# 6.390 Intro to Machine Learning

## Lecture 7: Convolutional Neural Networks

Shen Shen
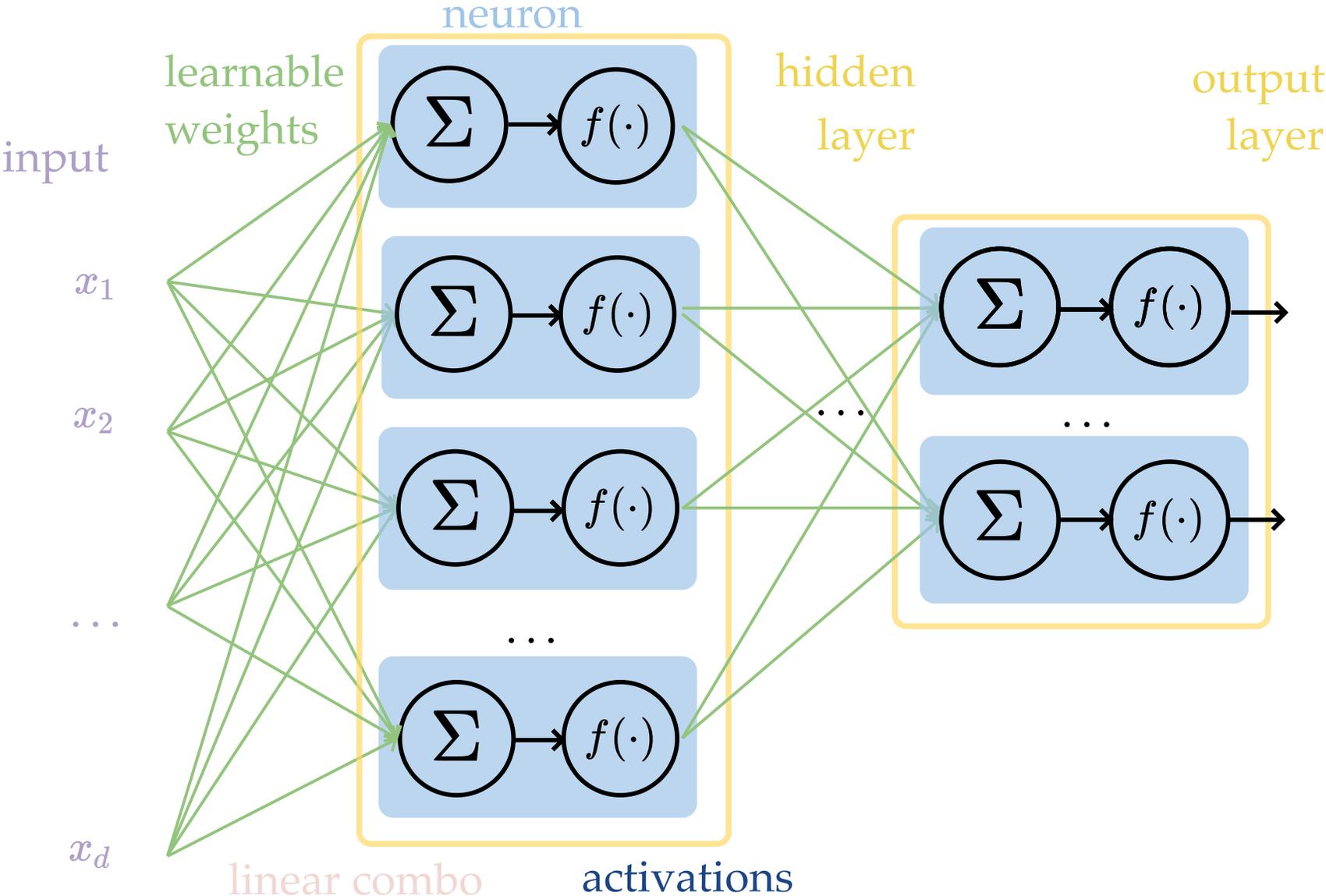
Mar 16, 2026

3pm, Room 10-250

Slides and Lecture Recording
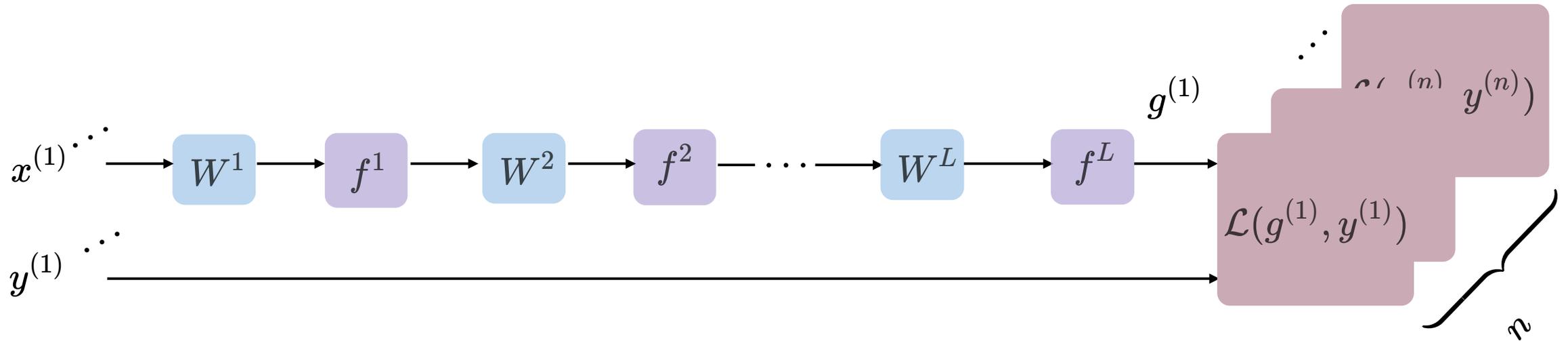
**Recap**    a (fully-connected, feed-forward) neural network

Forward pass: evaluate, *given* the current parameters



- the model outputs $g^{(i)} = f^L \left( \ldots f^2 \left( f^1(\mathbf{x}^{(i)}; \mathbf{W}^1); \mathbf{W}^2 \right) ; \ldots \mathbf{W}^L \right)$

- the loss incurred on the current data $\mathcal{L}(g^{(i)}, y^{(i)})$

- the training error $J = \frac{1}{n} \sum_{i=1}^{n} \mathcal{L}(g^{(i)}, y^{(i)})$
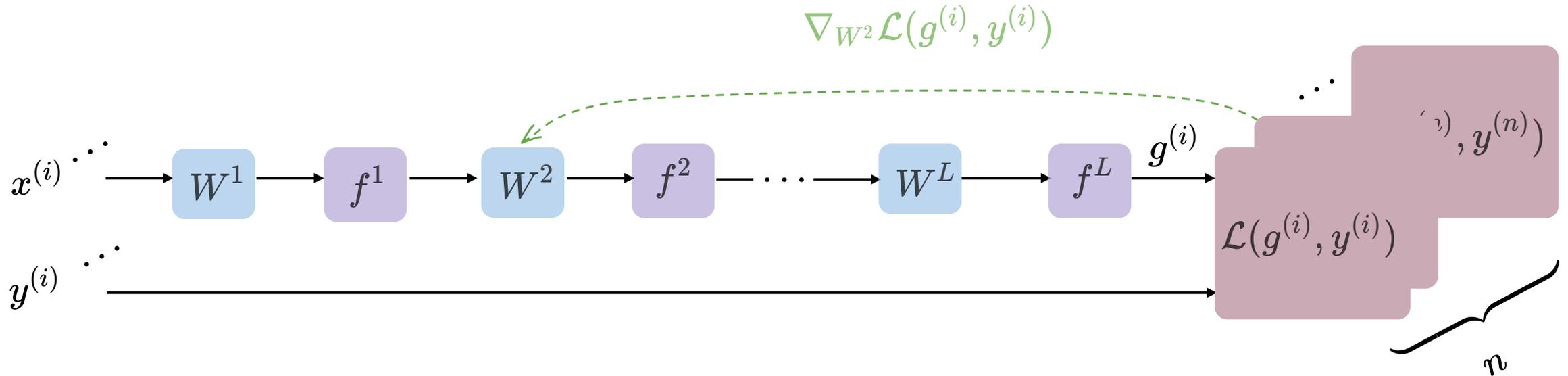
linear combination

(nonlinear) activation

loss function

Recap

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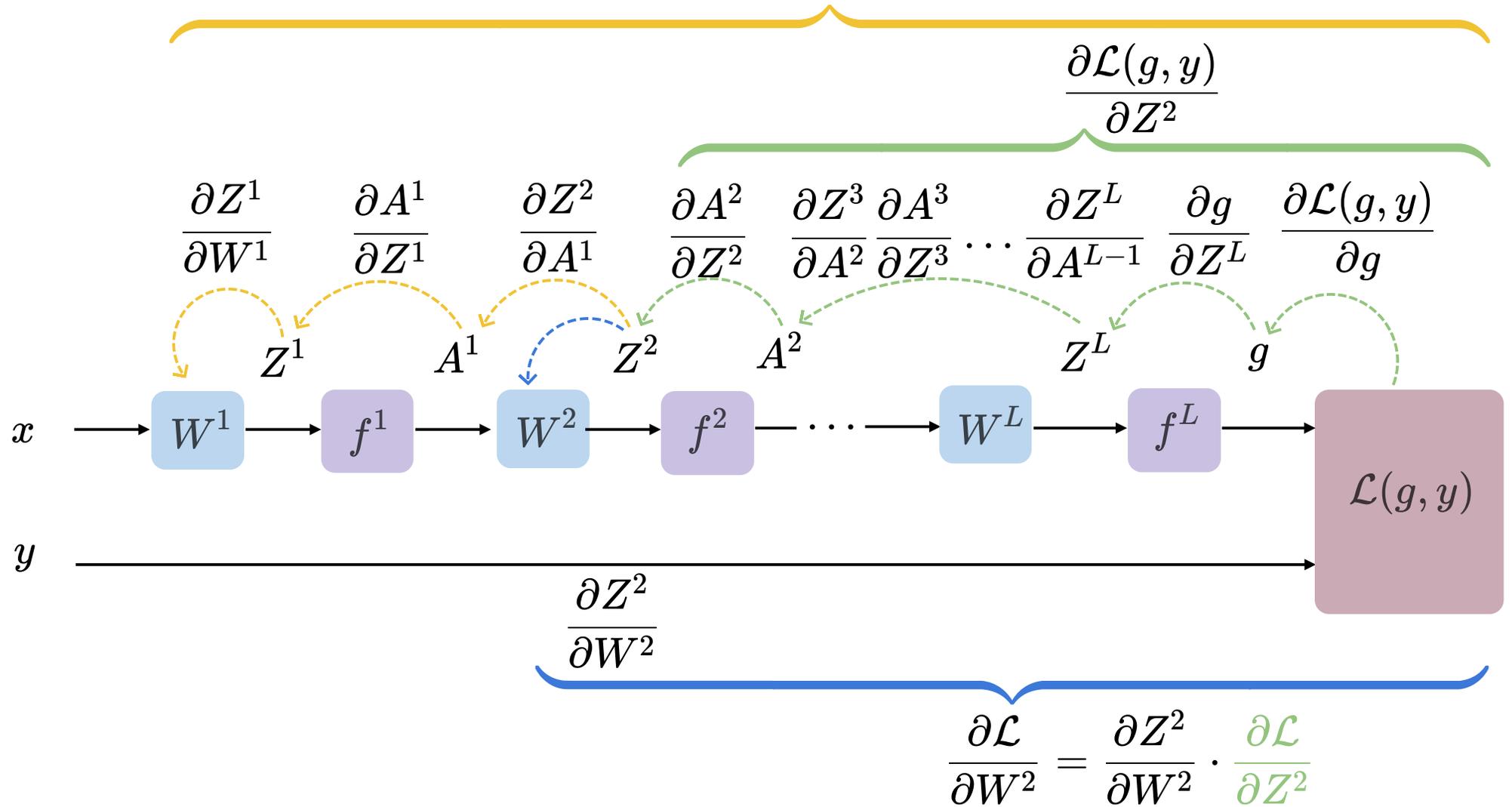
e.g. to update $W^2$

- Randomly pick a data point $(x^{(i)}, y^{(i)})$

- Evaluate the gradient $\nabla_{W^2} \mathcal{L}(g^{(i)}, y^{(i)})$

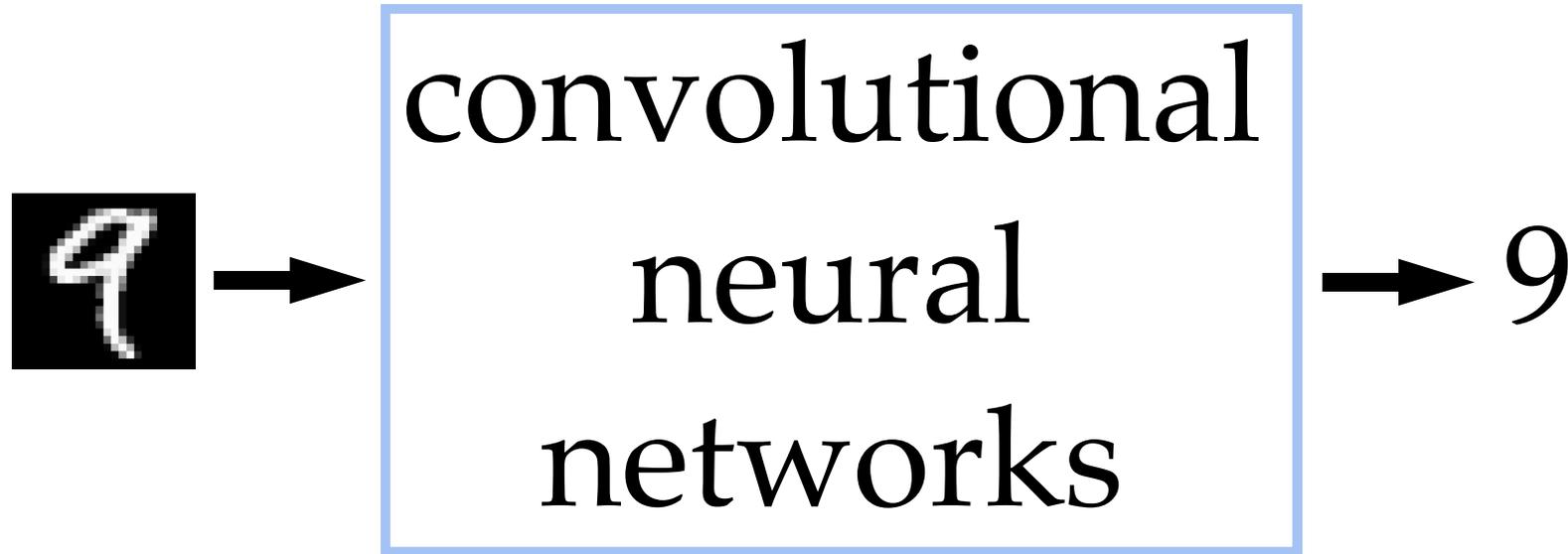- Update the weights $W^2 \leftarrow W^2 - \eta \nabla_{W^2} \mathcal{L}(g^{(i)}, y^{(i)})$

$\nabla_{W^2} \mathcal{L}(g^{(i)}, y^{(i)})$

$x^{(i)}$ $\cdots$ → $W^1$ → $f^1$ → $W^2$ → $f^2$ — $\cdots$ — $W^L$ → $f^L$ → $g^{(i)}$

$\quad ^{(i)}, y^{(n)})$

$\mathcal{L}(g^{(i)}, y^{(i)})$

$y^{(i)}$ $\cdots$

$n$

Recap

backpropagation: reuse of computation

$$\frac{\partial \mathcal{L}}{\partial W^1} = \frac{\partial Z^1}{\partial W^1} \cdot \frac{\partial A^1}{\partial Z^1} \cdot \frac{\partial Z^2}{\partial A^1} \cdot \textcolor{green}{\frac{\partial \mathcal{L}}{\partial Z^2}}$$
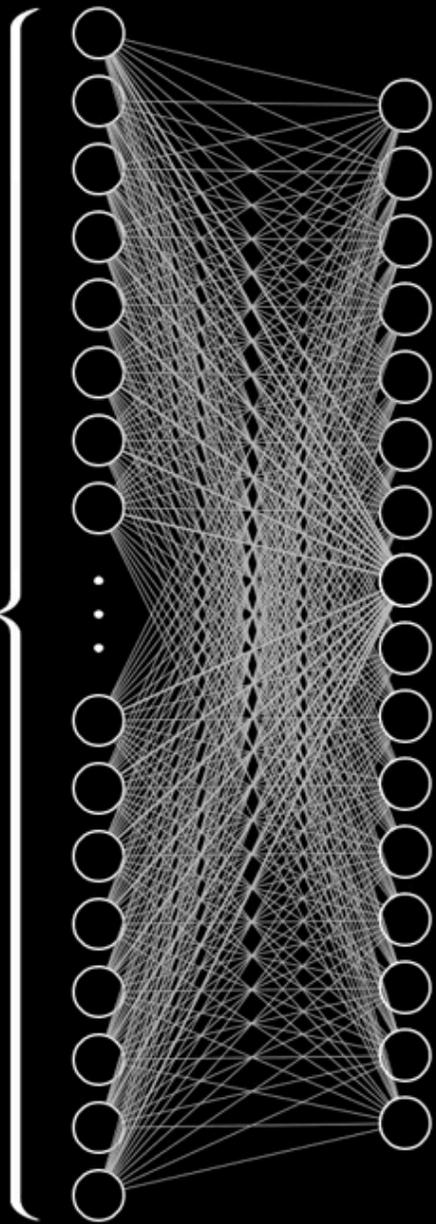
$$\frac{\partial \mathcal{L}(g,y)}{\partial Z^2}$$



$$\frac{\partial Z^1}{\partial W^1} \quad \frac{\partial A^1}{\partial Z^1} \quad \frac{\partial Z^2}{\partial A^1} \quad \frac{\partial A^2}{\partial Z^2} \quad \frac{\partial Z^3}{\partial A^2} \frac{\partial A^3}{\partial Z^3} \cdots \frac{\partial Z^L}{\partial A^{L-1}} \quad \frac{\partial g}{\partial Z^L} \quad \frac{\partial \mathcal{L}(g,y)}{\partial g}$$

$Z^1 \qquad A^1 \qquad Z^2 \qquad A^2 \qquad\qquad Z^L \qquad g$

$x \longrightarrow W^1 \rightarrow f^1 \rightarrow W^2 \rightarrow f^2 \rightarrow \cdots \rightarrow W^L \rightarrow f^L \rightarrow \mathcal{L}(g,y)$

$y \longrightarrow$

$$\frac{\partial Z^2}{\partial W^2}$$

$$\frac{\partial \mathcal{L}}{\partial W^2} = \frac{\partial Z^2}{\partial W^2} \cdot \textcolor{green}{\frac{\partial \mathcal{L}}{\partial Z^2}}$$

convolutional neural networks → 9

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)

784×16 weights        16 biases

784

426-by-426
grayscale image



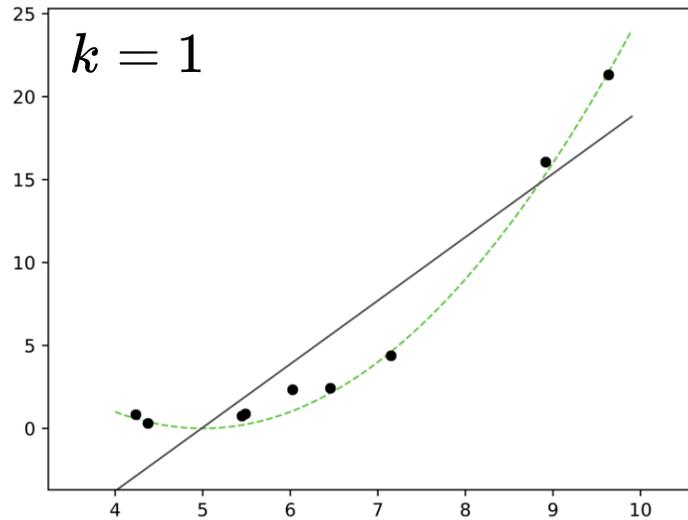Use the same 2 hidden-layer network to predict what top-10 engineering school seal this image is, need to learn ~3M parameters.

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

Underfitting

Appropriate

Overfitting



$k = 1$

$k = 2$

$k = 10$

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



Hubel & Weisel

featural 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 effectively.

typical CNN architecture for image classification

the same feedforward net as before



CNN

INPUT

FEATURE LEARNING

FLATTEN    FULLY CONNECTED    SOFTMAX

CAR
TRUCK
VAN

BICYCLE

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:

example: 1-dimensional convolution

input

| 0 | 1 | 0 | 1 | 1 |
|---|---|---|---|---|

filter

| -1 | 1 |
|----|---|

$$(0 * -1) + (1 * 1) = 1$$

convolved output

| 1 | | | |
|---|---|---|---|

example: 1-dimensional convolution

input

| 0 | 1 | 0 | 1 | 1 |
|---|---|---|---|---|

filter

| -1 | 1 |
|----|---|

$$(1 * -1) + (0 * 1) = -1$$

convolved output

| 1 | -1 | | |
|---|----|--|--|

example: 1-dimensional convolution

input

| 0 | 1 | 0 | 1 | 1 |

filter

| -1 | 1 |

$$(0 * -1) + (1 * 1) = 1$$

convolved output

| 1 | -1 | 1 | |

example: 1-dimensional convolution

input

| 0 | 1 | 0 | 1 | 1 |
|---|---|---|---|---|

filter

| -1 | 1 |
|----|---|

$$(1 * -1) + (1 * 1) = 0$$

convolved output

| 1 | -1 | 1 | 0 |
|---|----|---|---|

*convolution interpretation 1:*                    template matching

input

| 0 | 1 | -1 | 1 | 1 |

filter

| -1 | 1 |

convolved
output

| 1 | -2 | 2 | 0 |

*convolution interpretation 2:*

"look" *locally* through the filter

this local region = *receptive field*

input

| 0 | 1 | -1 | 1 | 1 |
|---|---|----|---|---|

filter

| -1 | 1 |
|----|---|

convolved
output

| 1 | -2 | 2 | 0 |
|---|----|---|---|

input

output

*convolution interpretation 3:*

sparse-connected layer
with parameter sharing

| 0 | 1 | -1 | 1 | 1 |
|---|---|---|---|---|

convolve with

| -1 | 1 |
|---|---|

dot product with

|   |   |   |   |
|---|---|---|---|
| -1 | 0 | 0 | 0 |
| 1 | -1 | 0 | 0 |
| 0 | 1 | -1 | 0 |
| 0 | 0 | 1 | -1 |
| 0 | 0 | 0 | 1 |

=

| 1 | -2 | 2 | 0 |
|---|---|---|---|

| 0 | 1 | 0 | 1 | 1 |
|---|---|---|---|---|

convolve with

| 1 |
|---|

dot product with

$I_{5\times 5}$

| 0 | 1 | 0 | 1 | 1 |
|---|---|---|---|---|

*convolution interpretation 4:*                    translational equivariance

input

| | | | | |
|---|---|---|---|---|
| 0 | 1 | 0 | 1 | 0 |

filter

| | |
|---|---|
| 0 | 1 |

convolved
output

| | | | |
|---|---|---|---|
| 1 | 0 | 1 | 0 |

example: 2-dimensional convolution

input

| 3 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|
| 0 | 0 | 1 | 3 | 1 |
| 3 | 1 | 2 | 2 | 3 |
| 2 | 0 | 0 | 2 | 2 |
| 2 | 0 | 0 | 0 | 1 |

filter

| 0 | 1 | 2 |
|---|---|---|
| 2 | 2 | 0 |
| 0 | 1 | 2 |

convolved output

| 12 | 12 | 17 |
|----|----|----|
| 10 | 17 | 19 |
| 9  | 6  | 14 |

stride of 2

input

| 3 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|
| 0 | 0 | 1 | 3 | 1 |
| 3 | 1 | 2 | 2 | 3 |
| 2 | 0 | 0 | 2 | 2 |
| 2 | 0 | 0 | 0 | 1 |

filter

| 0 | 1 | 2 |
|---|---|---|
| 2 | 2 | 0 |
| 0 | 1 | 2 |

output

| 12 | 17 |
|----|----|
| 9  | 14 |

stride of 2, with padding of size 1

input

| 3 | 3 | 2 | 1 | 0 |
|---|---|---|---|---|
| 0 | 0 | 1 | 3 | 1 |
| 3 | 1 | 2 | 2 | 3 |
| 2 | 0 | 0 | 2 | 2 |
| 2 | 0 | 0 | 0 | 1 |

filter

| 0 | 1 | 2 |
|---|---|---|
| 2 | 2 | 0 |
| 0 | 1 | 2 |

output

| 6 | 17 | 3 |
|---|----|---|
| 8 | 17 | 13 |
| 6 | 4 | 4 |

`quick summary: hyperparameters for 1d convolution`

- Zero-padding

| 0 | 0 | 1 | 0 | 1 | 1 | 0 |
|---|---|---|---|---|---|---|

- Stride (e.g. stride of 2)

| 0 | 0 | 1 | 0 | 1 | 1 | 0 |
|---|---|---|---|---|---|---|

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

| -1 | 1 |
|----|---|

| 1 |
|---|

these weights are what
CNN learn eventually

# quick summary: hyperparameters for 2d convolution

Input (5, 5)
After-padding (5, 5)
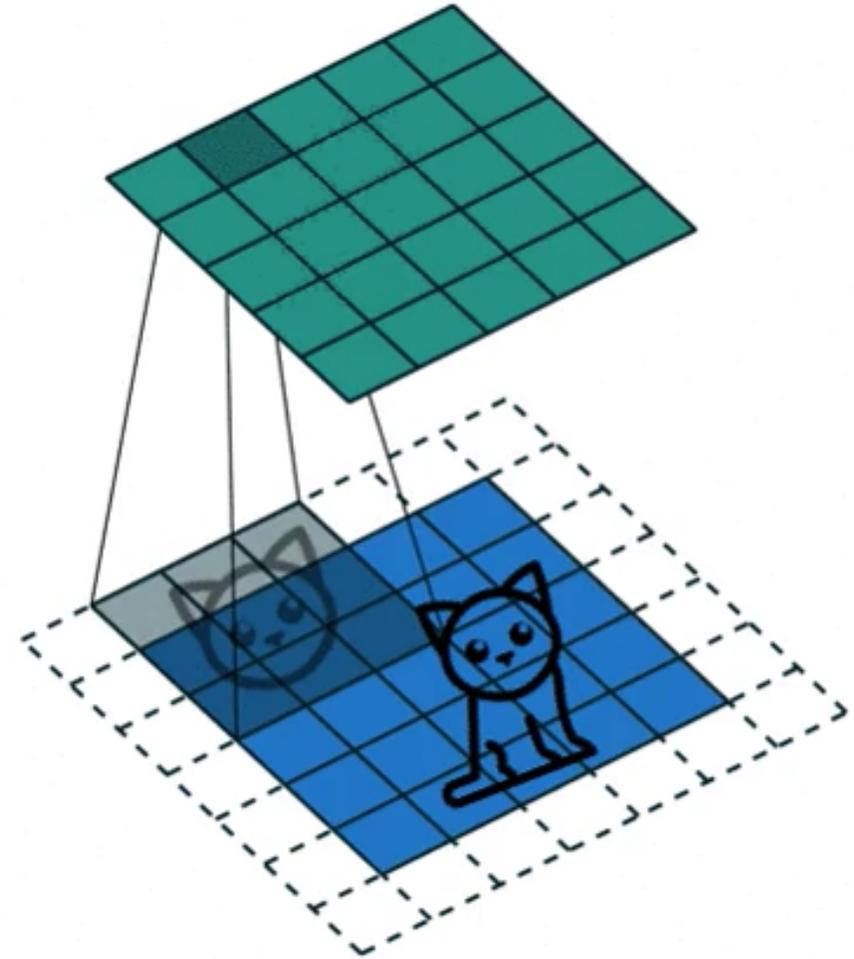
Output (4, 4)

Input Size:

5

## Hyperparameters

Padding:

0

Filter Size:

2

Stride:

1

ⓘ *Hover over* to change focus. *Click outside* to resume.

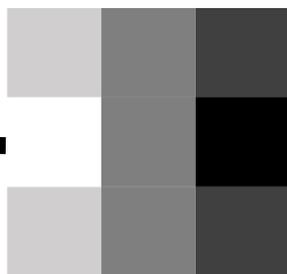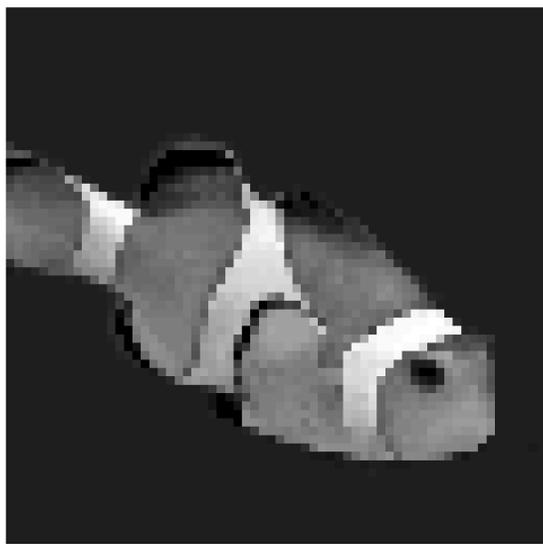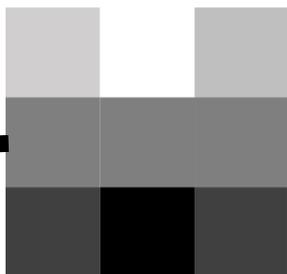Adapted from CNN Explainer with ❤️ by Shen² | Bug Report

# quick summary: convolution interpretation

- Look locally (sparse connections)

- Parameter sharing

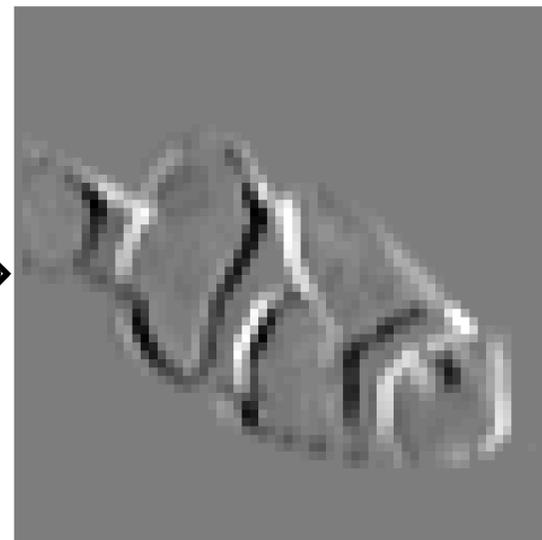- Template matching

- Translational equivariance
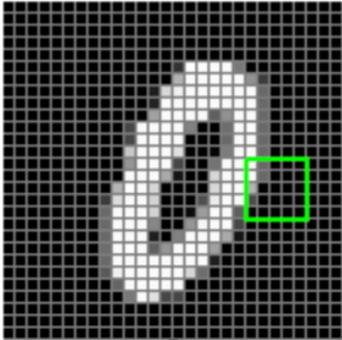
hand-designed filters (e.g. Sobel)

filter 1

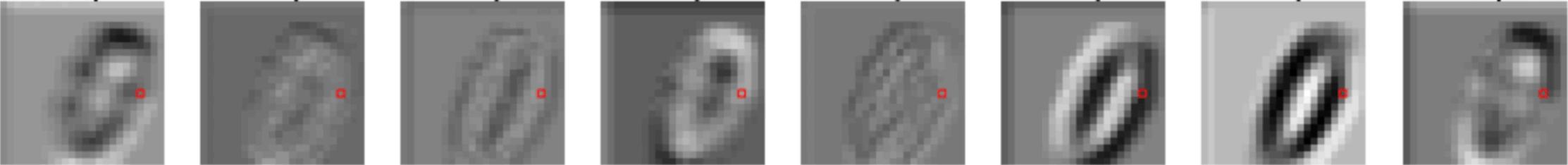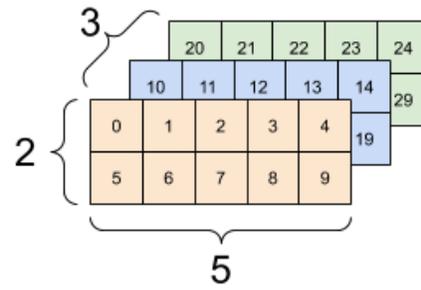filter 2

learned filters detect many patterns
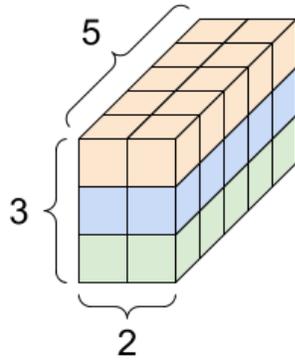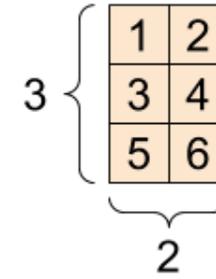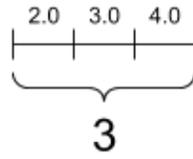
input

filters

conv'd output

# Outline

- Vision problem structure

- Convolution

  - 1-dimensional and 2-dimensional *convolution*

  - **3-dimensional** *tensors*

- Max pooling

- (Case studies)

# A tender intro to tensor:
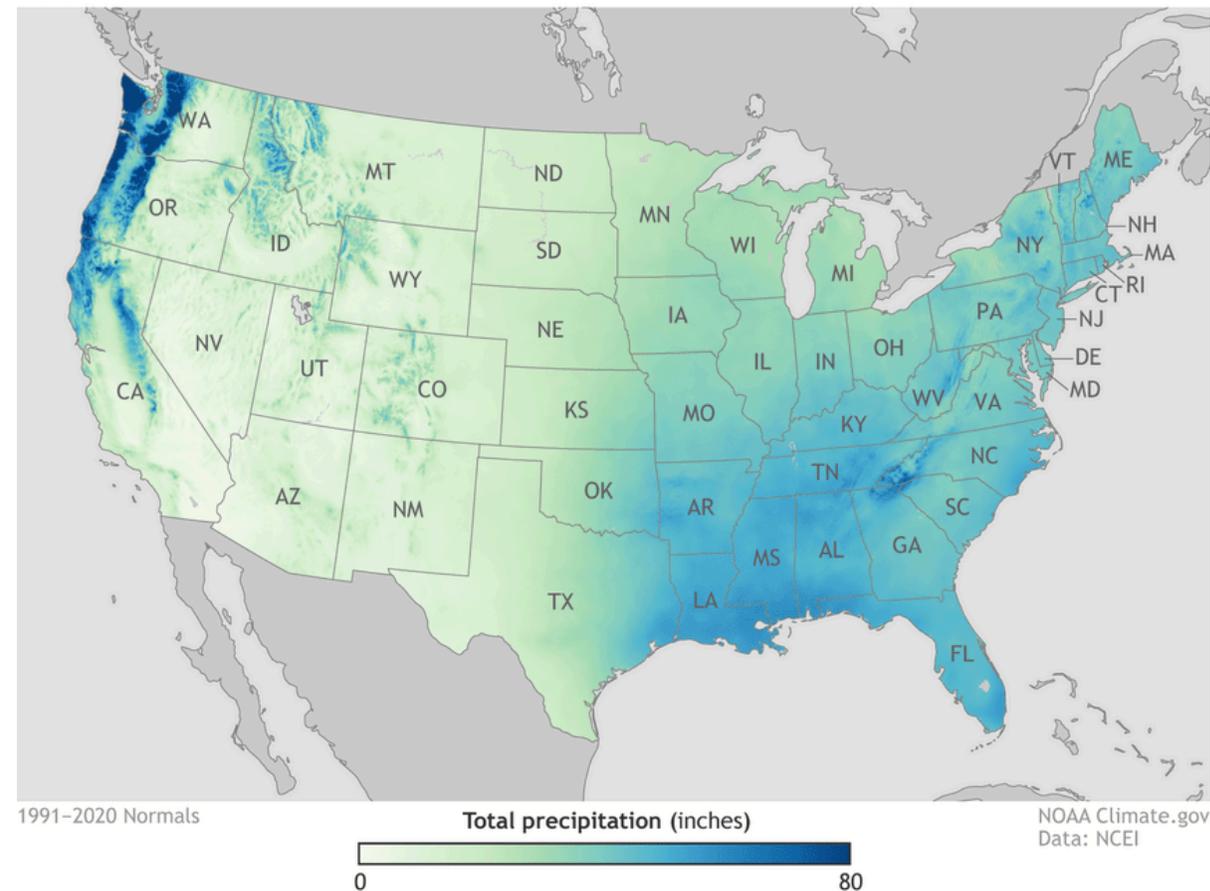
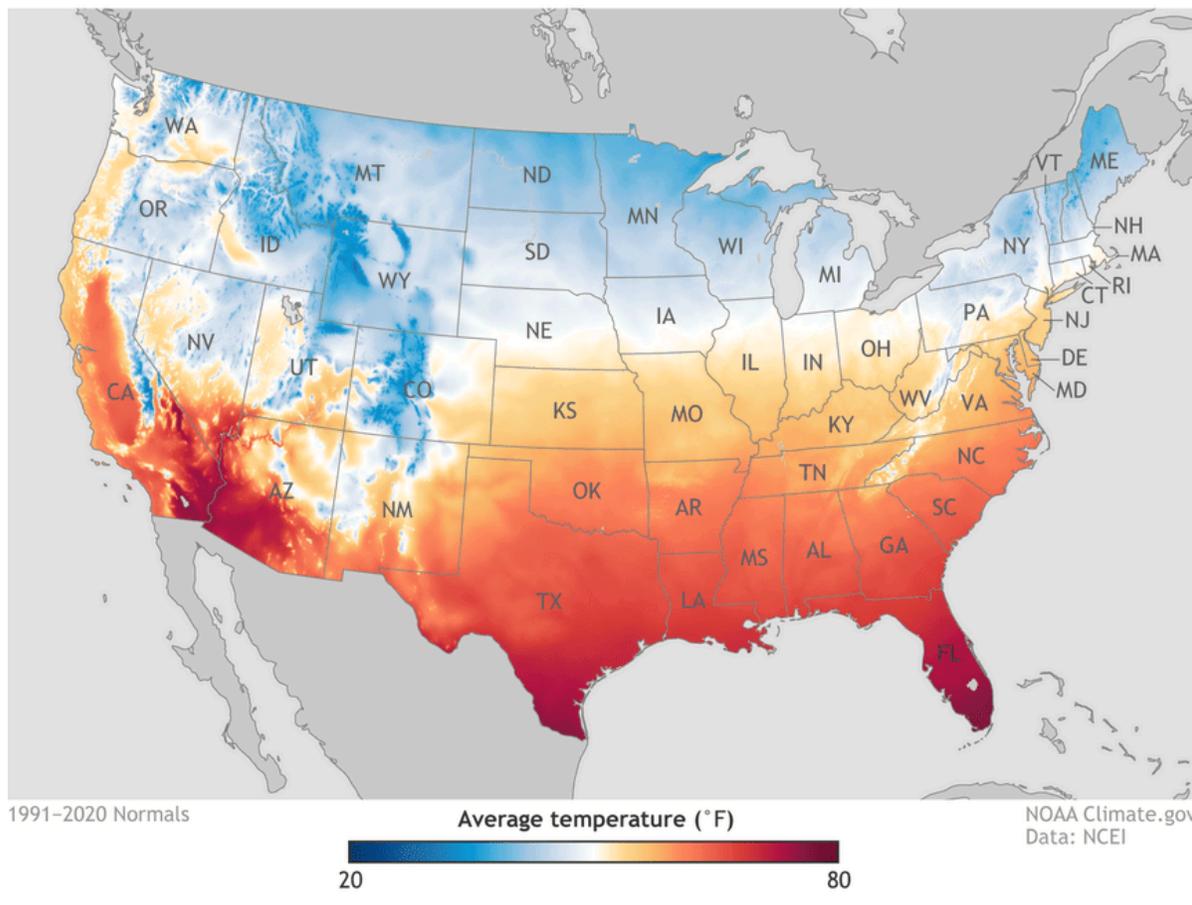# color images and channels



red

green

blue

each channel is a *complete but independent* view of the same scene

like when we think of weather:



so channels are often referred to as *feature maps*
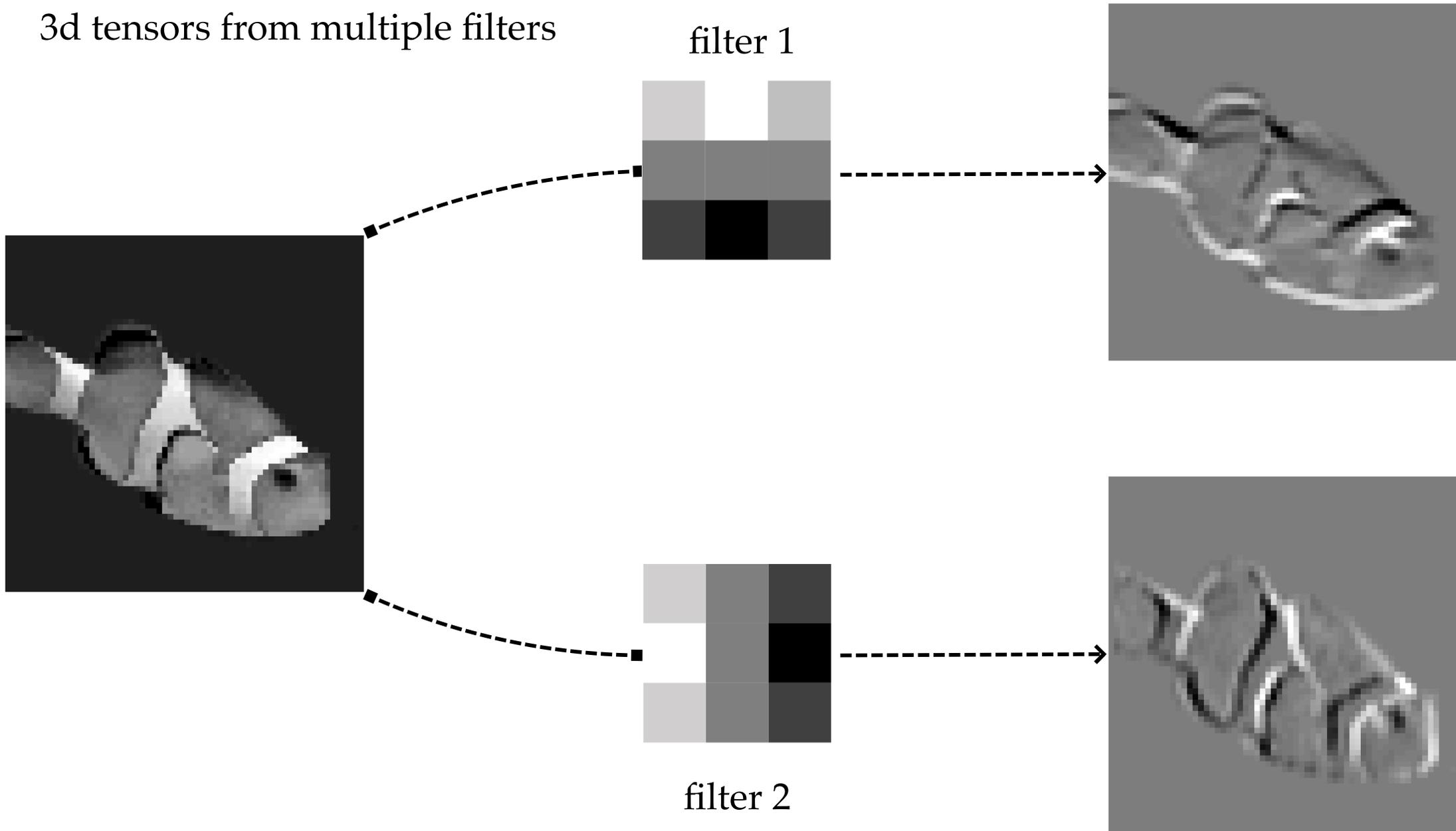
# 3d tensors from color channels



image height

image width

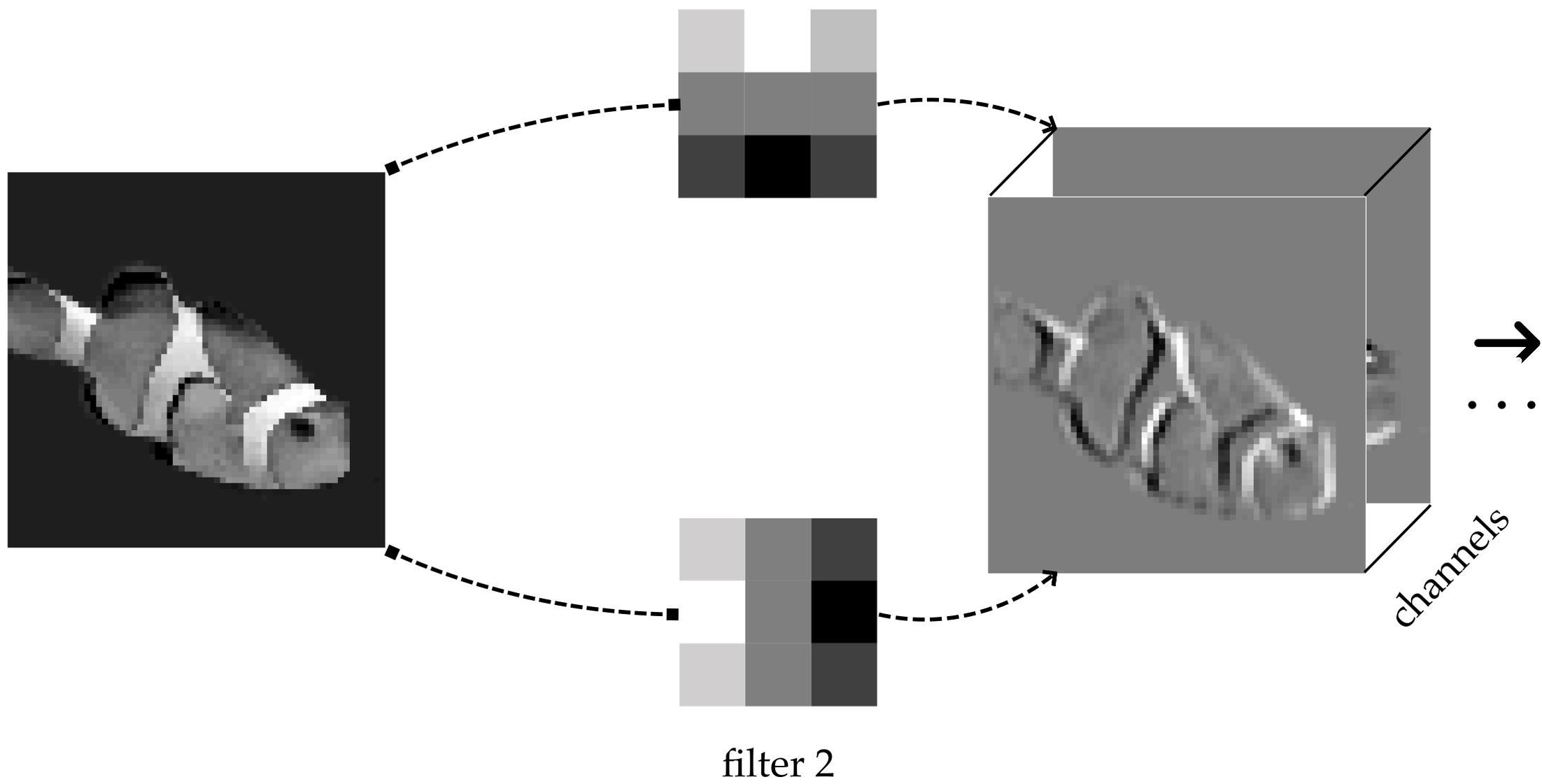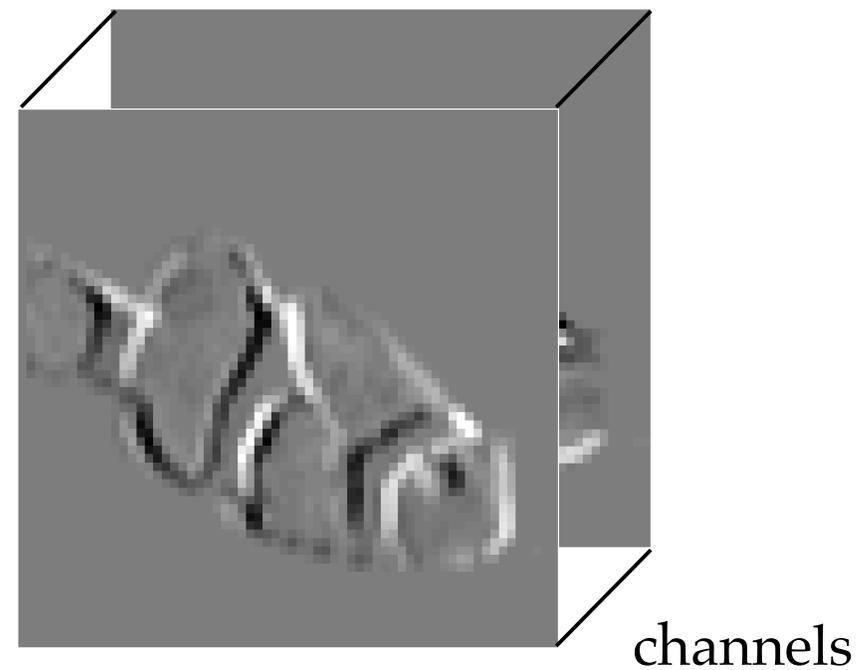image channels

3d tensors from multiple filters

filter 1

filter 2

3d tensors from multiple filters

filter 1

filter 2

channels

where do channels come from?

1. color input

height

width

channels
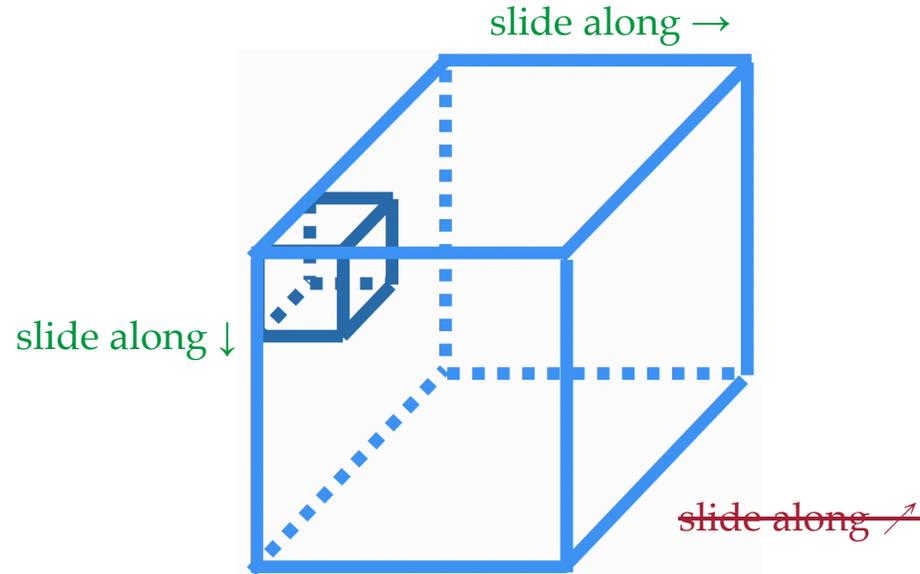
2. multiple filters → multiple channels
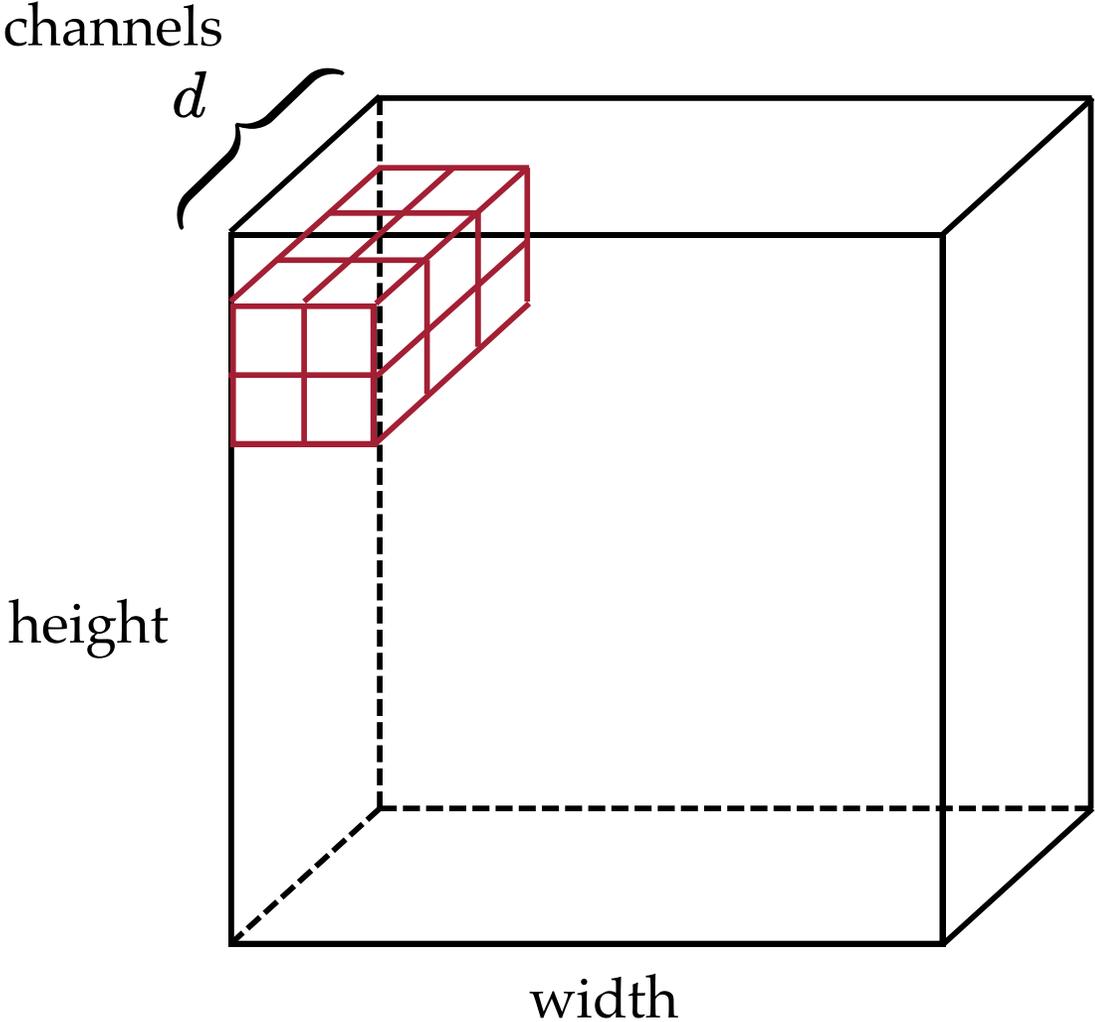
channels

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



Convolution **shares** weights across **shifted** positions

shifting makes sense spatially (a cat can be anywhere)

but not across channels (red ≠ shifted green)

3D conv is used when the third axis is spatial/temporal (MRI, video)
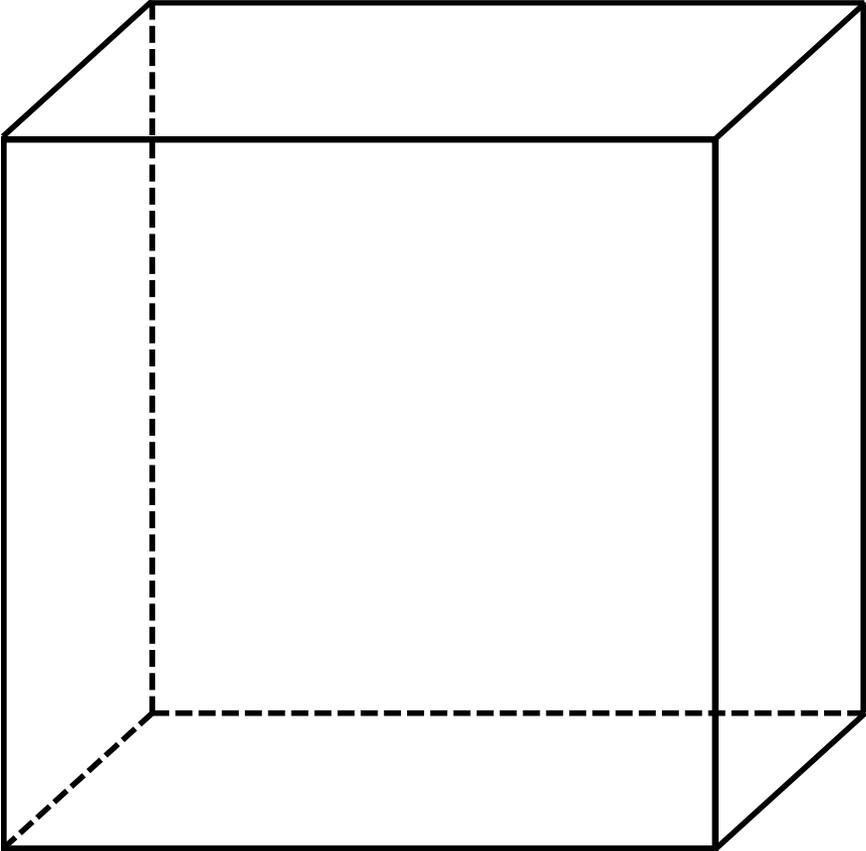
full-depth 2D convolution

channels
$d$

height

width

- 3d tensor input, channel $d$
- 3d tensor filter, channel $d$

- 2d convolution, 2d output

output

| | | | ... | |
|---|---|---|---|---|
| | | | ... | |
| | | | | |
| ... | | | | |
| | | | | |

input tensor

multiple filters

multiple output matrices

input tensor $\qquad$ $k$ filters $\qquad$ output tensor

$d$

$d$

$k$

$k$

# Every convolutional layer works with 3d tensors:

## 1. color input



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



input tensor     $k$ filters     output tensor

# Feature Visualization — Appendix

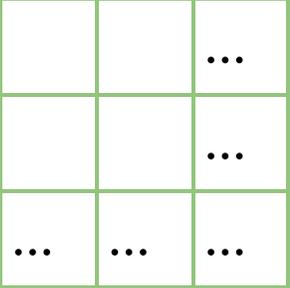After reading "Feature Visualization" you may be curious what other channels of GoogLeNet look like.



This appendix contains layers 3a through 5b of GoogLeNet.

*Below you can click on the boxes to see all units of that layer. You may al...*

# 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

convolution                    ReLU        max pooling



slide w. stride               slide w. stride

learnable filter weights       no learnable parameter

detects pattern                summarizes strongest response

# 2d max pooling

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

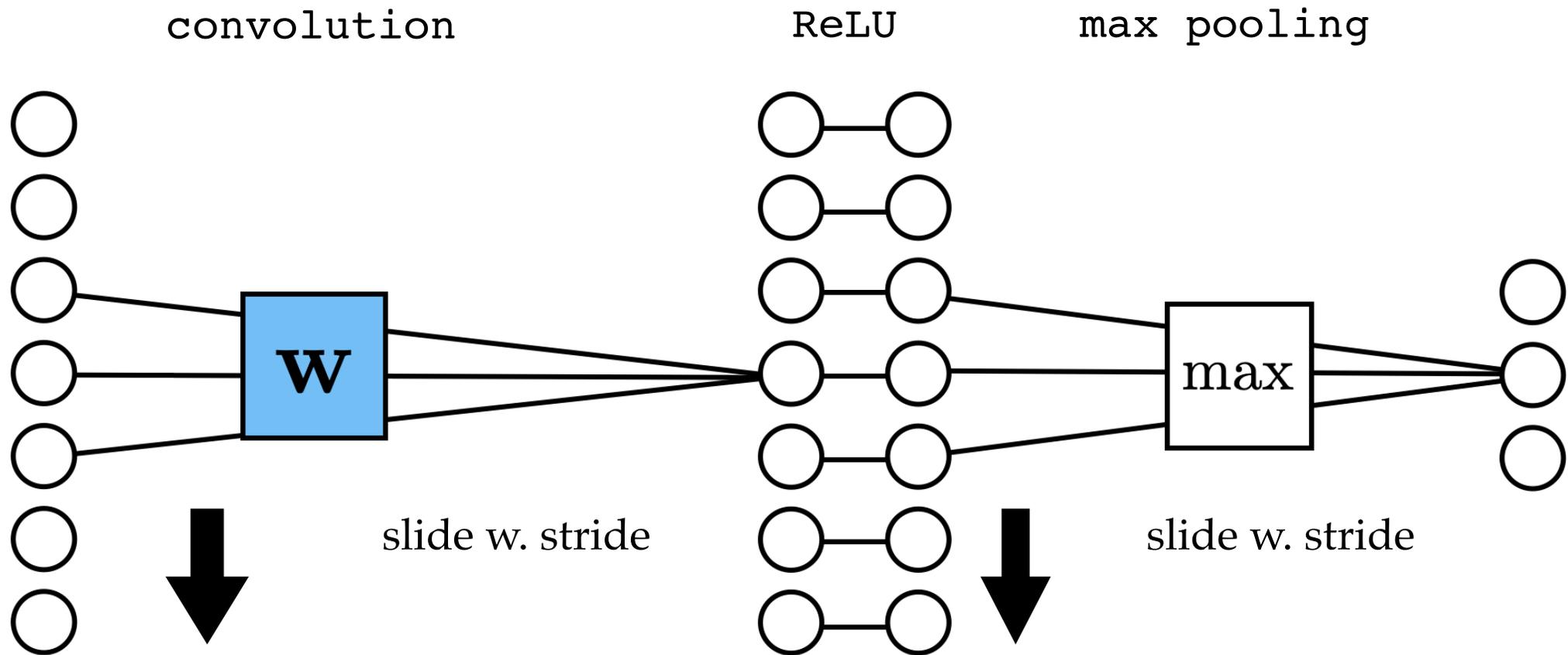  - 1-dimensional and 2-dimensional *convolution*

  - 3-dimensional *tensors*

- Max pooling

- (Case studies)

ImageNet Classification Error (Top 5)

CNN renaissance

[image credit Philip Isola]

filter weights

fully-connected neuron weights

$\nabla \mathcal{L}$

image

224

224

3

11

11

Stride of 4

55

55

96

Max pooling

27

27

256

Max pooling

5

5

13

13

384

3

3

13

13

384

3

3

13

13

256

Max pooling

dense

dense

4096

4096

1000

$\mathcal{L}$

label

[AlexNet paper]

[all max pooling are 3×3 filter, stride 2; pooled outputs not explicitly shown on diagram: 27×27×96, 13×13×256, 6×6×256]

AlexNet
5 conv. layers

VGG
16 conv. layers

GoogLeNet
22 conv. layers

ResNet
>100 conv. layers

2012    2013    2014    2015    2016

[image credit Philip Isola]

# VGG '14

## Main developments:

- small convolutional filters: only 3x3
- increased depth: about 16 or 19 layers of the same modules

**two 3×3 convolutions**

$3 \times 3$

$3 \times 3$

$5 \times 5$ input

### AlexNet '12

| |
|---|
| 11x11 conv, 96, /4, pool/2 |
| 5x5 conv, 256, pool/2 |
| 3x3 conv, 384 |
| 3x3 conv, 384 |
| 3x3 conv, 256, pool/2 |
| fc, 4096 |
| fc, 4096 |
| fc, 1000 |

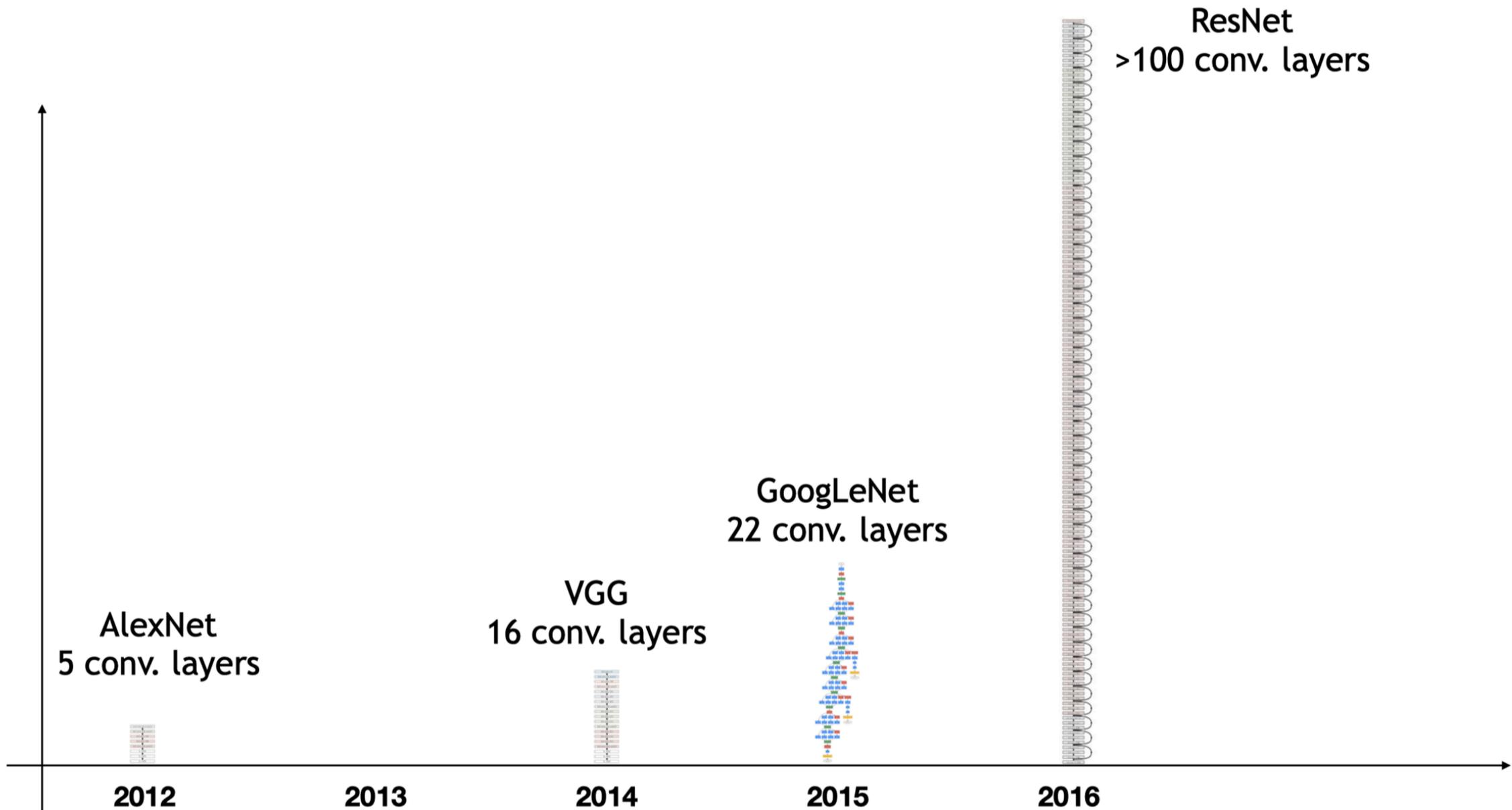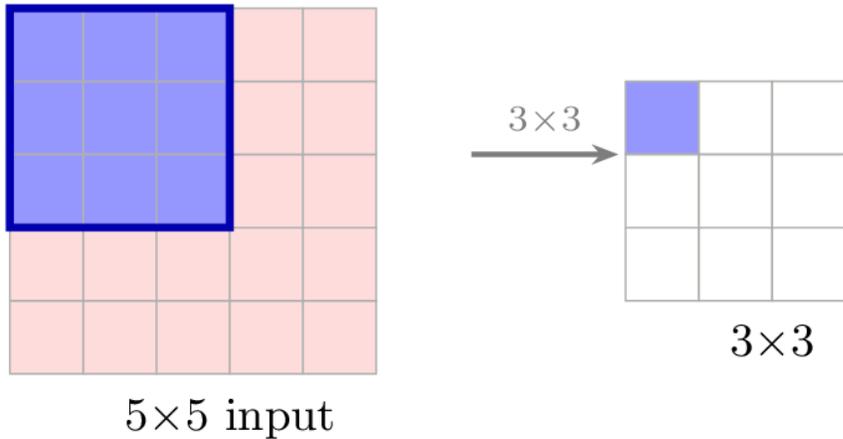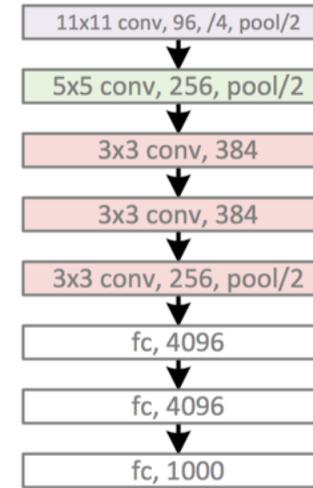### VGG '14

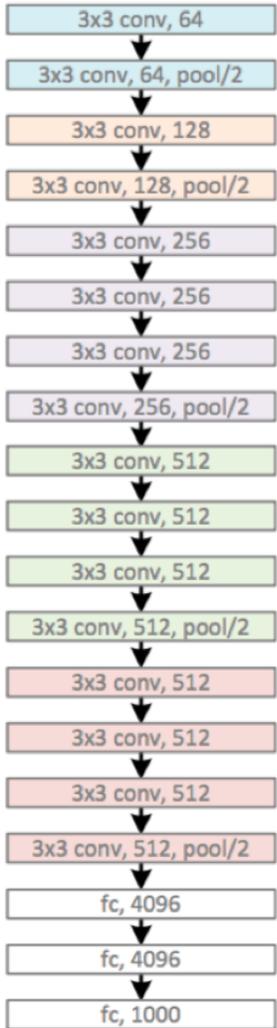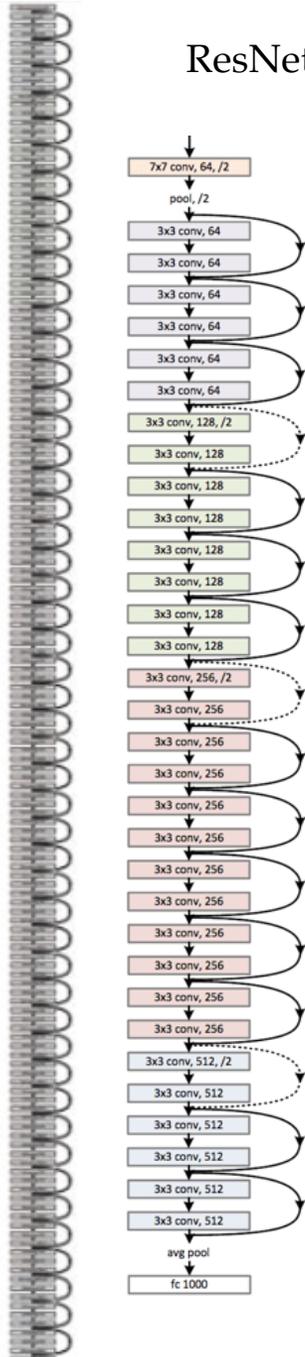| |
|---|
| 3x3 conv, 64 |
| 3x3 conv, 64, pool/2 |
| 3x3 conv, 128 |
| 3x3 conv, 128, pool/2 |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| 3x3 conv, 256 |
| 3x3 conv, 256, pool/2 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512, pool/2 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512 |
| 3x3 conv, 512, pool/2 |
| fc, 4096 |
| fc, 4096 |
| fc, 1000 |

"Very Deep Convolutional Networks for Large-Scale Image Recognition", Simonyan & Zisserman. ICLR 2015

[image credit Philip Isola and Kaiming He]

VGG '14

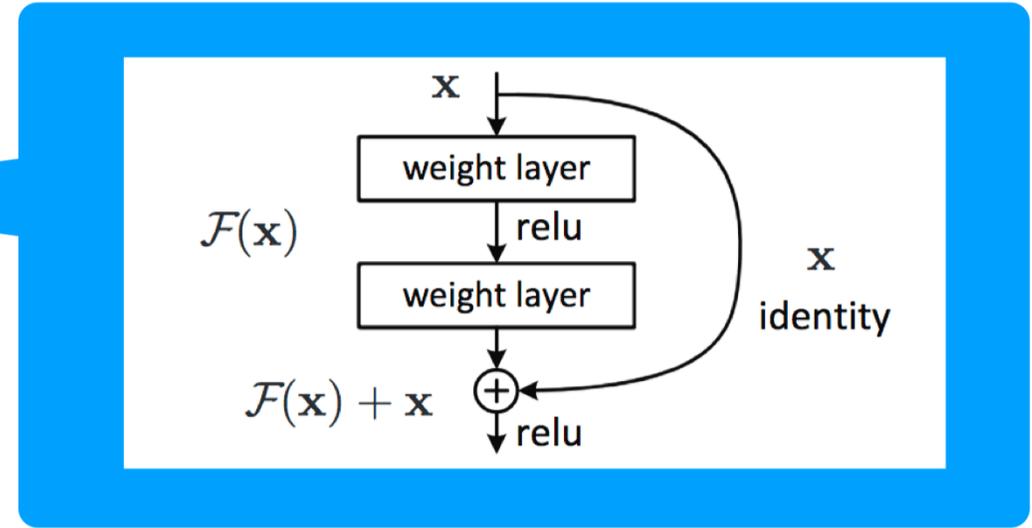ResNet '16

Main developments:

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

$\mathcal{F}(\mathbf{x})$

$\mathbf{x}$

weight layer

relu

weight layer

$\mathbf{x}$ identity

$\mathcal{F}(\mathbf{x}) + \mathbf{x}$

relu

[He et al: Deep Residual Learning for Image Recognition, CVPR 2016]

[image credit Philip Isola and Kaiming He]

# Summary

- Even though NNs are universal approximators, matching the architecture to problem structure — visual hierarchy, locality, translational invariance — improves generalization and efficiency.

- Convolution slides a small learned filter across the input, detecting local patterns with shared weights — sparse and efficient.

- Max pooling summarizes spatial information: "did a pattern occur?" rather than "where exactly?"

- Filter weights are learned end-to-end; convolutional layers extract features, fully connected layers classify.