



6.390 Intro to Machine Learning

Lecture 7: Convolutional Neural Networks

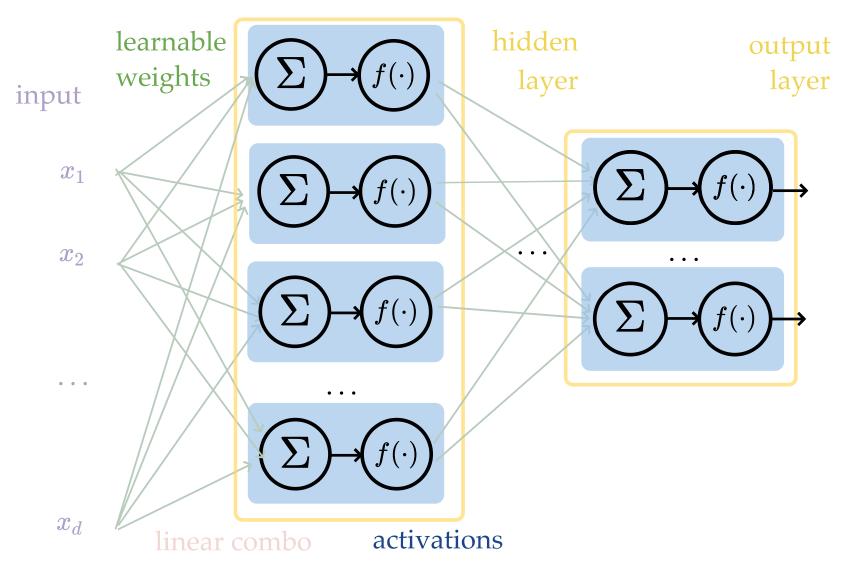
Shen Shen

Oct 16, 2025

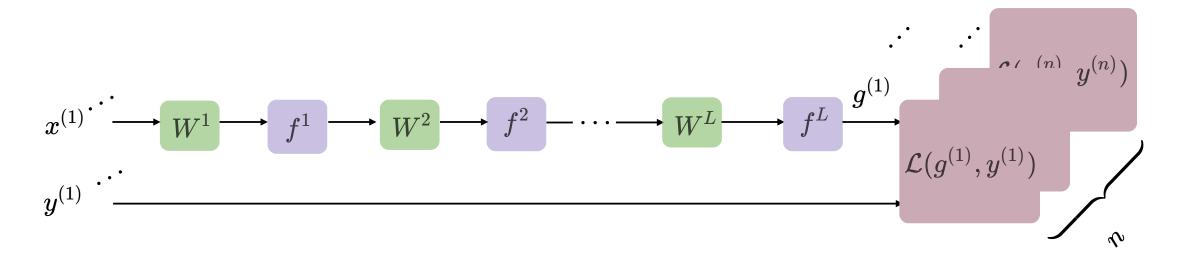
11am, Room 10-250

Interactive Slides and Lecture Recording

Recap: A (fully-connected, feed-forward) neural network neuron



Forward pass: evaluate, given the current parameters



- the model outputs $g^{(i)} = f^L\left(\dots f^2\left(f^1(\mathbf{x}^{(i)};\mathbf{W}^1);\mathbf{W}^2\right);\dots \mathbf{W}^L\right)$
- ullet the loss incurred on the current data $\mathcal{L}(g^{(i)}, y^{(i)})$

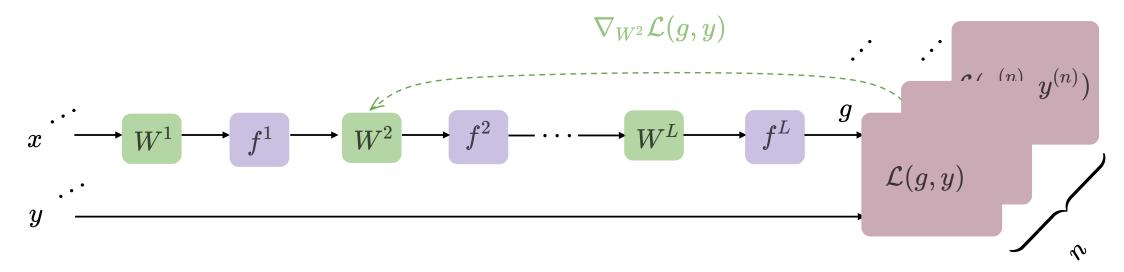
• the training error $J = rac{1}{n} \sum_{i=1}^n \mathcal{L}(g^{(i)}, y^{(i)})$

linear combination

(nonlinear) activation

loss function

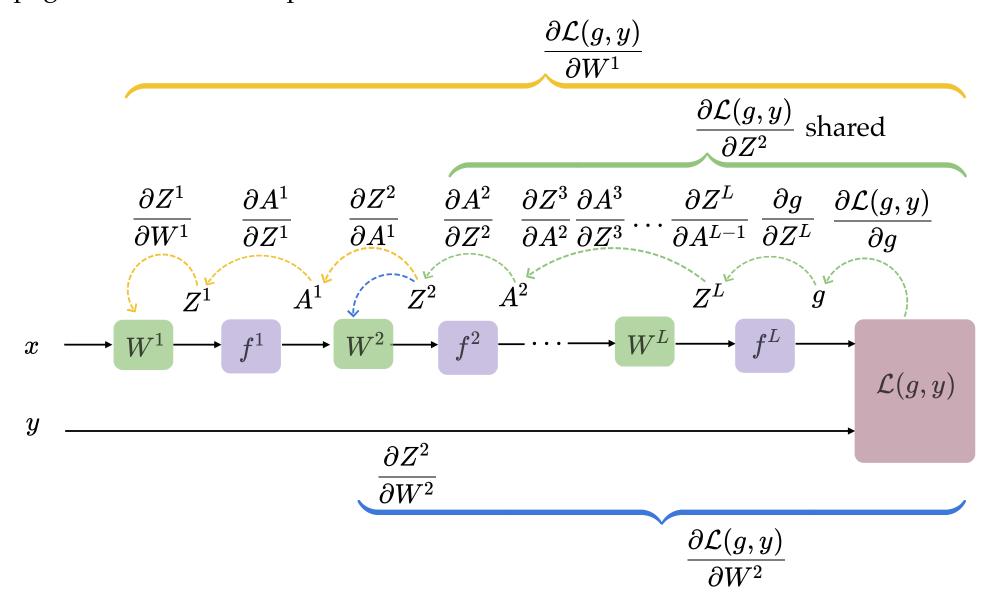
Backward pass: run SGD to learn all parameters e.g. to update W^2



- Randomly pick a data point (x, y)
- Evaluate the gradient $abla_{W^2}\mathcal{L}(g,y)$
- Update the weights $W^2 \leftarrow W^2 \eta \nabla_{W^2} \mathcal{L}(g,y)$

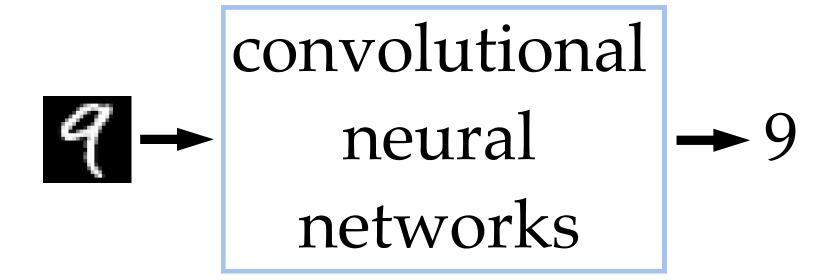
suppose we sampled a particular (x, y), how to find $\partial \mathcal{L}(g,y)$ ∂W^2 $\partial \mathcal{L}(g,y)$ ∂Z^2 $\partial A^2 \quad \partial Z^3 \; \partial A^3 \qquad \partial Z^L \quad \underline{\partial g}$ ∂Z^2 $\partial \mathcal{L}(g,y)$ $\overline{\partial Z^2} \quad \overline{\partial A^2} \ \overline{\partial Z^3} \cdots \overline{\partial A^{L-1}} \ \overline{\partial Z^L}$ $\overline{\partial W^2}$ ∂g A^1 $\mathcal{L}(g,y)$ y

back propagation: reuse of computation





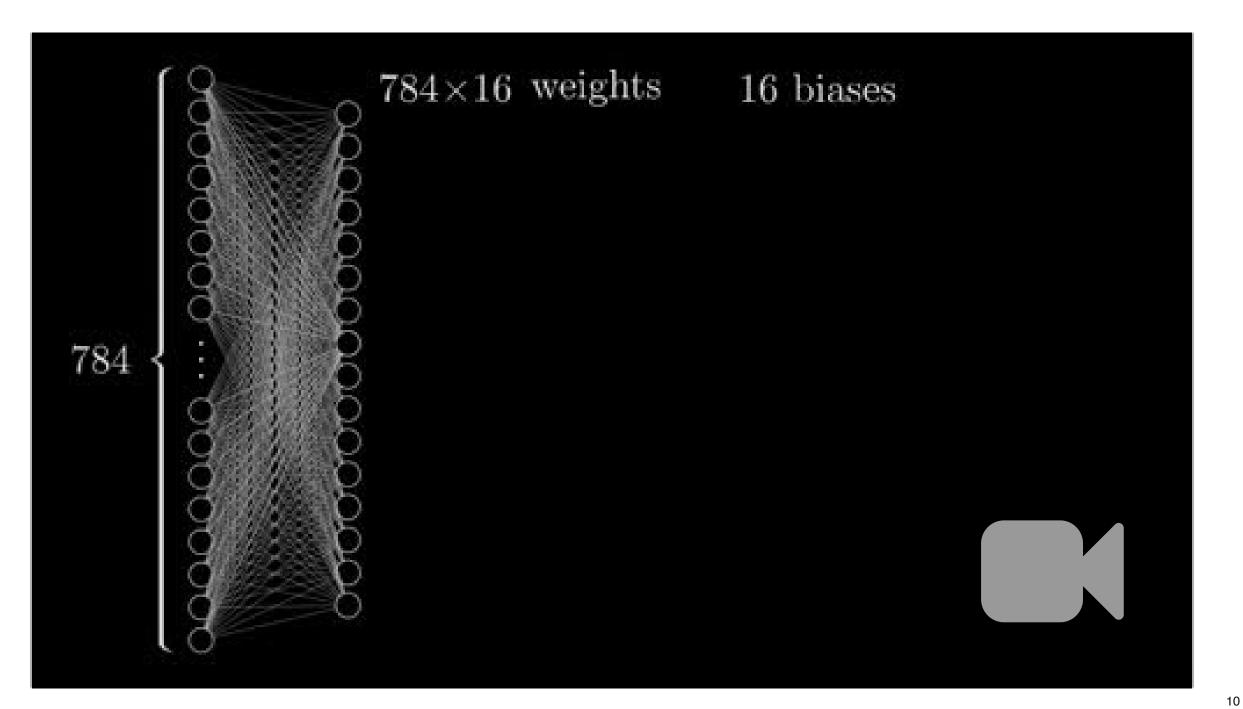
https://shenshen.mit.edu/demos/nn/playground/



- 1. Why do we need a special network for images?
- 2. Why is CNN (the) special network for images?

Outline

- Vision problem structure
- Convolution
 - 1-dimensional and 2-dimensional *convolution*
 - 3-dimensional *tensors*
- Max pooling
- (Case studies)



426-by-426 grayscale image



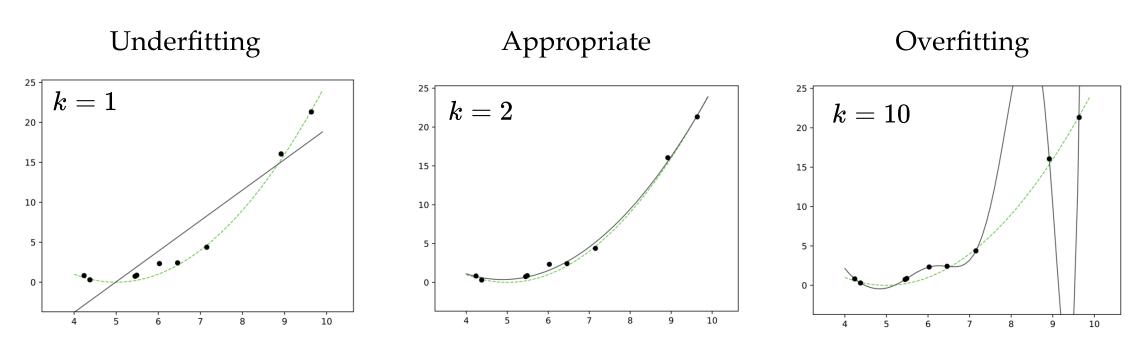
Use the same 2 hidden-layer network to predict what top-10 engineering school seal this image is, need to learn ~3M parameters.

For higher-resolution images, or more complex tasks, or larger networks, the number of parameters can grow very fast.

Why do we need a specialized network (hypothesis class)?

- Partly, fully-connected nets don't scale well for vision tasks
- More importantly, a carefully chosen hypothesis class helps fight overfitting

Recall, models with needless parameters tend to overfit



If we know the data is generated by the green curve, it's easy to choose the appropriate quadratic hypothesis class.

so... do we know anything about vision problems?

Why do we humans think

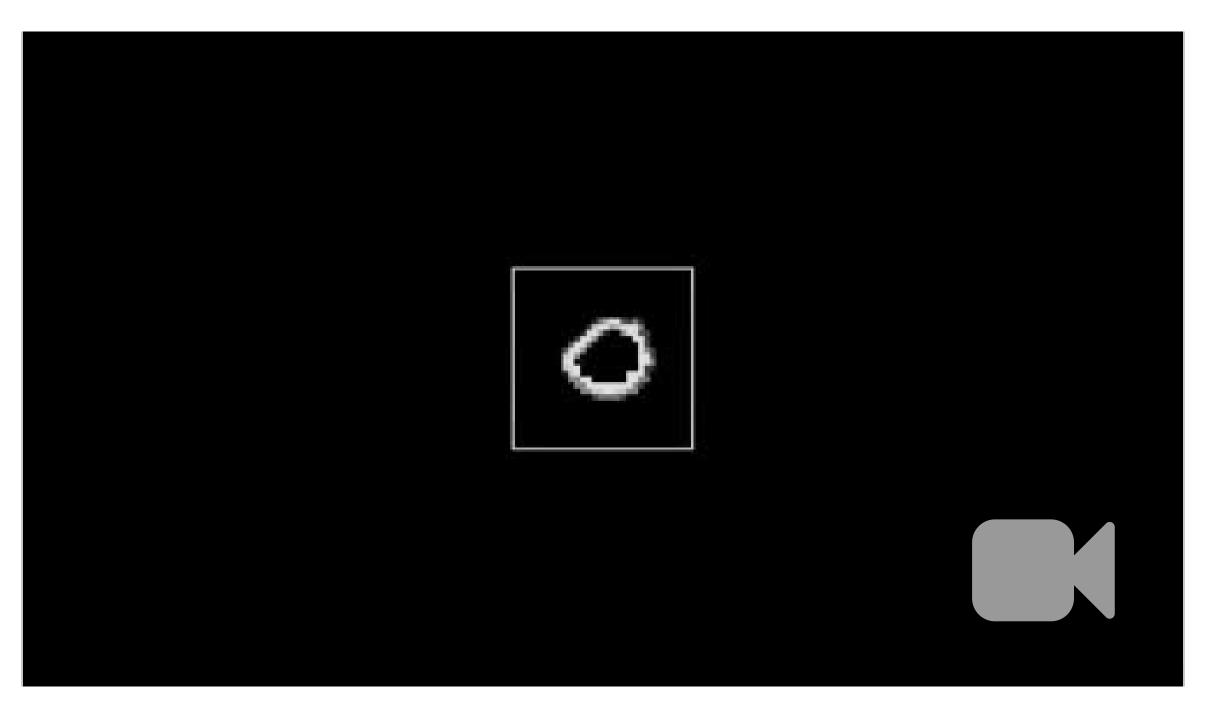


is a 9?

Why do we think any of

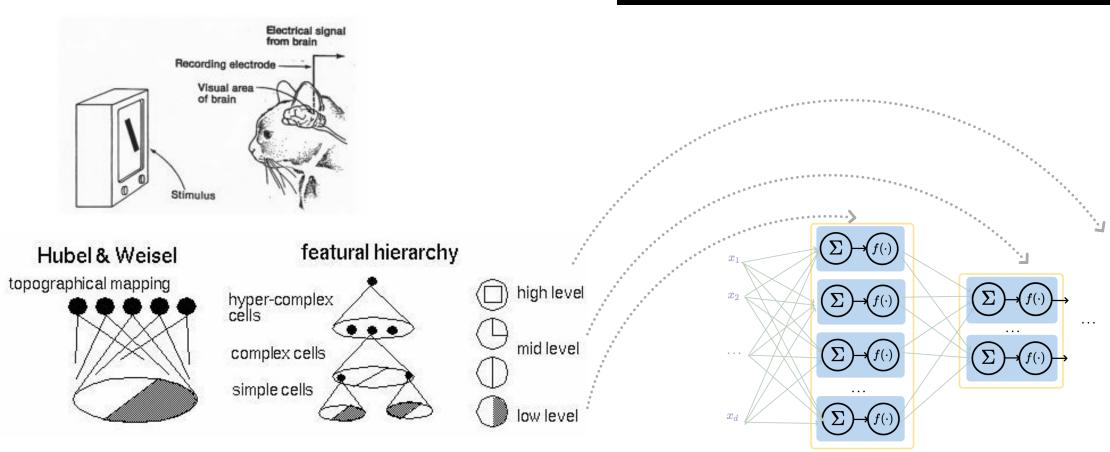


is a 9?



• Visual hierarchy





Layered structure are well-suited to model this hierarchical processing.

• Visual hierarchy



• Spatial locality



• Translational invariance

CNN exploits

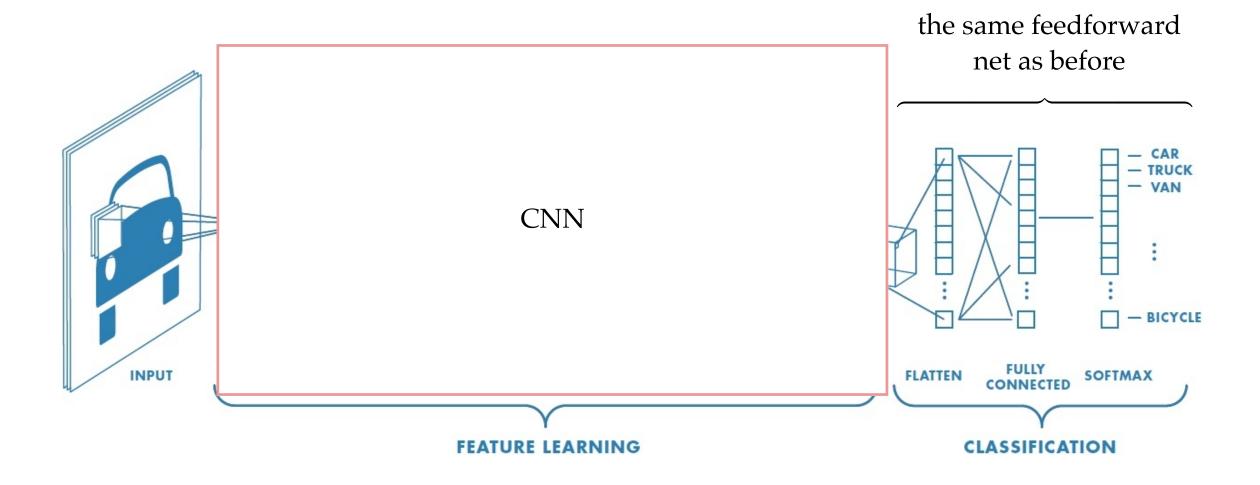
- Visual hierarchy
- Spatial locality
- Translational invariance

via

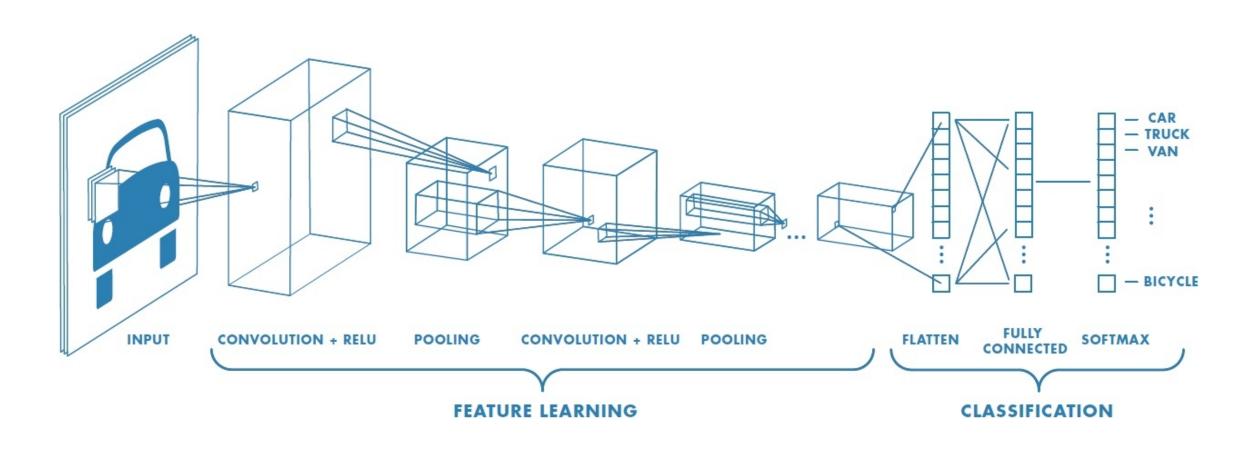
- layered structure
- convolution
- pooling

to handle images efficiently and sensibly.

typical CNN architecture for image classification



typical CNN structure for image classification



Outline

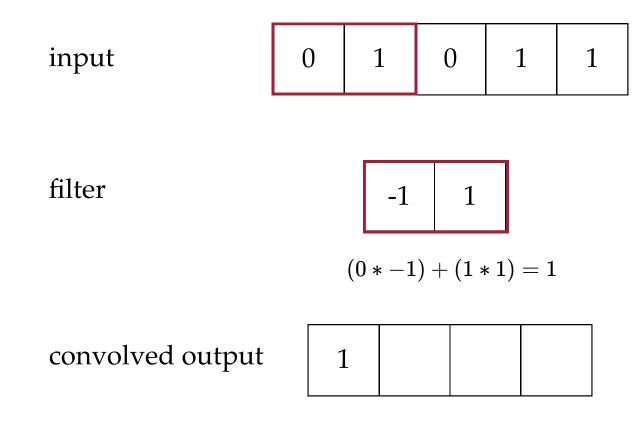
- Vision problem structure
- Convolution
 - 1-dimensional and 2-dimensional *convolution*
 - 3-dimensional *tensors*
- Max pooling
- (Case studies)

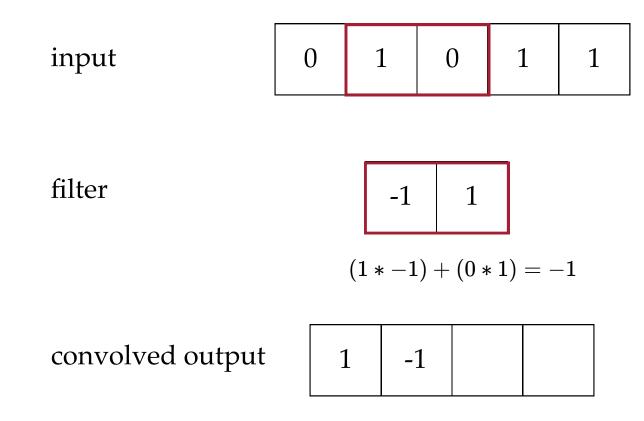
Convolutional layer might sound foreign, but it's very similar to a fully-connected layer

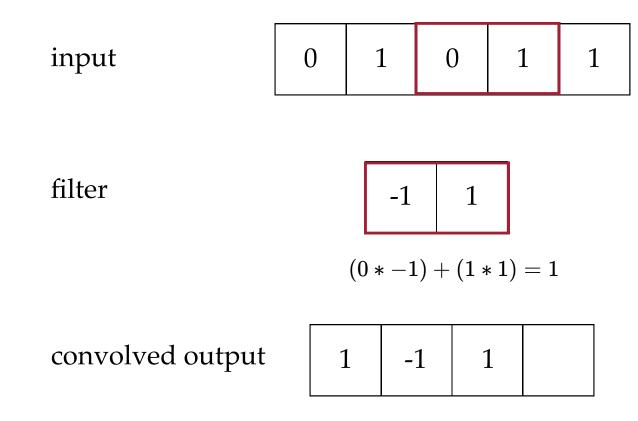
Layer	Forward pass, do	Backward pass, learn	Design choices
fully-connected	dot-product, activation	neuron weights	neuron count, etc.
convolutional	convolution, activation	filter weights	conv specs, etc.

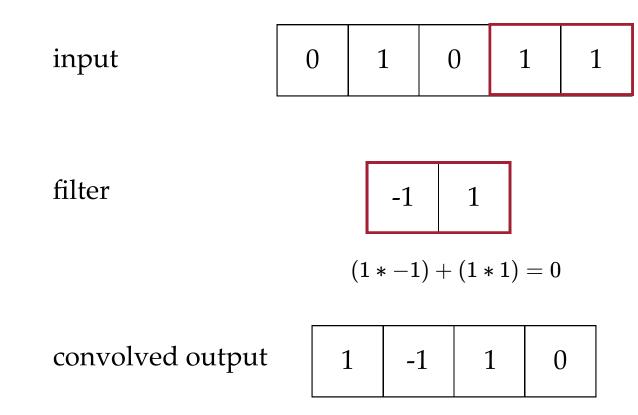
Convolution result:



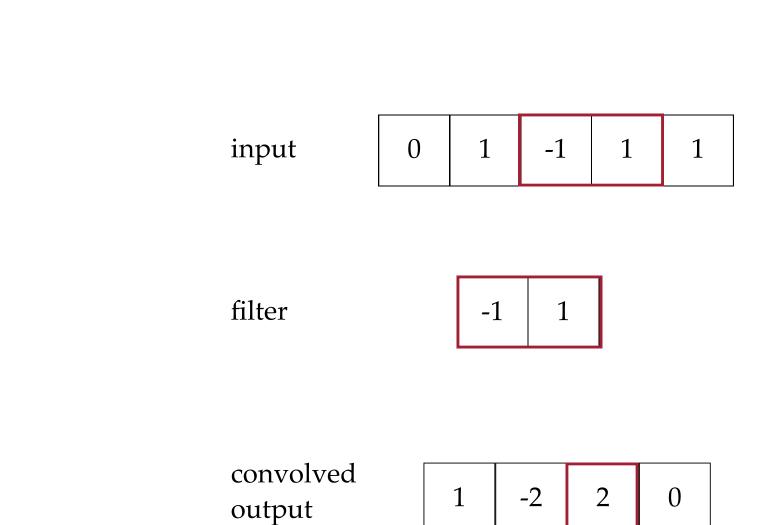




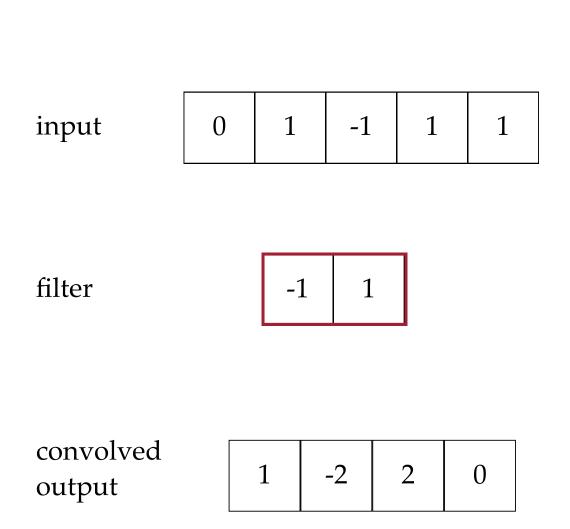


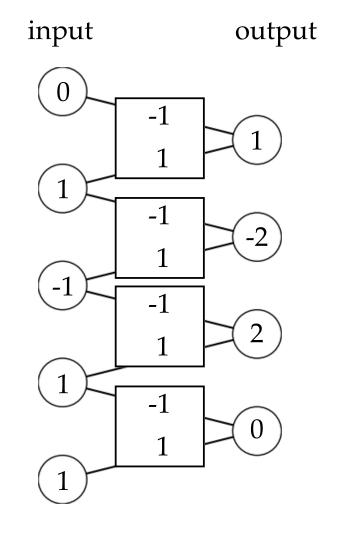


convolution interpretation 1: template matching

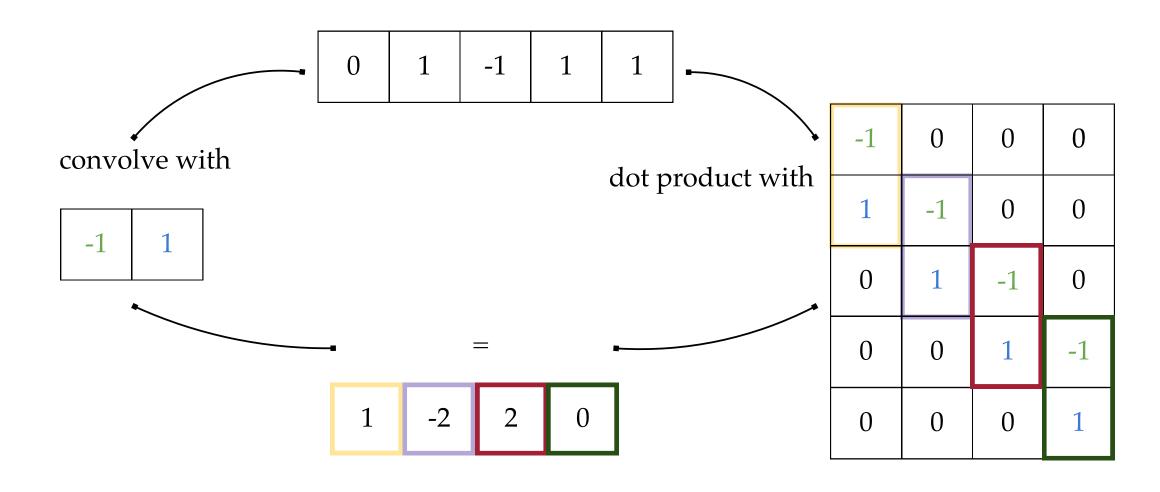


convolution interpretation 2: "look" locally through the filter (sparse-connected layer)





convolution interpretation 3: sparse-connected layer with parameter sharing





0	1	0	1	1
---	---	---	---	---

convolve with

1

1

0 1 0 1 1

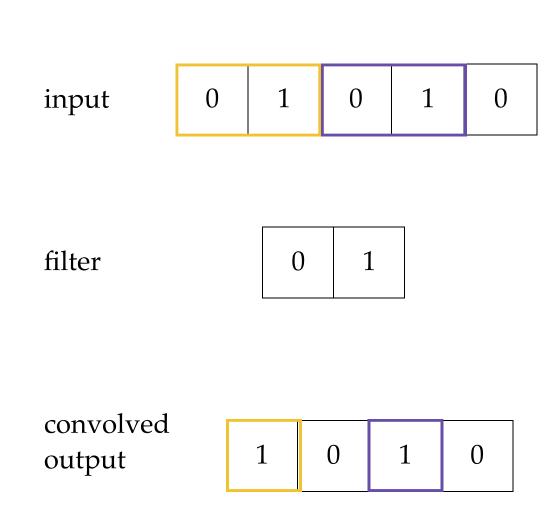
7

dot product with

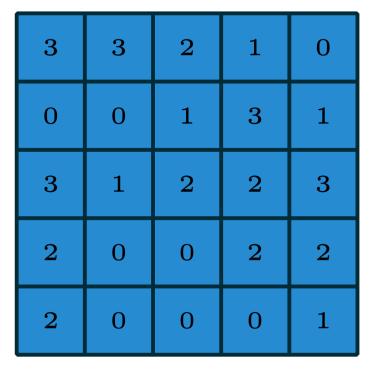
 $I_{5 imes 5}$



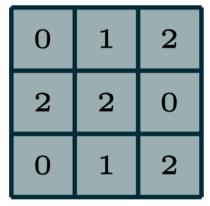
convolution interpretation 4: translational equivariance





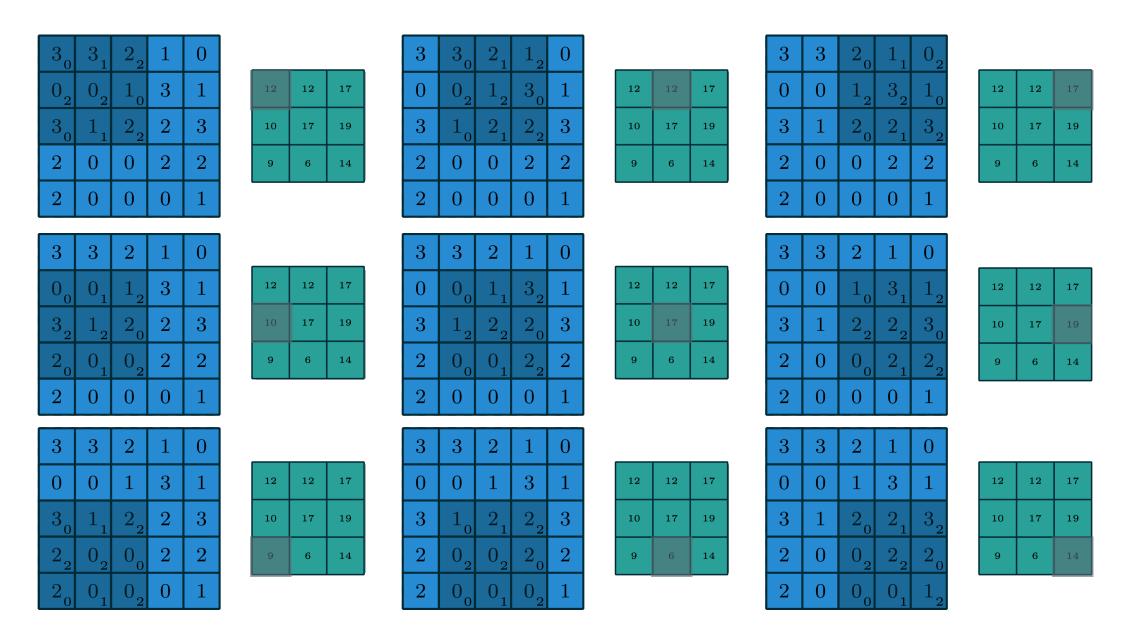


filter

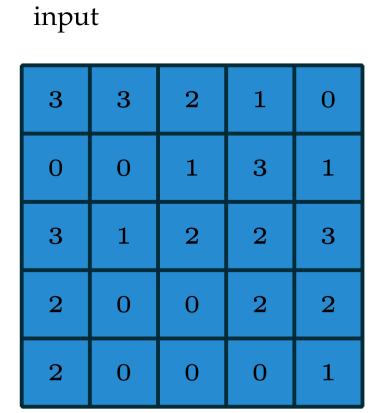


convolved output

12	12	17
10	17	19
9	6	14



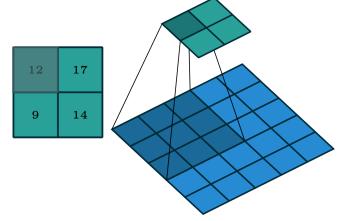
stride of 2



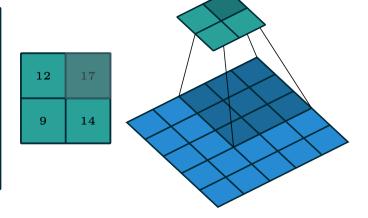
	filter		out	put
0	1	2	12	17
2	2	0	9	14
0	1	2		

stride of 2

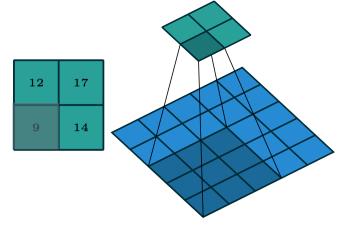
30	3	2_{2}	1	0
02	0_2	1_{0}	3	1
30	1,	2_{2}	2	3
2	0	0	2	2
2	0	0	0	1

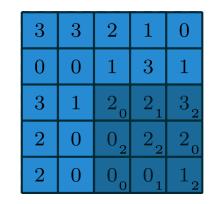


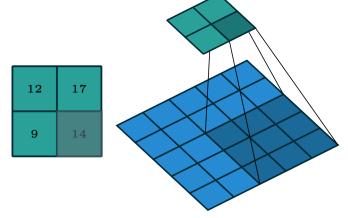
3	2_0	1_{1}	0_2
0	1_2	32	10
1	2_0	2_{1}	3_2
0	0	2	2
0	0	0	1
	0	0 1 ₂ 1 2 ₀	$\begin{array}{c cccc} 0 & 1_2 & 3_2 \\ 1 & 2_0 & 2_1 \\ 0 & 0 & 2 \\ \end{array}$



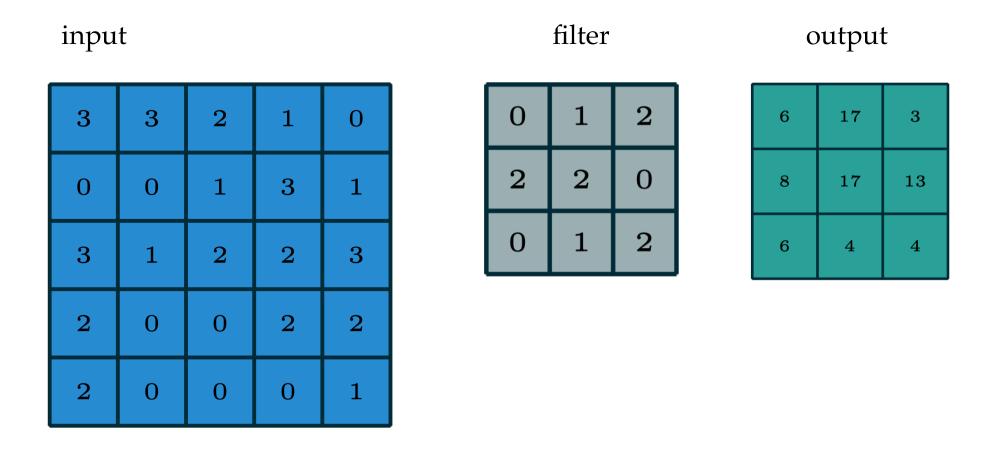
3	3	2	1	0
0	0	1	3	1
30	1,	2_2	2	3
2_{2}	0_2	00	2	2
2_{0}	01	0_2	0	1

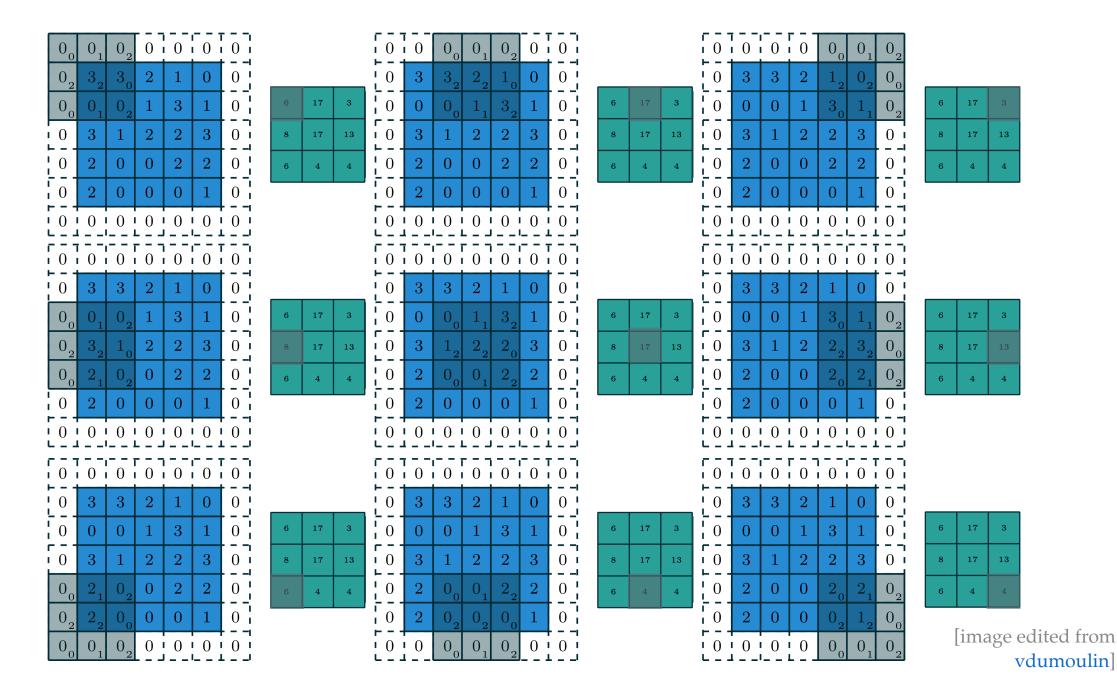






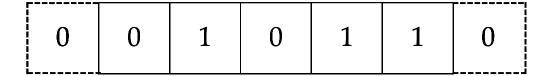
stride of 2, with padding of size 1





quick summary: hyperparameters for 1d convolution

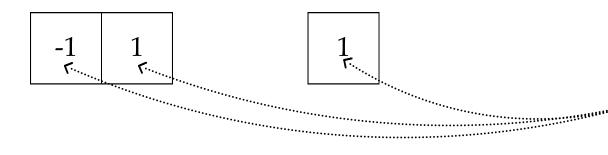
Zero-padding



• Stride (e.g. stride of 2)

0	0	1	0	1	1	0
i						

• Filter size (e.g. we saw these two in 1-d)



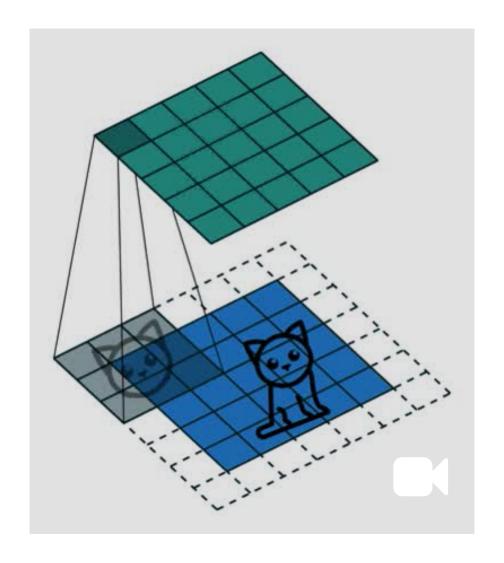
these weights are what CNN learn eventually

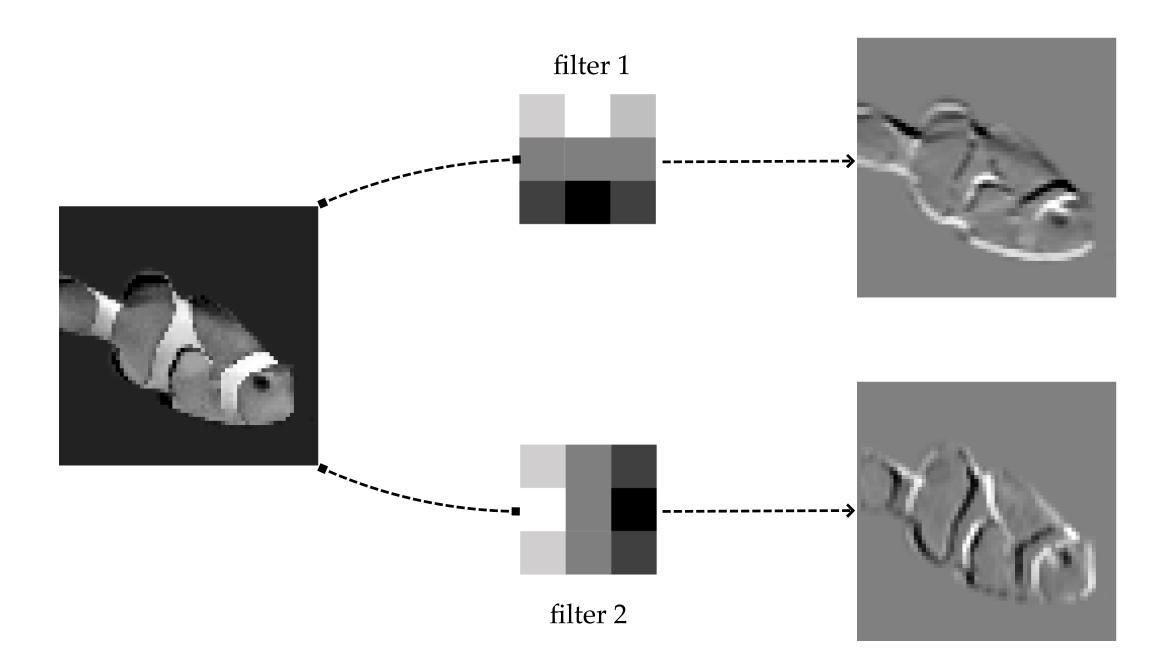
quick summary: hyperparameters for 2d convolution

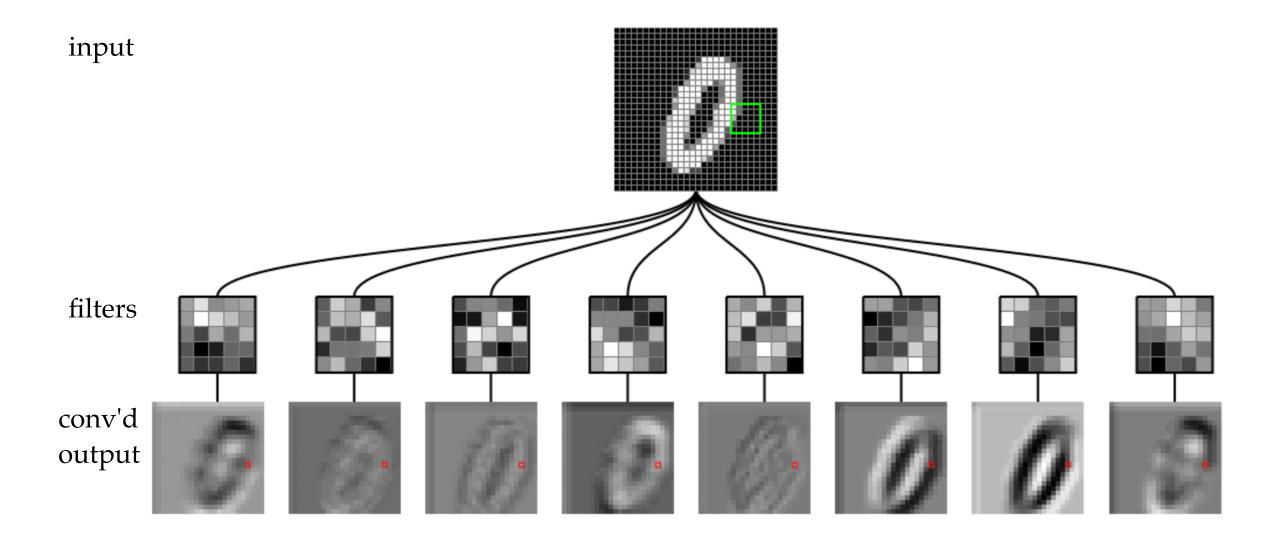
nttps://shenshen.mit.edu/demos/cnn/hyperparameters.html

quick summary: convolution interpretation

- Look locally (sparse connections)
- Parameter sharing
- Template matching
- Translational equivariance

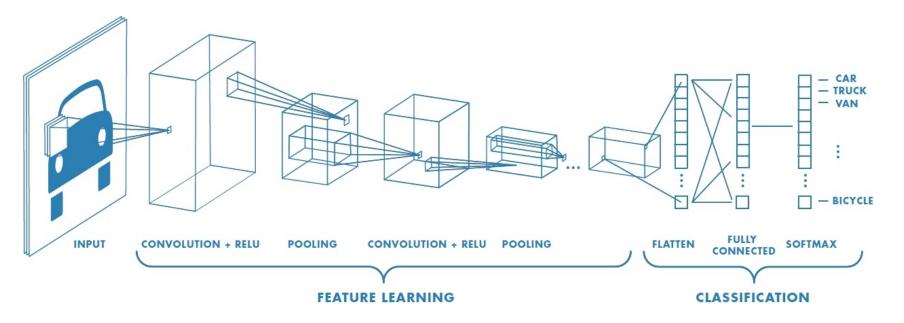






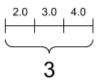
Outline

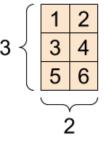
- Vision problem structure
- Convolution
 - 1-dimensional and 2-dimensional *convolution*
 - 3-dimensional *tensors*
- Max pooling
- (Case studies)

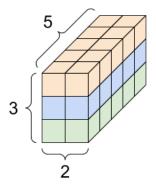


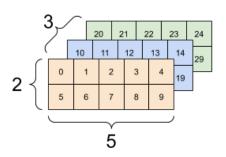
A tender intro to tensor:

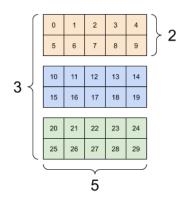
4











color images and channels







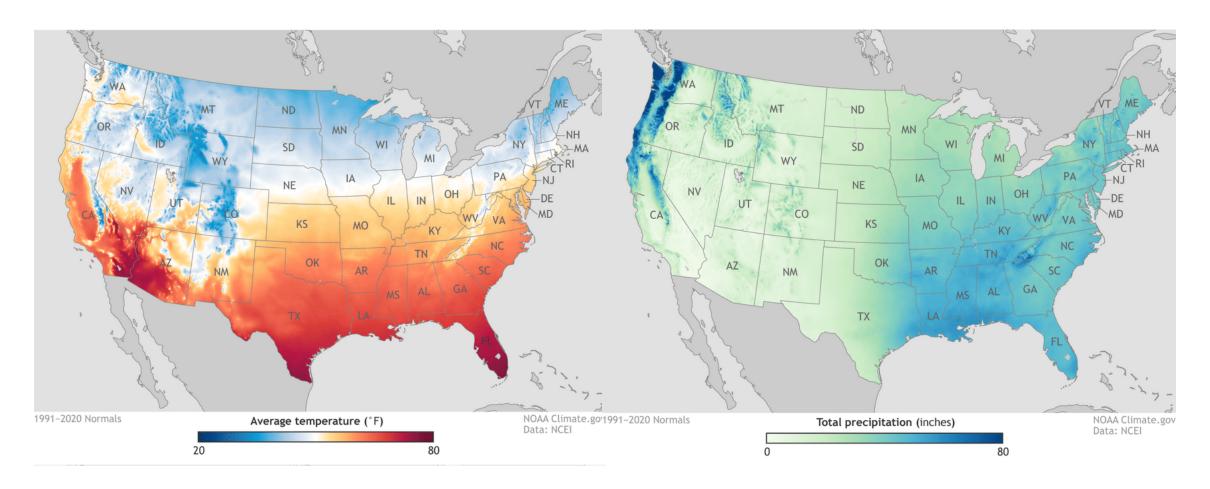


green



blue

each channel encodes a *holistic but independent* perspective of the same image, similar to:



so channels are often referred to as feature maps

3d tensors from color channels

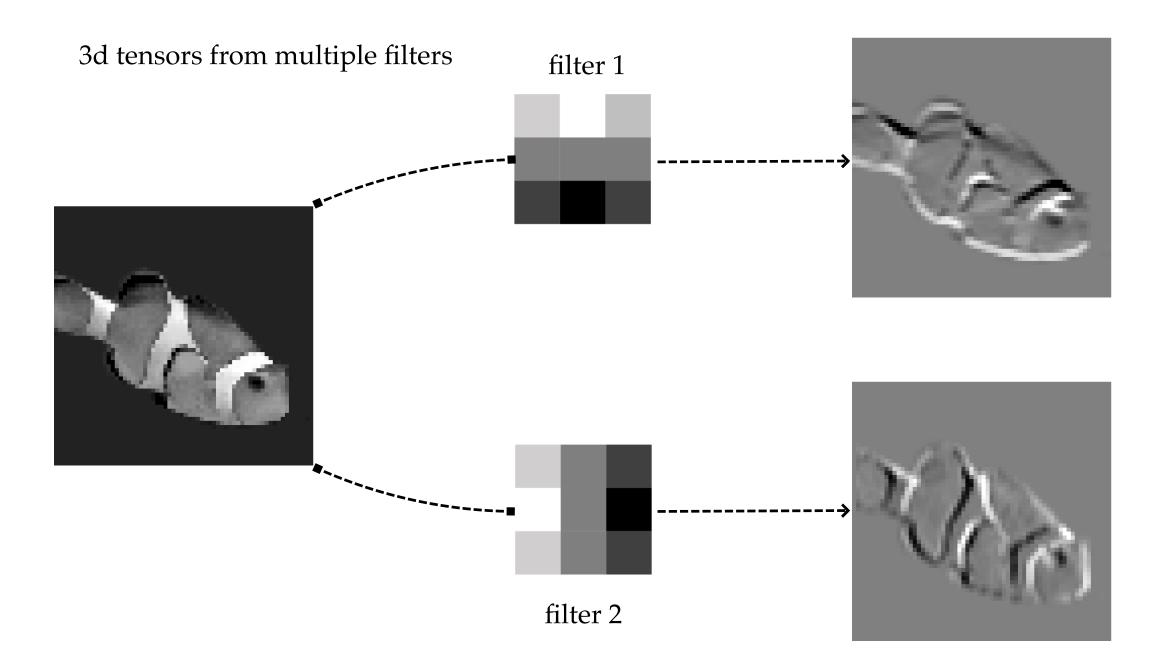


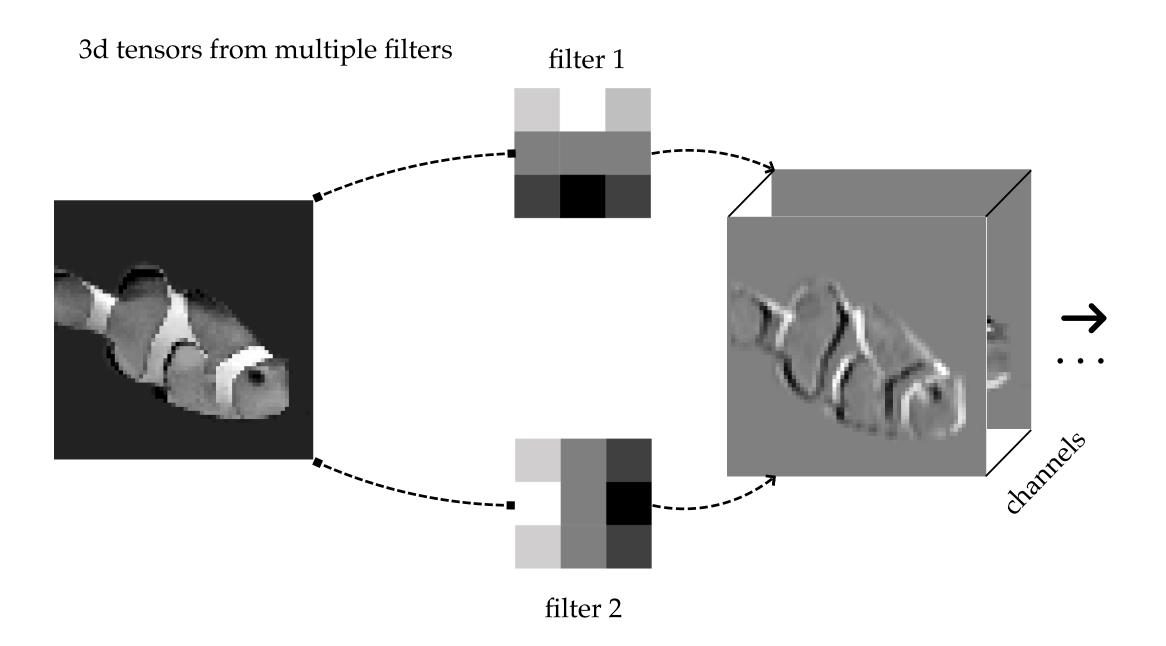
image height



image channels

image width



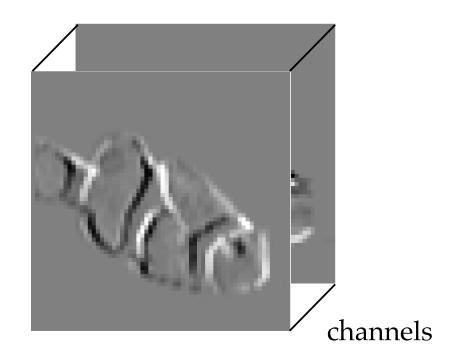


Why 3d tensors:

1. color input

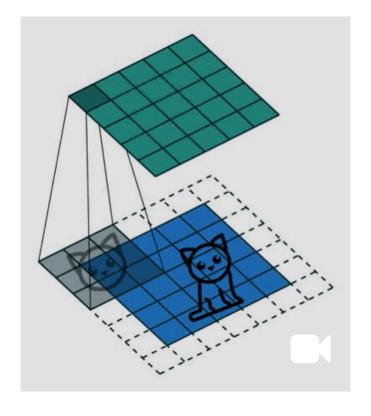


2. using multiple filters

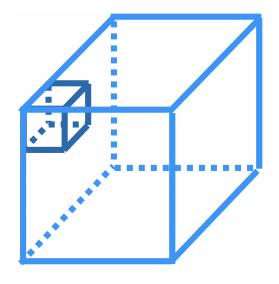


Why we *don't* typically do 3d convolution

2d convolution

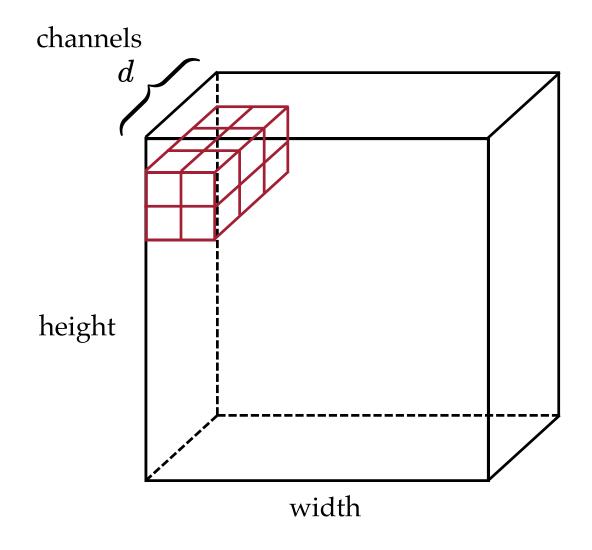


3d convolution



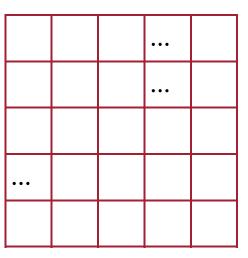


We *don't* typically do 3-dimensional convolution. Instead:



- 3d tensor input, channel *d*
- 3d tensor filter, channel *d*
- 2d convolution, 2d output





multiple filters multiple output matrices input tensor •

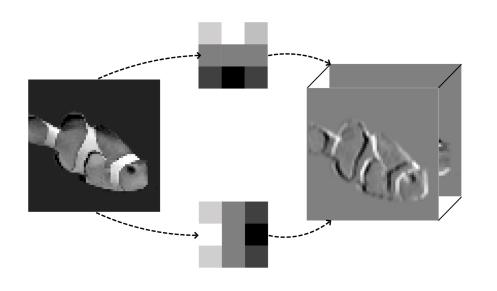
input tensor output tensor *k* filters

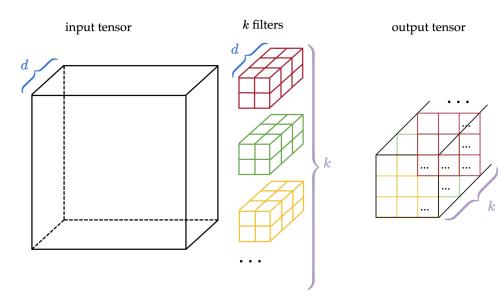
Every convolutional layer works with 3d tensors:

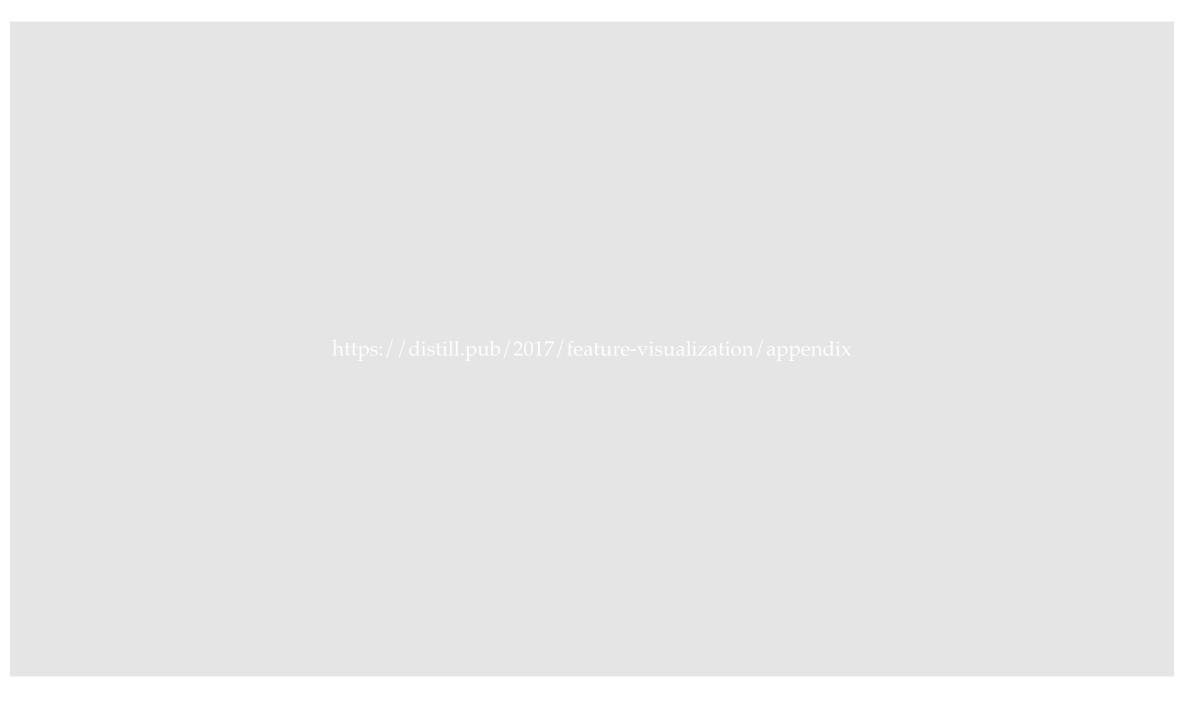
1. color input



2. the use of multiple filters in doing 2d convolution



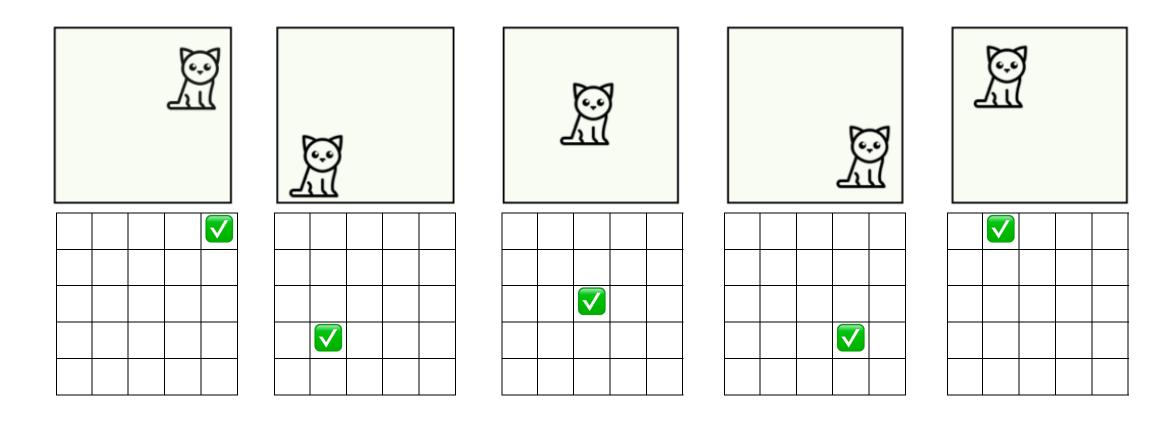




Outline

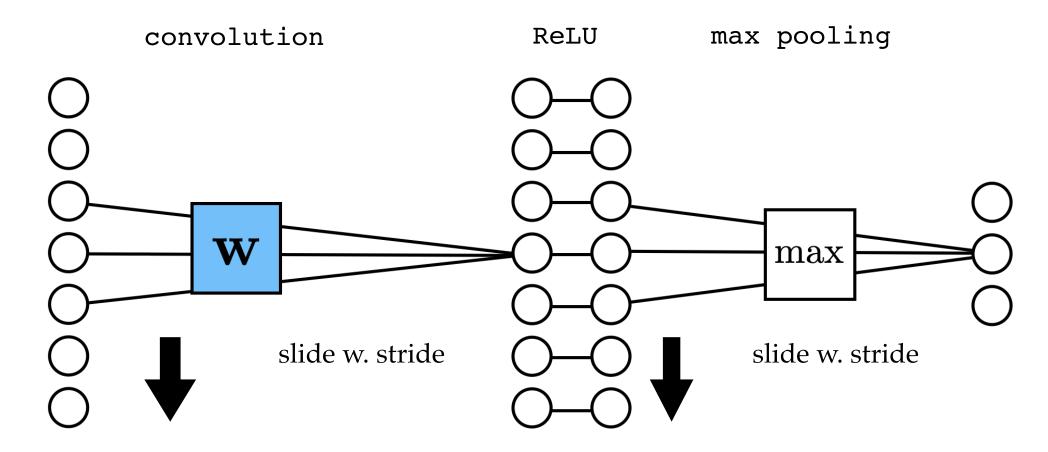
- Vision problem structure
- Convolution
 - 1-dimensional and 2-dimensional *convolution*
 - 3-dimensional *tensors*
- Max pooling
- (Case studies)

convolution helps detect pattern, but ...



cat moves, detection moves

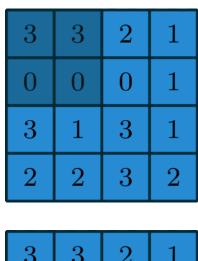
1d max pooling

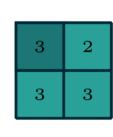


learnable filter weights detects pattern

no learnable parameter summarizes strongest response

2d max pooling



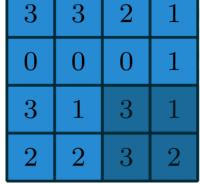


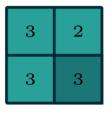
3	3	2	1
0	0	0	1
3	1	3	1
2	2	3	2

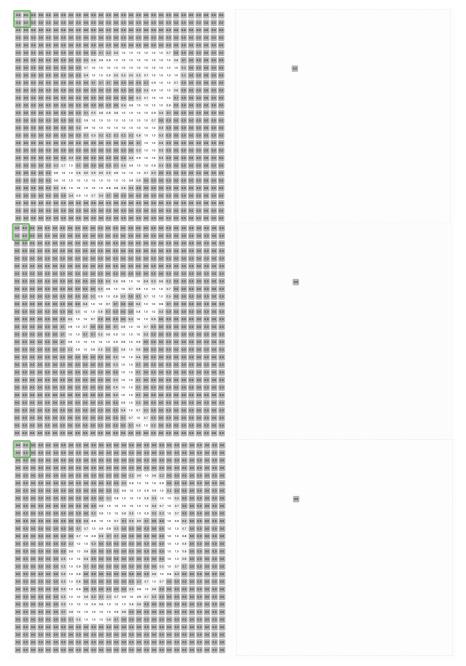


3	3	2	1
0	0	0	1
3	1	3	1
2	2	3	2



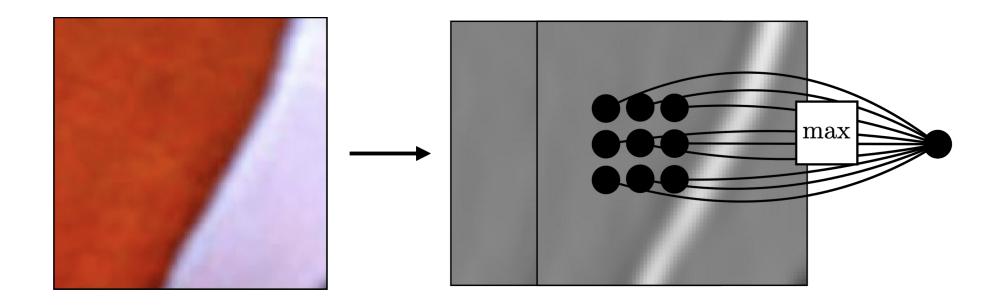






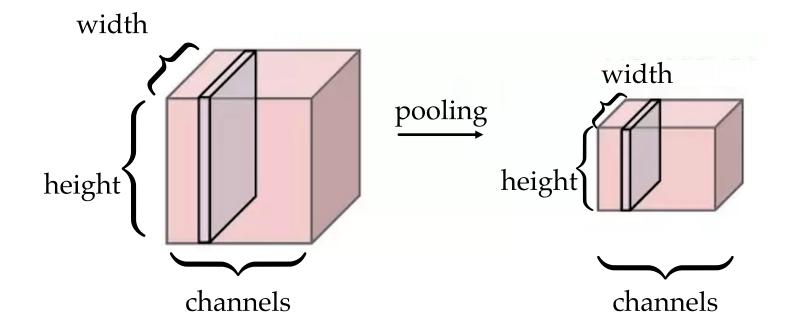
[image edited from vdumoulin gif adapted from demo source]

Pooling across *spatial* locations achieves invariance w.r.t. small translations:



large response regardless of exact position of edge

Pooling across *spatial* locations achieves invariance w.r.t. small translations:



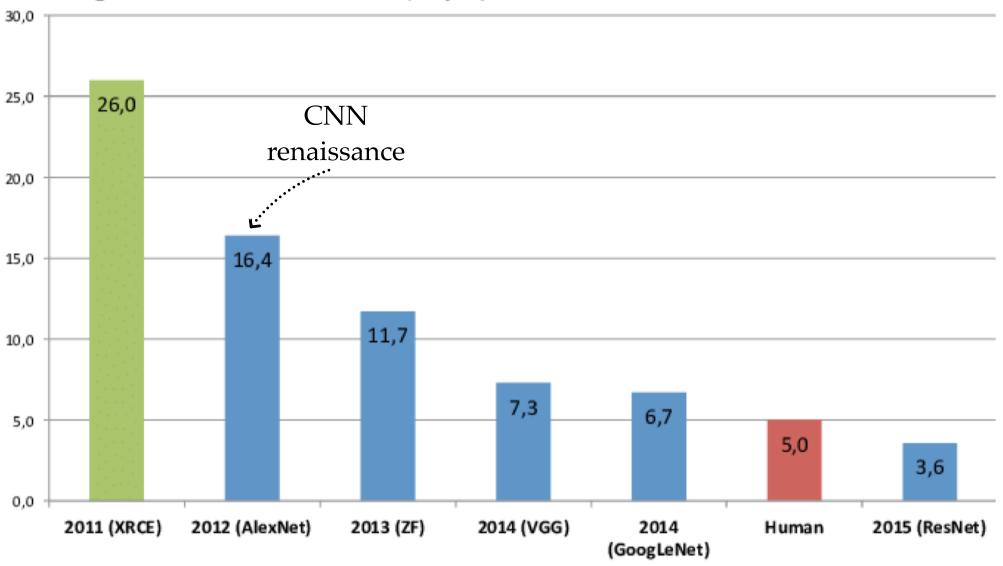
applied independently across all channels

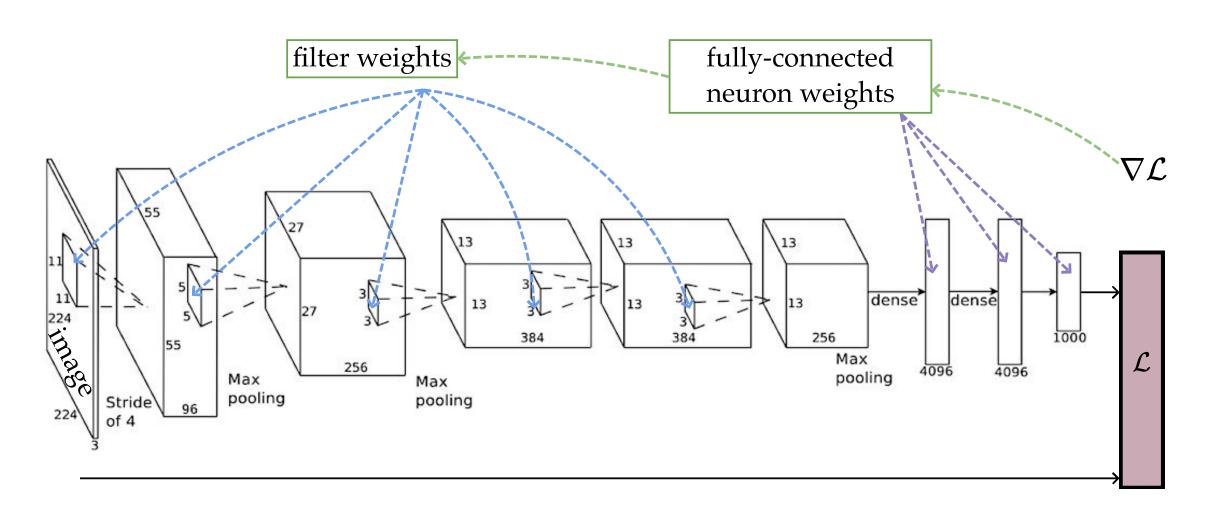
so the *channel* dimension remains *unchanged*

Outline

- Vision problem structure
- Convolution
 - 1-dimensional and 2-dimensional *convolution*
 - 3-dimensional *tensors*
- Max pooling
- (Case studies)

ImageNet Classification Error (Top 5)

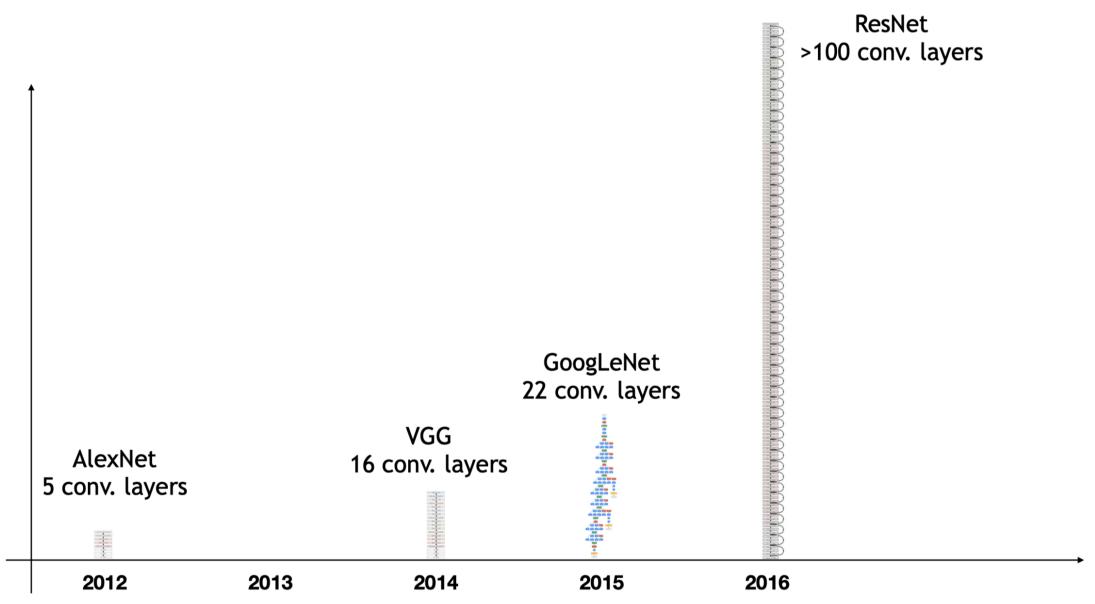




label

[AlexNet paper]

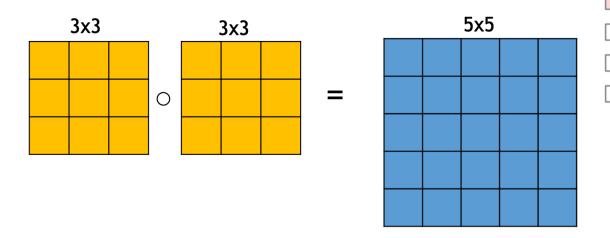
[all max pooling are via 3-by-3 filter, stride of 2]



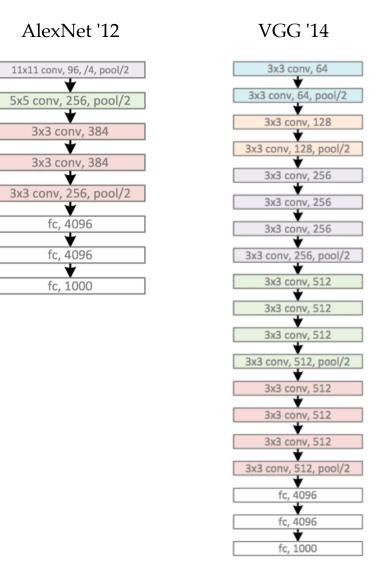
VGG '14

Main developments:

• small convolutional filters: only 3x3

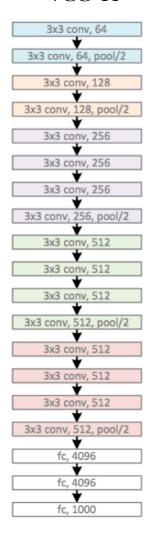


- increased depth: about 16 or 19 layers
- stack the same modules

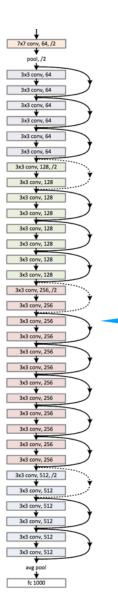


"Very Deep Convolutional Networks for Large-Scale Image Recognition", Simonyan & Zisserman. ICLR 2015 [image credit Philip Isola and Kaiming He]

VGG '14

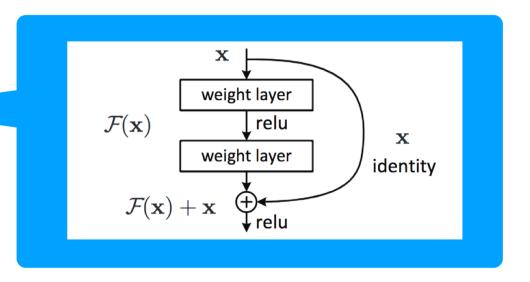


ResNet '16



Main developments:

- Residual block -- gradients can propagate faster (via the identity mapping)
- increased depth: > 100 layers



[He et al: Deep Residual Learning for Image Recognition, CVPR 2016] [image credit Philip Isola and Kaiming He]

Summary

- Though NN are technically "universal approximators", designing the NN structure so that it matches what we know about the underlying structure of the problem can substantially improve generalization ability and computational efficiency.
- Images are a very important input type and they have important properties that we can take advantage of: visual hierarchy, translation invariance, spatial locality.
- Convolution is an important image-processing technique that builds on these ideas. It can be interpreted as locally connected network, with weight-sharing.
- Pooling layer helps aggregate local info effectively, achieving bigger receptive field.
- We can train the parameters in a convolutional filtering function using backprop and combine convolutional filtering operations with other neural-network layers.

https://forms.gle/36SX9pqCTWpp323N8

We'd love to hear your thoughts.

Thanks!