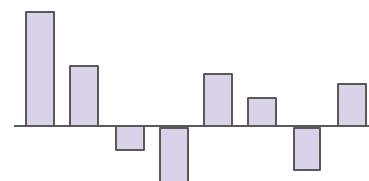
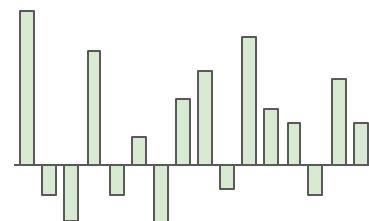
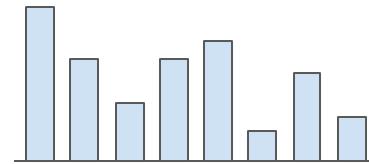
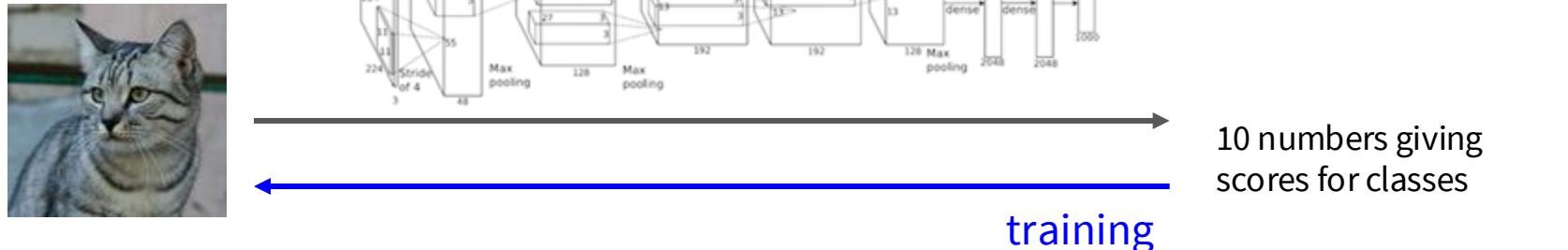
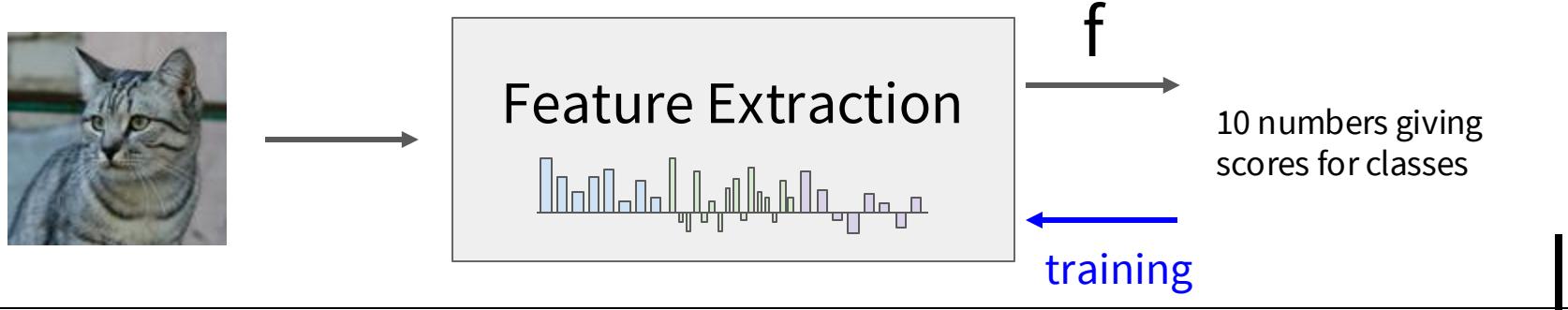


Image features vs. ConvNets



10 numbers giving
scores for classes

training

Last Time: Neural Networks

Linear score function:

$$f = Wx$$

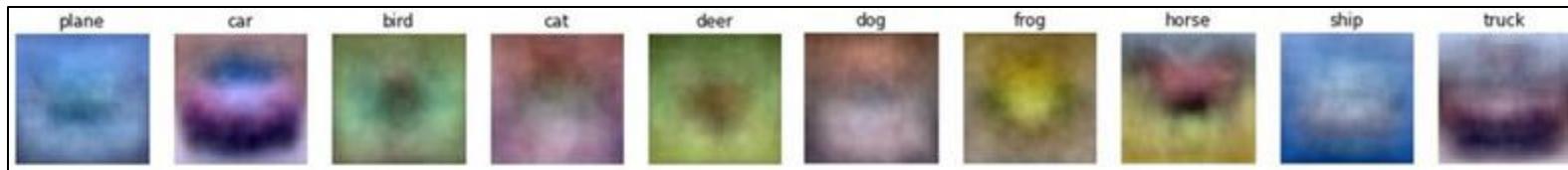
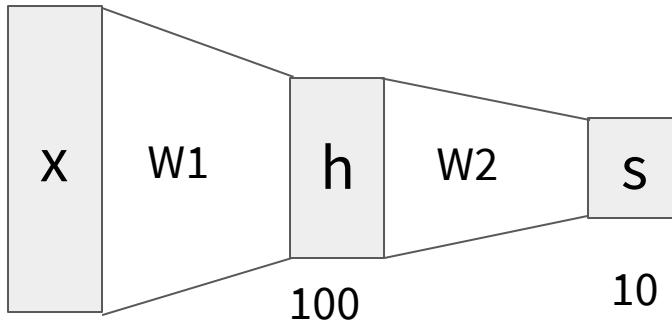
2-layer Neural Network

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

Problem: The spatial
structure of images
is destroyed!

32x32x3

3072



Next: Convolutional Neural Networks

Convolution and **pooling** operators extract features while respecting 2D image structure

Fully-Connected layers form an MLP at the end to predict scores

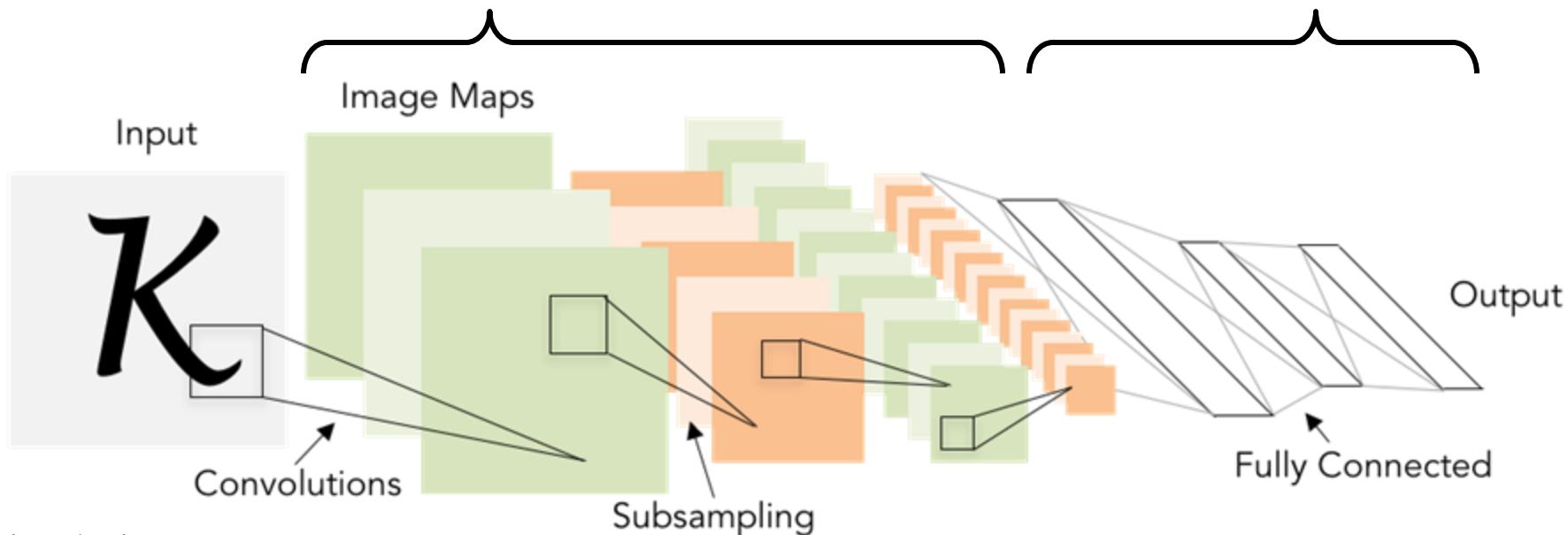


Illustration of LeCun et al. 1998 from CS231n 2017 Lecture 1

Next: Convolutional Neural Networks

Trained end-to-end with backprop + gradient descent

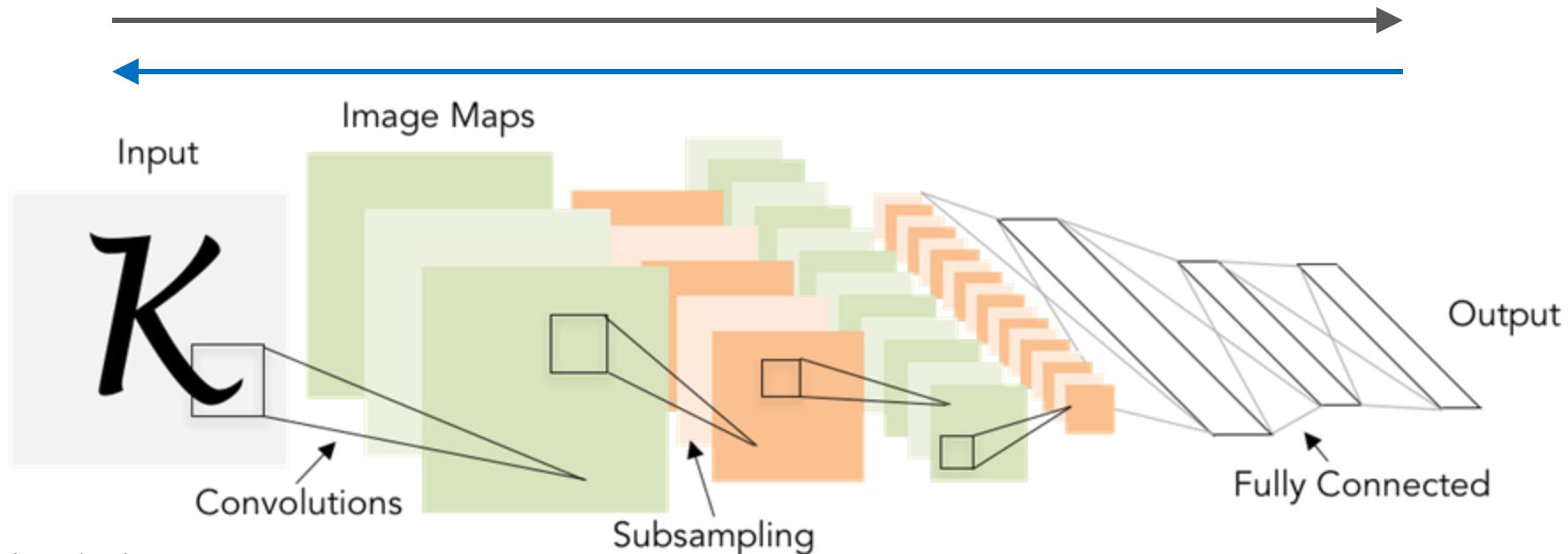
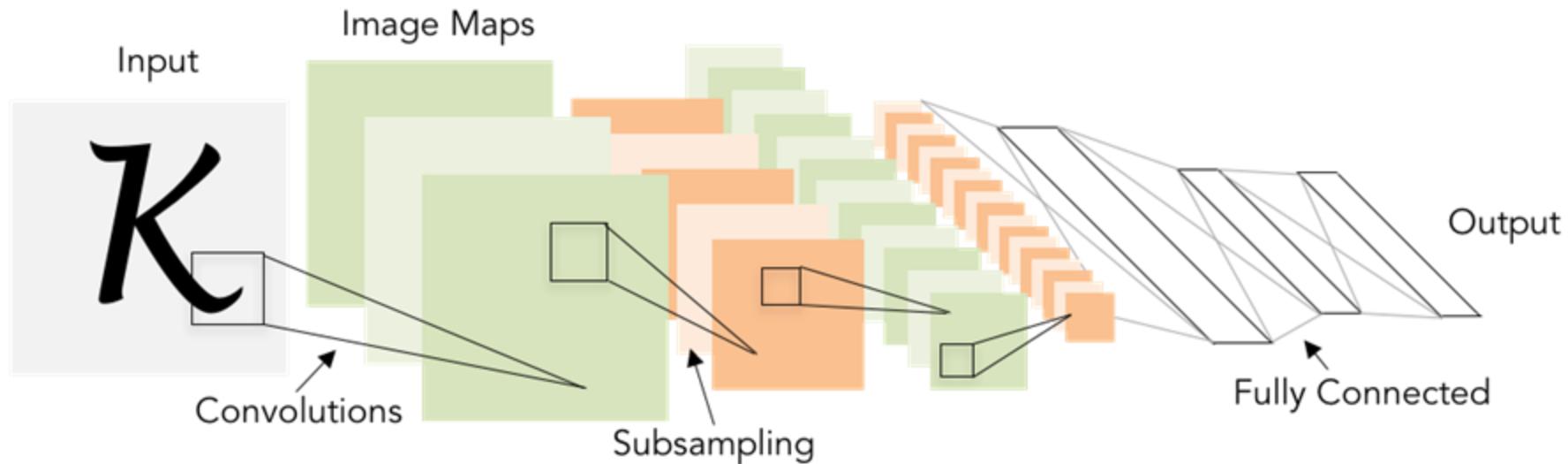


Illustration of LeCun et al. 1998 from CS231n 2017 Lecture 1

A bit of history:

Gradient-based learning applied to
document recognition

[LeCun, Bottou, Bengio, Haffner 1998]



A bit of history:

ImageNet Classification with Deep Convolutional Neural Networks

[Krizhevsky, Sutskever, Hinton, 2012]

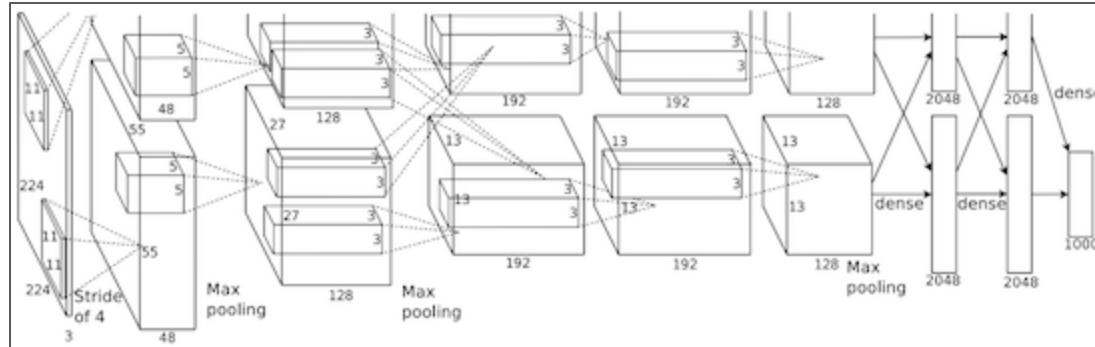
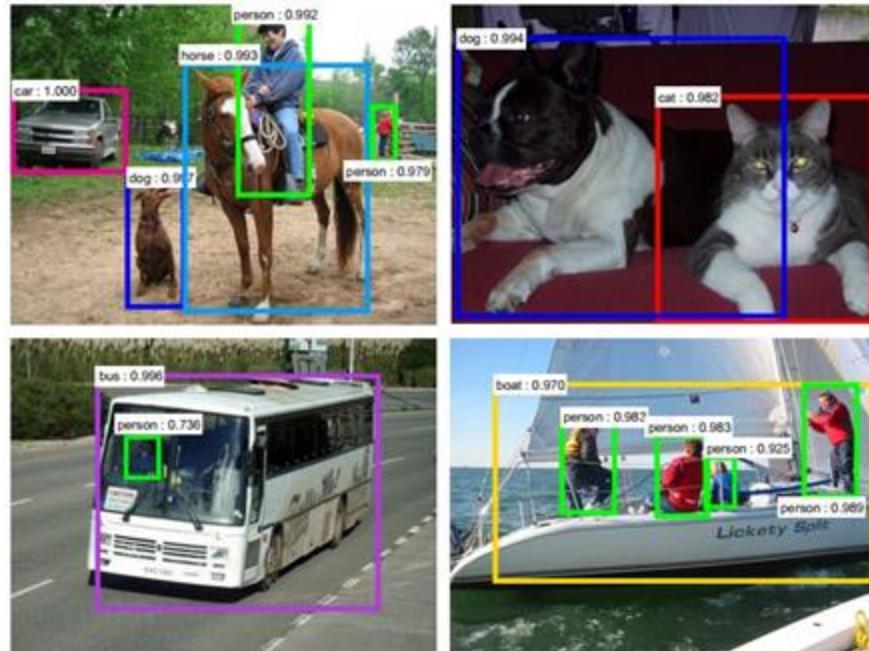


Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

“AlexNet”

~2012 – 2020: ConvNets dominate all vision tasks

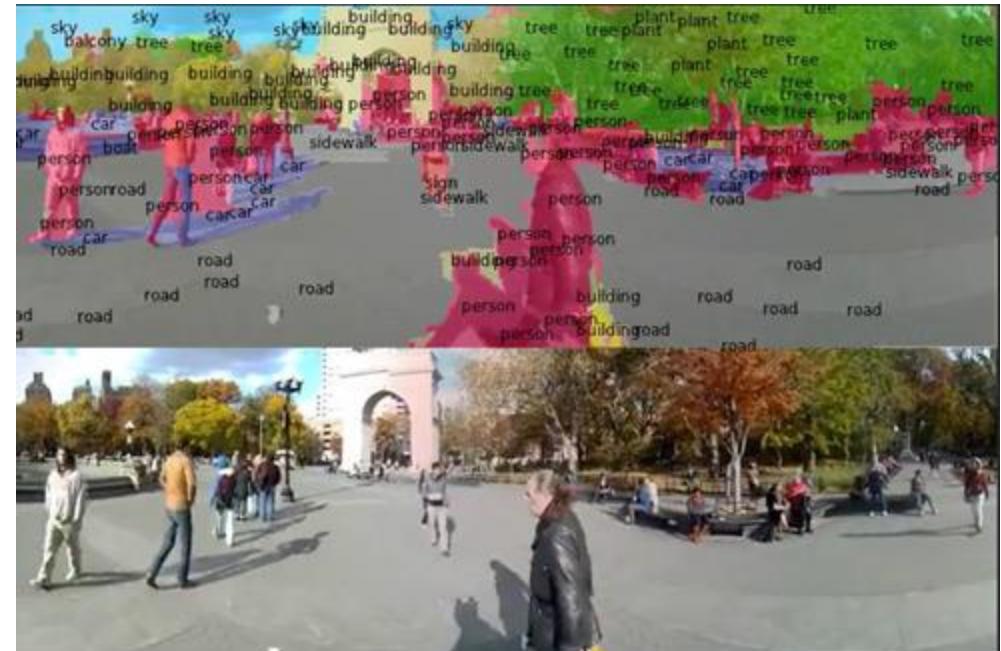
Detection



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

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

Detection



Figures copyright Clement Farabet, 2012. Reproduced with permission.

[Farabet et al., 2012]

~2012 – 2020: ConvNets dominate all vision tasks

Image Captioning



A white teddy bear sitting in the grass



A man in a baseball uniform throwing a ball



A woman is holding a cat in her hand



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](#)

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

~2012 – 2020: ConvNets dominate all vision tasks

Text-to-Image Generation

Rombach et al, “High-Resolution Image Synthesis with Latent Diffusion Models”, CVPR 2022



A zombie in the style of Picasso

An image of a half mouse half octopus

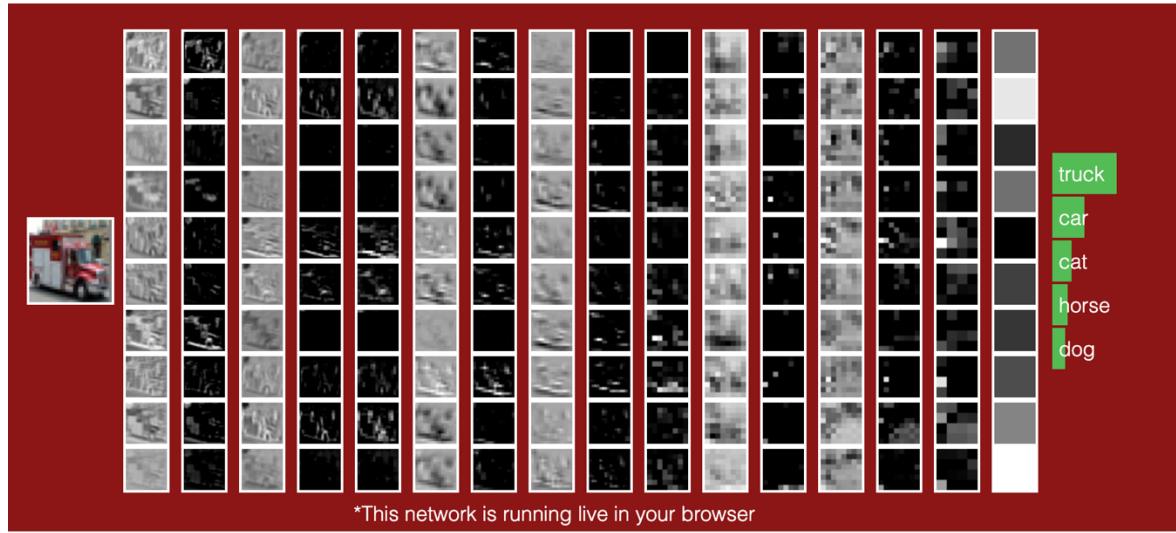
A painting of a squirrel eating a burger

A watercolor painting of a chair that looks like an octopus

A shirt with the inscription: “I love generative models!”

~2012 – 2020: ConvNets dominate all vision tasks

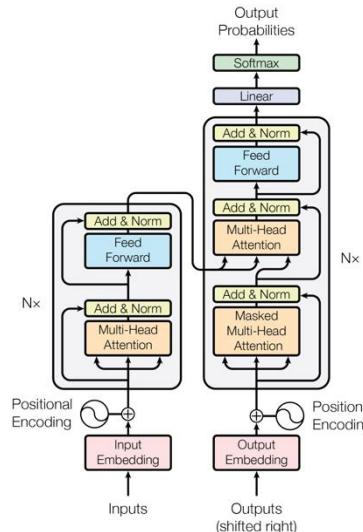
CS231n: Convolutional Neural Networks for Visual Recognition



This class used to be focused on ConvNets!

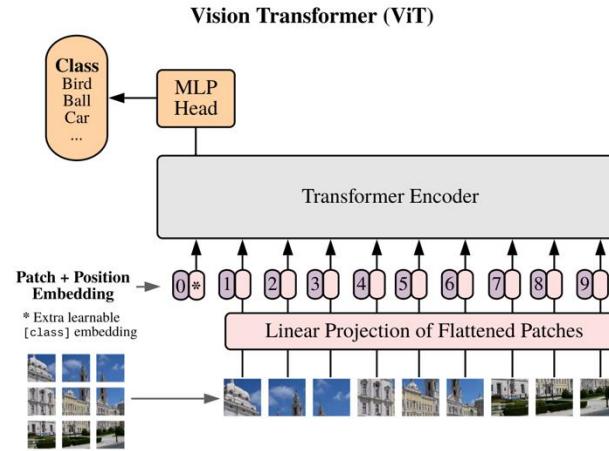
2021 - Present: Transformers have taken over

2017: Transformers
for language tasks

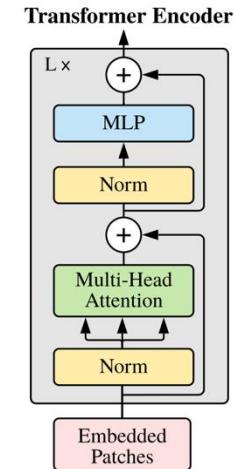


Vaswani et al, "Attention is all you need", NeurIPS 2017

2021: Transformers
for vision tasks



Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021

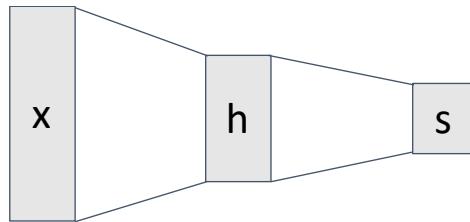


Convolutional Neural Networks

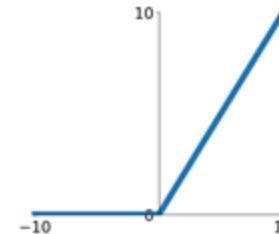
Today: Convolutional Networks

We have
already
seen these

Fully-Connected Layer

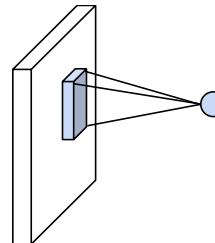


Activation Function

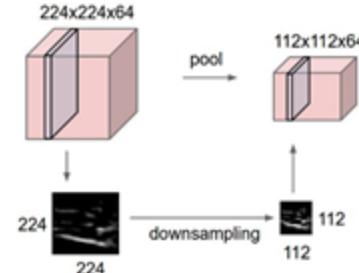


Convolution Layer

Today: Image-specific operators



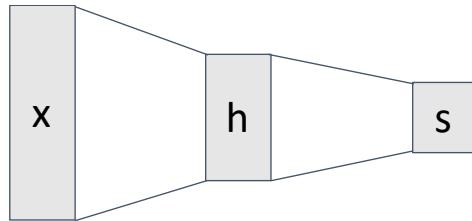
Pooling Layer



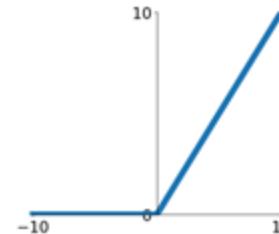
Today: Convolutional Networks

We have
already
seen these

Fully-Connected Layer

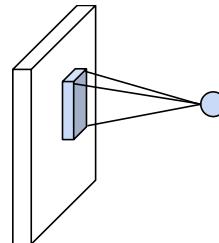


Activation Function

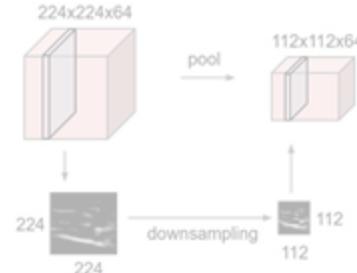


Convolution Layer

Today: Image-specific operators

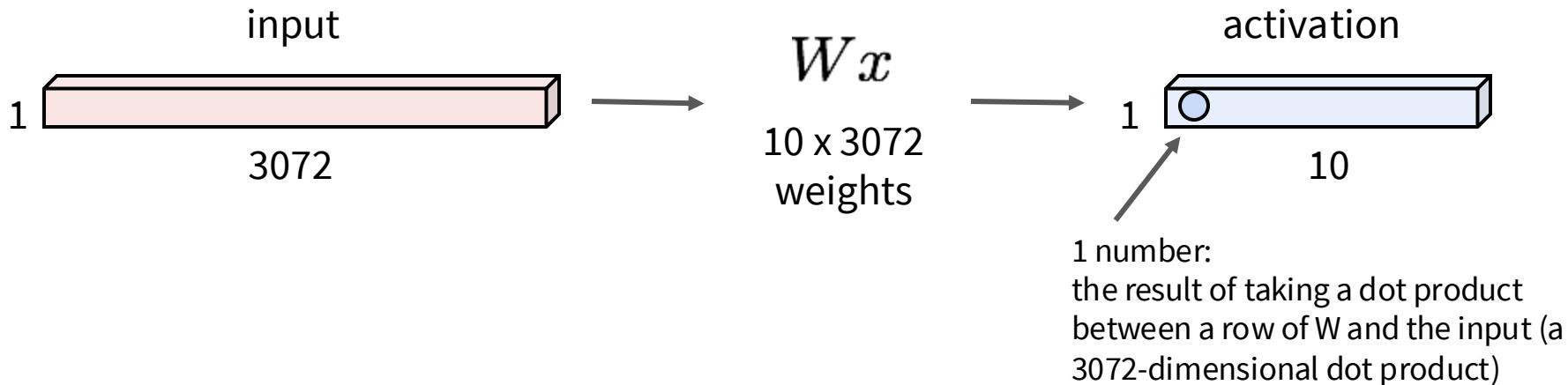


Pooling Layer



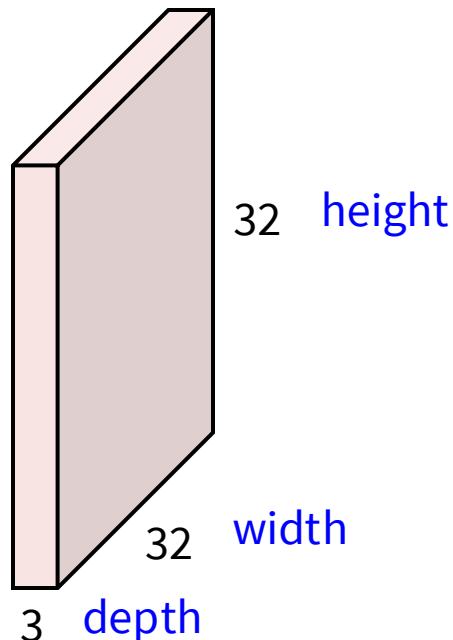
Fully Connected Layer

32x32x3 image -> stretch to 3072 x 1



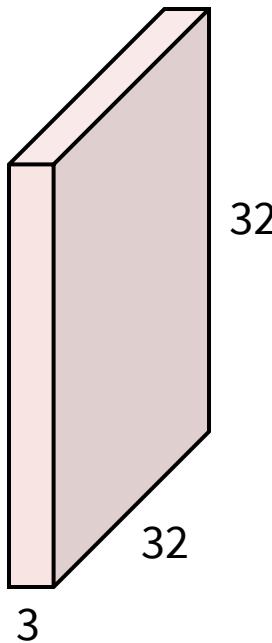
Convolution Layer

32x32x3 image -> preserve spatial structure



Convolution Layer

32x32x3 image



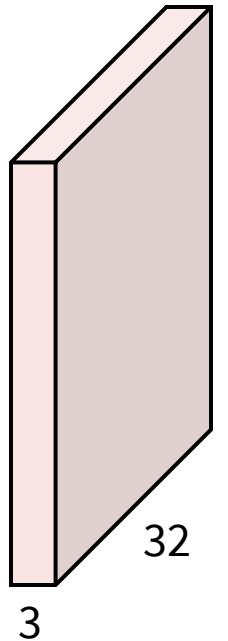
5x5x3 filter



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

Convolution Layer

32x32x3 image



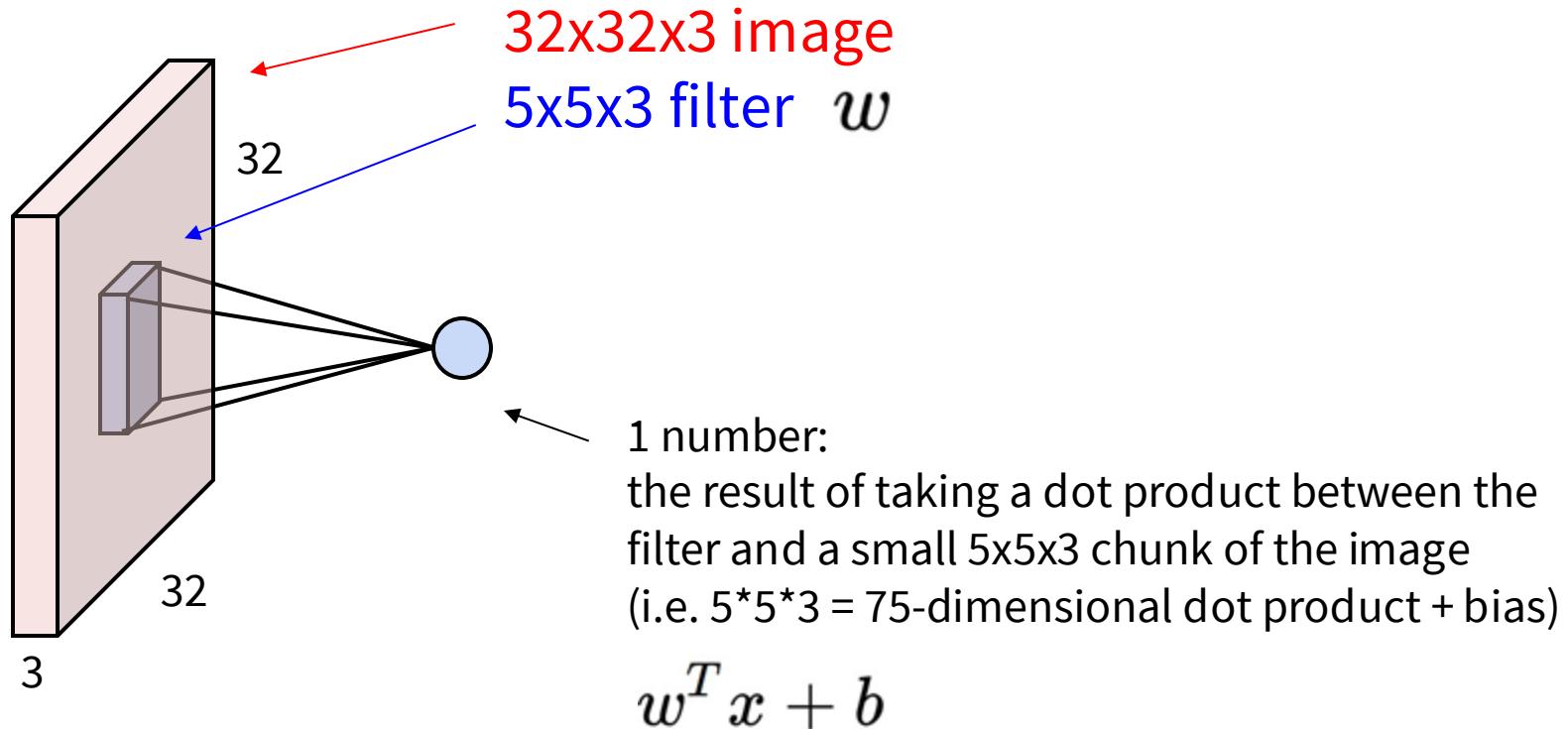
5x5x3 filter



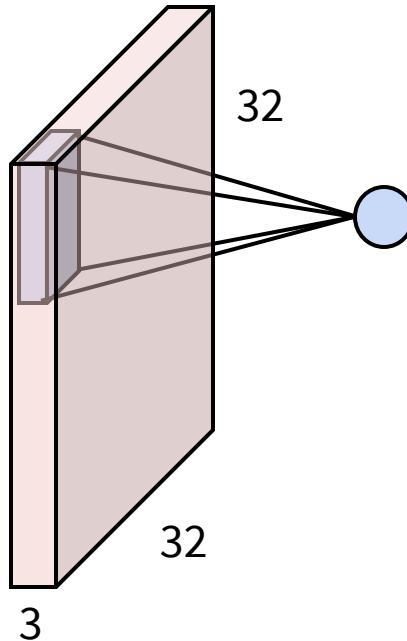
Filters always extend the full depth of the input volume

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

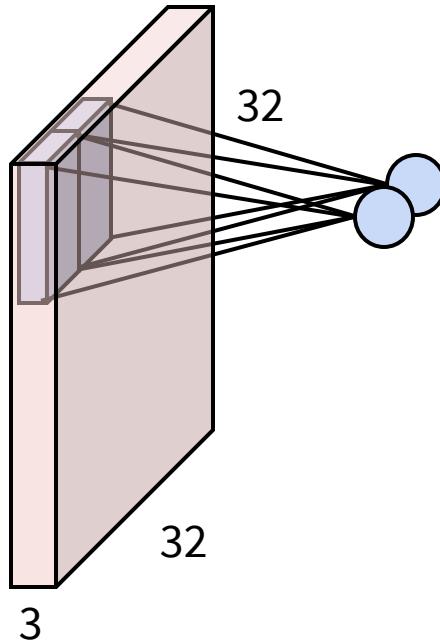
Convolution Layer



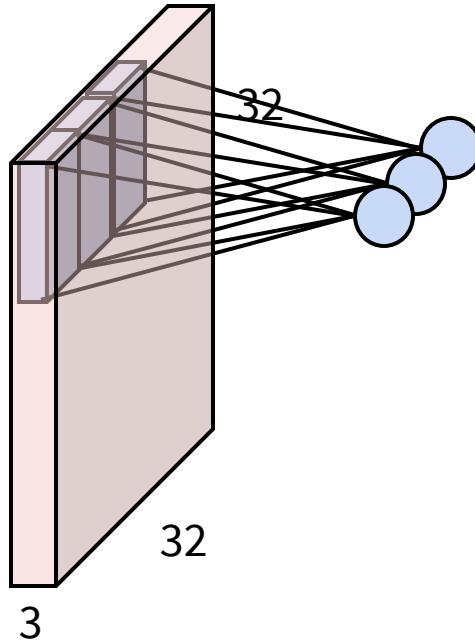
Convolution Layer



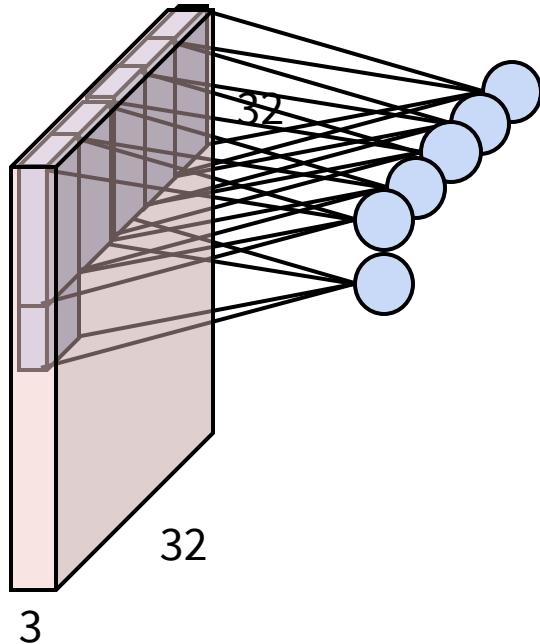
Convolution Layer



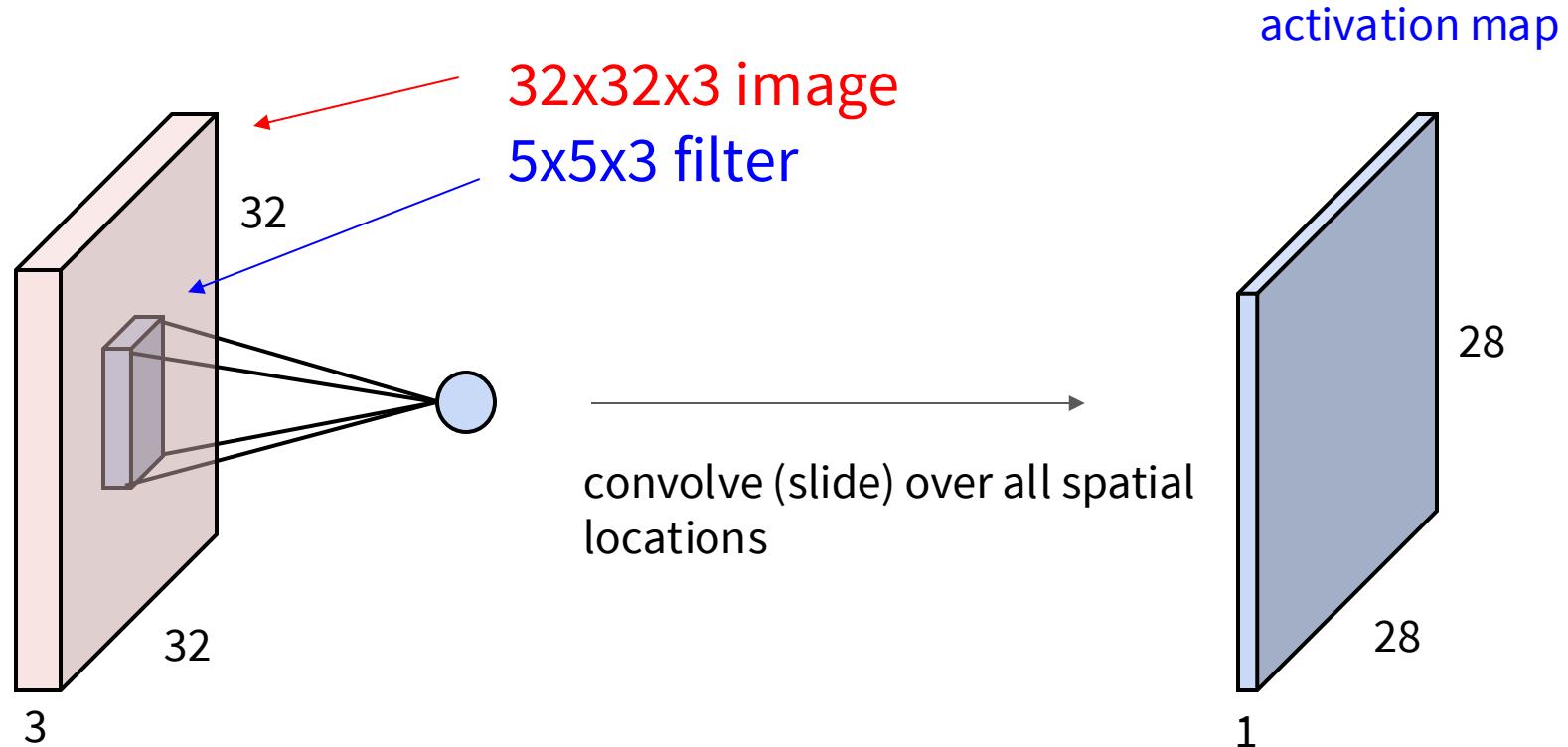
Convolution Layer



Convolution Layer

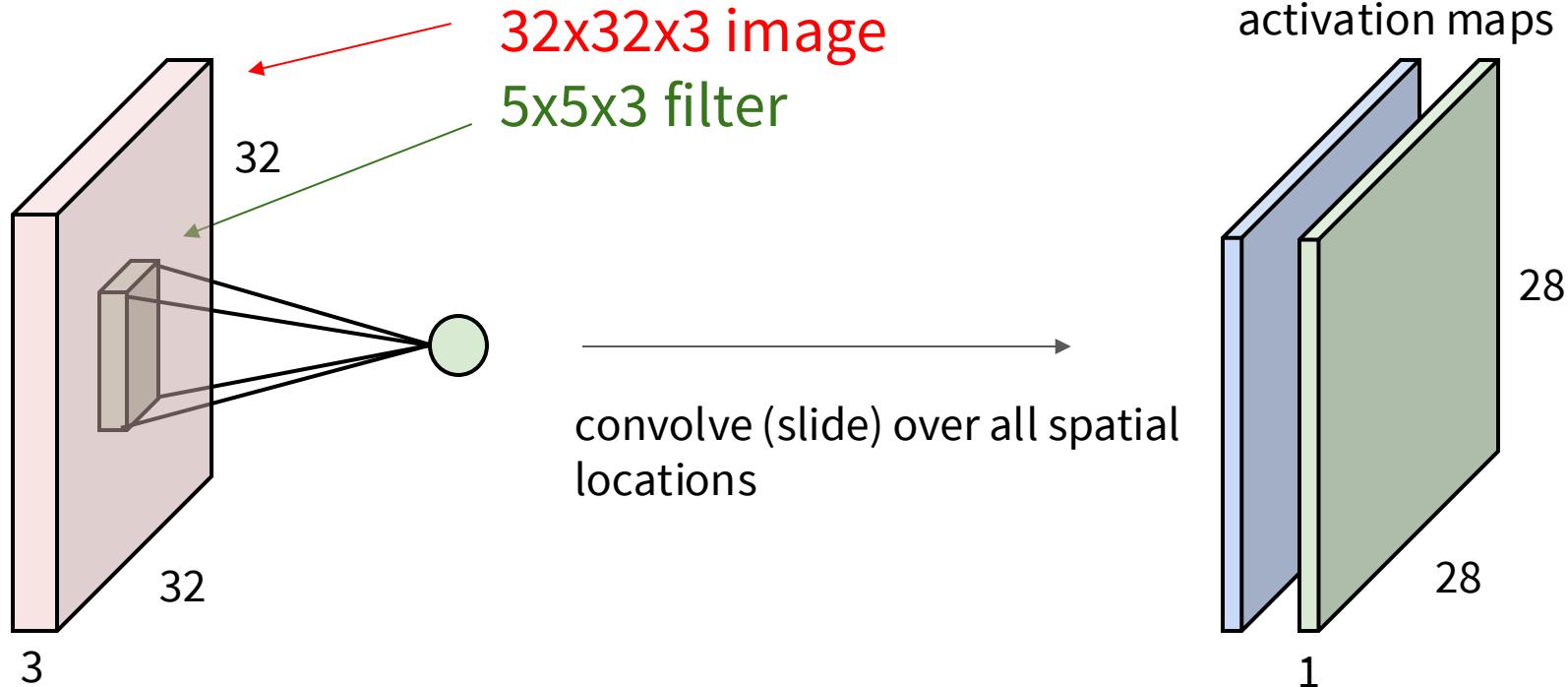


Convolution Layer



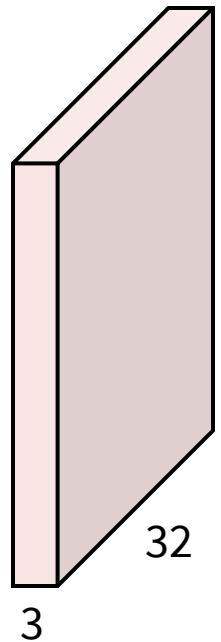
Convolution Layer

consider a second, green filter

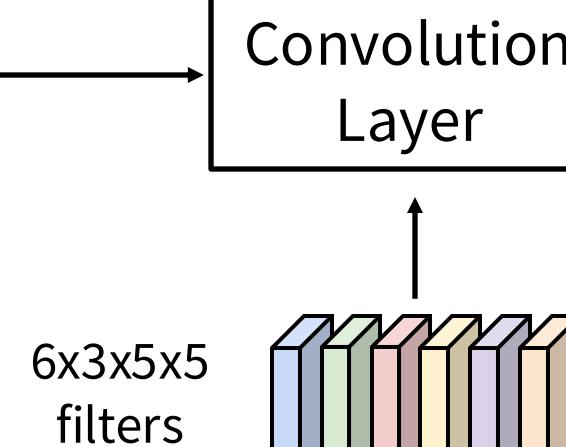


Convolution Layer

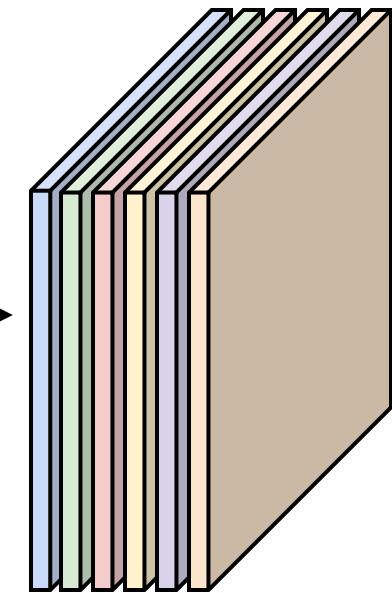
3x32x32 image



Consider 6 filters,
each 3x5x5



6 activation maps,
each 1x28x28

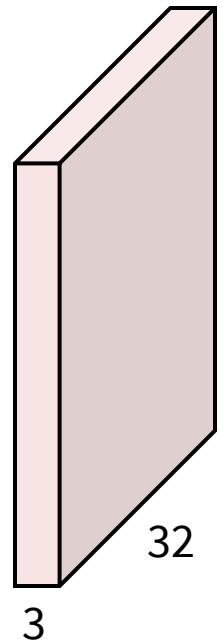


Stack activations to get a
6x28x28 output image!

Slide inspiration: Justin Johnson

Convolution Layer

3x32x32 image

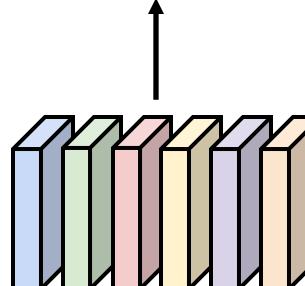


Also 6-dim bias vector:

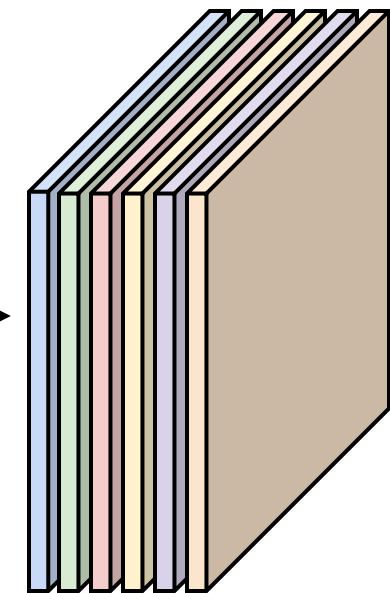


Convolution
Layer

6x3x5x5
filters



6 activation maps,
each 1x28x28



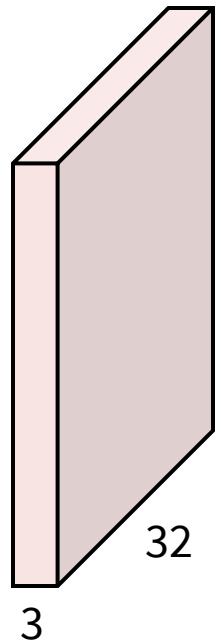
Stack activations to get a
6x28x28 output image!

Slide inspiration: Justin Johnson

Convolution Layer

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

3x32x32 image

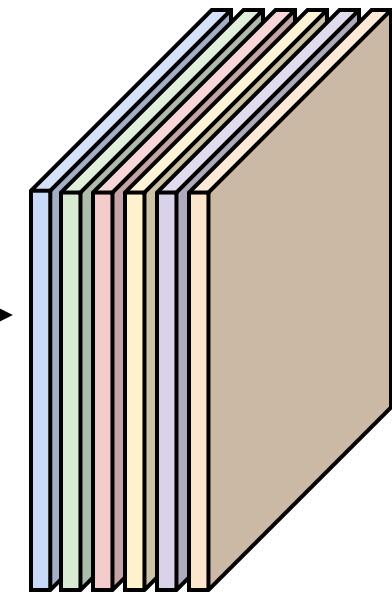
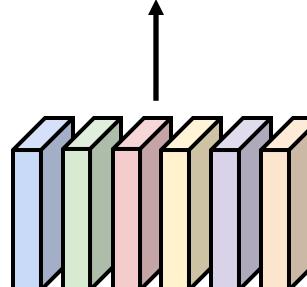


Also 6-dim bias vector:



Convolution
Layer

6x3x5x5
filters

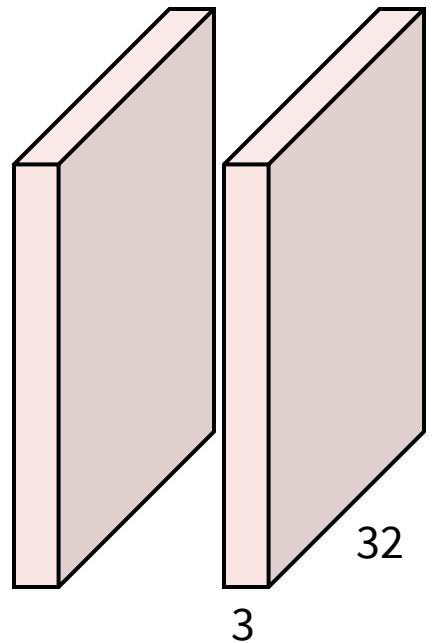


Stack activations to get a 6x28x28 output image!

Slide inspiration: Justin Johnson

Convolution Layer

$2 \times 3 \times 32 \times 32$
Batch of images

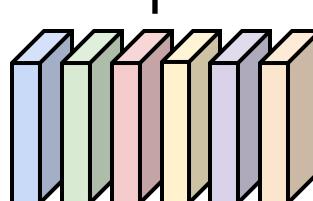


Also 6-dim bias vector:



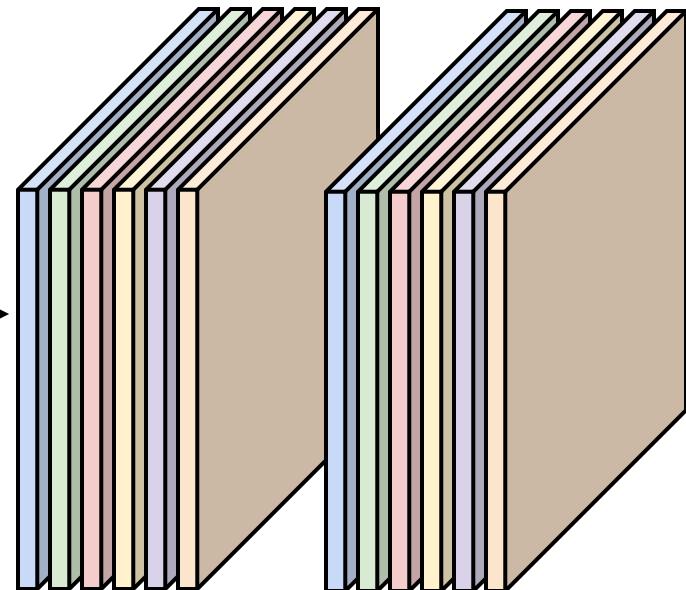
Convolution
Layer

$6 \times 3 \times 5 \times 5$
filters



$2 \times 6 \times 28 \times 28$

Batch of outputs

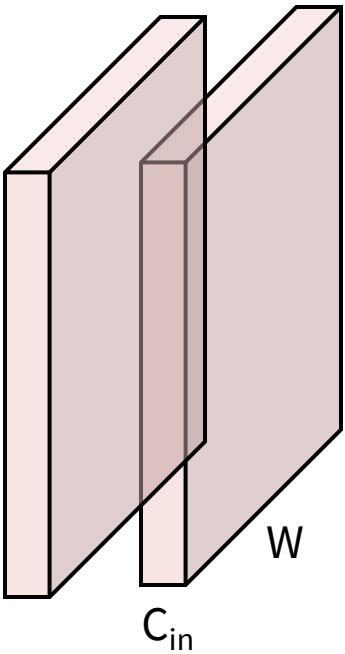


Slide inspiration: Justin Johnson

Convolution Layer

$N \times C_{in} \times H \times W$

Batch of images

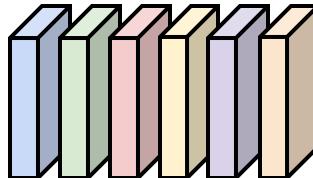


Also C_{out} -dim bias vector:

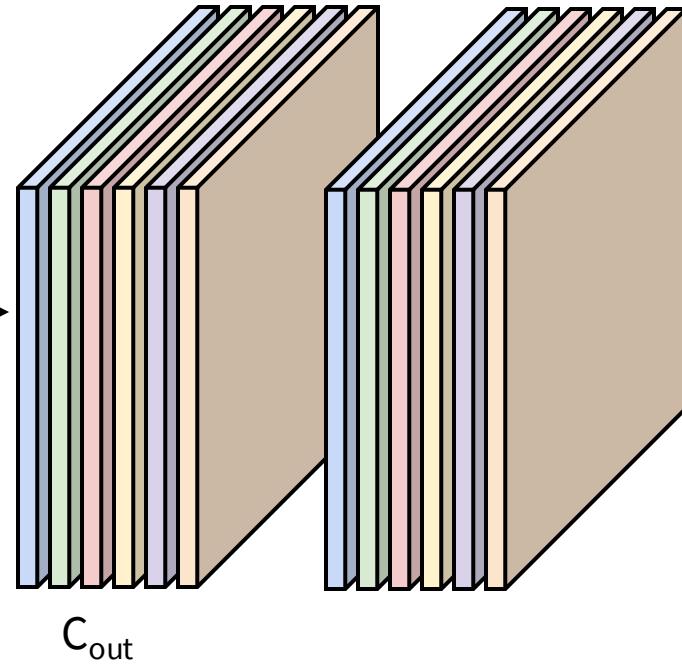


Convolution
Layer

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

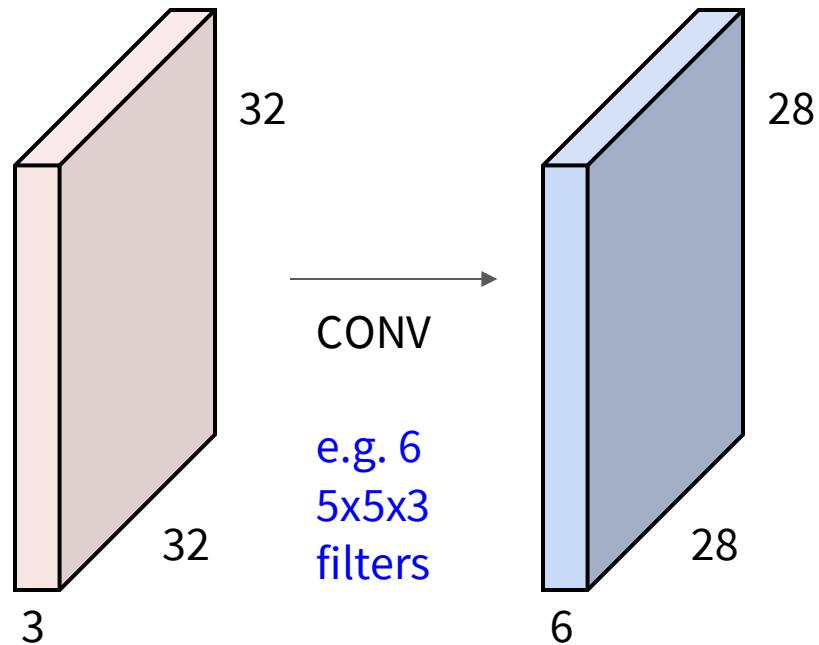


$N \times C_{out} \times H' \times W'$
Batch of outputs

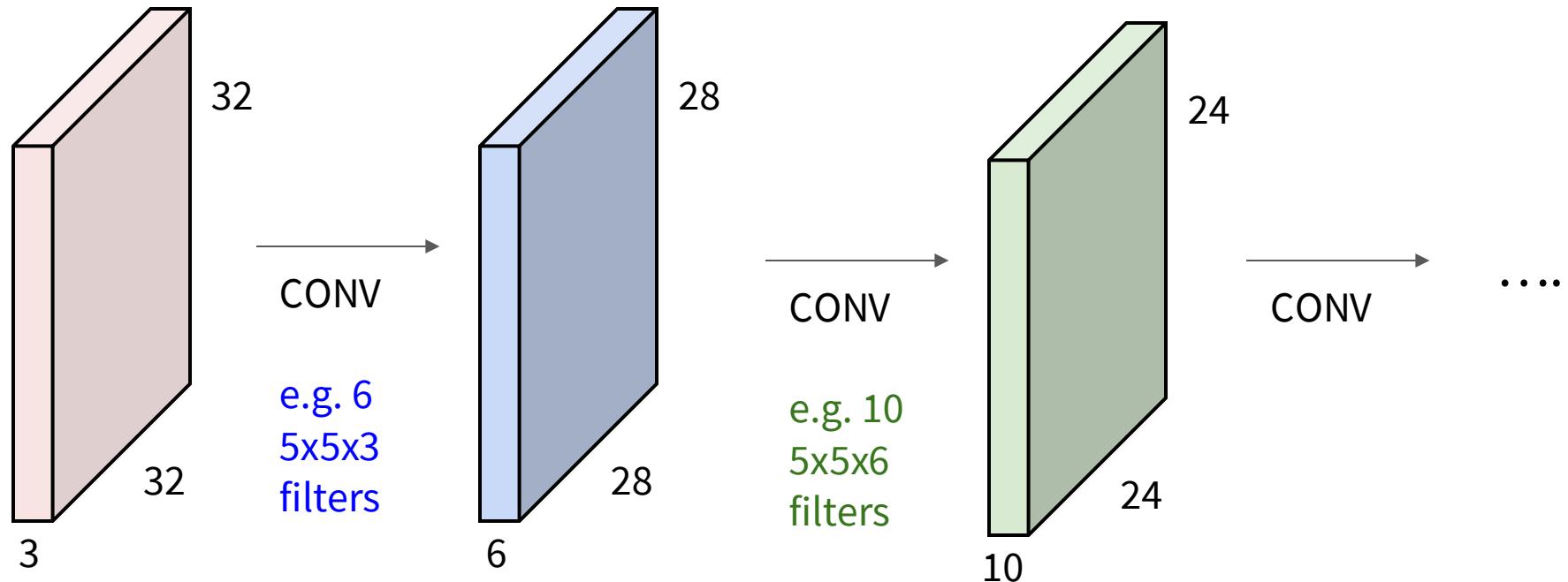


Slide inspiration: Justin Johnson

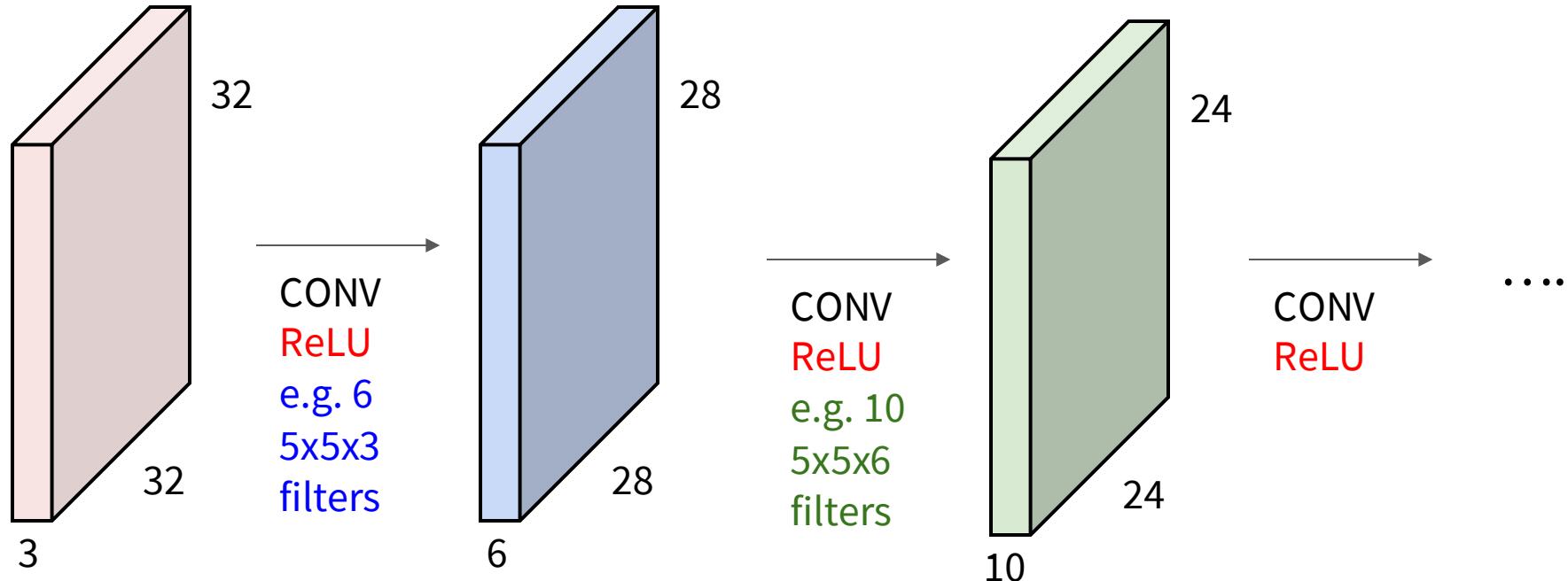
A ConvNet is a neural network with Conv layers



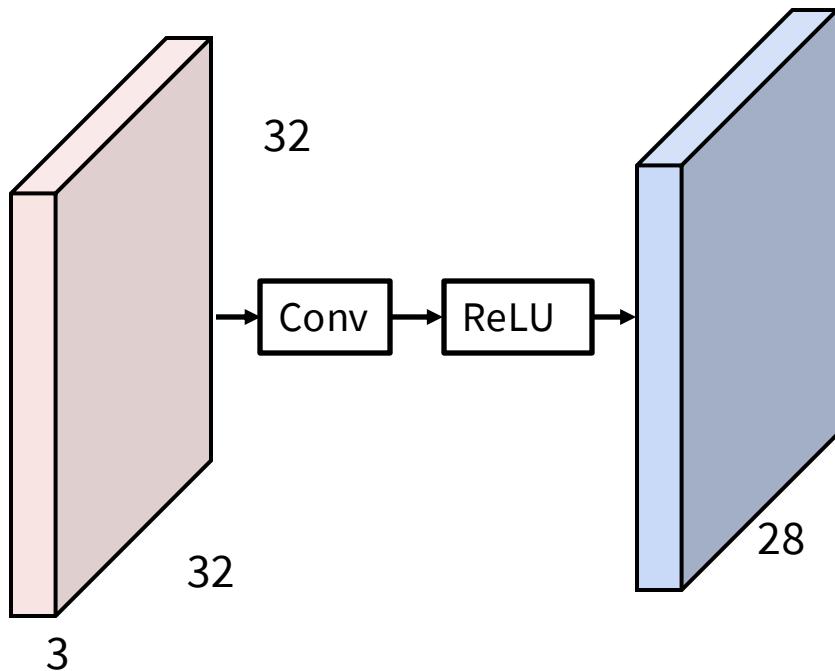
A ConvNet is a neural network with Conv layers



A ConvNet is a neural network with Conv layers with activation functions!



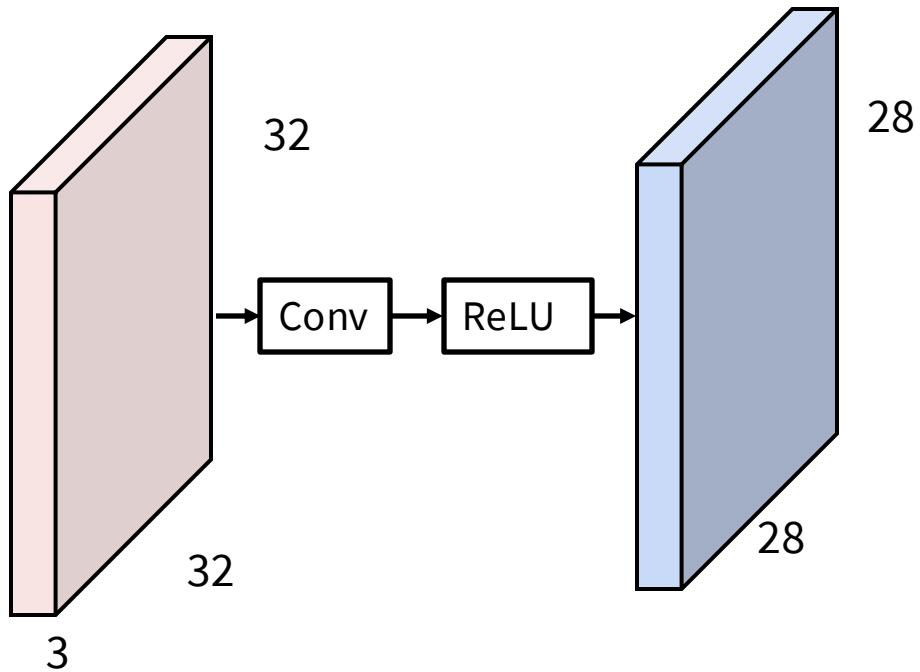
What do Conv filters learn?



Linear classifier: One template per class



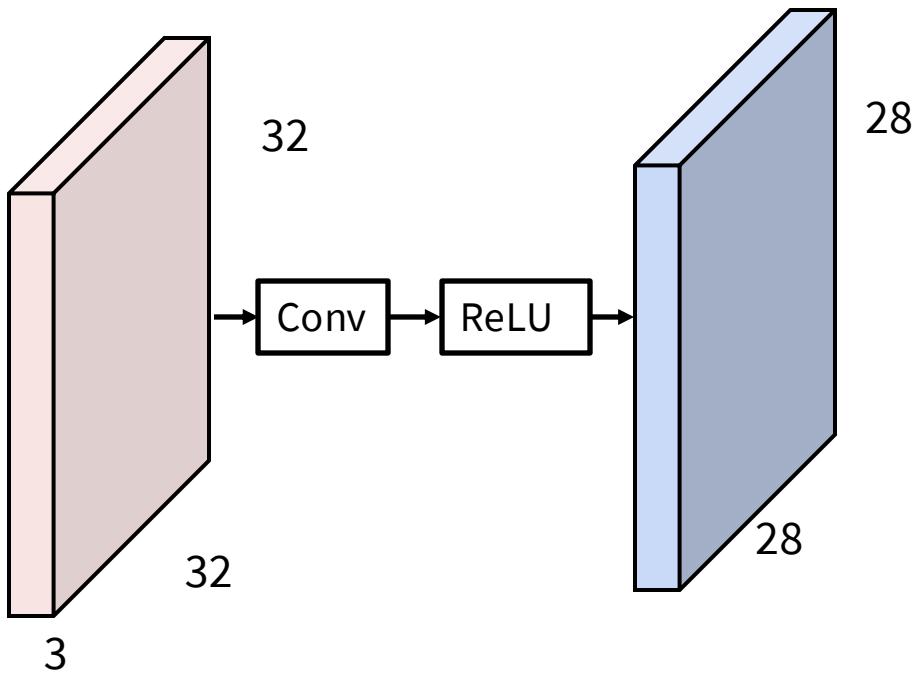
What do Conv filters learn?



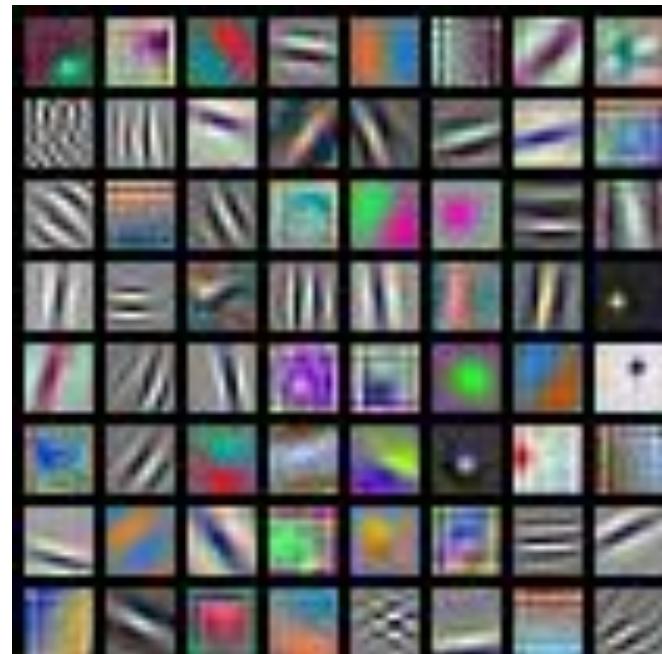
MLP: Bank of whole-image templates



What do Conv filters learn?

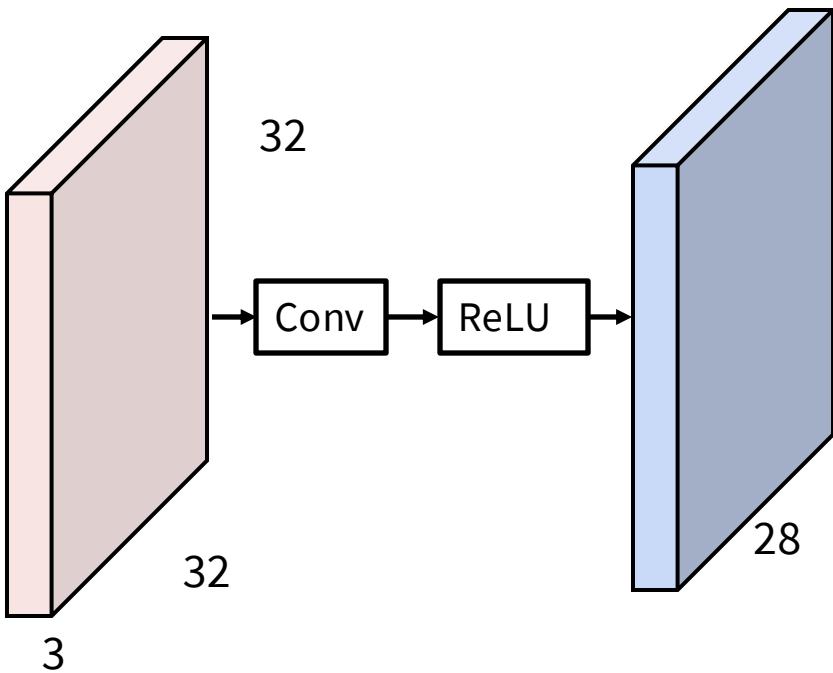


First-layer conv filters: local image templates
(Often learns oriented edges, opposing colors)

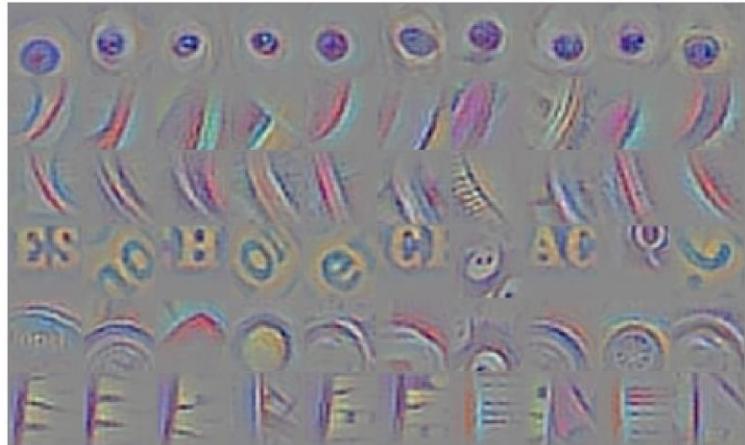


AlexNet: 64 filters, each 3x11x11

What do Conv filters learn?



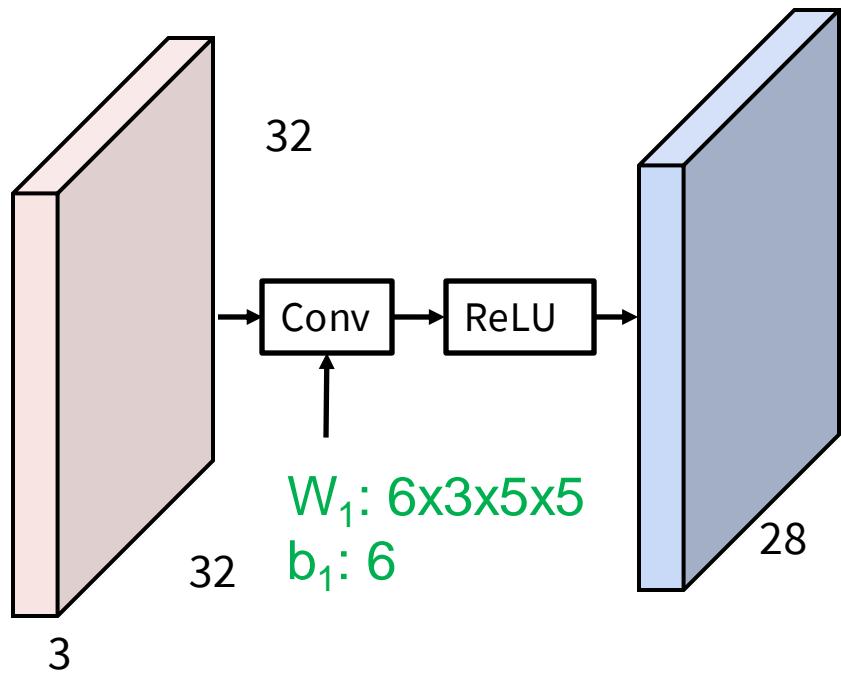
Deeper conv layers: Harder to visualize
Tend to learn larger structures e.g. eyes, letters



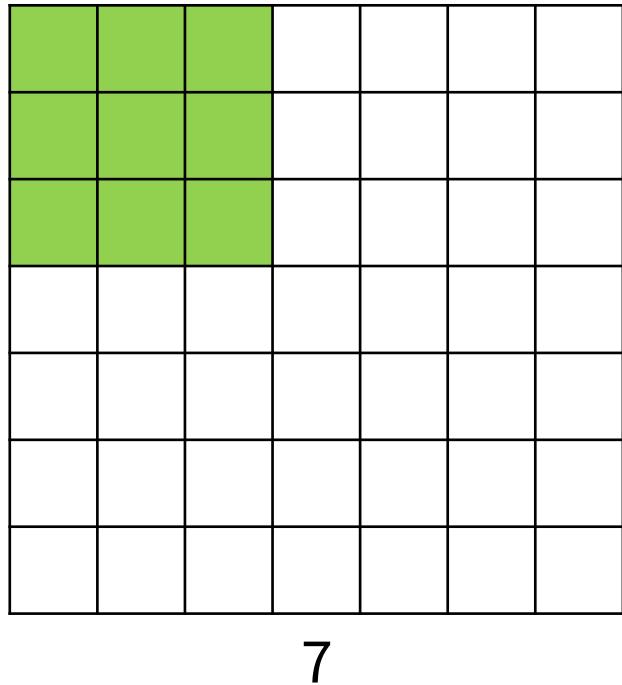
6th layer conv layer from an ImageNet model

Visualization from [Springenberg et al, ICLR 2015]

Convolution: Spatial Dimensions

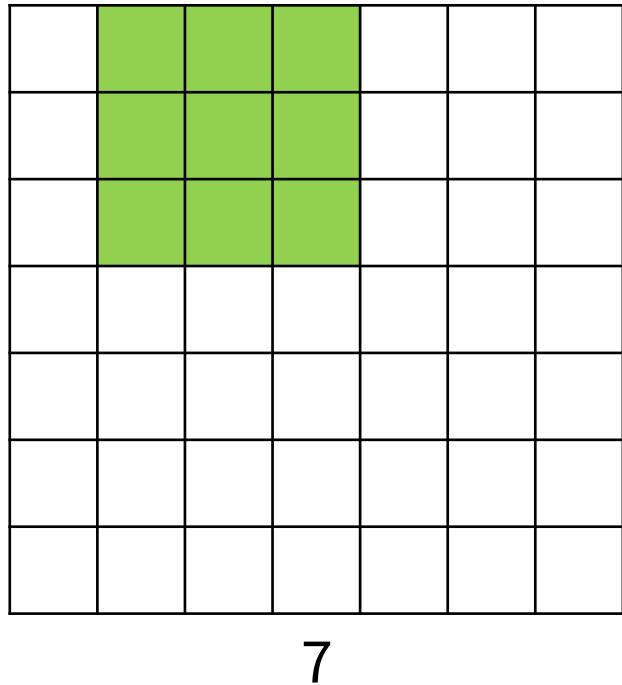


Convolution: Spatial Dimensions



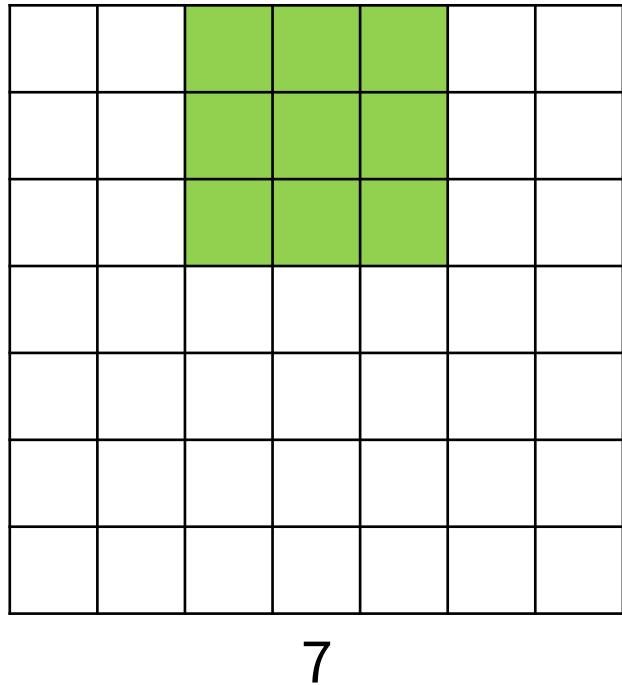
Input: 7x7
Filter: 3x3

Convolution: Spatial Dimensions



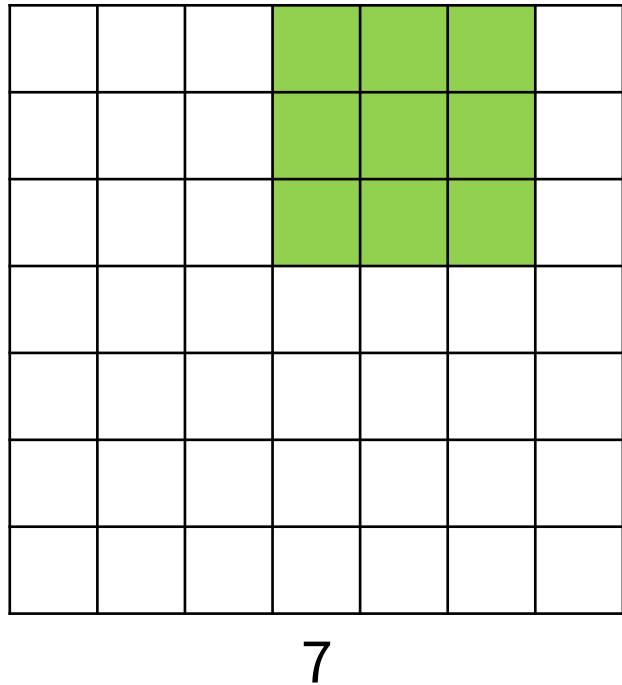
Input: 7x7
Filter: 3x3

Convolution: Spatial Dimensions



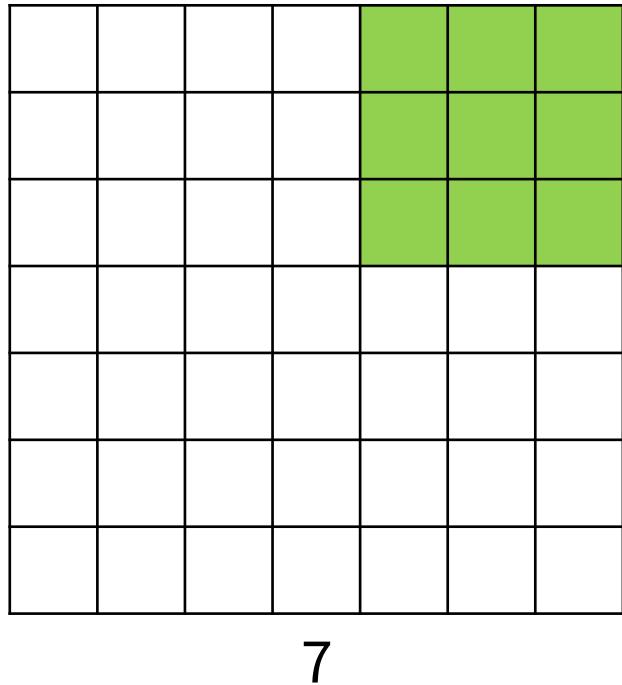
Input: 7x7
Filter: 3x3

Convolution: Spatial Dimensions



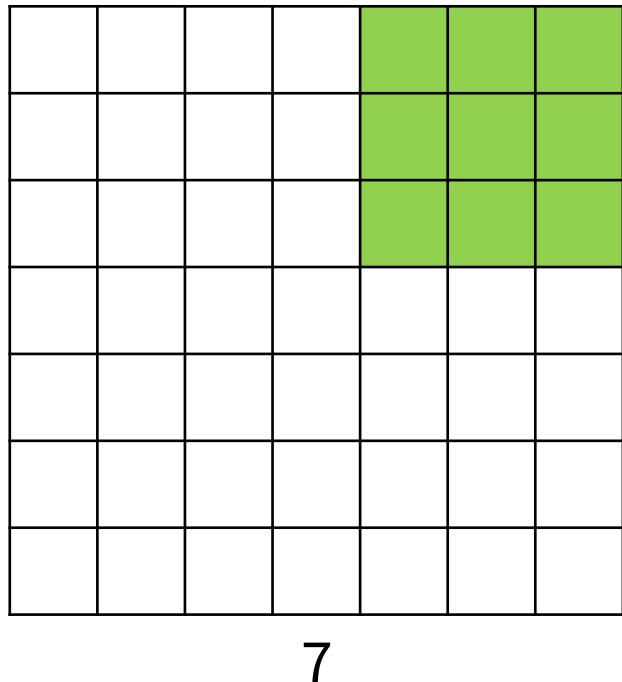
Input: 7x7
Filter: 3x3

Convolution: Spatial Dimensions



Input: 7x7
Filter: 3x3
Output: 5x5

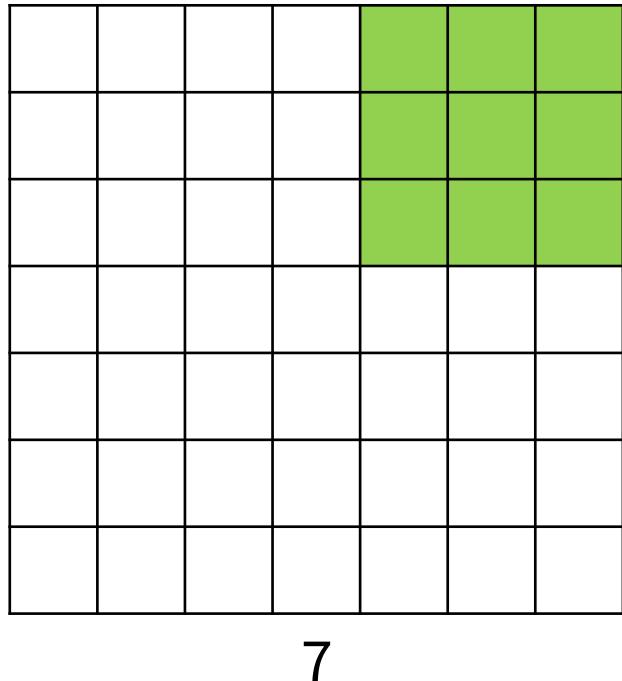
Convolution: Spatial Dimensions



Input: 7x7
Filter: 3x3
Output: 5x5

In general
Input: W
Filter: K
Output: $W - K + 1$

Convolution: Spatial Dimensions



Input: 7x7
Filter: 3x3
Output: 5x5

Problem: Feature maps shrink with each layer!

In general
Input: W
Filter: K
Output: $W - K + 1$

Convolution: Spatial Dimensions

0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

Input: 7x7

Filter: 3x3

Output: 5x5

In general

Input: W

Filter: K

Padding: P

Output: $W - K + 1 + 2P$

Problem: Feature maps shrink with each layer!

Solution: Add **padding** around the input before sliding the filter

Convolution: Spatial Dimensions

0	0	0	0	0	0	0	0	0
0								0
0								0
0								0
0								0
0								0
0								0
0								0
0	0	0	0	0	0	0	0	0

Input: 7x7

Filter: 3x3

Output: 5x5

In general

Input: W

Filter: K

Padding: P

Output: $W - K + 1 + 2P$

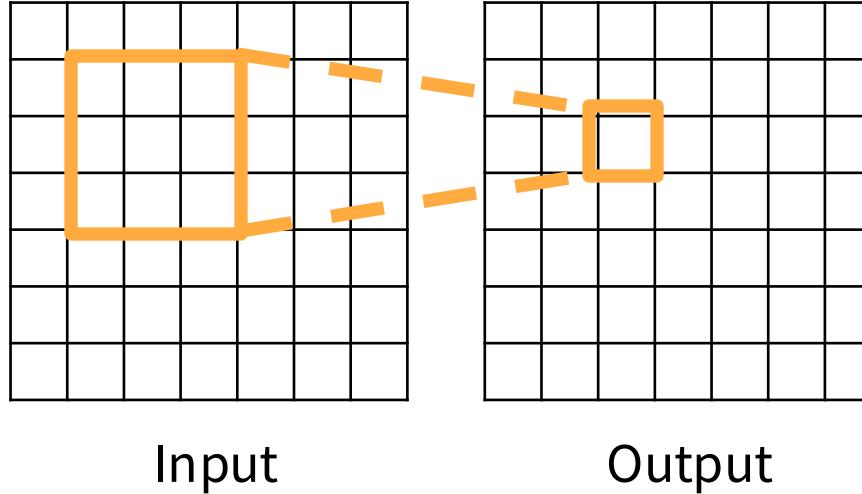
Common setting:

$$P = (K - 1) / 2$$

Means output has
same size as input

Receptive Fields

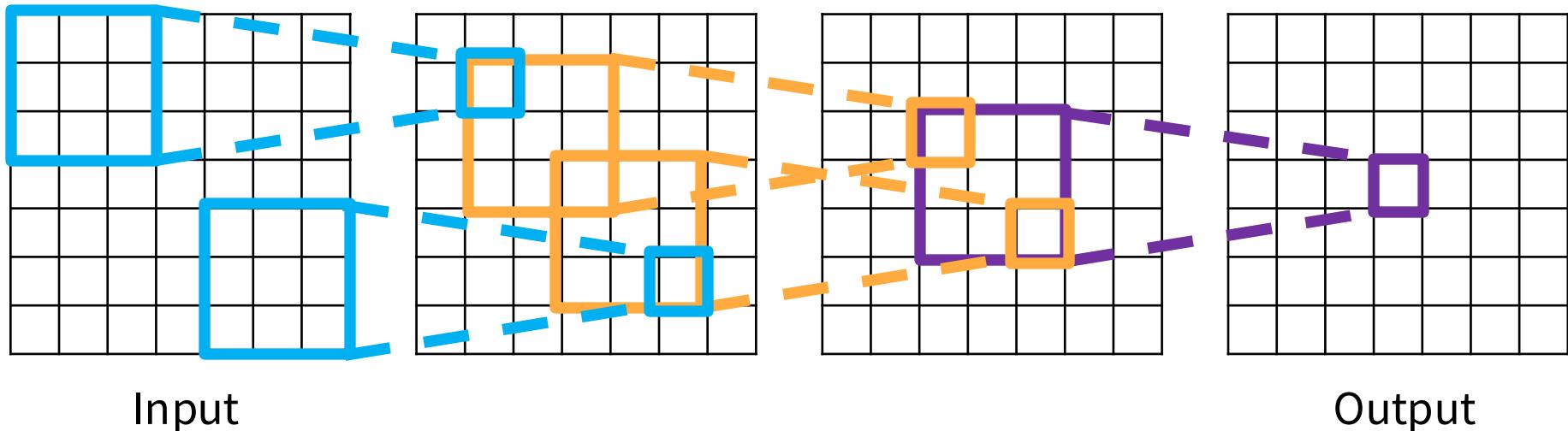
For convolution with **kernel size K**, each element in the output depends on a $K \times K$ receptive field in the input



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

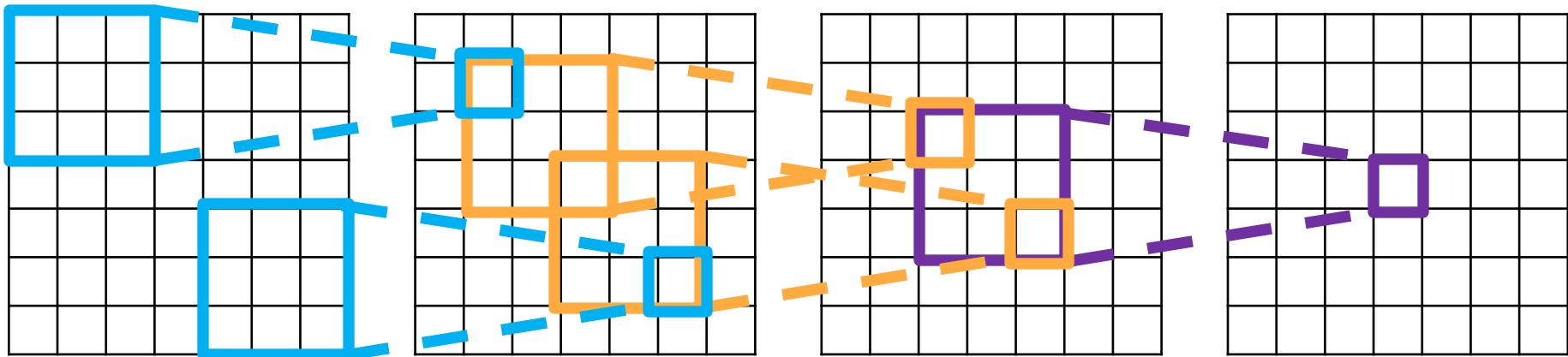
Output

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

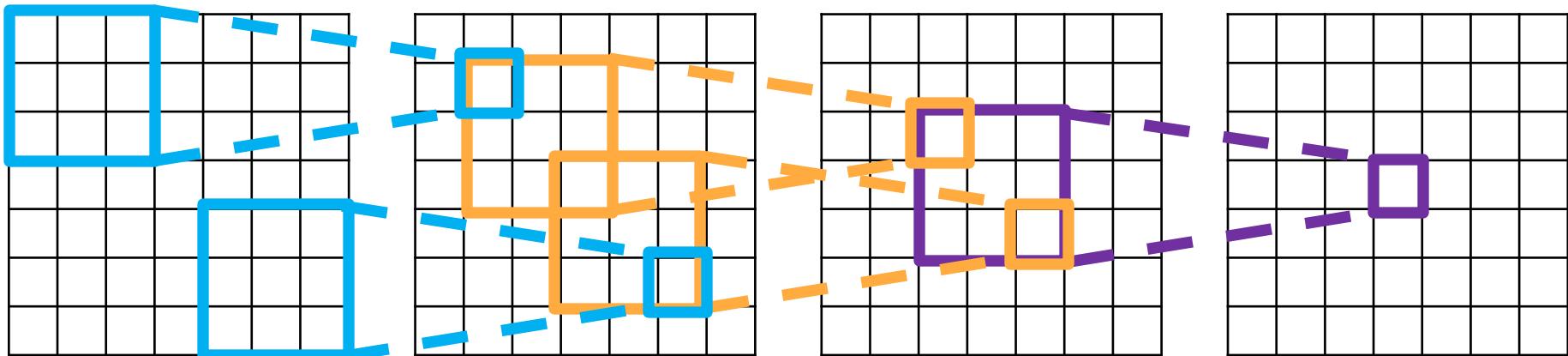
Problem: For large images we need many layers for each output to “see” the whole image

Output

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

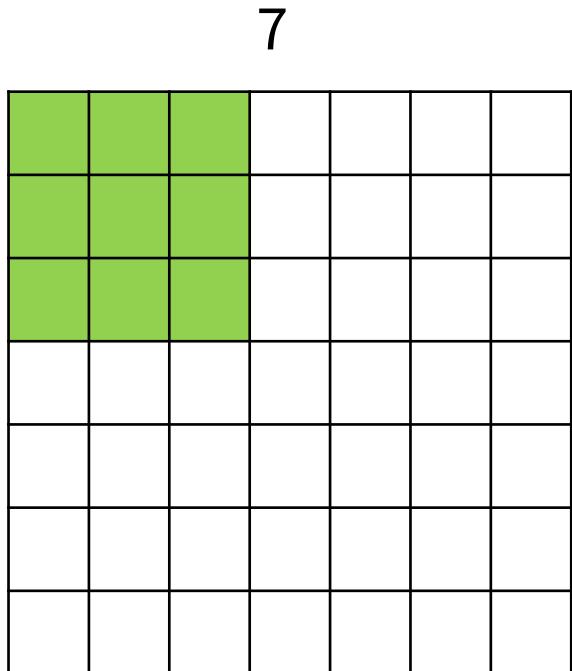
Problem: For large images we need many layers for each output to “see” the whole image

Solution: Downsample inside the network

Output

Slide inspiration: Justin Johnson

Strided Convolution



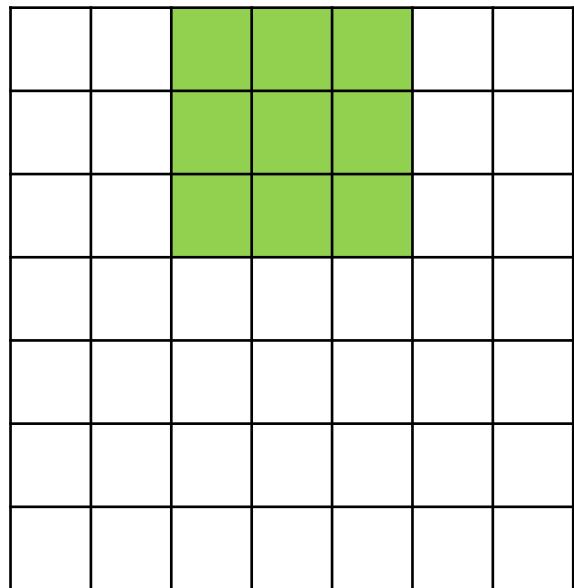
Input: 7x7

Filter: 3x3

Stride: 2

Strided Convolution

7



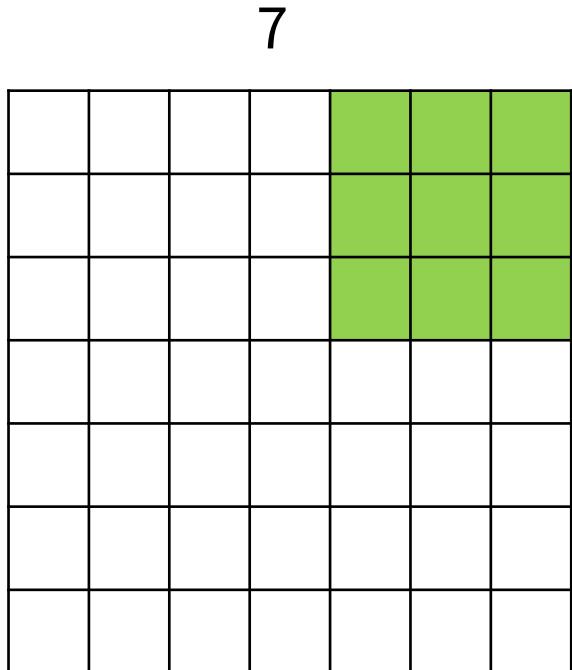
Input: 7x7

Filter: 3x3

Stride: 2

7

Strided Convolution



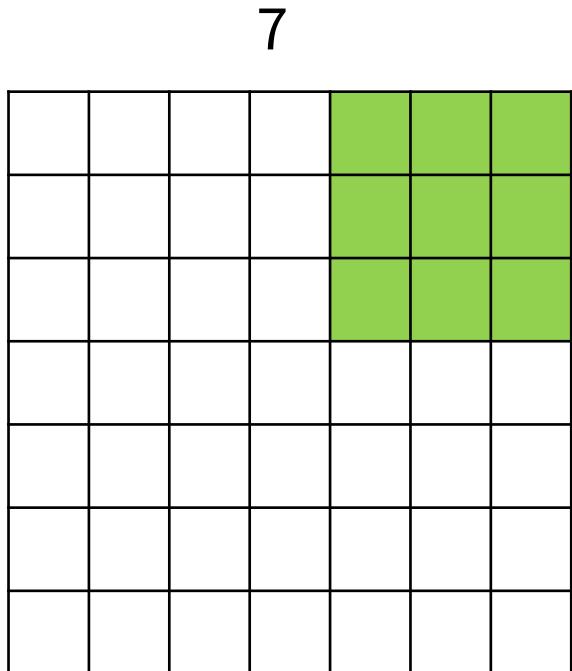
Input: 7x7

Filter: 3x3

Stride: 2

Output: 3x3

Strided Convolution



Input: 7x7

Filter: 3x3

Stride: 2

Output: 3x3

In general:

Input: W

Filter: K

Padding: P

Stride: S

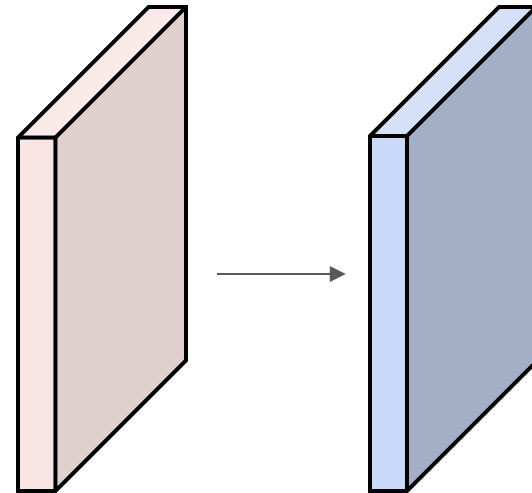
→ Output:
 $(W - K + 2P) / S + 1$

Convolution Example

Input volume: $3 \times 32 \times 32$

10 5x5 filters with stride 1, pad 2

Output volume size: ?



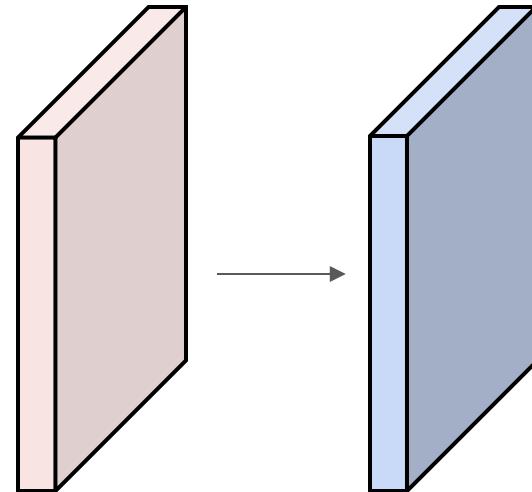
Convolution Example

Input volume: $3 \times 32 \times 32$

$10 \text{ } 5 \times 5$ filters with stride 1 , pad 2

Output volume size: $10 \times 32 \times 32$

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



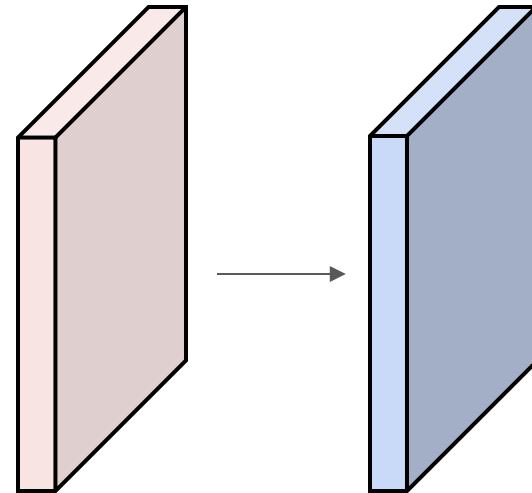
Convolution Example

Input volume: $3 \times 32 \times 32$

10 5x5 filters with stride 1, pad 2

Output volume size: $10 \times 32 \times 32$

Number of learnable parameters: ?



Convolution Example

Input volume: **3** x 32 x 32

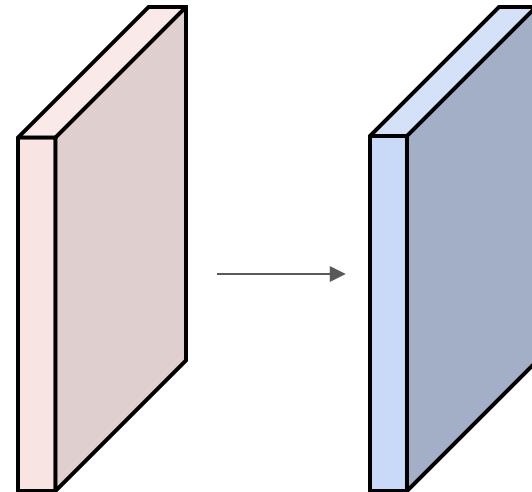
10 **5x5** filters with stride 1, pad 2

Output volume size: 10 x 32 x 32

Number of learnable parameters: 760

Parameters per filter: **3*5*5 + 1** (for bias) = **76**

10 filters, so total is **10 * 76 = 760**



Convolution Example

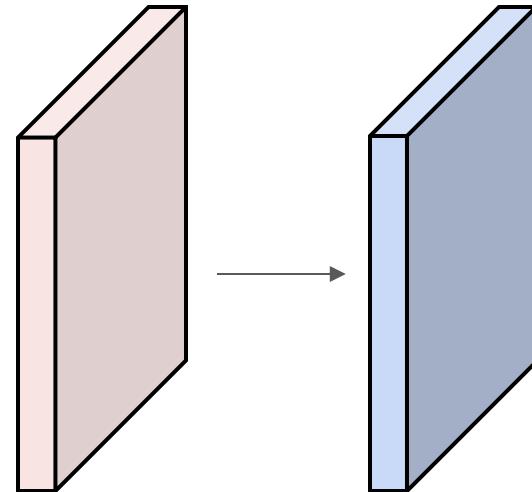
Input volume: $3 \times 32 \times 32$

10 5x5 filters with stride 1, pad 2

Output volume size: $10 \times 32 \times 32$

Number of learnable parameters: 760

Number of multiply-add operations?



Convolution Example

Input volume: **3 x 32 x 32**

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

Output volume size: **10 x 32 x 32**

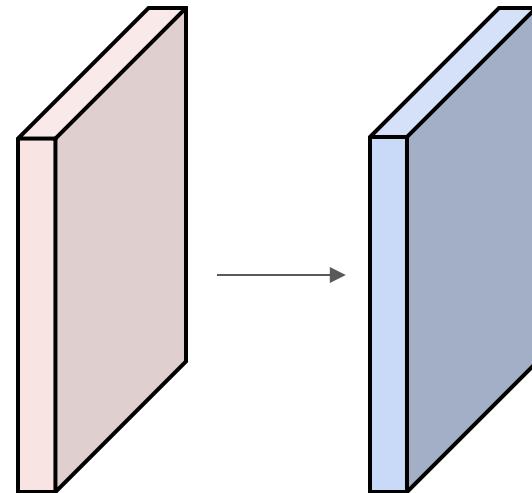
Number of learnable parameters: 760

Number of multiply-add operations: **768,000**

$10 * 32 * 32 = 10,240$ outputs

Each output is the inner product of two **3x5x5** tensors (75 elems)

Total = **$75 * 10240 = 768K$**



Convolution Summary

Input: $C_{in} \times H \times W$

Hyperparameters:

- **Kernel size:** $K_H \times K_W$
- **Number filters:** C_{out}
- **Padding:** P
- **Stride:** S

Weight matrix: $C_{out} \times C_{in} \times K_H \times K_W$

giving C_{out} filters of size $C_{in} \times K_H \times K_W$

Bias vector: C_{out}

Output size: $C_{out} \times H' \times W'$ where:

- $H' = (H - K + 2P) / S + 1$
- $W' = (W - K + 2P) / S + 1$

Common settings:

$K_H = K_W$ (Small square filters)

$P = (K - 1) / 2$ ("Same" padding)

$C_{in}, C_{out} = 32, 64, 128, 256$ (powers of 2)

$K = 3, P = 1, S = 1$ (3x3 conv)

$K = 5, P = 2, S = 1$ (5x5 conv)

$K = 1, P = 0, S = 1$ (1x1 conv)

$K = 3, P = 1, S = 2$ (Downsample by 2)

PyTorch Convolution Layer

```
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.

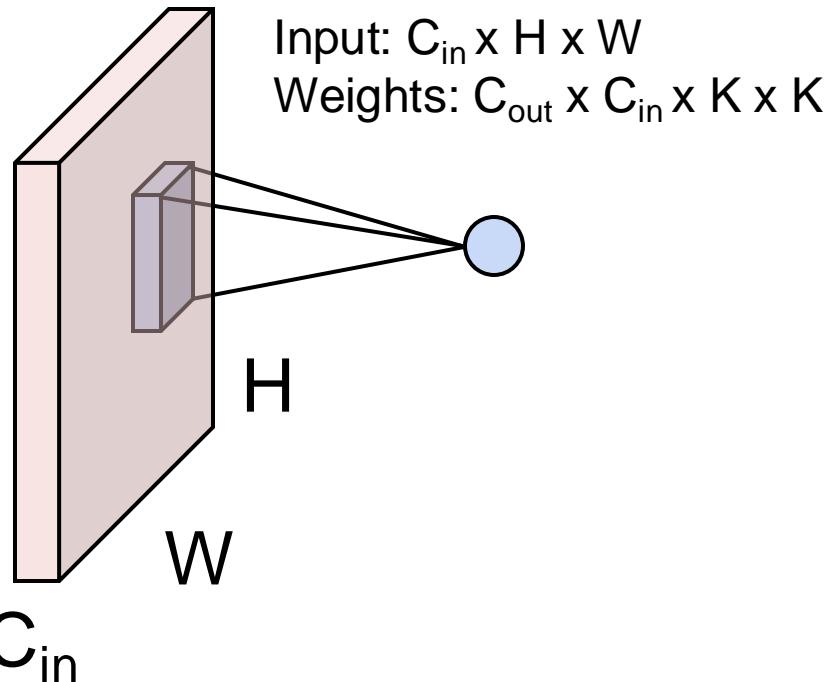
In the simplest case, the output value of the layer with input size (N, C_{in}, H, W) and output $(N, C_{\text{out}}, H_{\text{out}}, W_{\text{out}})$ can be precisely described as:

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

We didn't talk about groups or dilation...

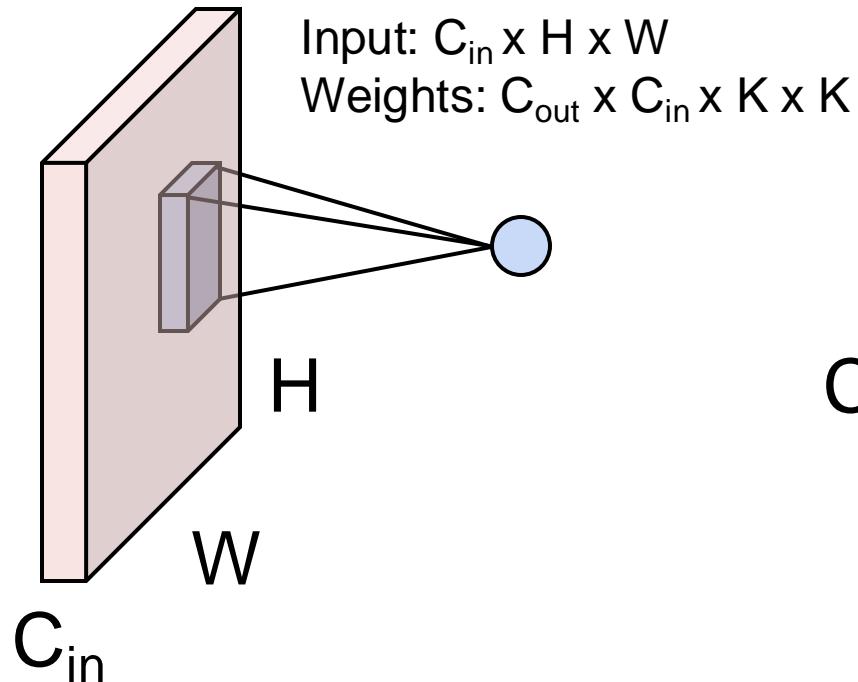
Other Types of Convolution

So far: 2D Convolution

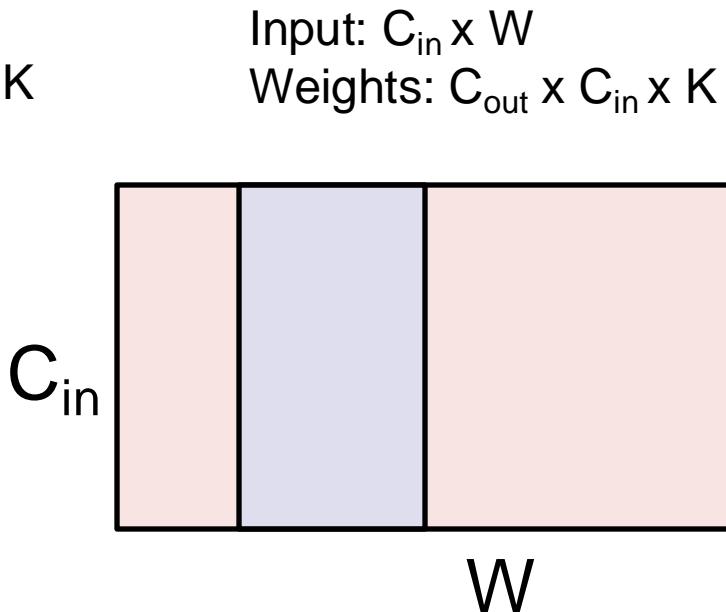


Other Types of Convolution

So far: 2D Convolution

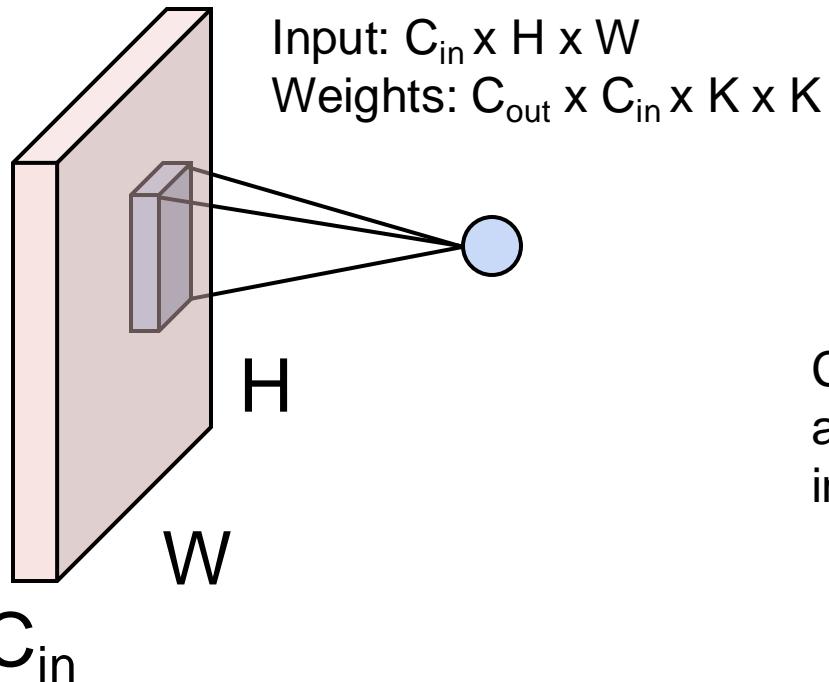


1D Convolution



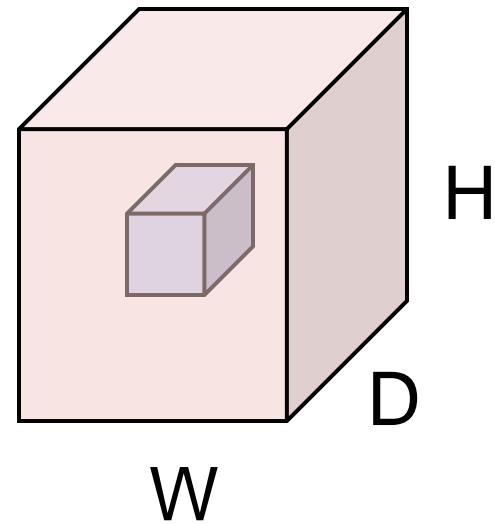
Other Types of Convolution

So far: 2D Convolution



3D Convolution

Input: $C_{in} \times H \times W \times D$
Weights: $C_{out} \times C_{in} \times K \times K \times K$

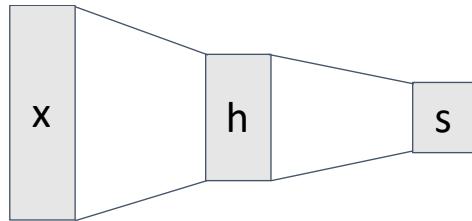


C_{in} -dim vector
at each point
in the volume

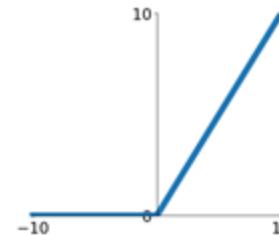
Convolutional Networks

We have
already
seen these

Fully-Connected Layer

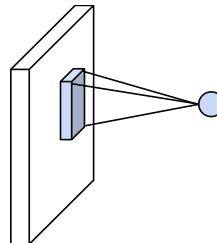


Activation Function

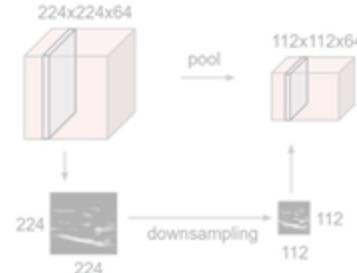


Convolution Layer

Today: Image-specific operators



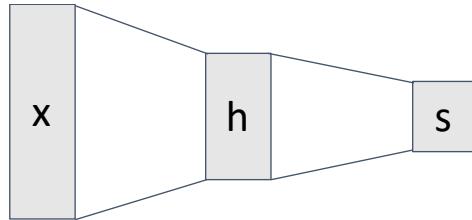
Pooling Layer



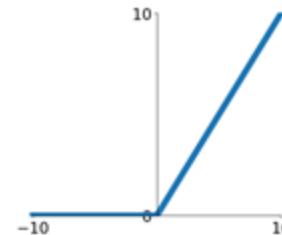
Convolutional Networks

We have
already
seen these

Fully-Connected Layer

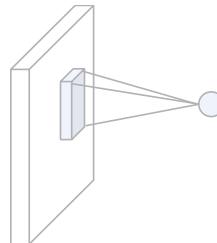


Activation Function

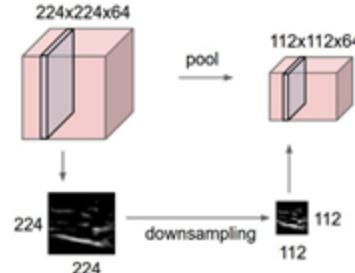


Convolution Layer

Today: Image-specific operators

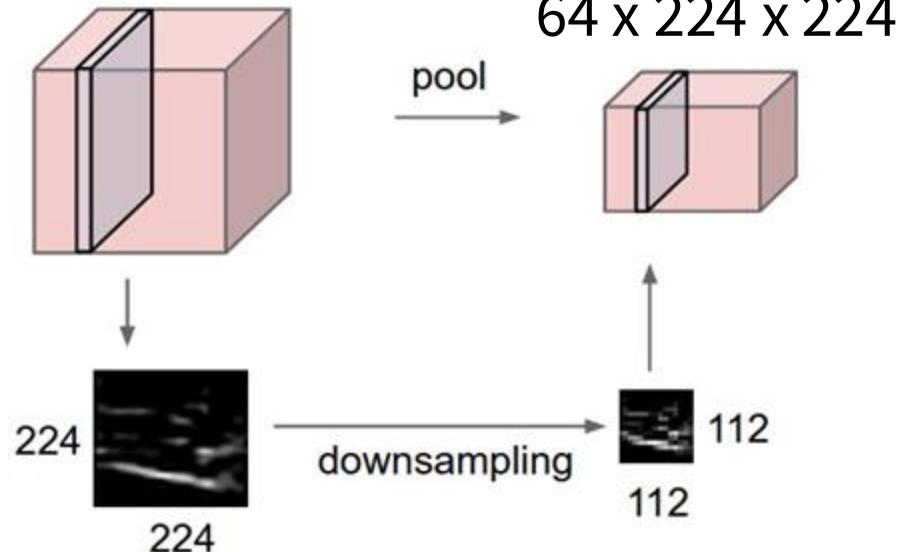


Pooling Layer



Pooling Layers: Another way to downsample

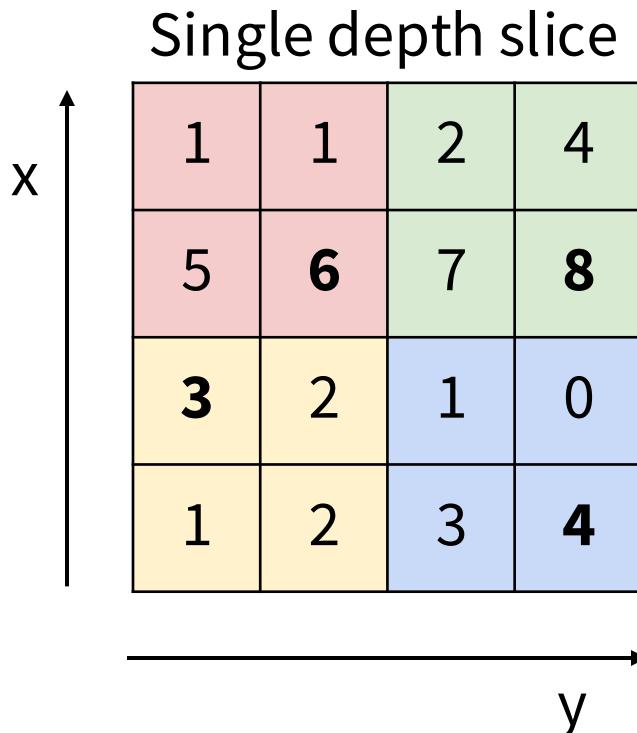
$64 \times 112 \times 112$



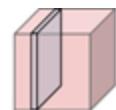
Given an input $C \times H \times W$,
downsample each $1 \times H \times W$ plane

Hyperparameters:
Kernel Size
Stride
Pooling function

Pooling Layers: Another way to downsample



64 x 224 x 224



Max pooling with 2x2 kernel size and stride 2

6	8
3	4

Gives **invariance** to small spatial shifts. No learnable parameters.

Pooling Summary

Input: $C \times H \times W$

Hyperparameters:

- **Kernel size:** K
- **Stride:** S
- **Pooling function:** max, avg

Output size: $C \times H' \times W'$ where:

- $H' = (H - K) / S + 1$
- $W' = (W - K) / S + 1$

No learnable parameters

Common setting:

max, $K=2$, $S=2 \Rightarrow$ Gives 2x downsampling

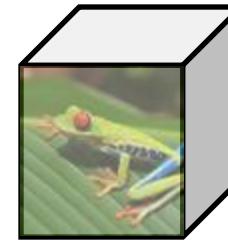
Convolution and Pooling: Translation Equivariance

$H \times W \times C$

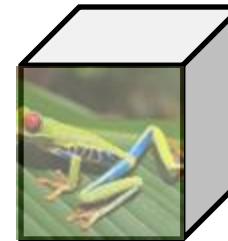


Conv or Pool

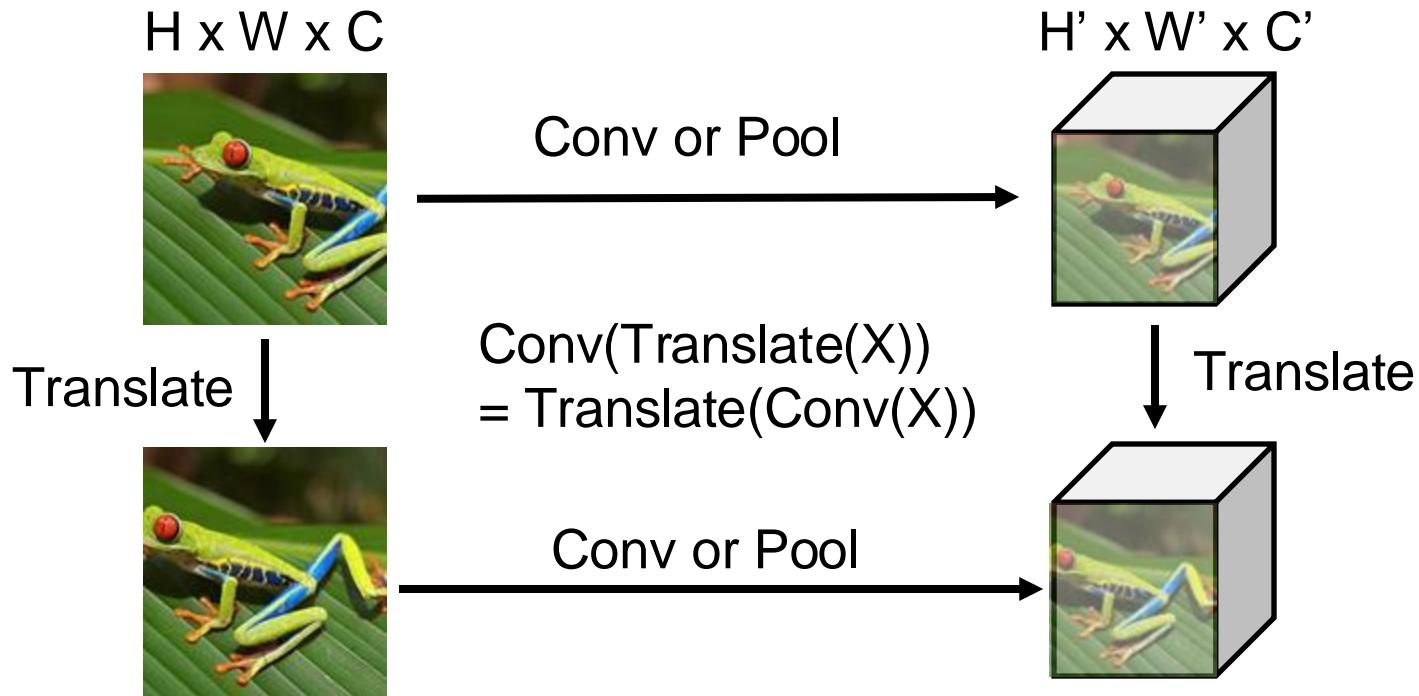
$H' \times W' \times C'$



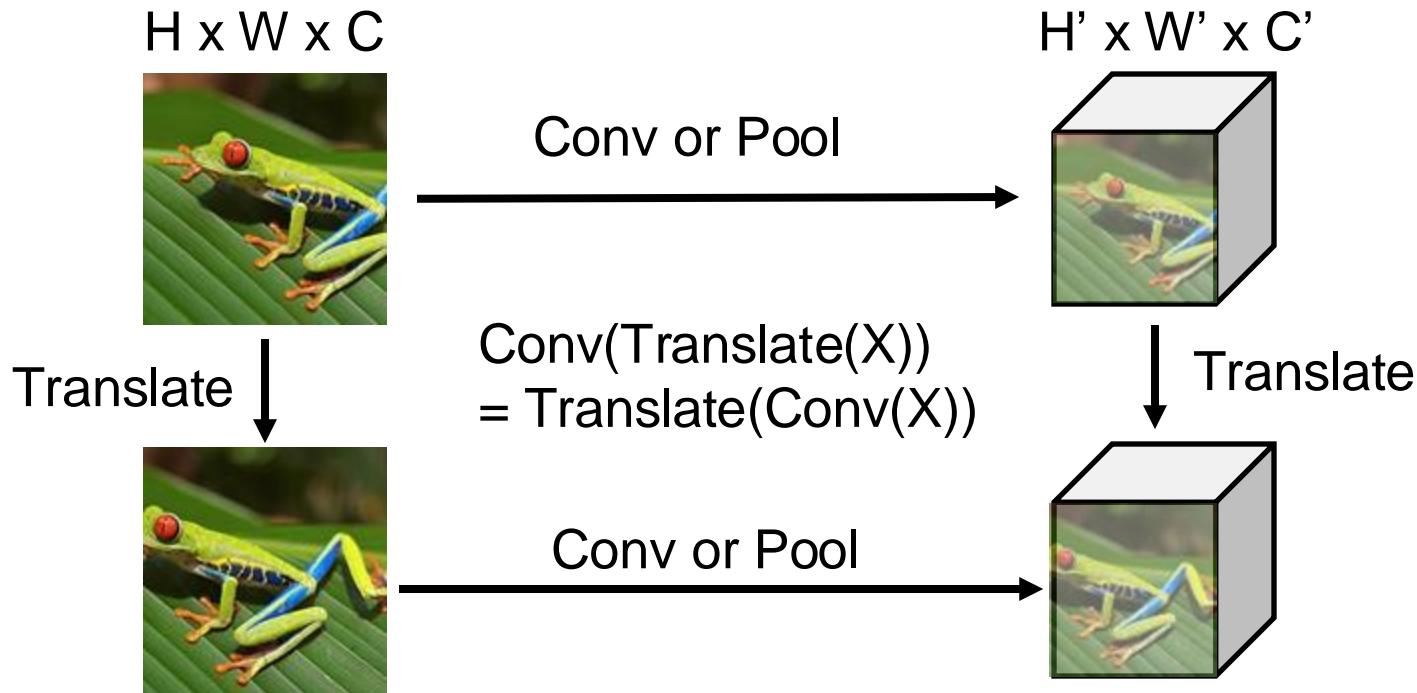
Translate



Convolution and Pooling: Translation Equivariance



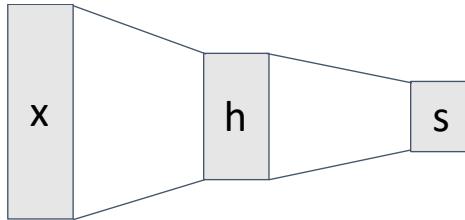
Convolution and Pooling: Translation Equivariance



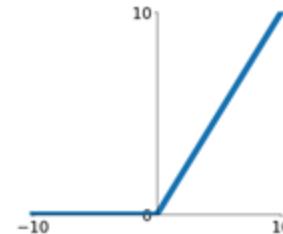
Intuition: Features of images don't depend on their location in the image

Summary: Convolutional Networks

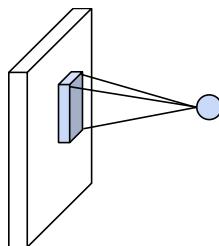
Fully-Connected Layer



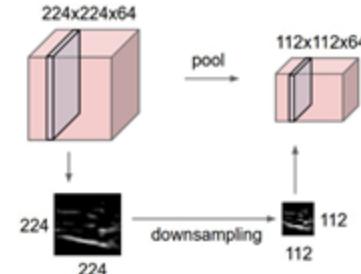
Activation Function



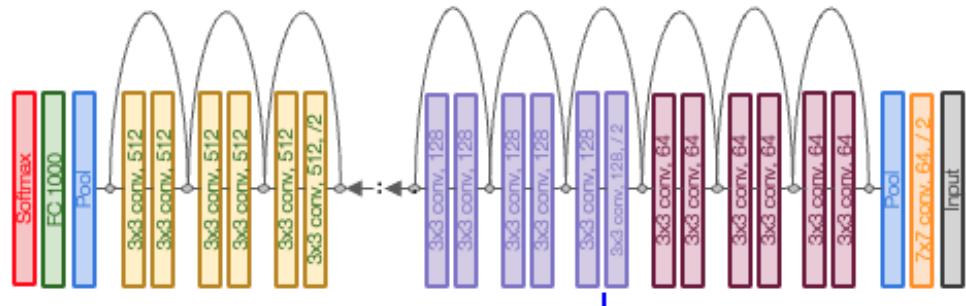
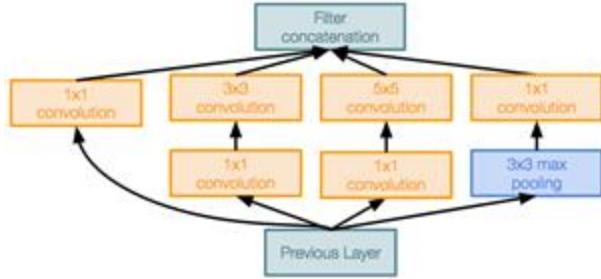
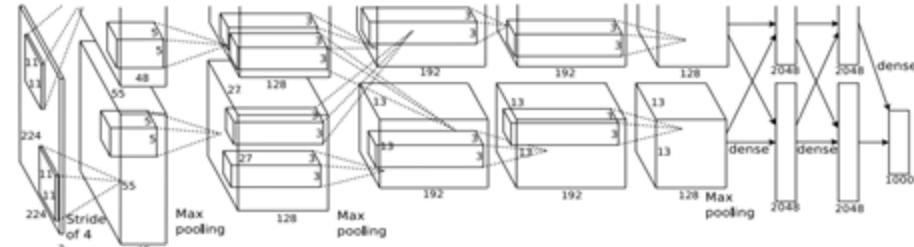
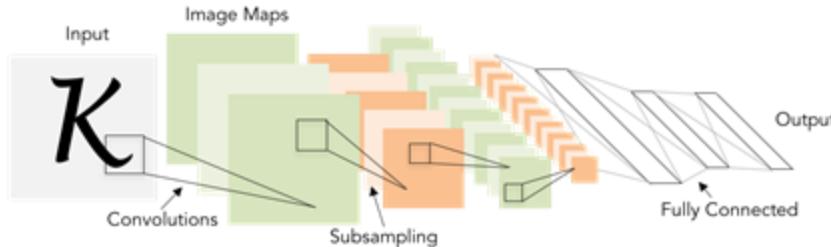
Convolution Layer

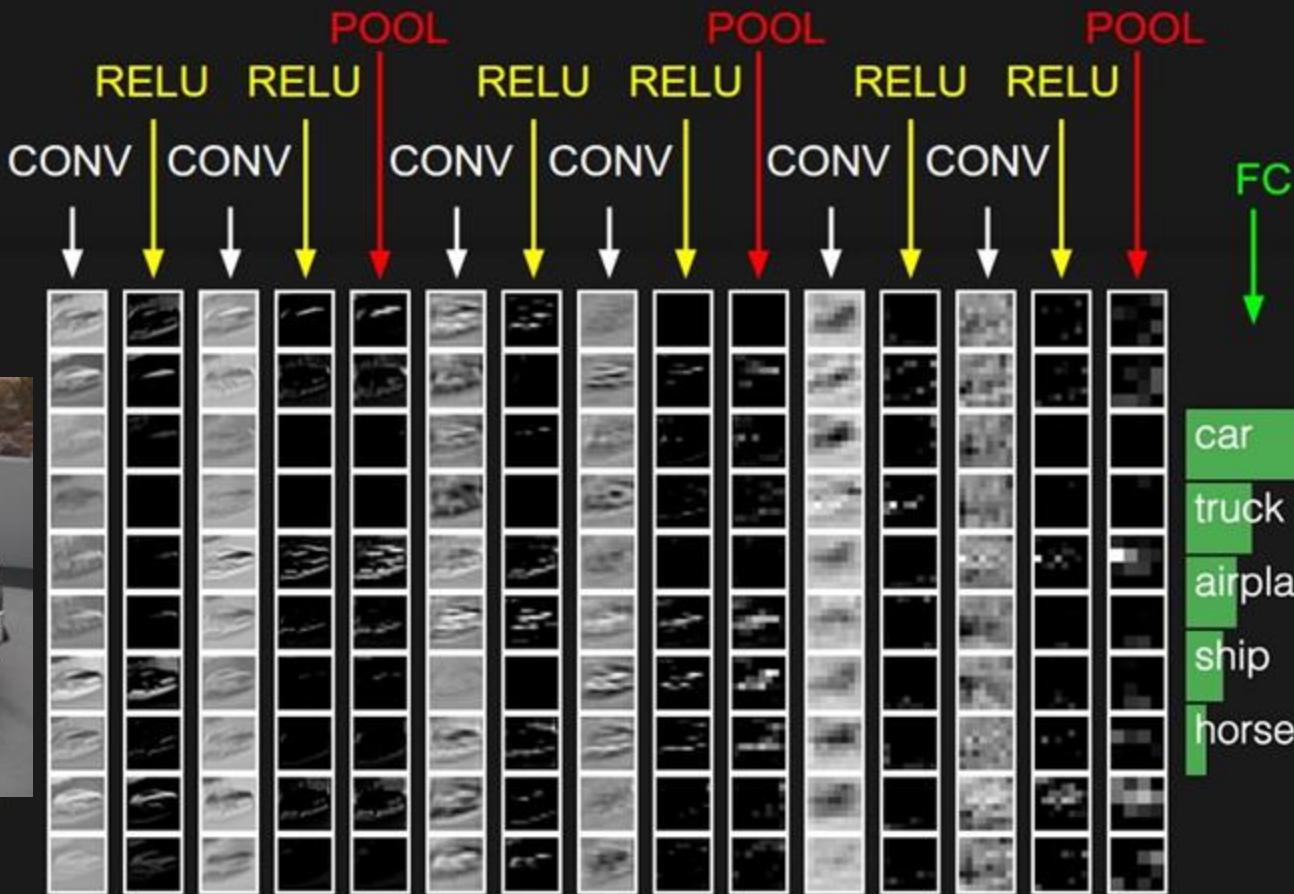


Pooling Layer



Next time: CNN Architectures





[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](#).



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