

Lecture 9: Convolutional networks

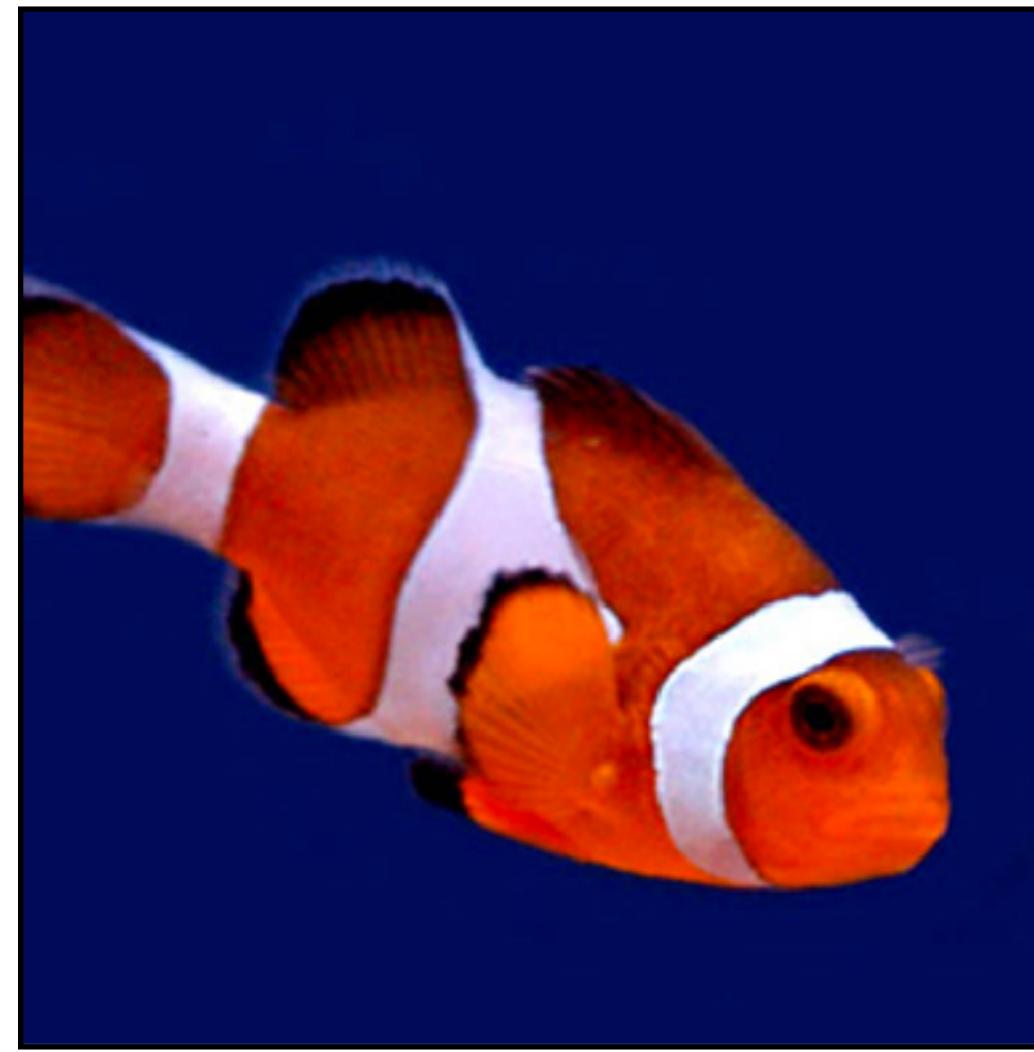
CS 5670: Introduction to Computer Vision



Announcements

- PS2 due on Monday
- PyTorch tutorial link (from Justin Johnson)
- PS1 grades back in ~2 weeks
- Thoughts on PS1?

Image classification



"Fish"

A rectangular box with a black border containing the word "Fish" in black, sans-serif capital letters.

image **x**

label **y**

Image classification

What should these be?

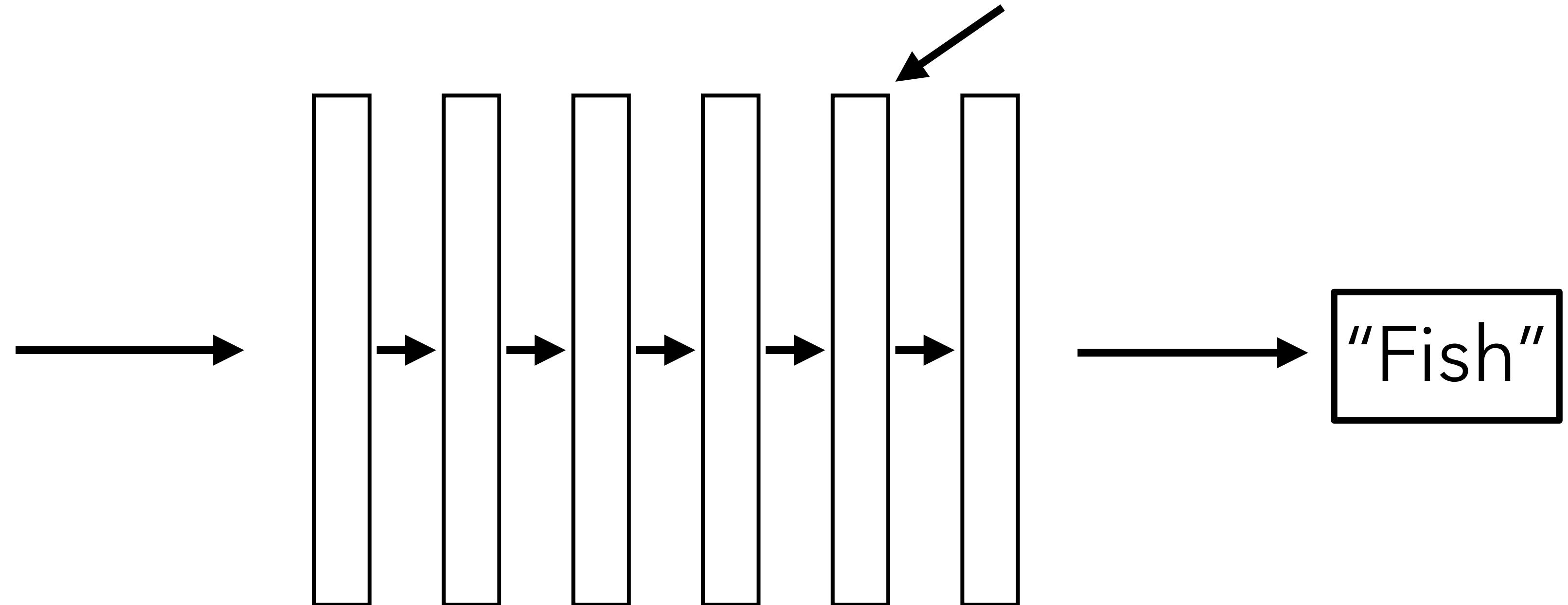
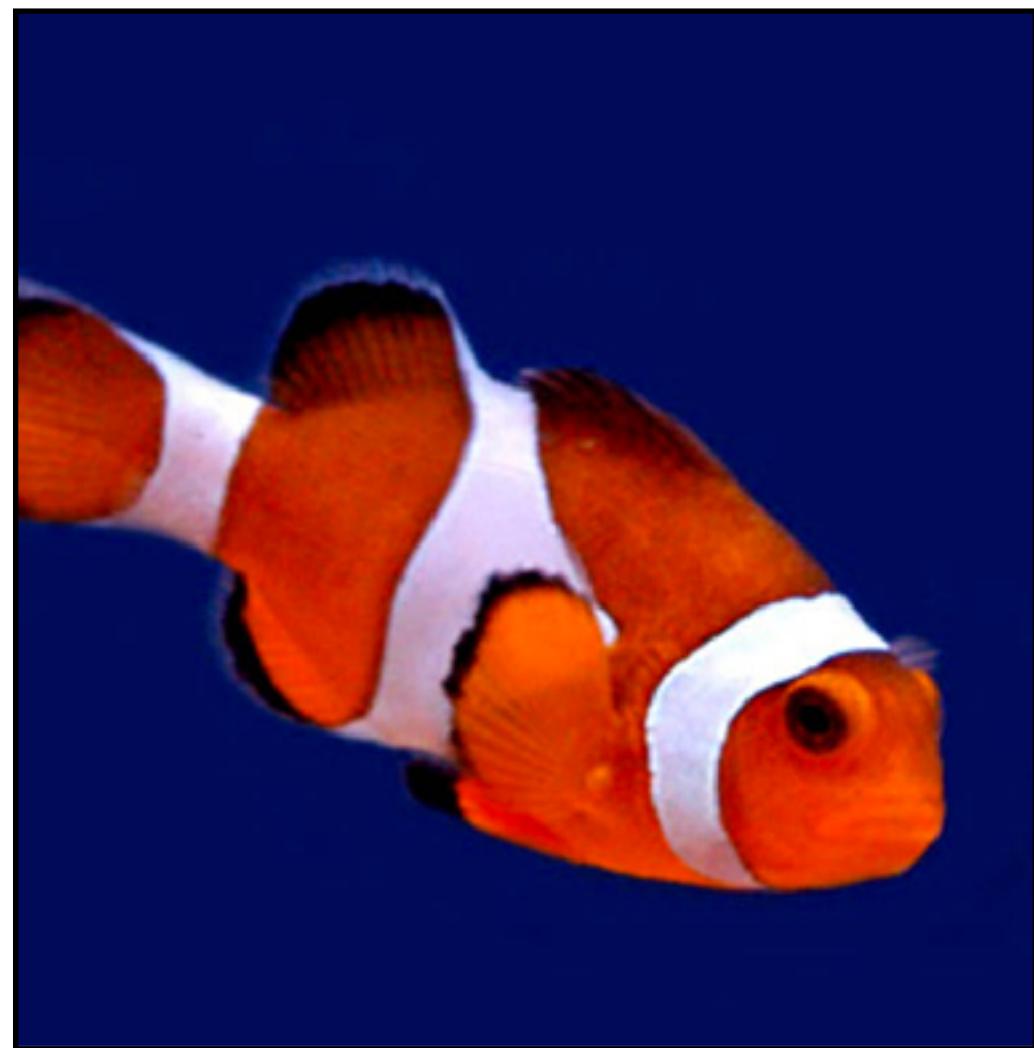
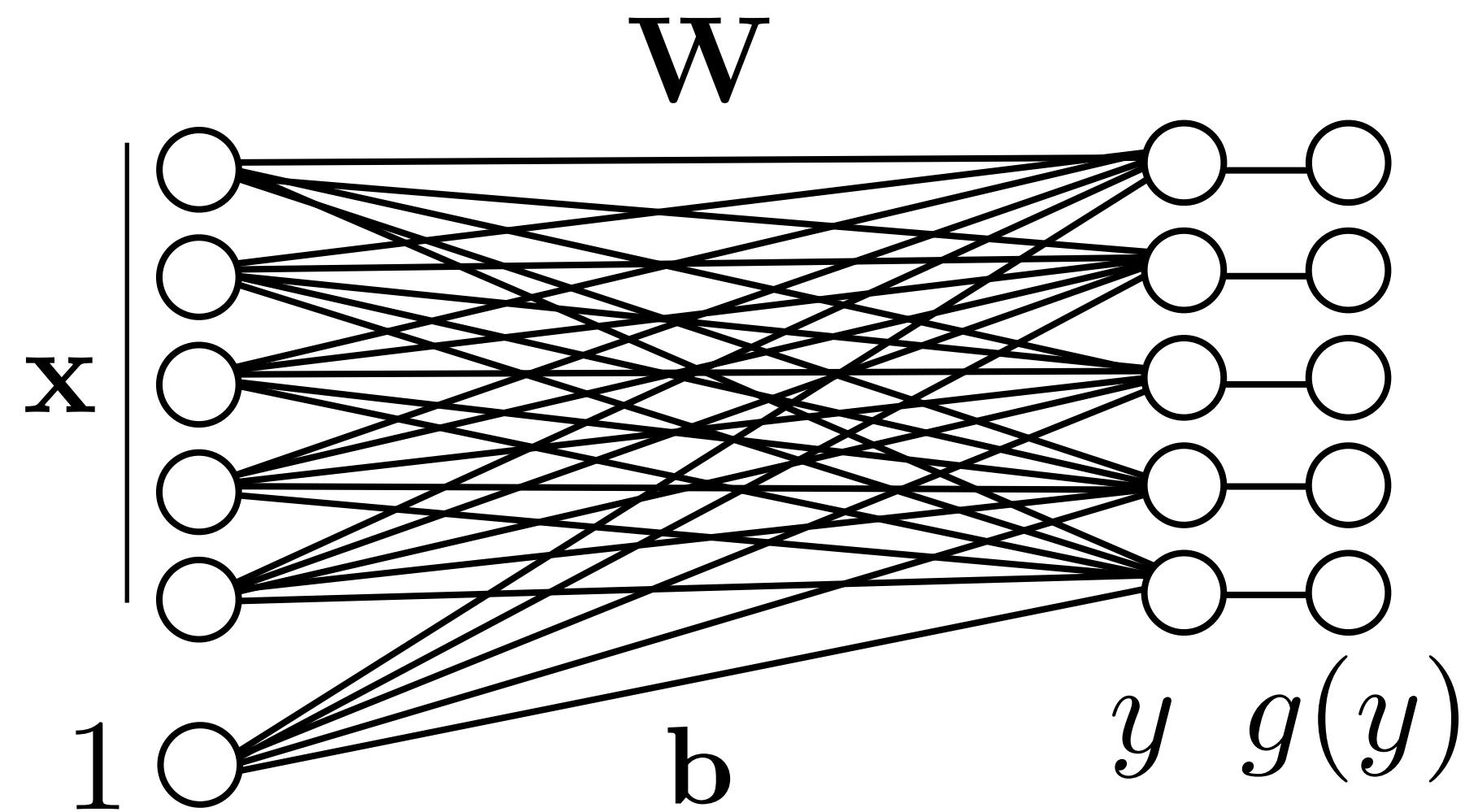


image **x**

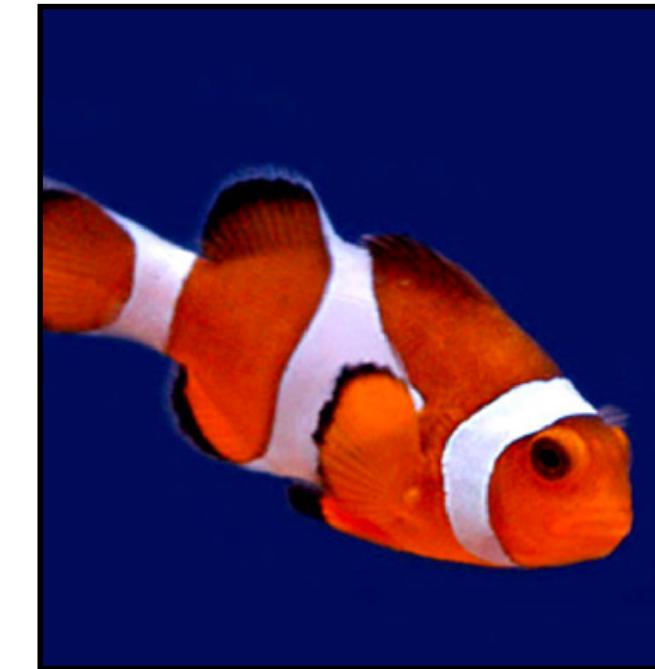
label **y**

Idea #1: Fully-connected network

Fully-connected (a.k.a. linear) layer



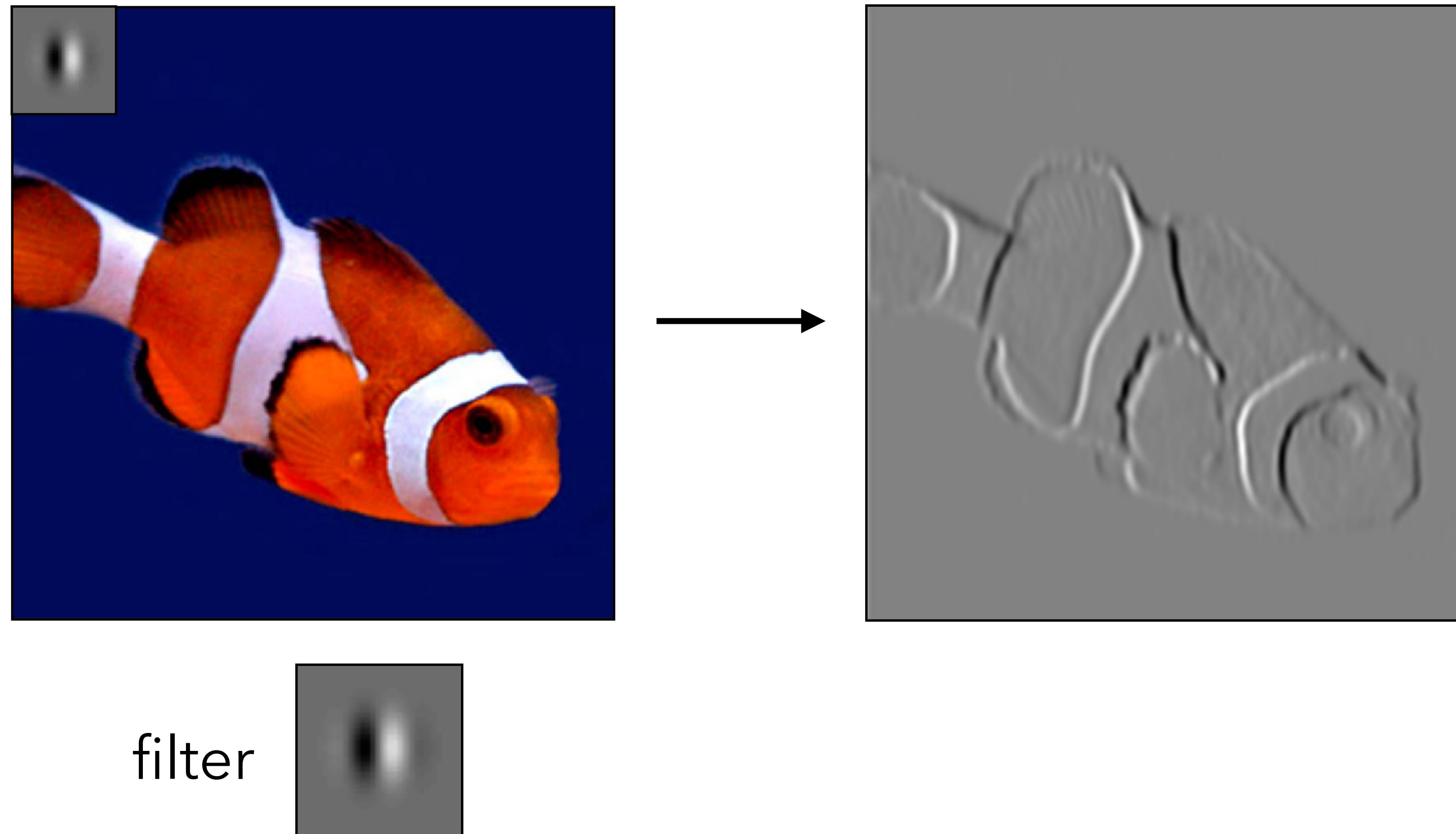
$\mathbf{x} =$



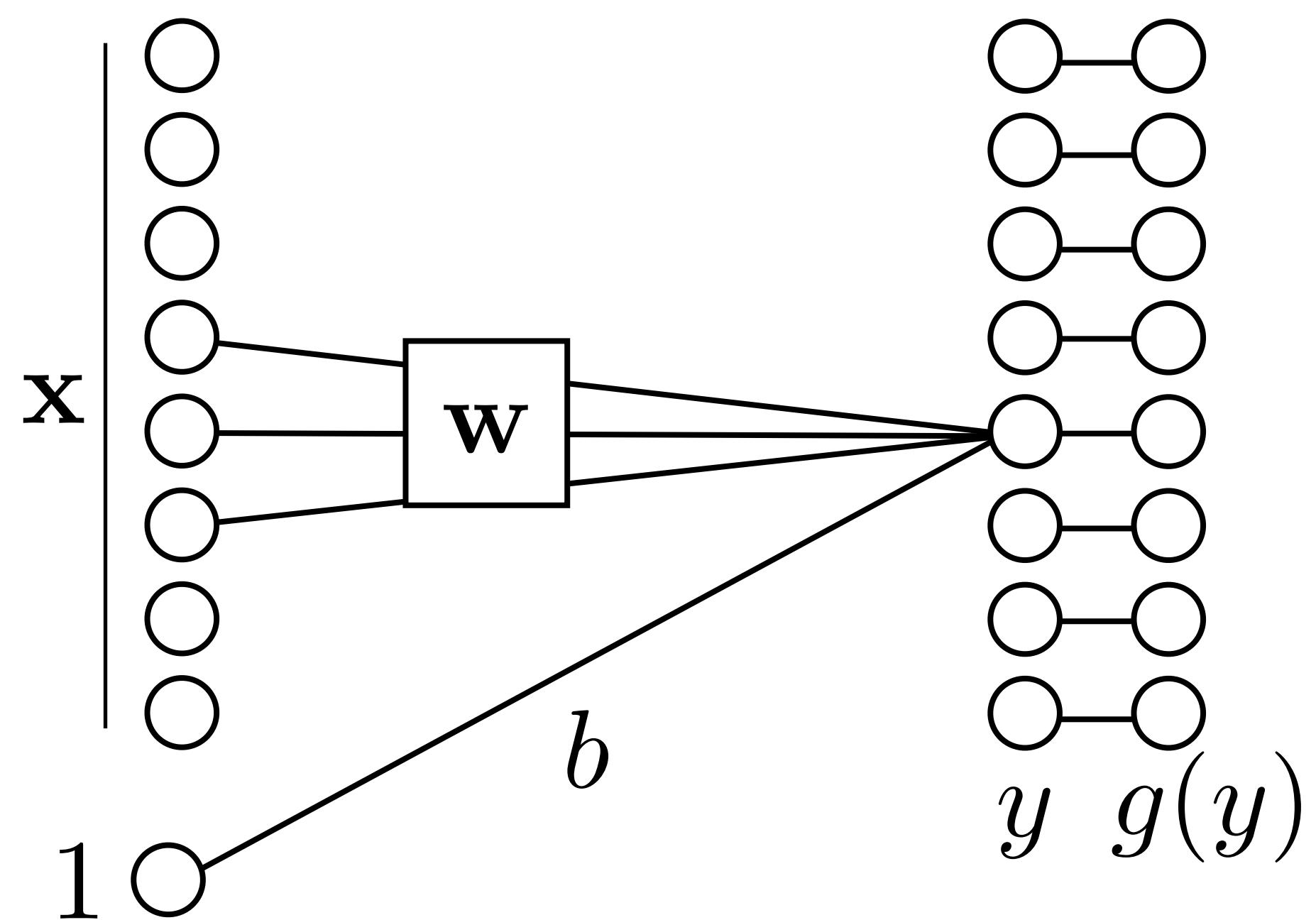
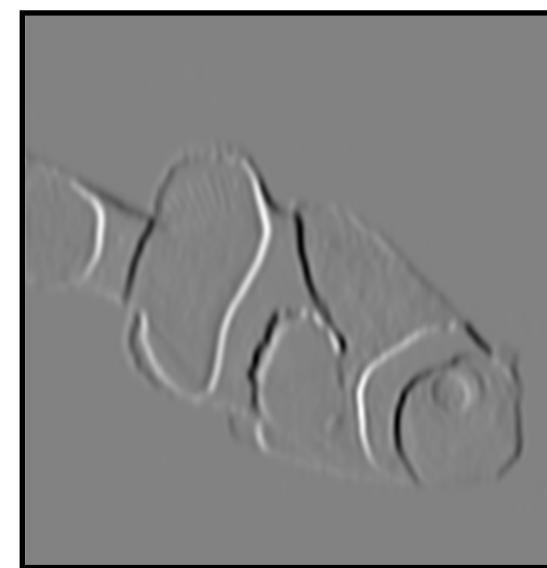
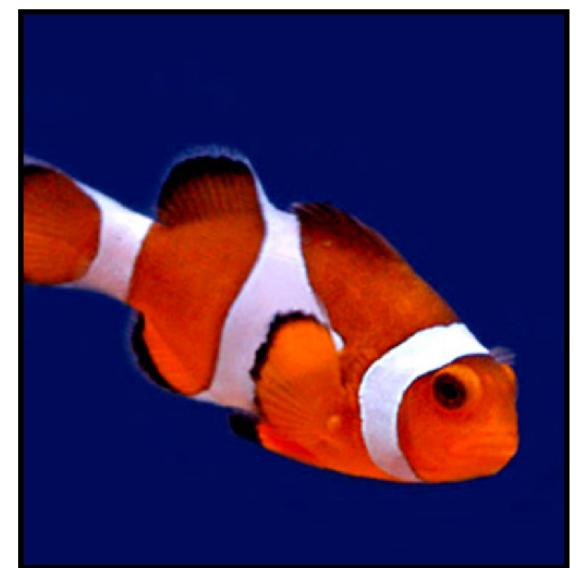
But \mathbf{x} is really big!

Say, $256 \times 256 \times 3 = 197k$

Can we use convolution in a neural network?



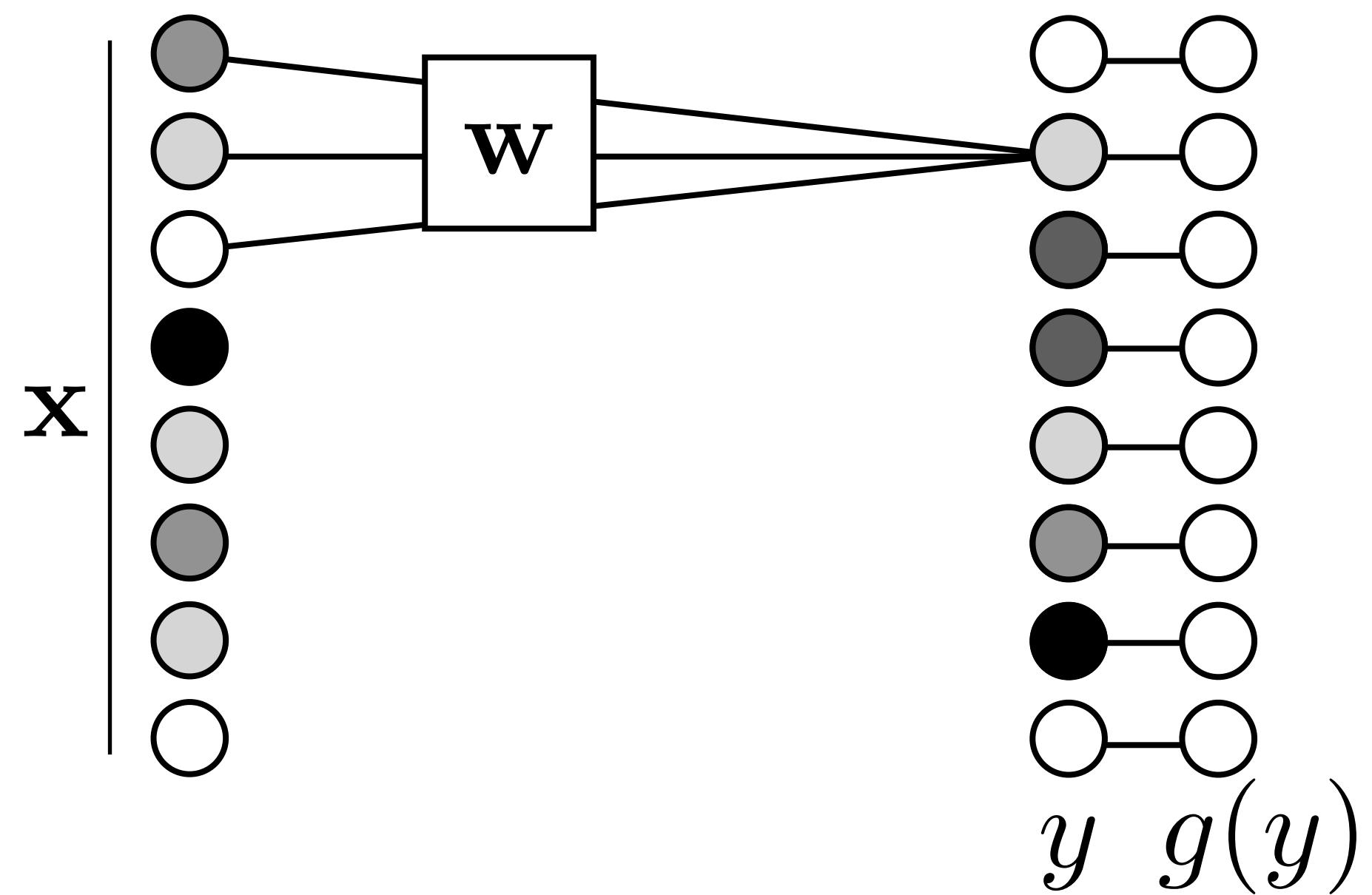
Sparsely connected network



Each unit is connected to a subset of the units in the previous layer.

Convolutional neural network

Conv layer

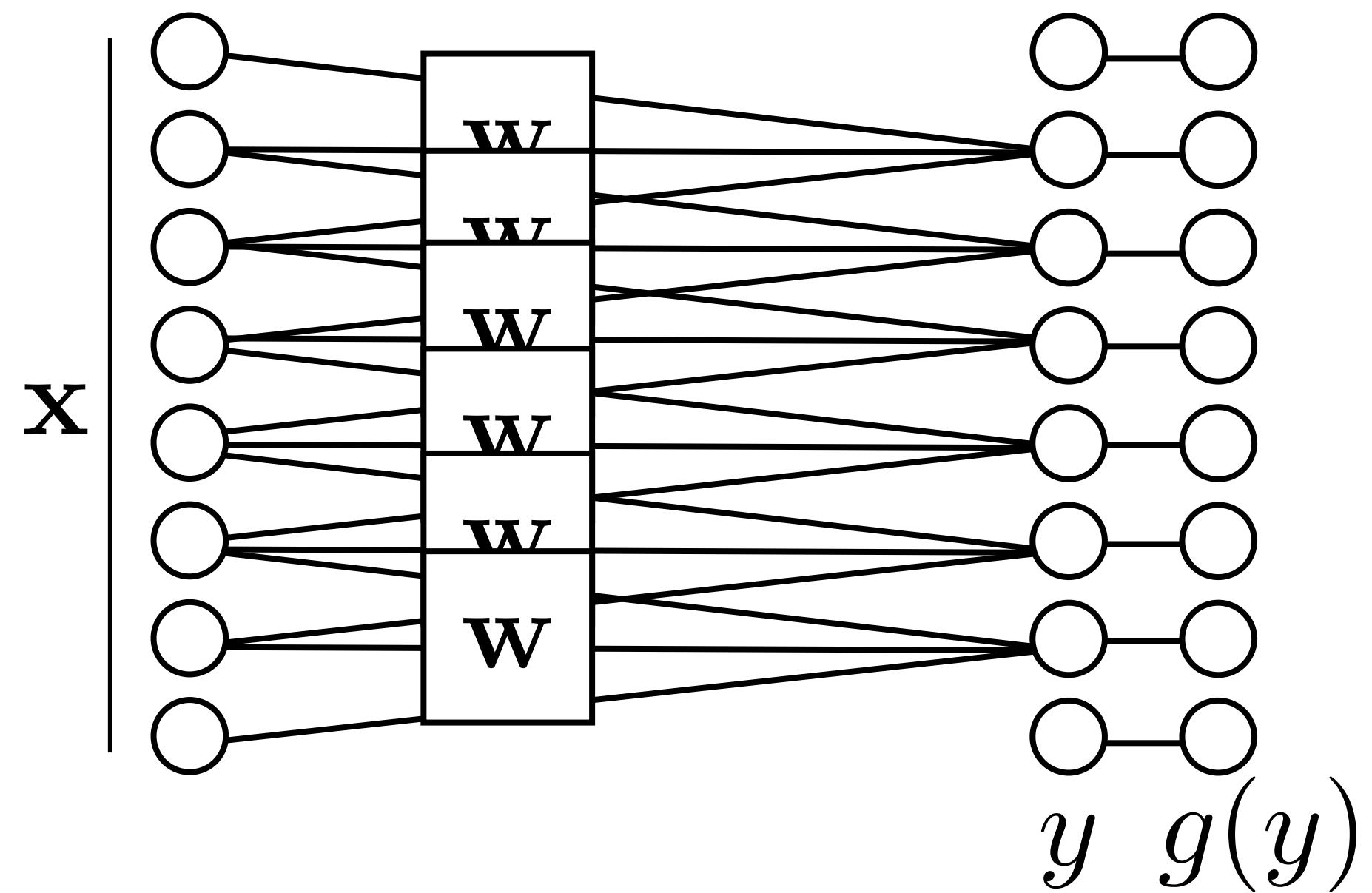


Each output unit is computed from an image patch.

a.k.a. CNN or ConvNet

Weight sharing

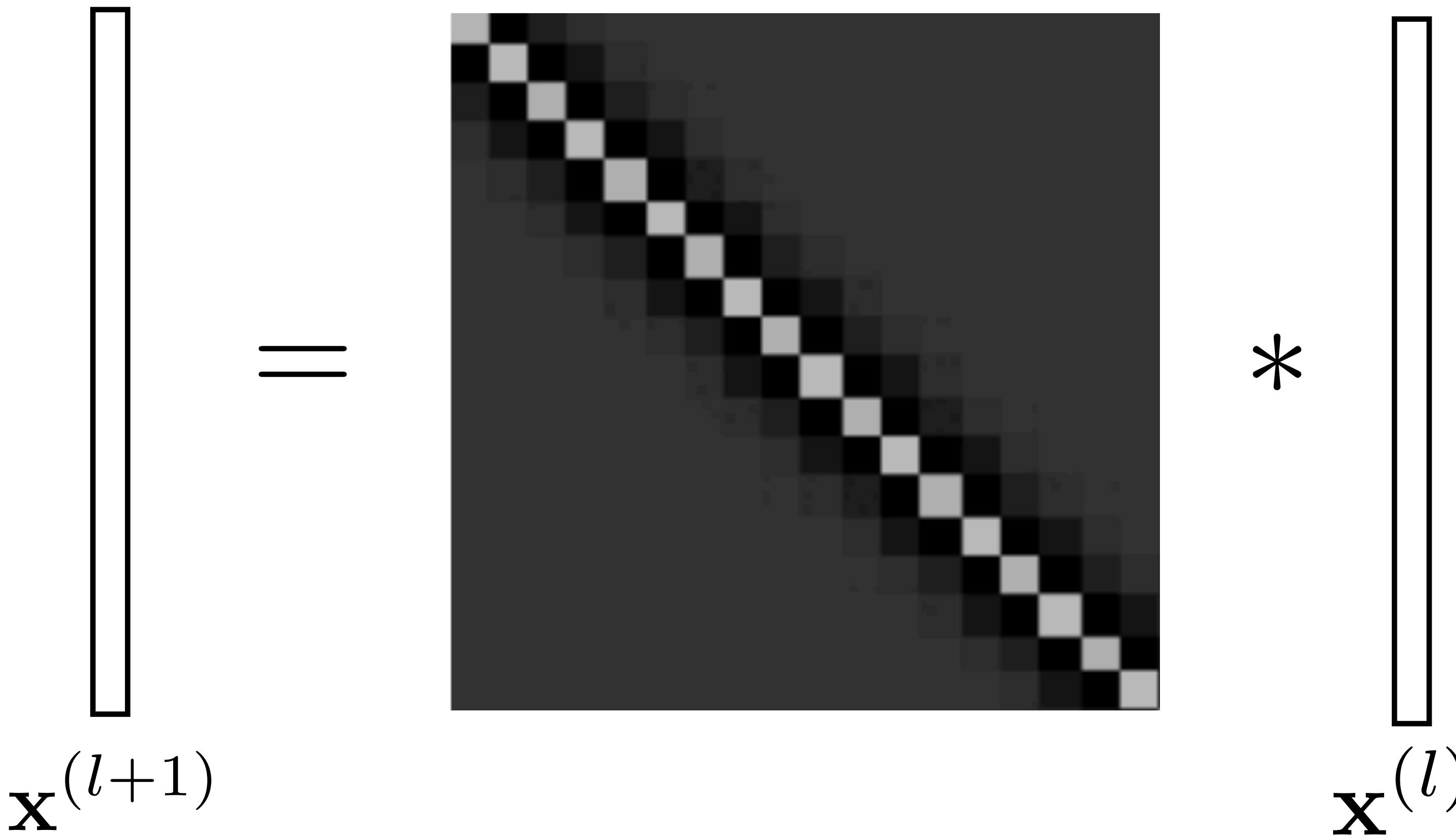
Conv layer



We “share” weights for each patch.

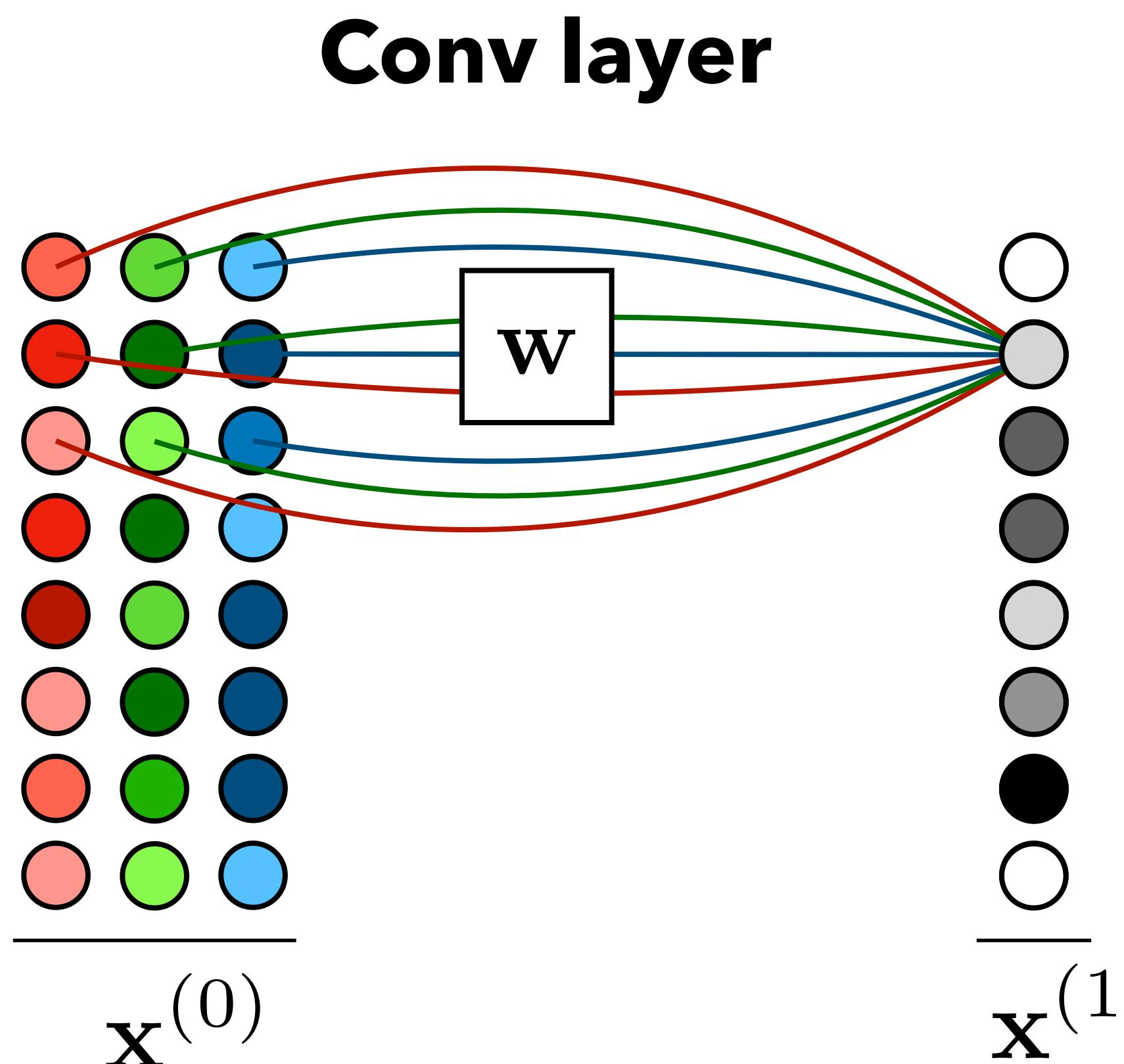
If a feature is useful in one position, it should be useful in others, too.

Convolution is a linear function



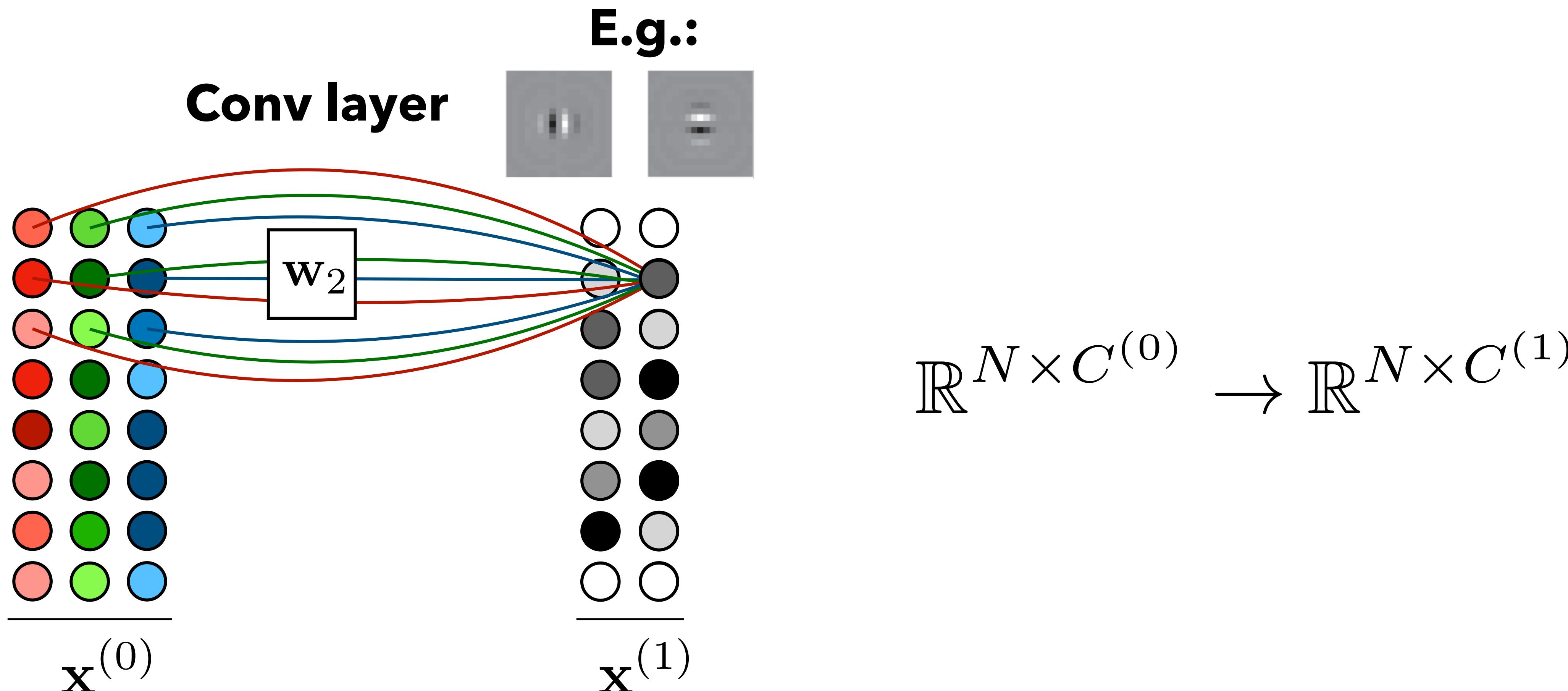
- Constrained linear layer e.g., image
- Fewer parameters: easier to learn, less overfitting
- Usually use zero padding

Multiple input channels



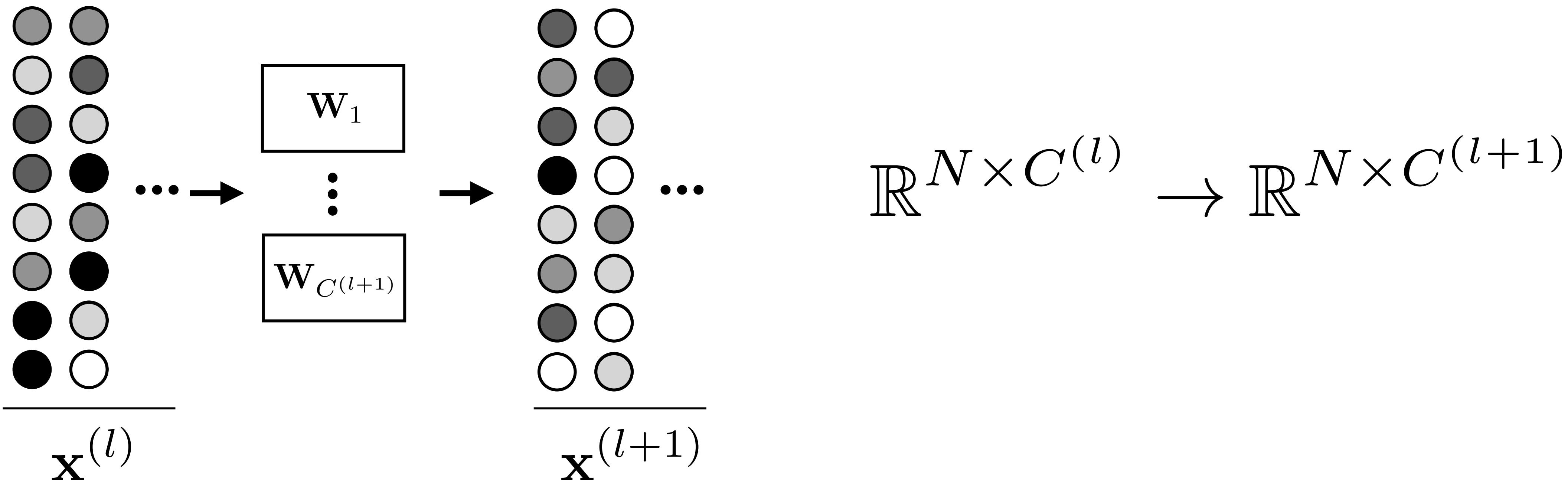
$$\mathbb{R}^{N \times C} \rightarrow \mathbb{R}^{N \times 1}$$

Multiple output channels

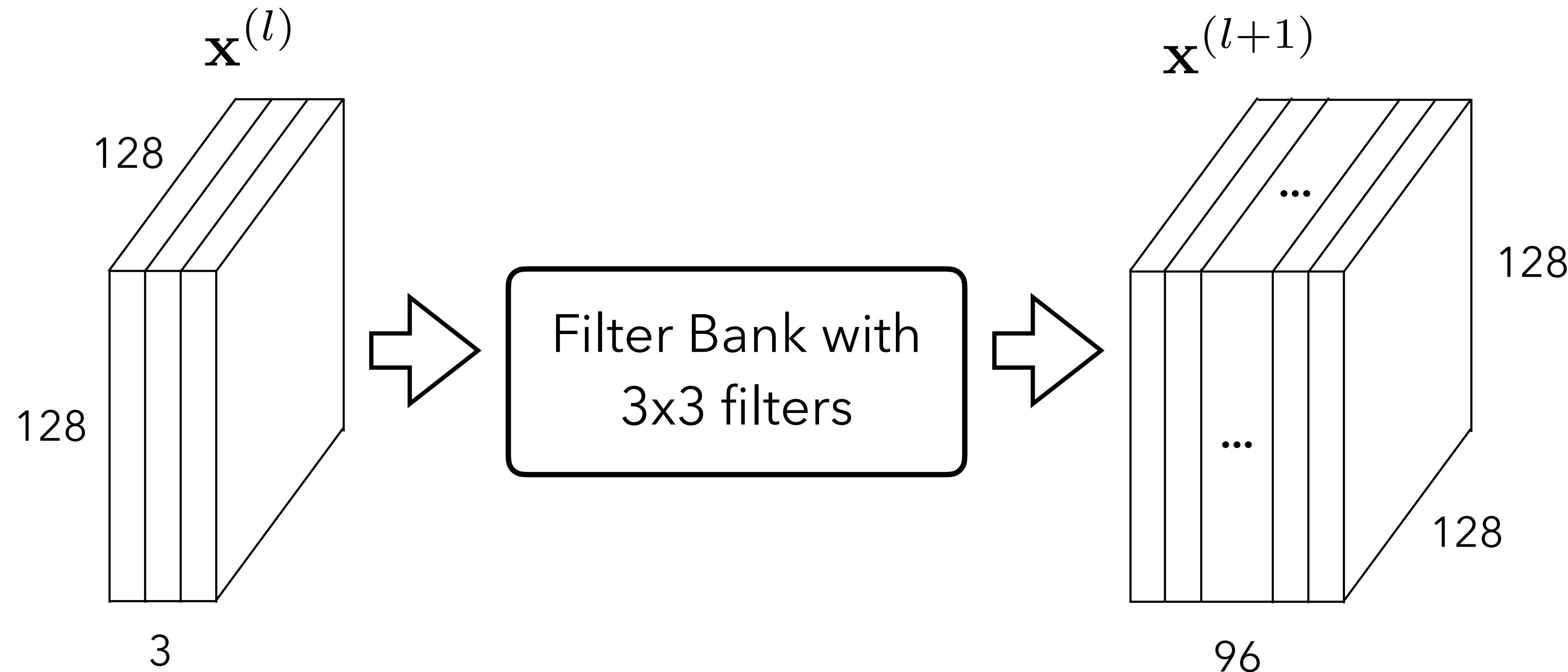


Multiple input and output channels

Conv layer



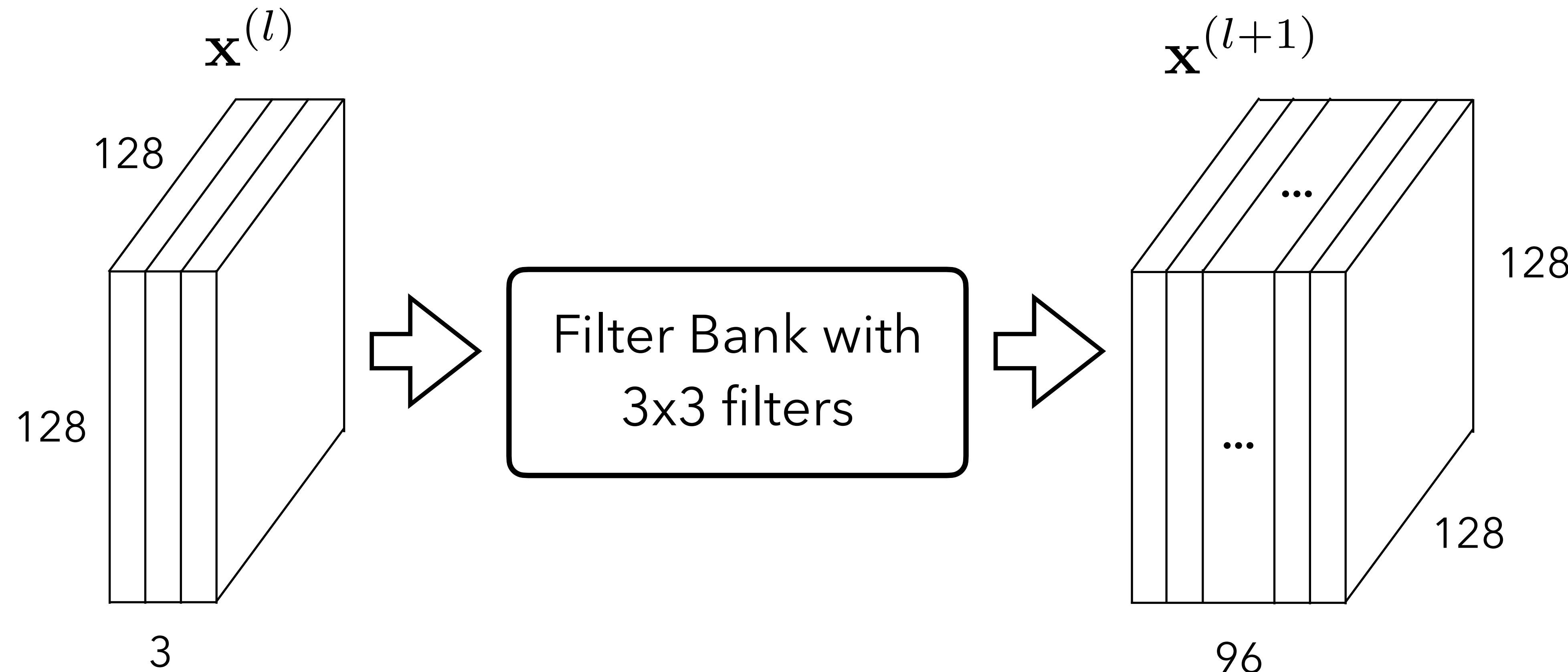
Multiple channels: Example



How many parameters does each *filter* have?

- (a) 9
- (b) 27
- (c) 96
- (d) 2592

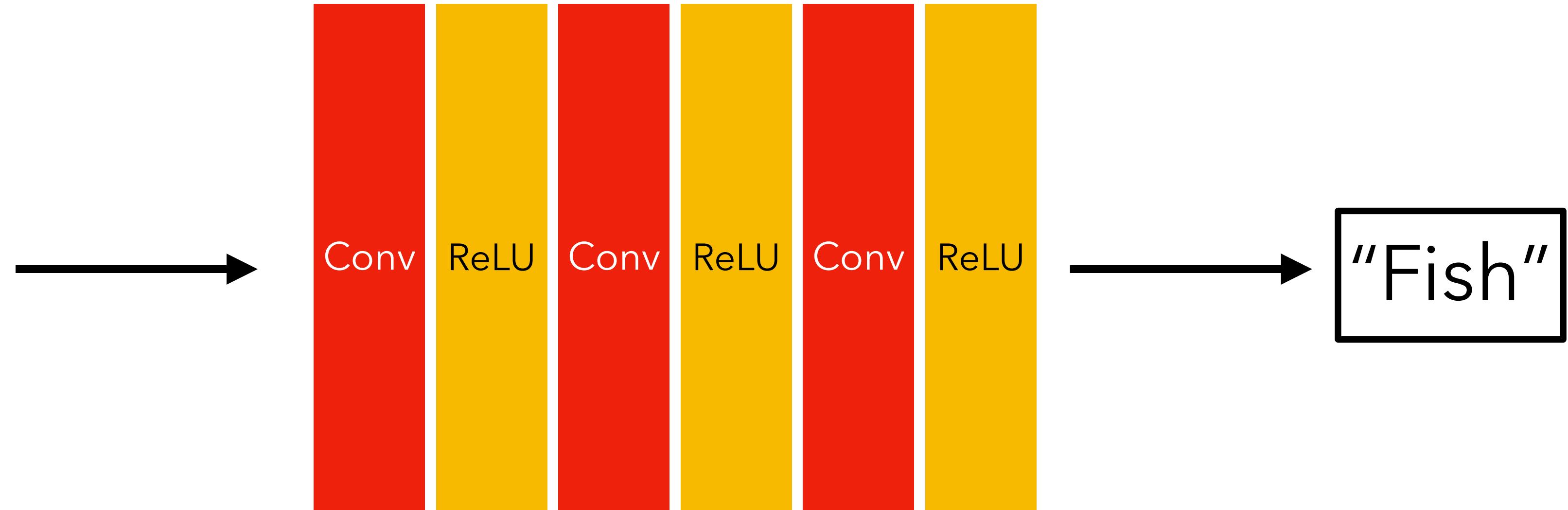
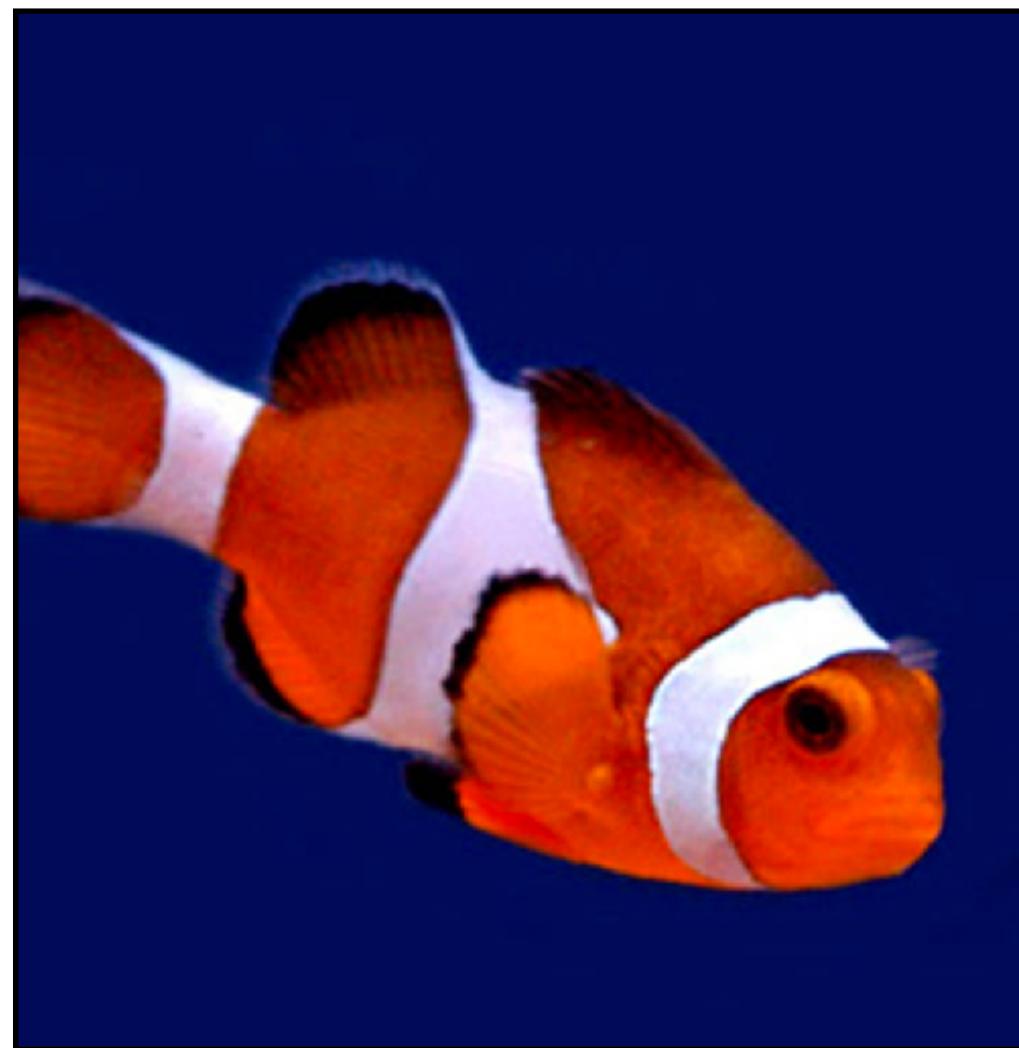
Multiple channels: Example



How many parameters *total* does this layer have?

- (a) 9
- (b) 27
- (c) 96
- (d) 2592

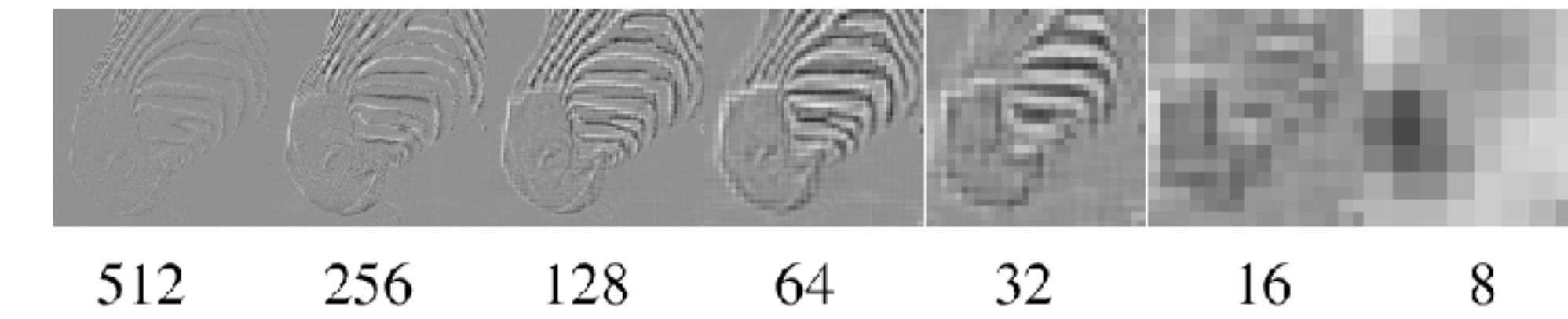
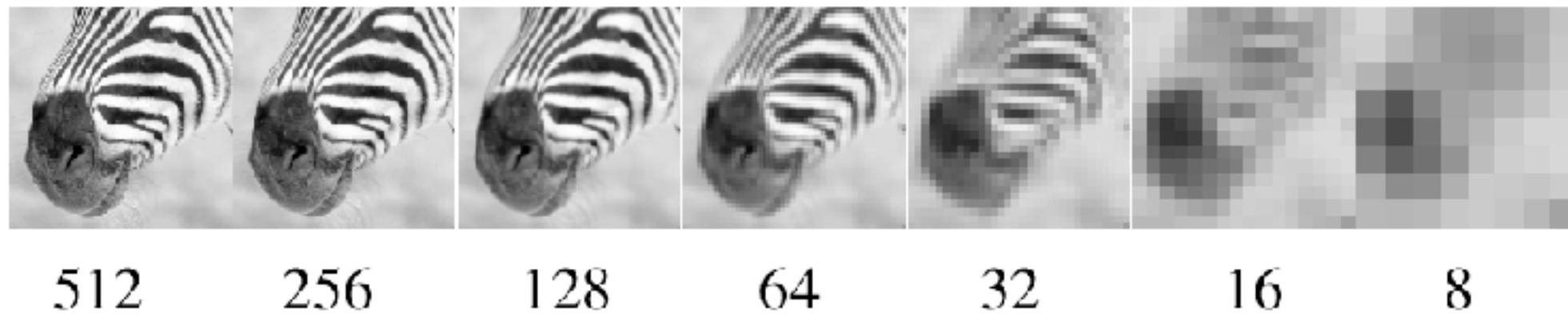
Image classification



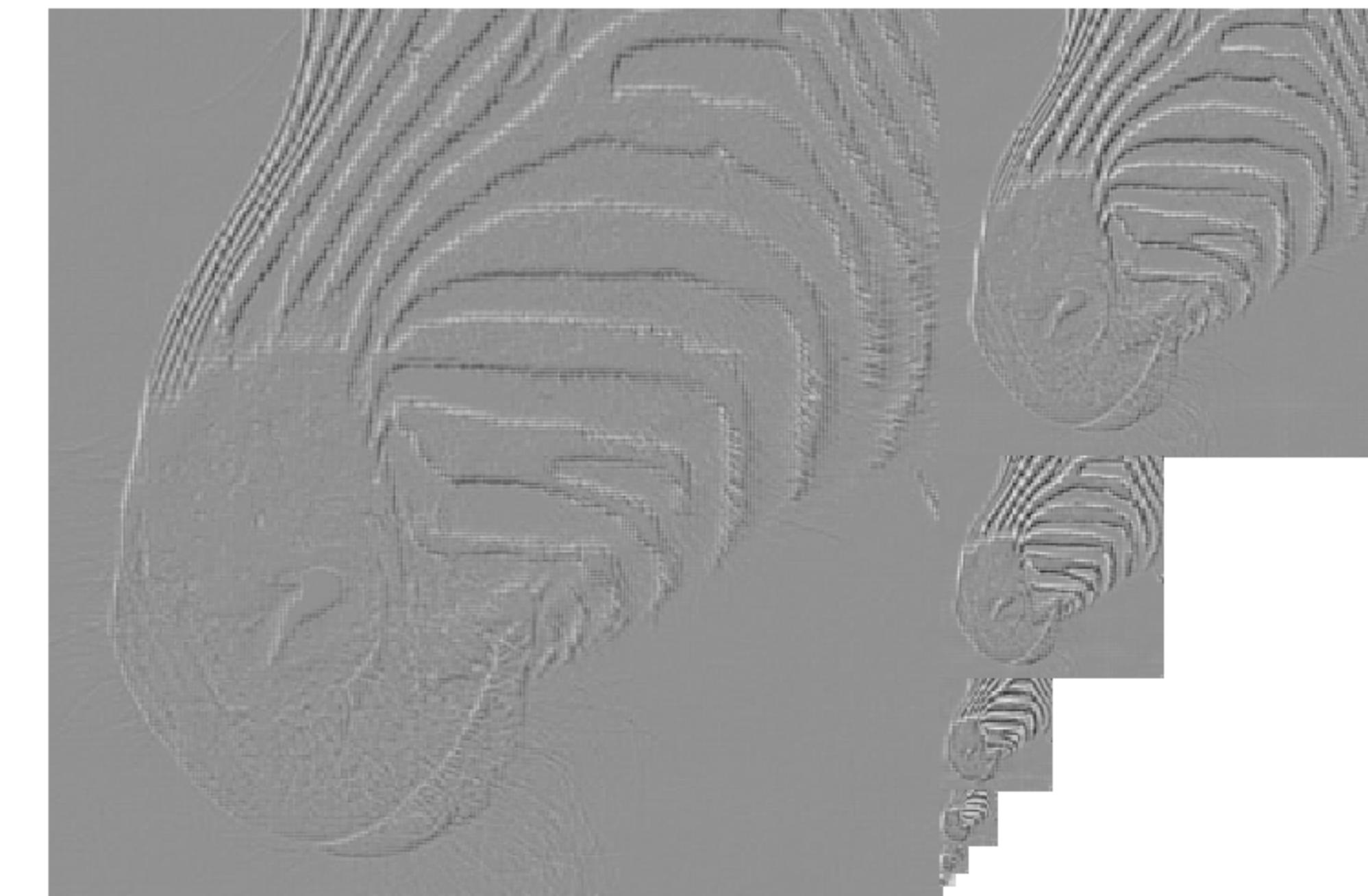
Problems with this idea:

1. No “global” processing.
2. How do you get the final label?

Recall: pyramid representations



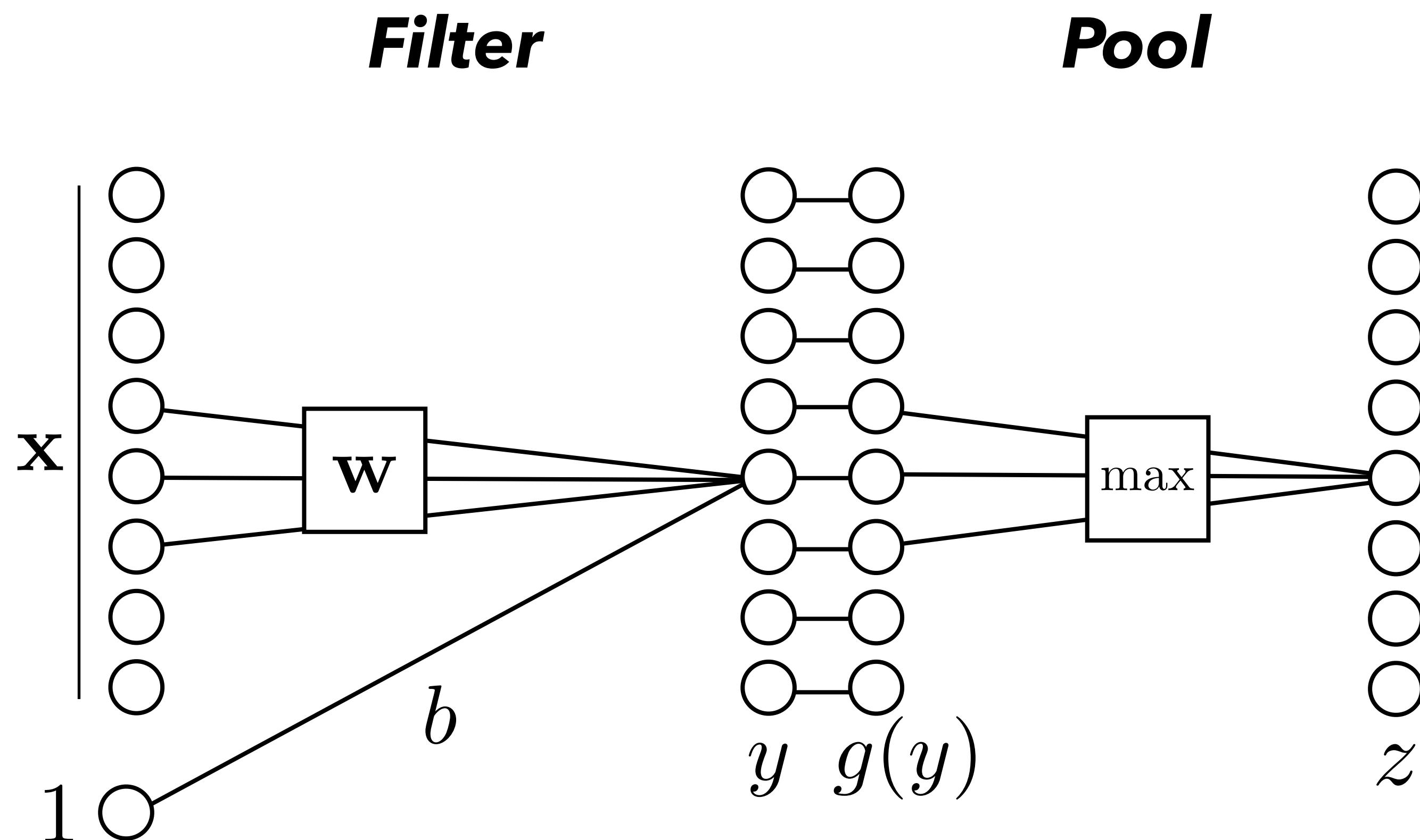
Gaussian Pyramid



Laplacian Pyramid

Can we use a similar idea in CNNs?

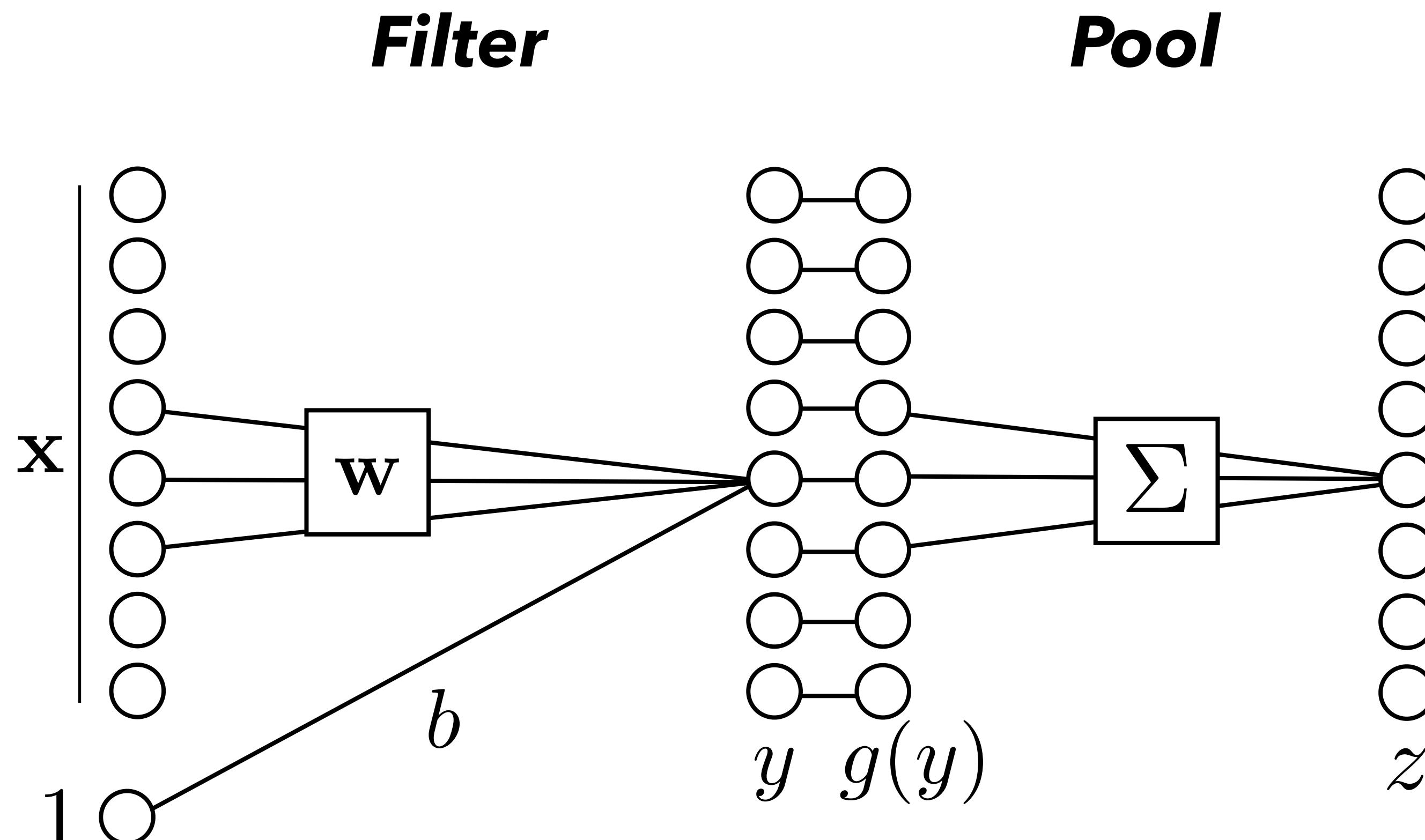
Pooling



Max pooling

$$z_k = \max_{j \in \mathcal{N}(j)} g(y_j)$$

Pooling



Max pooling

$$z_k = \max_{j \in \mathcal{N}(j)} g(y_j)$$

Mean pooling

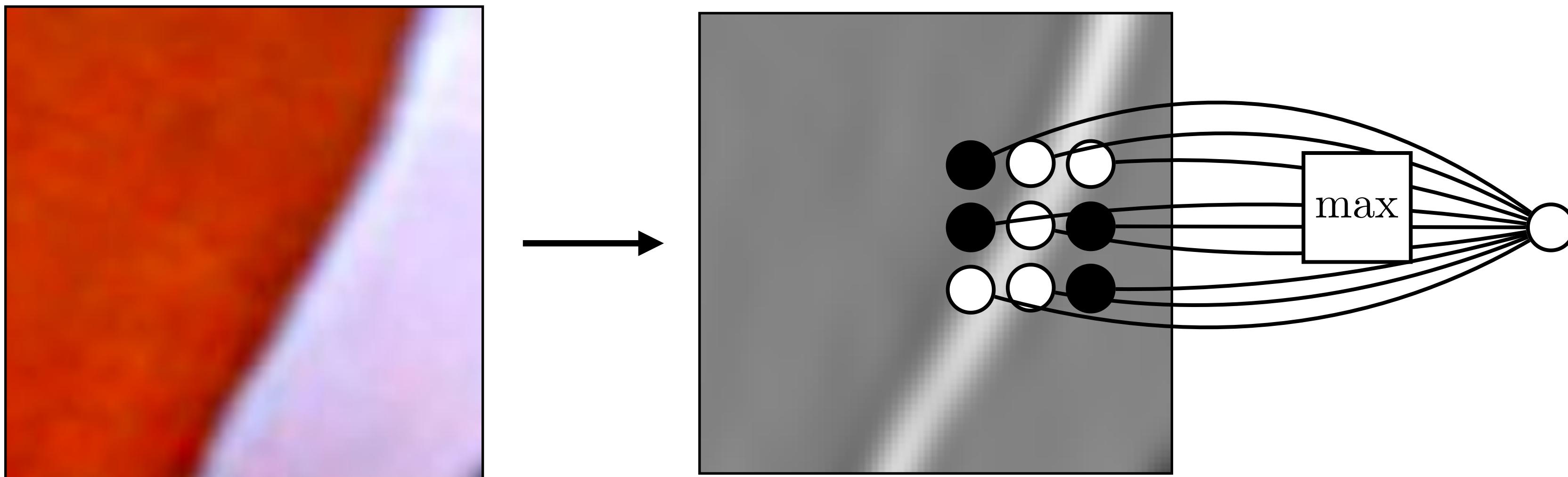
$$z_k = \frac{1}{|\mathcal{N}|} \sum_{j \in \mathcal{N}(j)} g(y_j)$$

Blurring [Zhang 2019]

$$z = \text{conv}(y, \frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix})$$

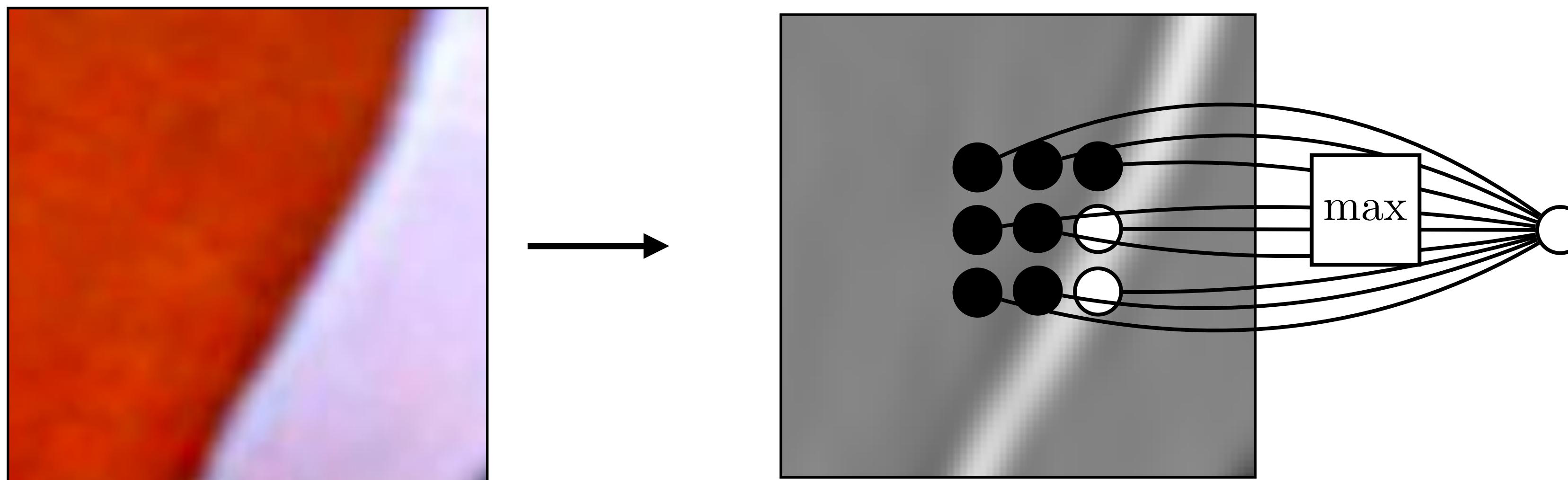
Pooling – Why?

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



Pooling – Why?

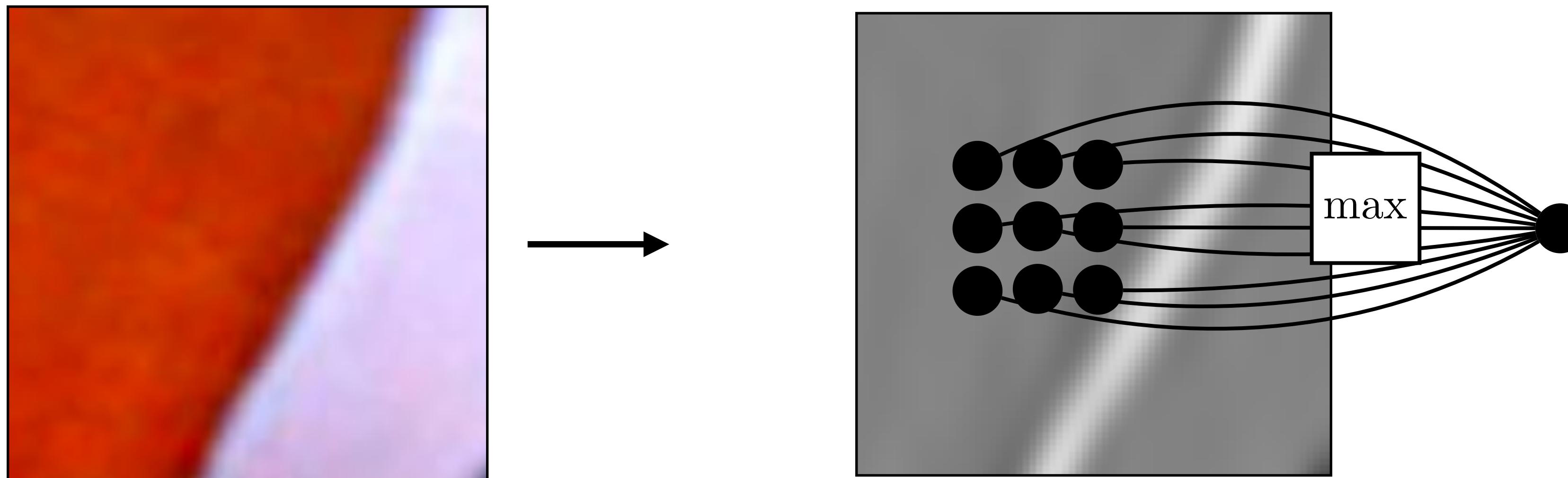
Pooling across spatial locations achieves stability w.r.t. small translations:



same large response
regardless of exact
position of edge

Pooling – Why?

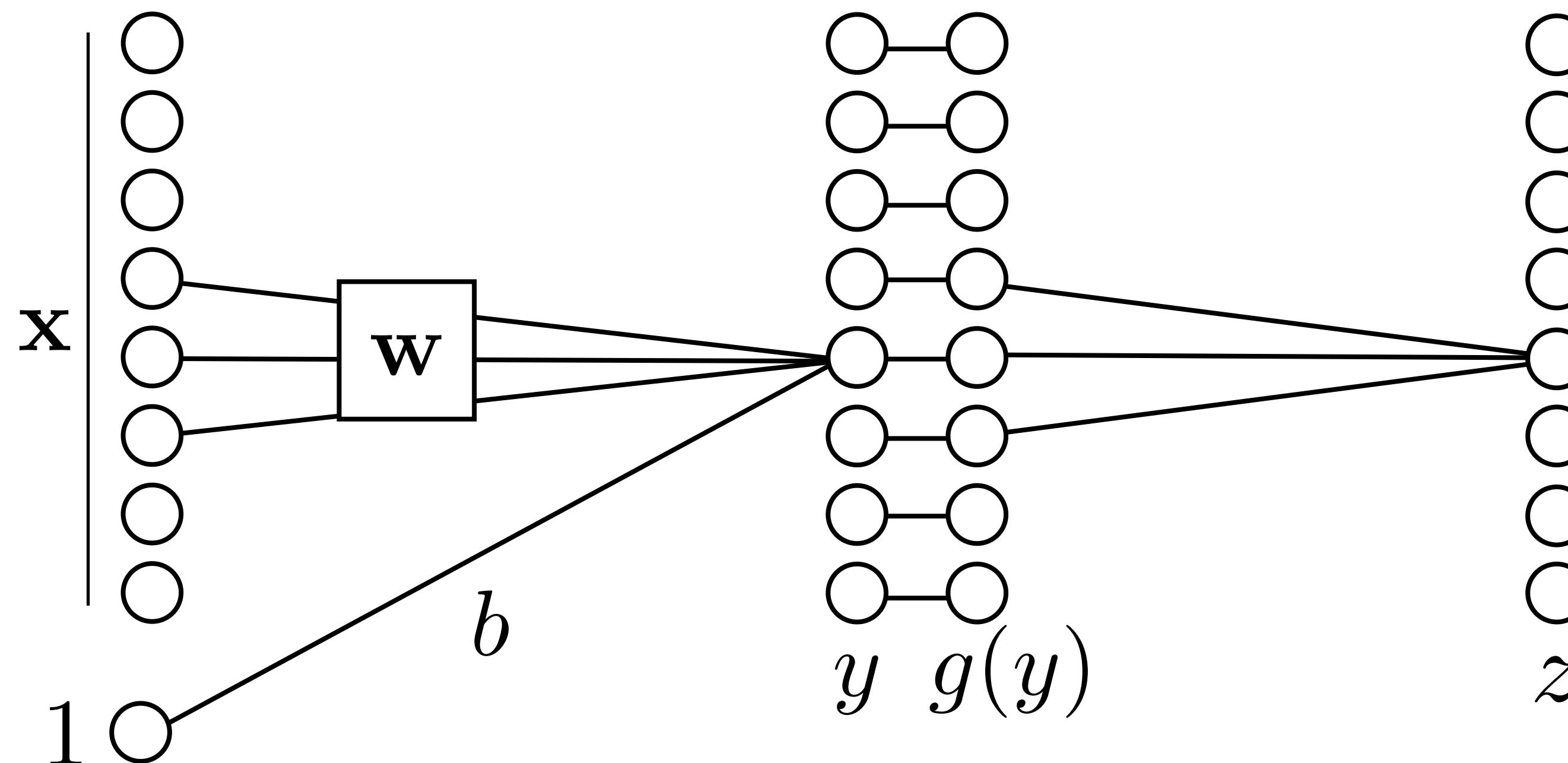
Pooling across spatial locations achieves stability w.r.t. small translations:



however, if the
image is translated a
lot...

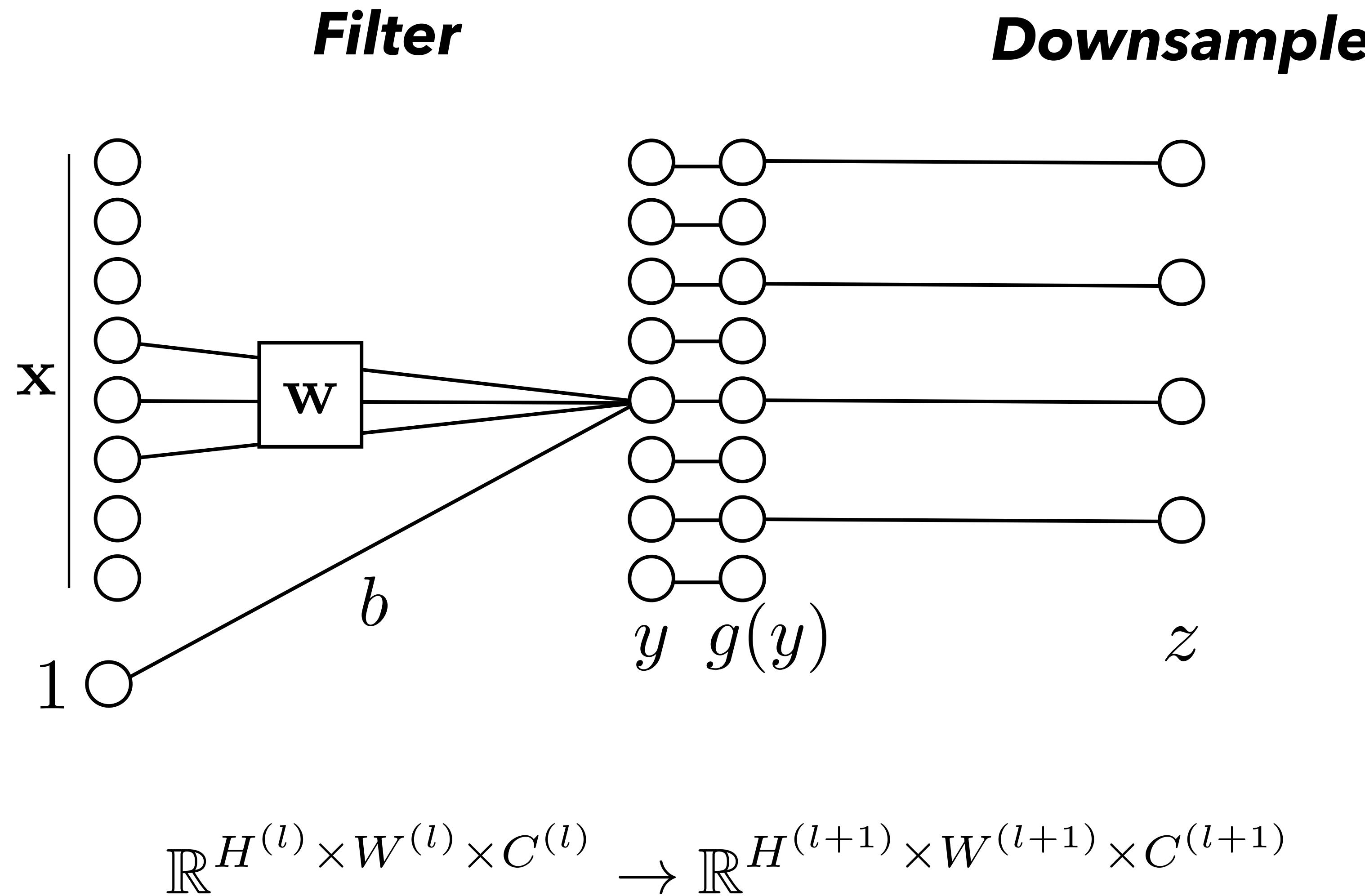
Downsampling

Filter

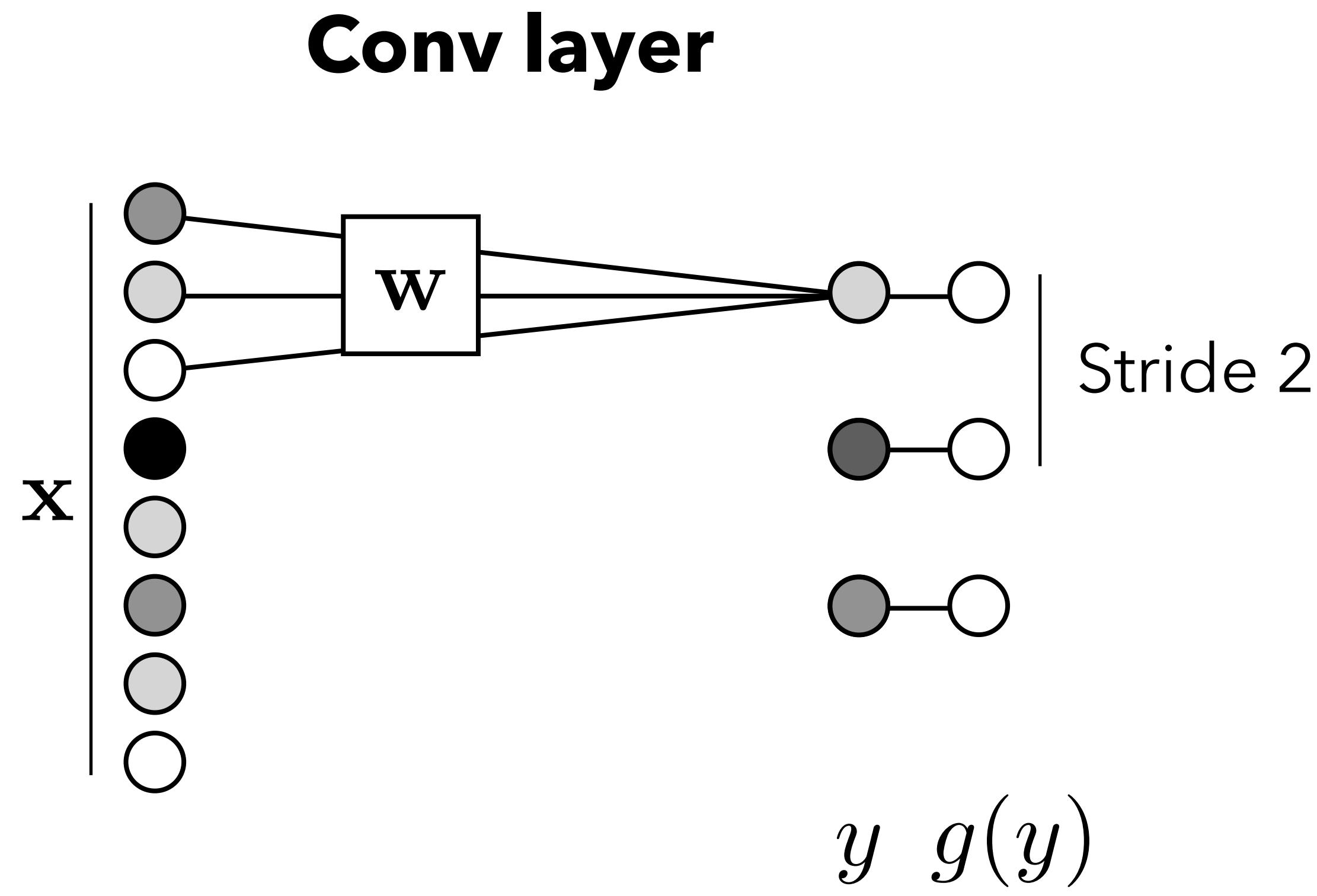


Pool and downsample

Downsampling



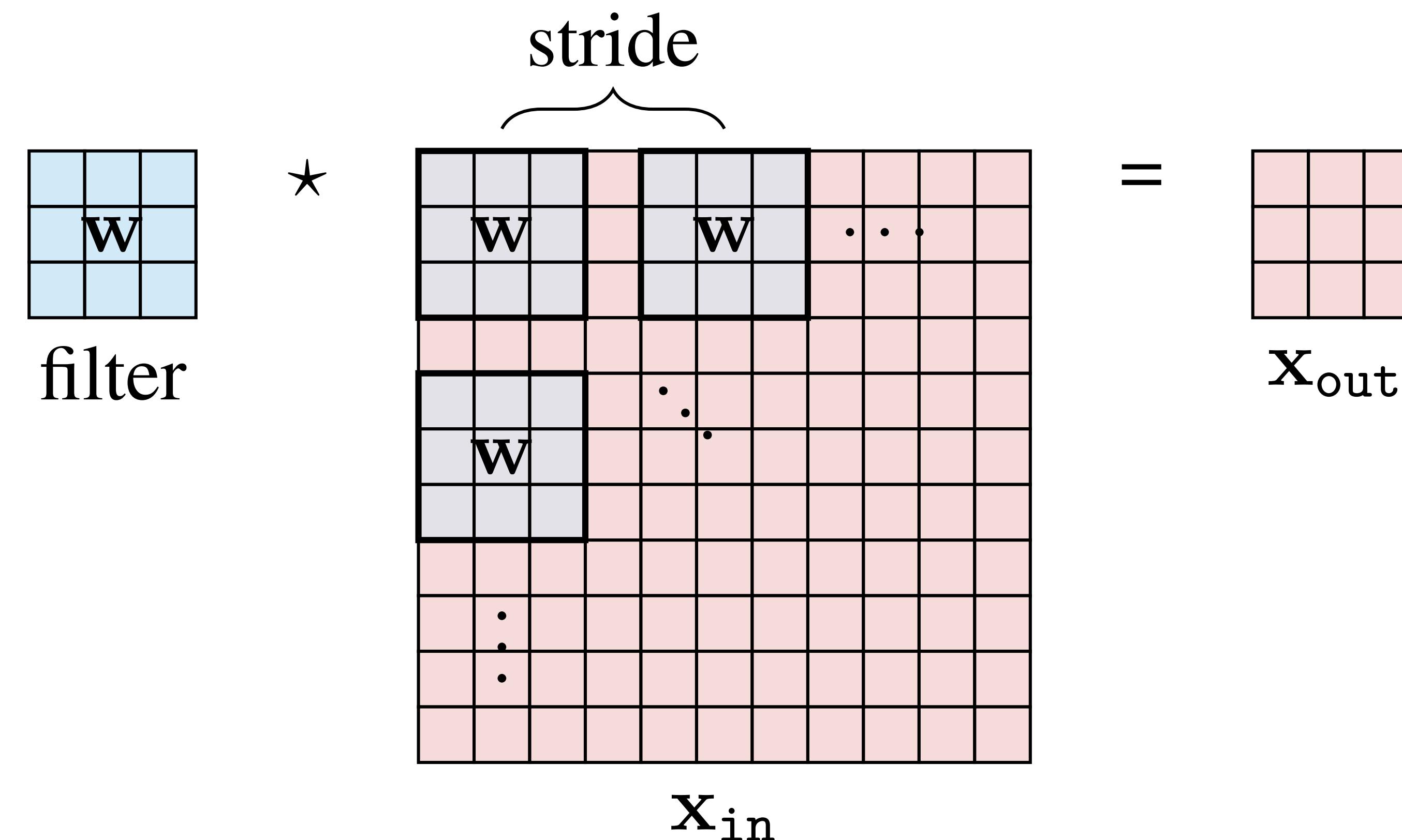
Strided operations



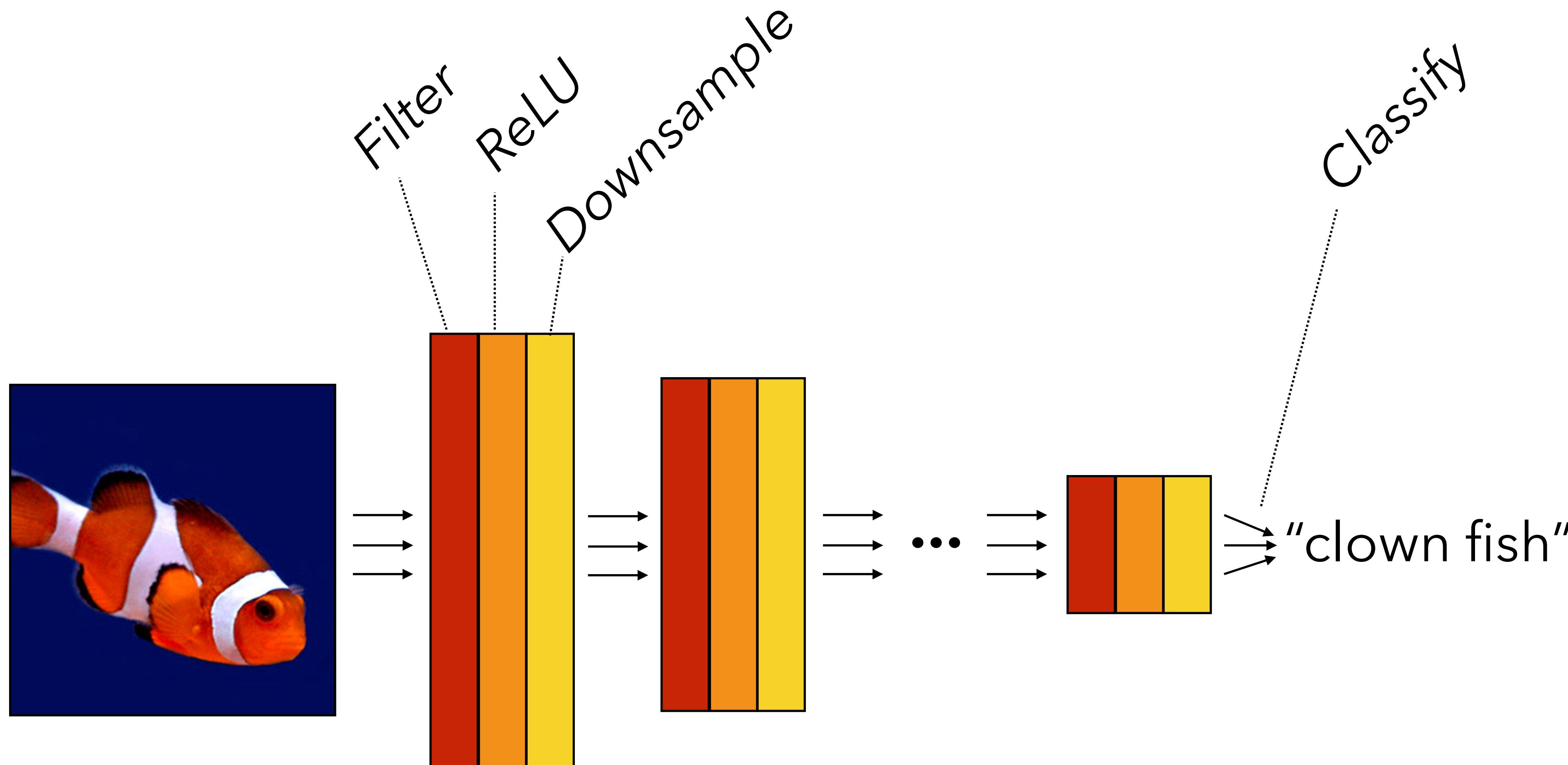
Strided operations combine a given operation (convolution or pooling) and downsampling into a single operation.

Strided convolution is an alternative to pooling layers: just do a strided convolution!

Strided operations (2D)

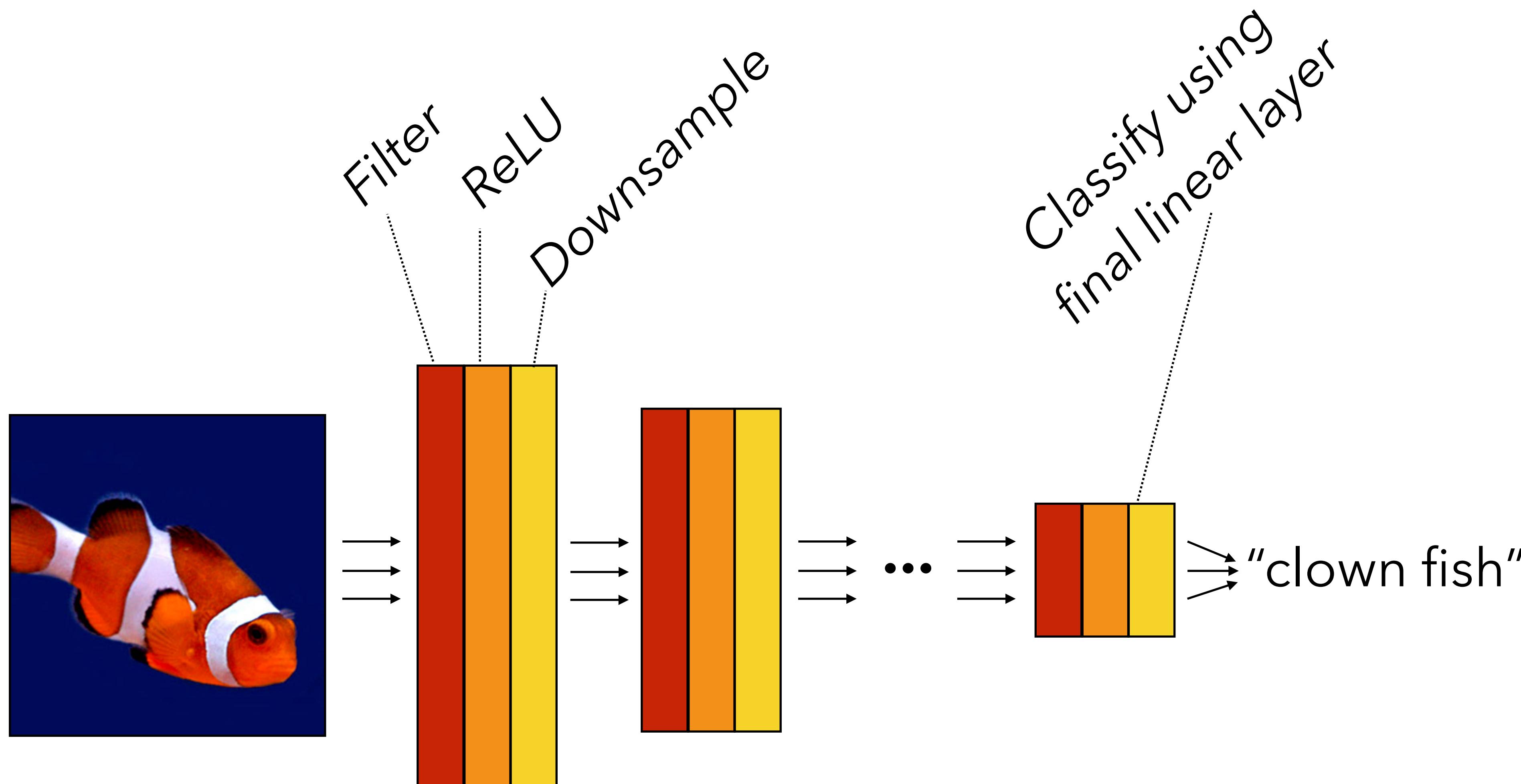


Computation in a neural net



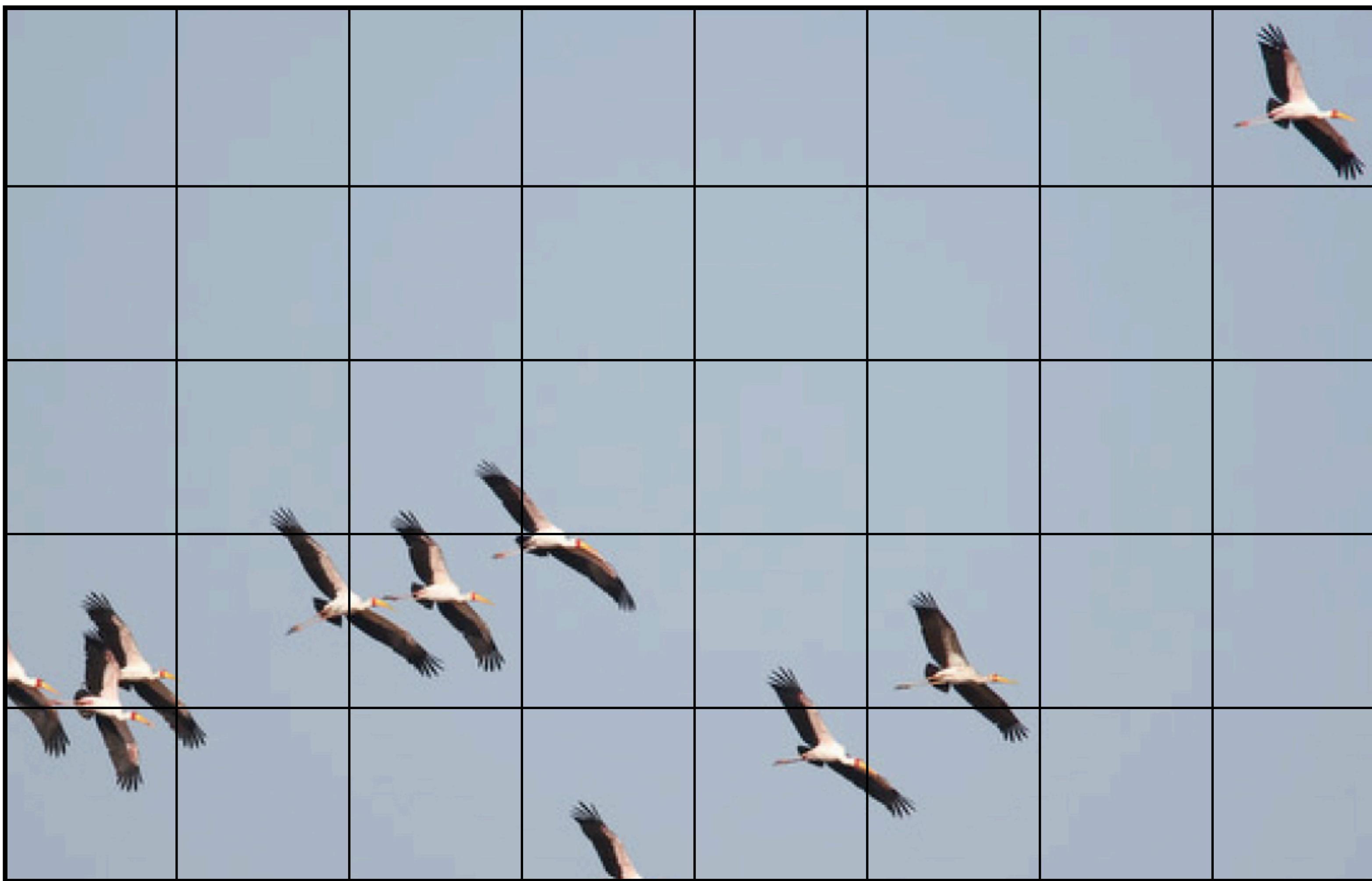
$$f(\mathbf{x}) = f_L(\dots f_2(f_1(\mathbf{x})))$$

Computation in a neural net

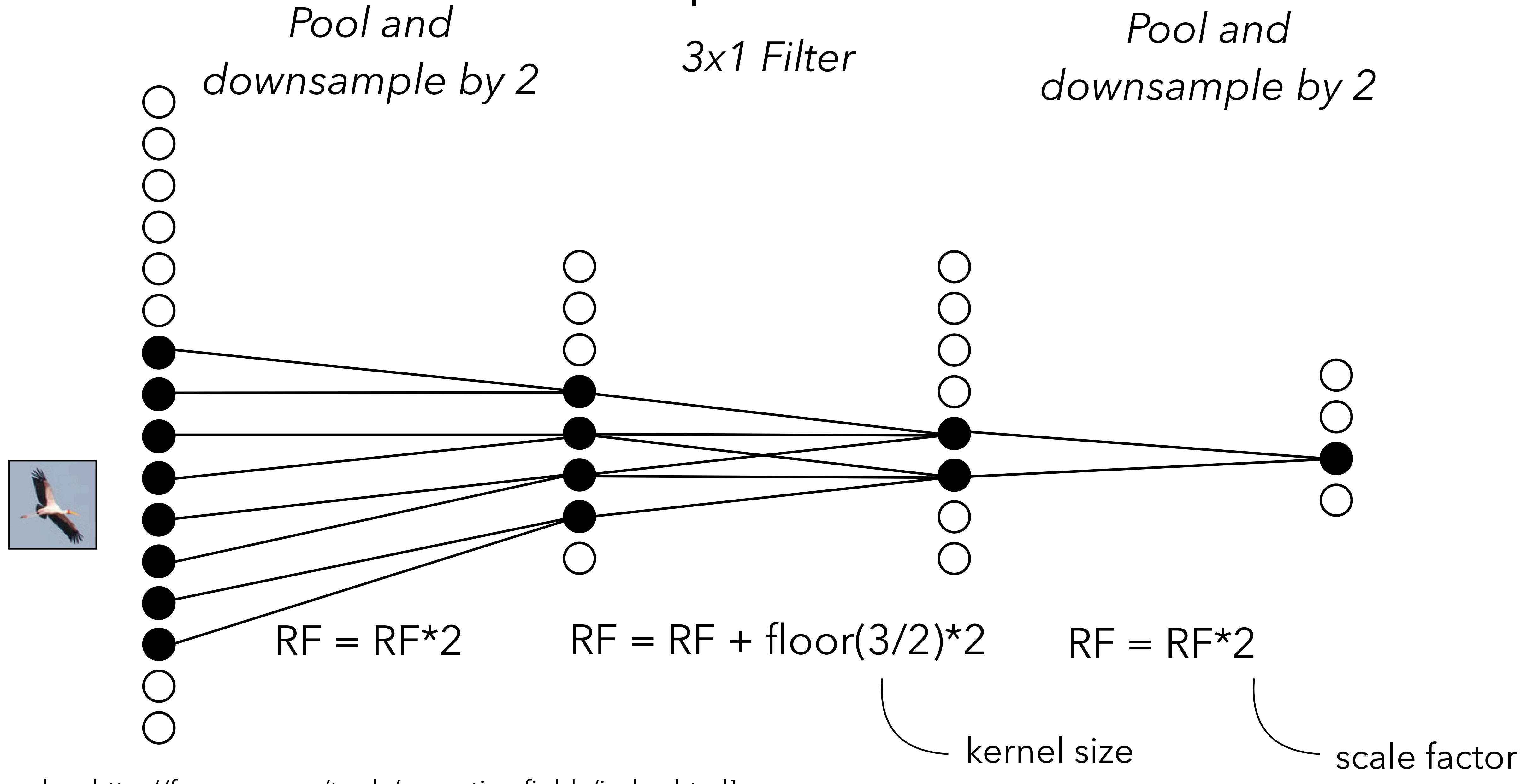


$$f(\mathbf{x}) = f_L(\dots f_2(f_1(\mathbf{x})))$$

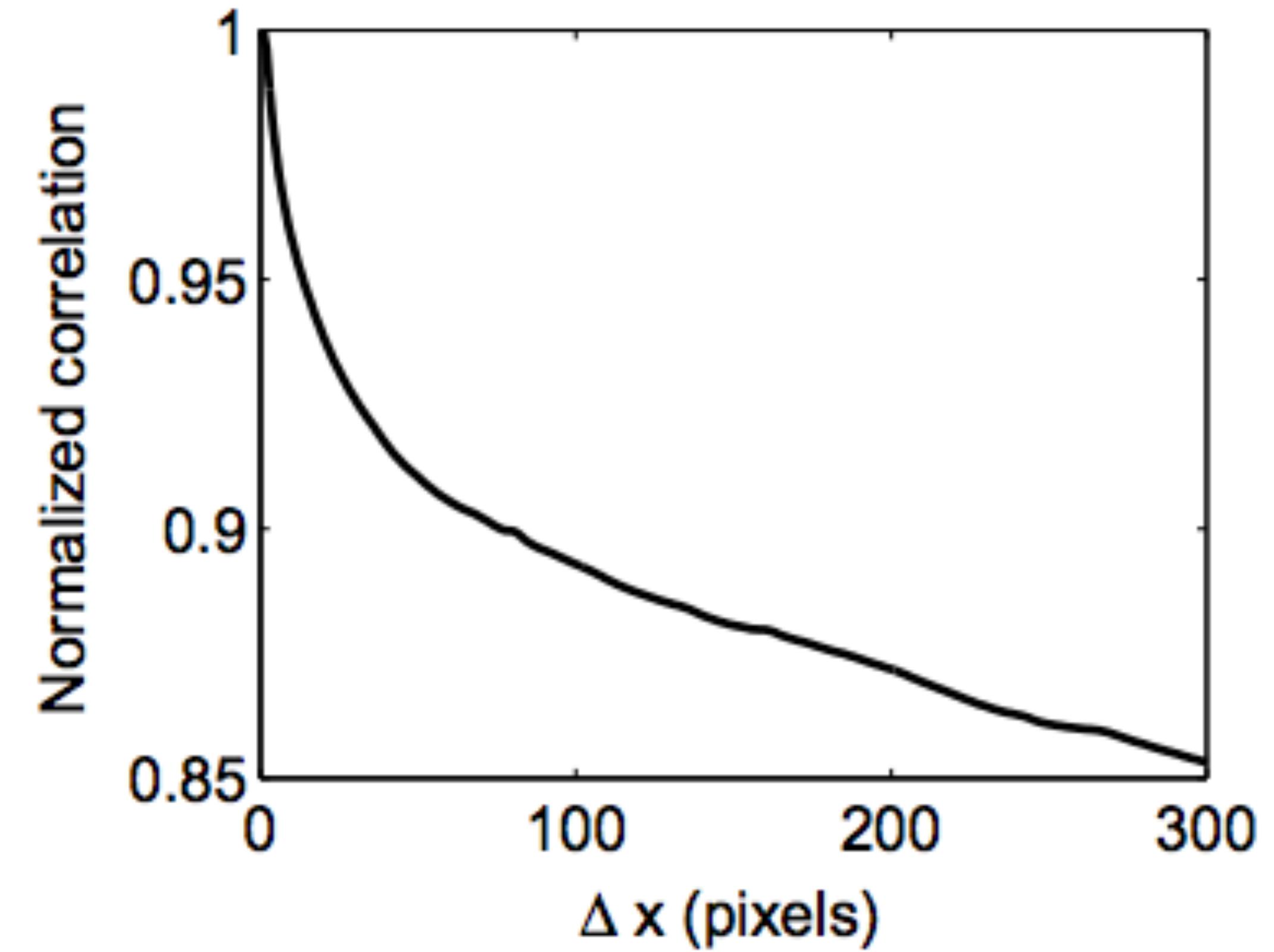
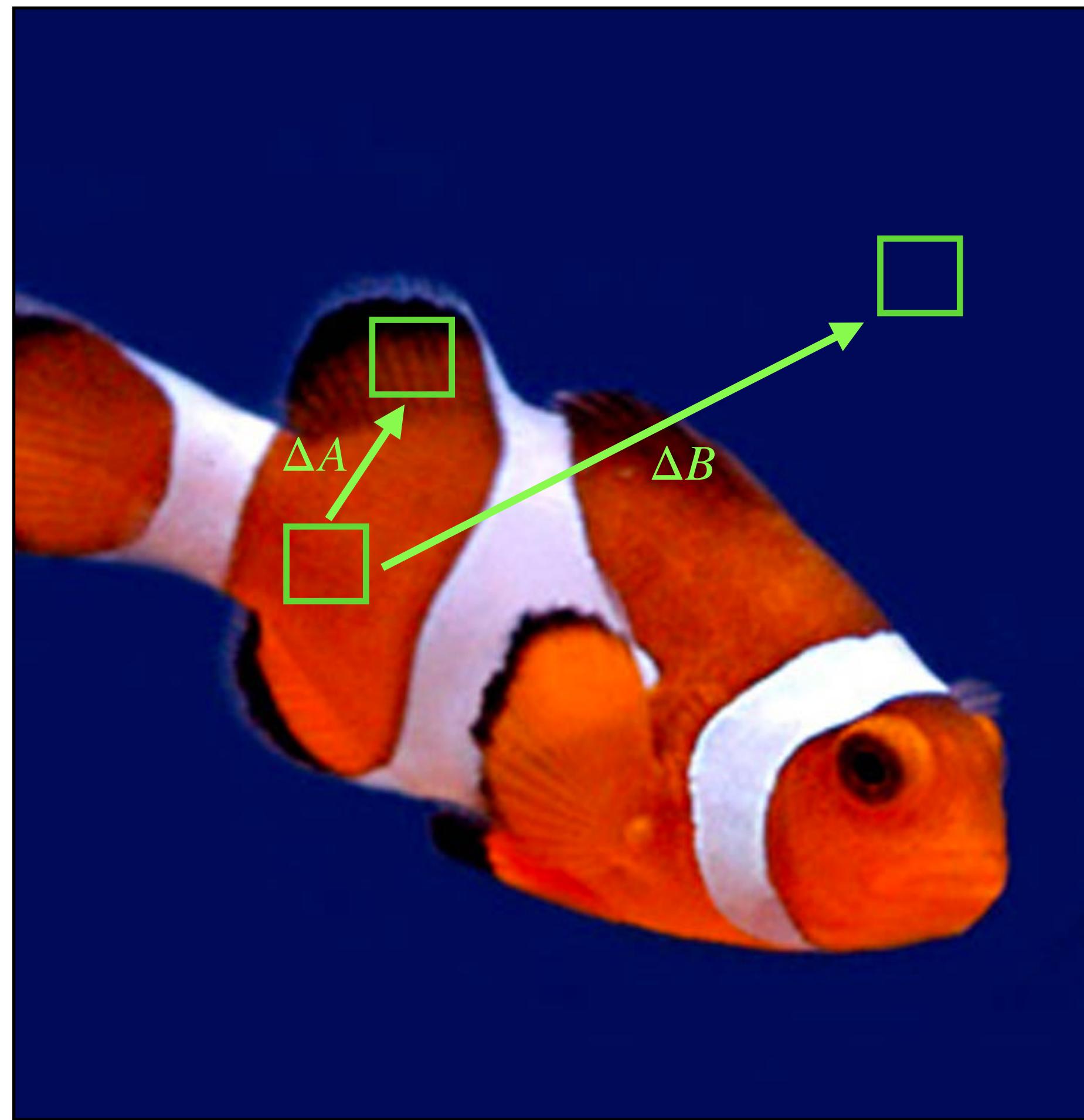
Receptive fields



Receptive fields



Why have coefficients only for nearby pixels?



[Simoncelli: Statistical Modelling of Photographic Images, 2005]

Implementation details

Convolutional layers

Option 1: Just make the giant convolution matrix and use the fully backprop equations for linear layers!
Early deep learning frameworks did this.



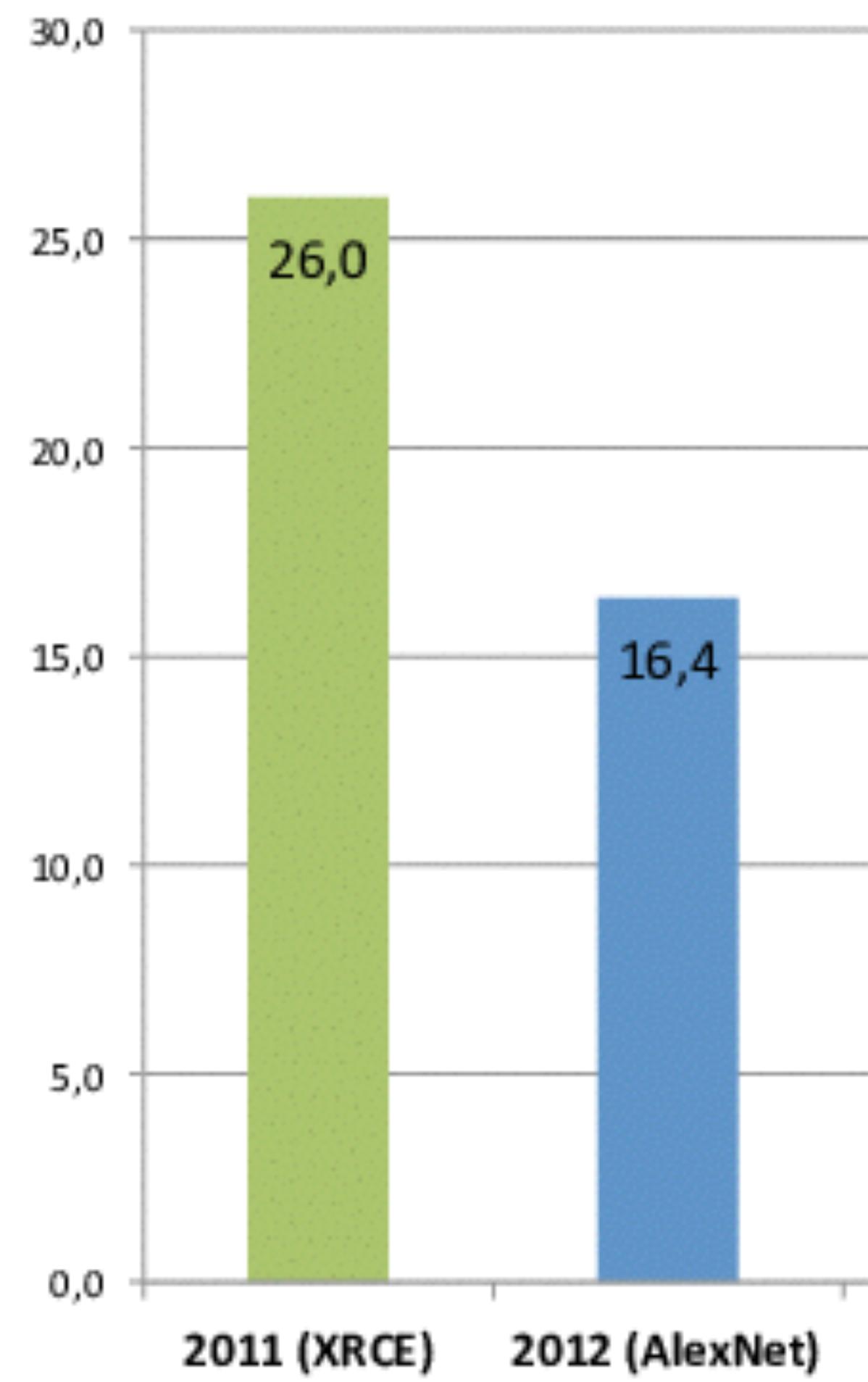
Option 2: Use a more optimized implementation.
There are optimized GPU implementations. For very large filters, these use Fourier Transform.

Network designs

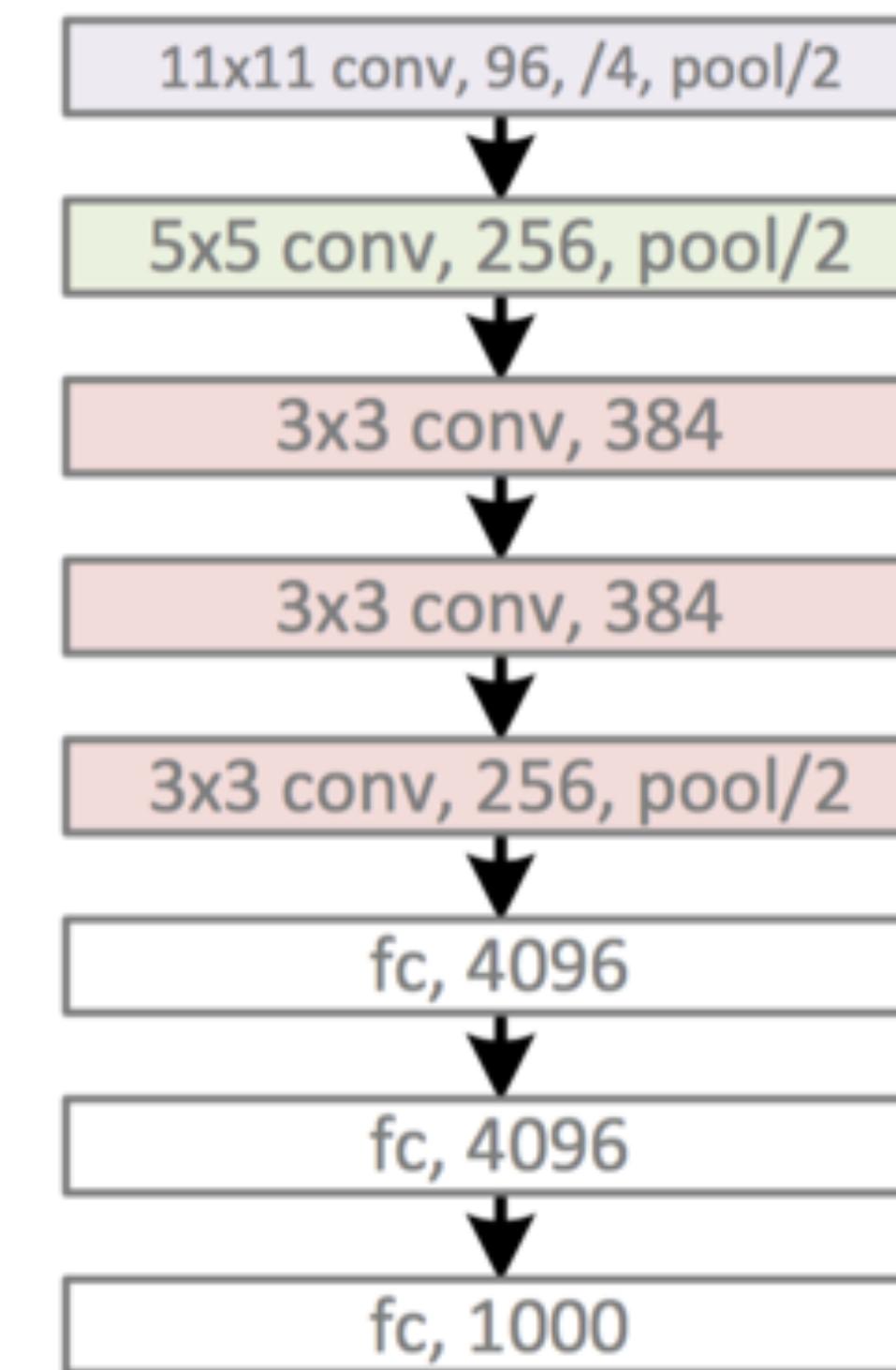
ImageNet classification (top 5)



ImageNet classification (top 5)



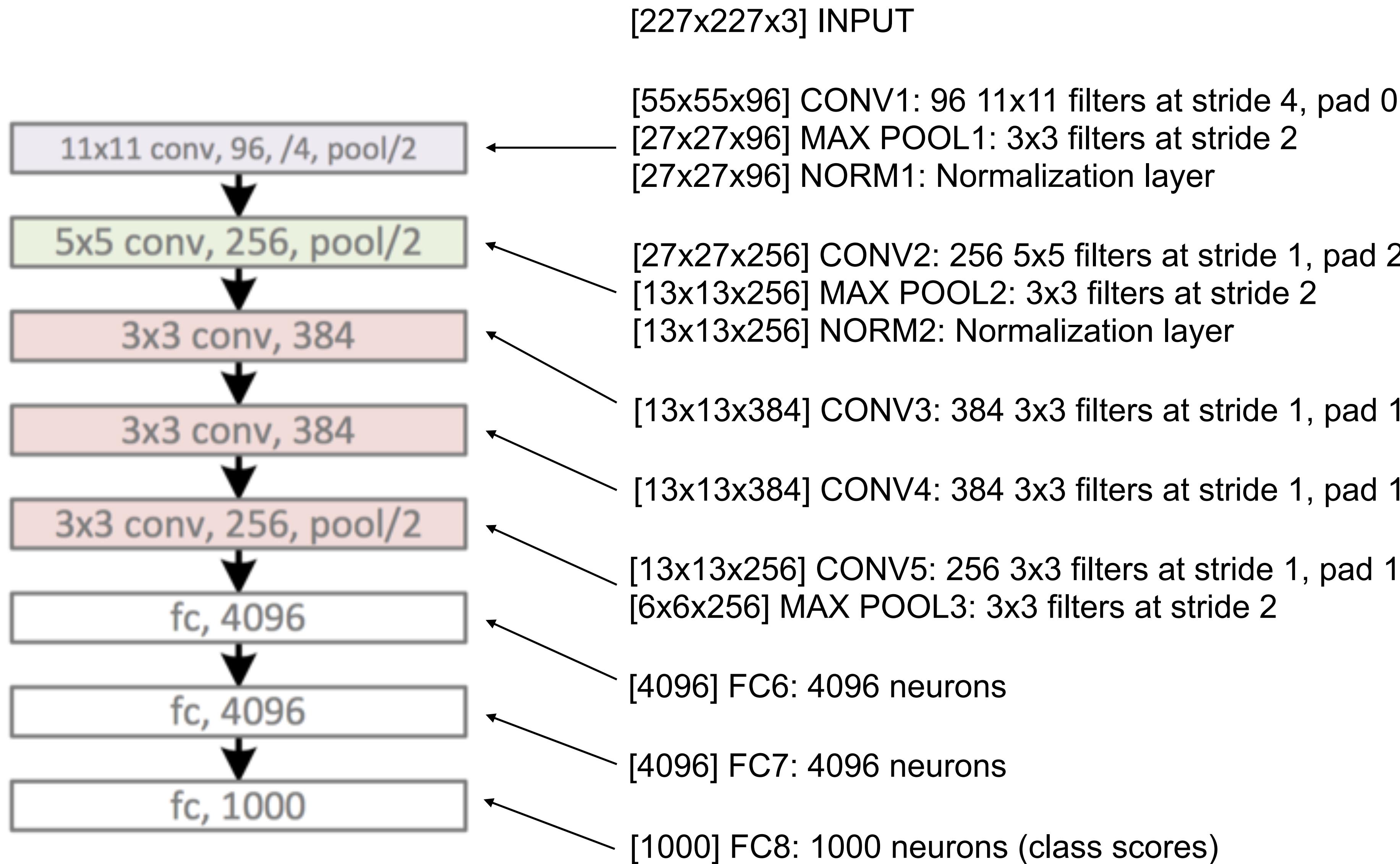
2012: AlexNet
5 conv. layers

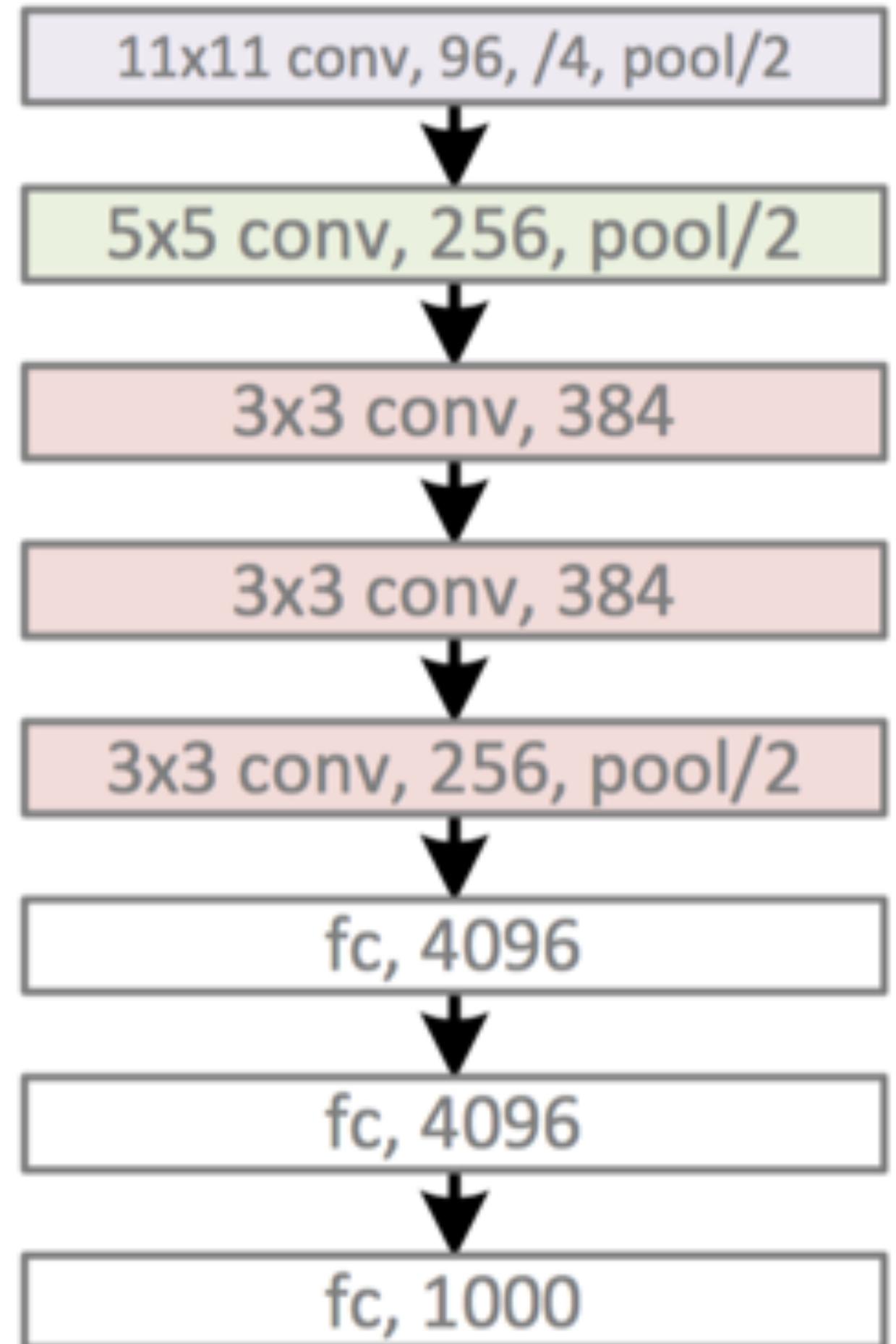


Error: 16.4%

[Krizhevsky et al: ImageNet Classification with Deep Convolutional Neural Networks, NeurIPS 2012]

Alexnet – [Krizhevsky et al. NIPS 2012]

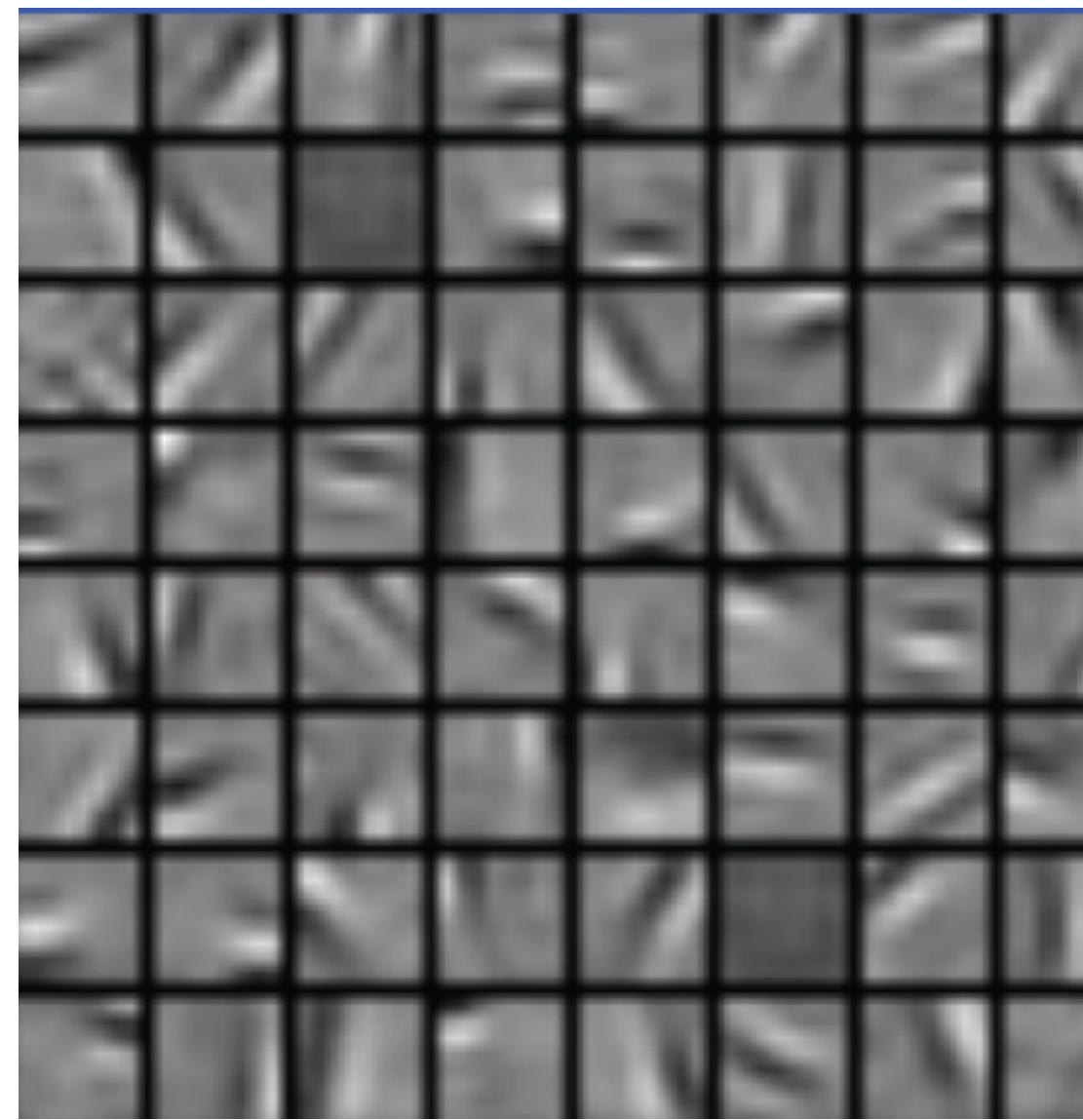




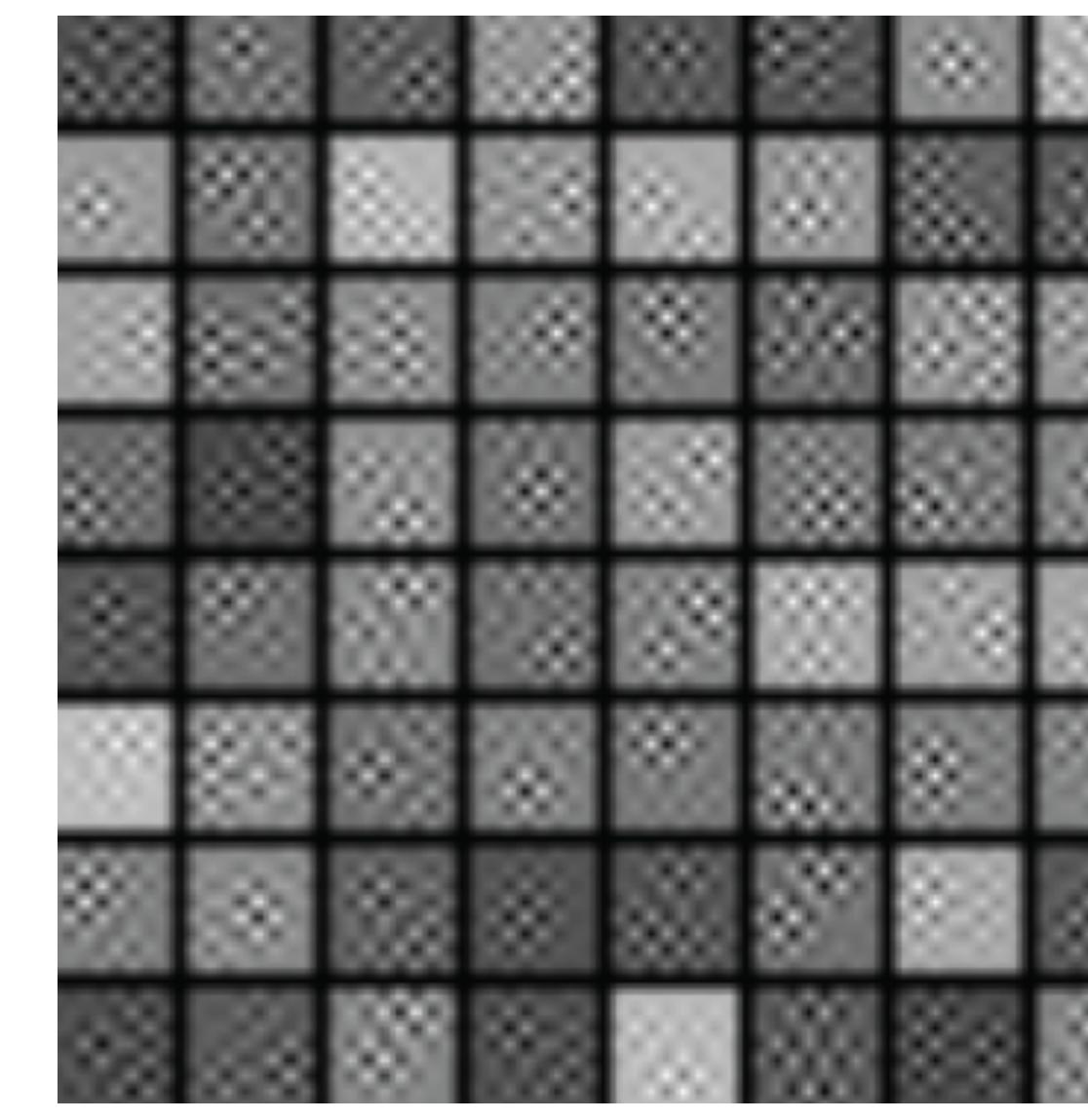
What filters are learned?

What filters are learned?

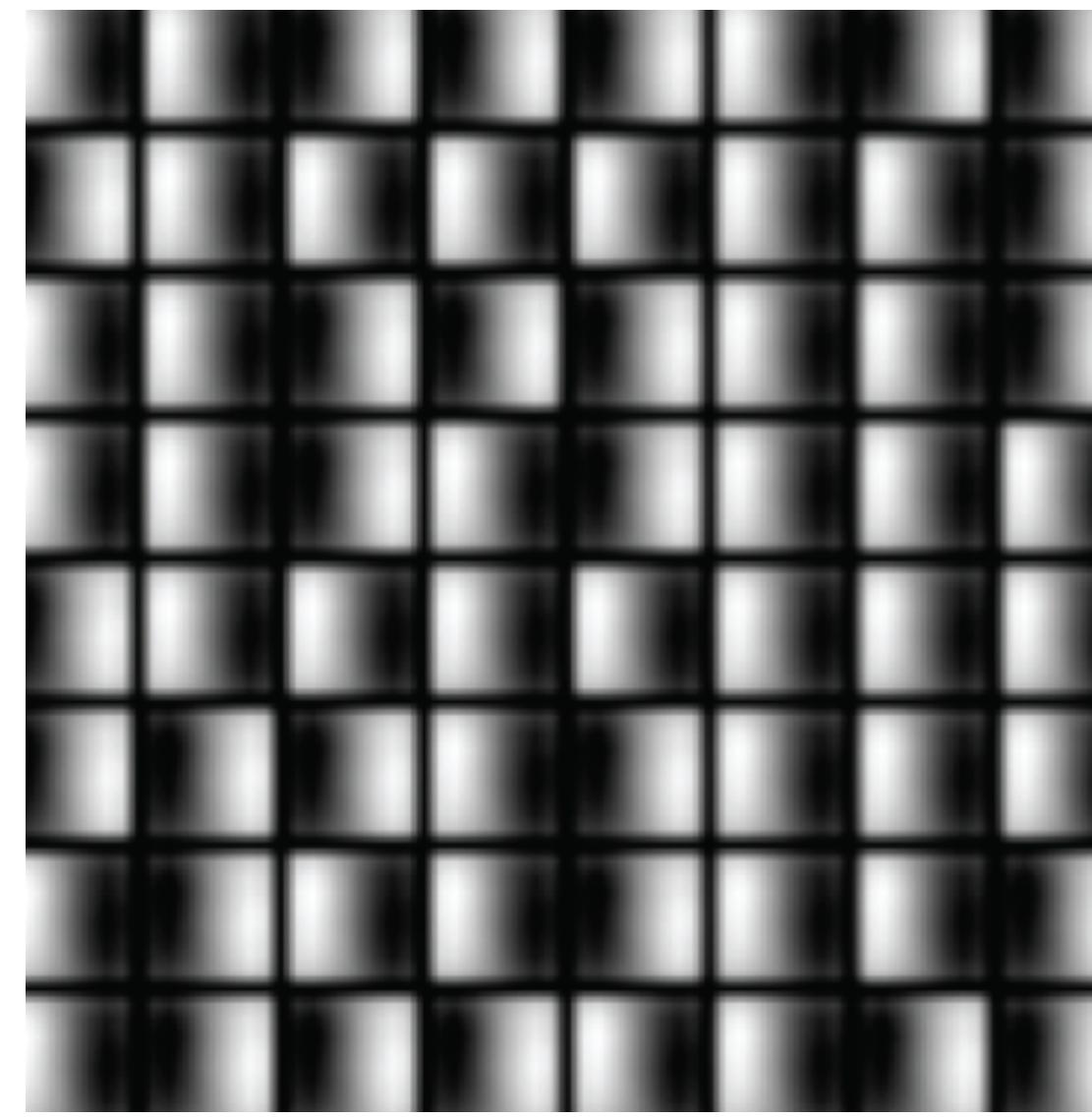
A



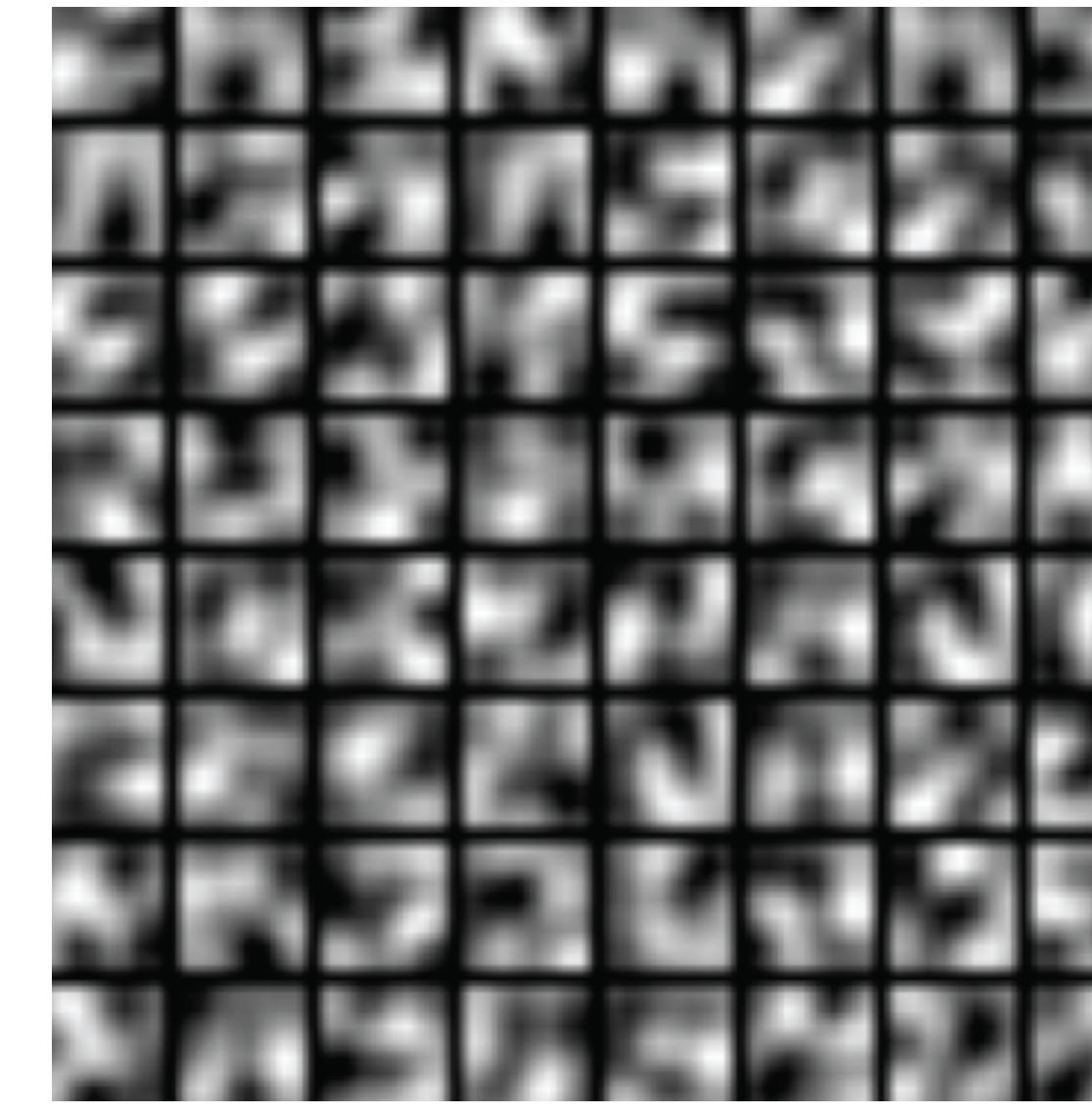
B



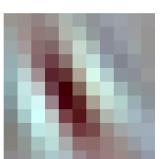
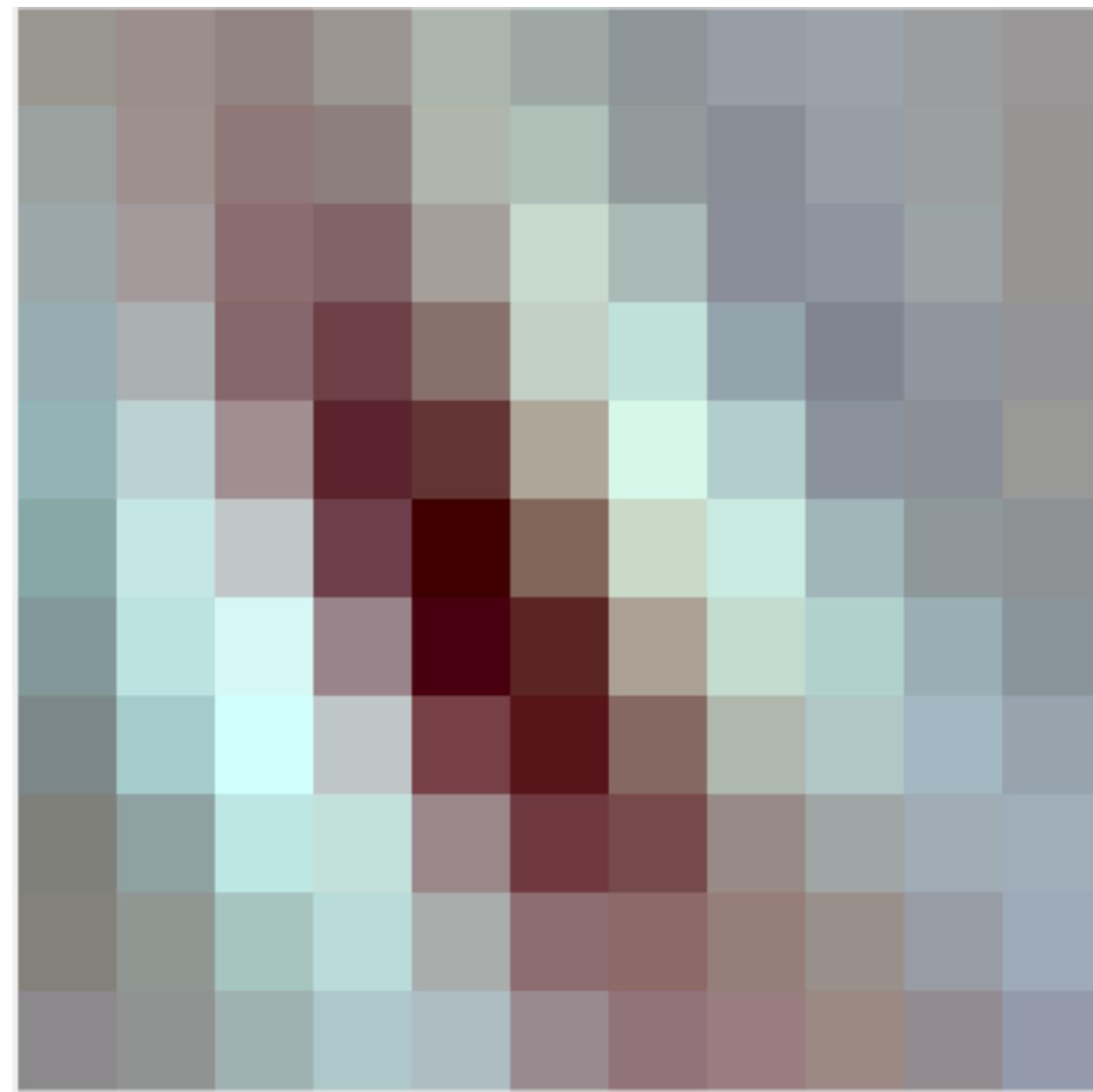
C



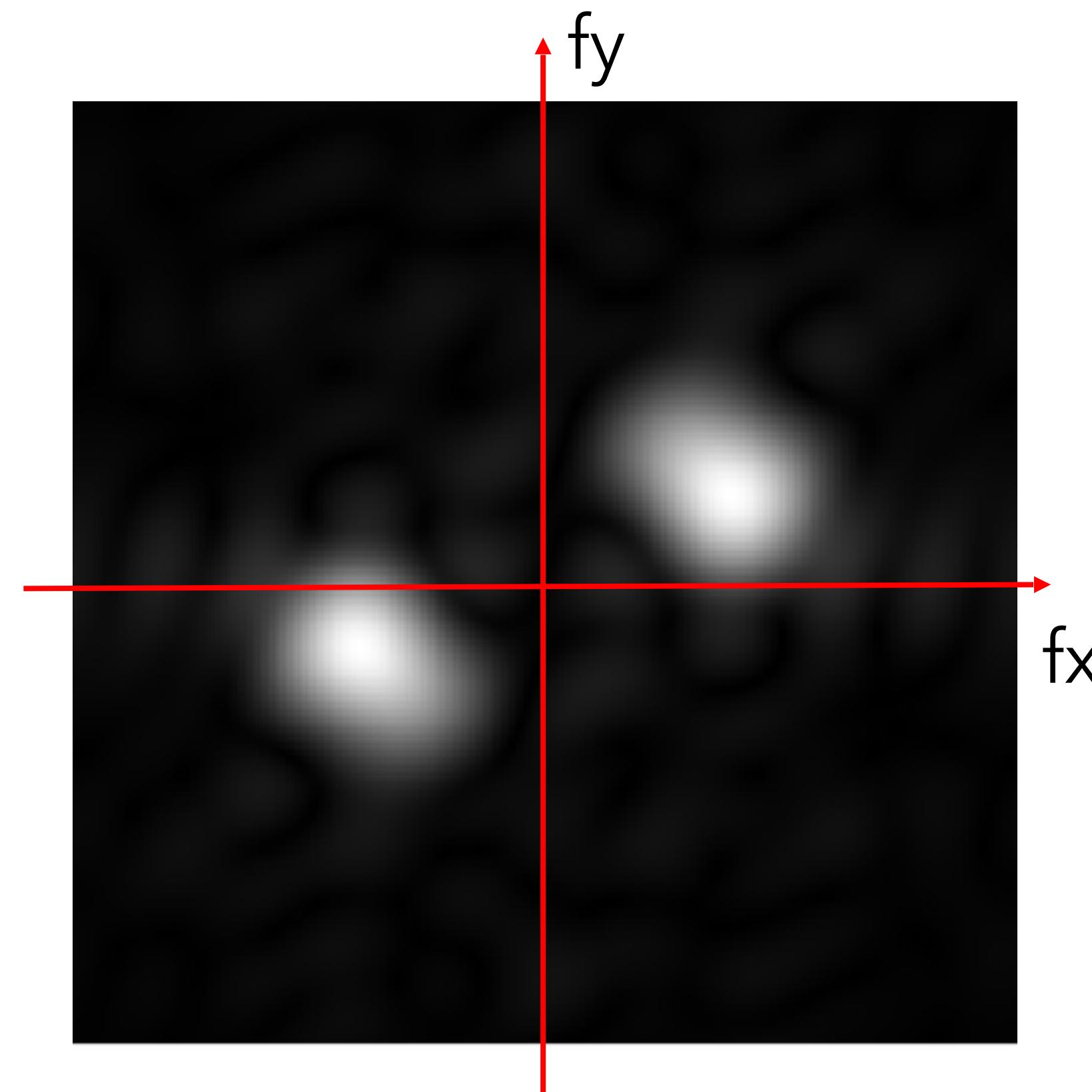
D



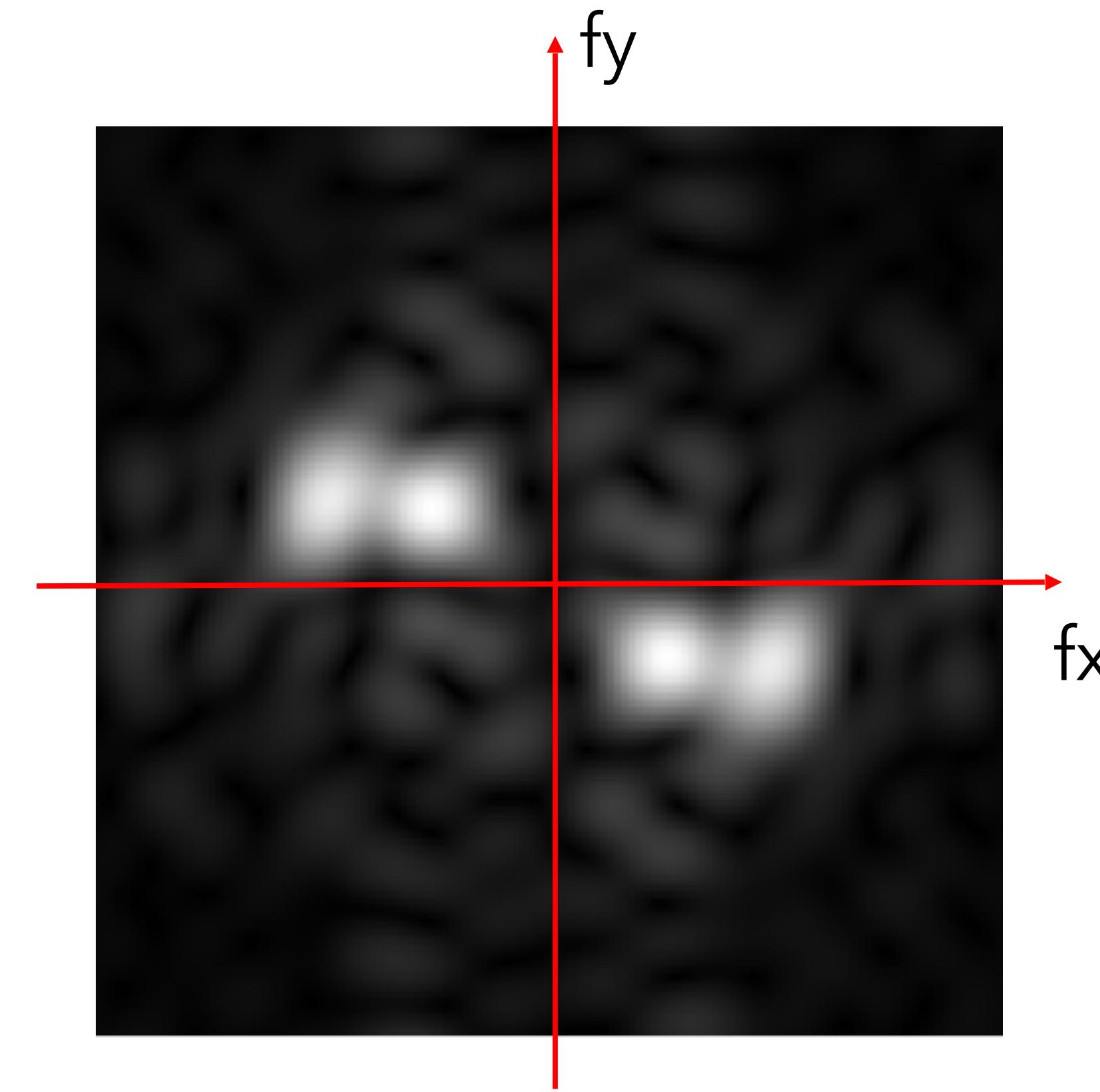
Filters in first layer



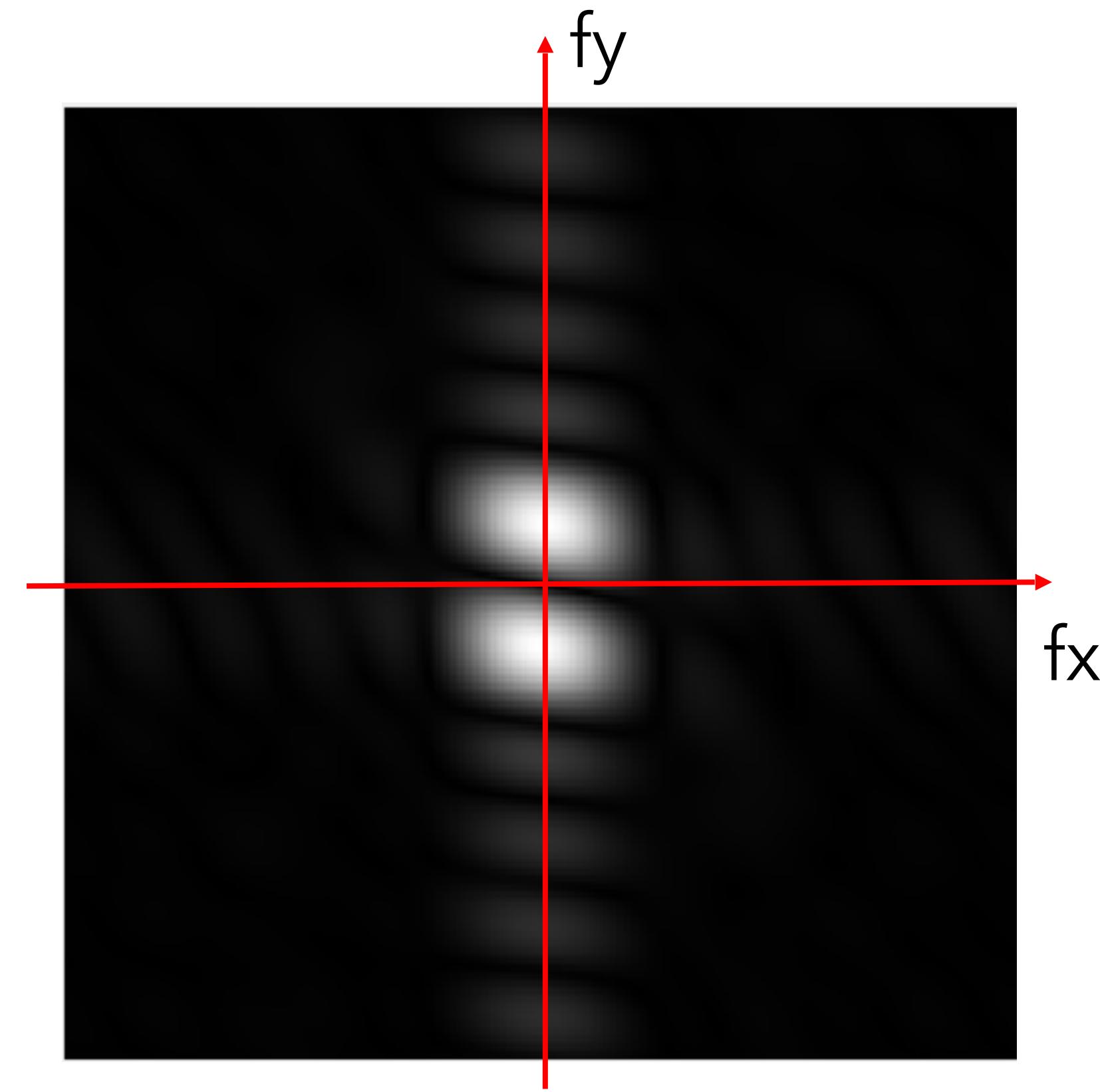
11x11 convolution kernel
(3 color channels)



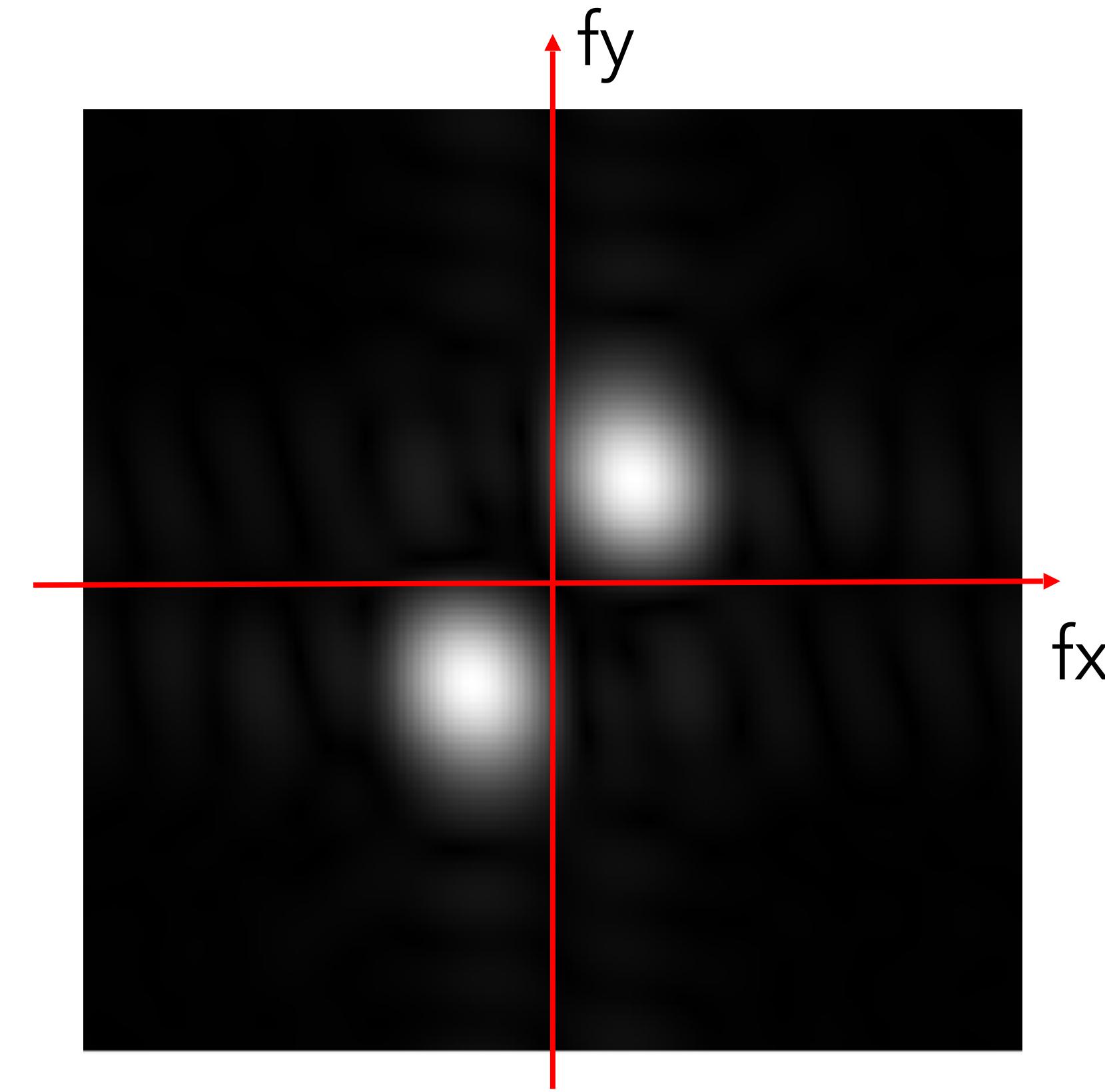
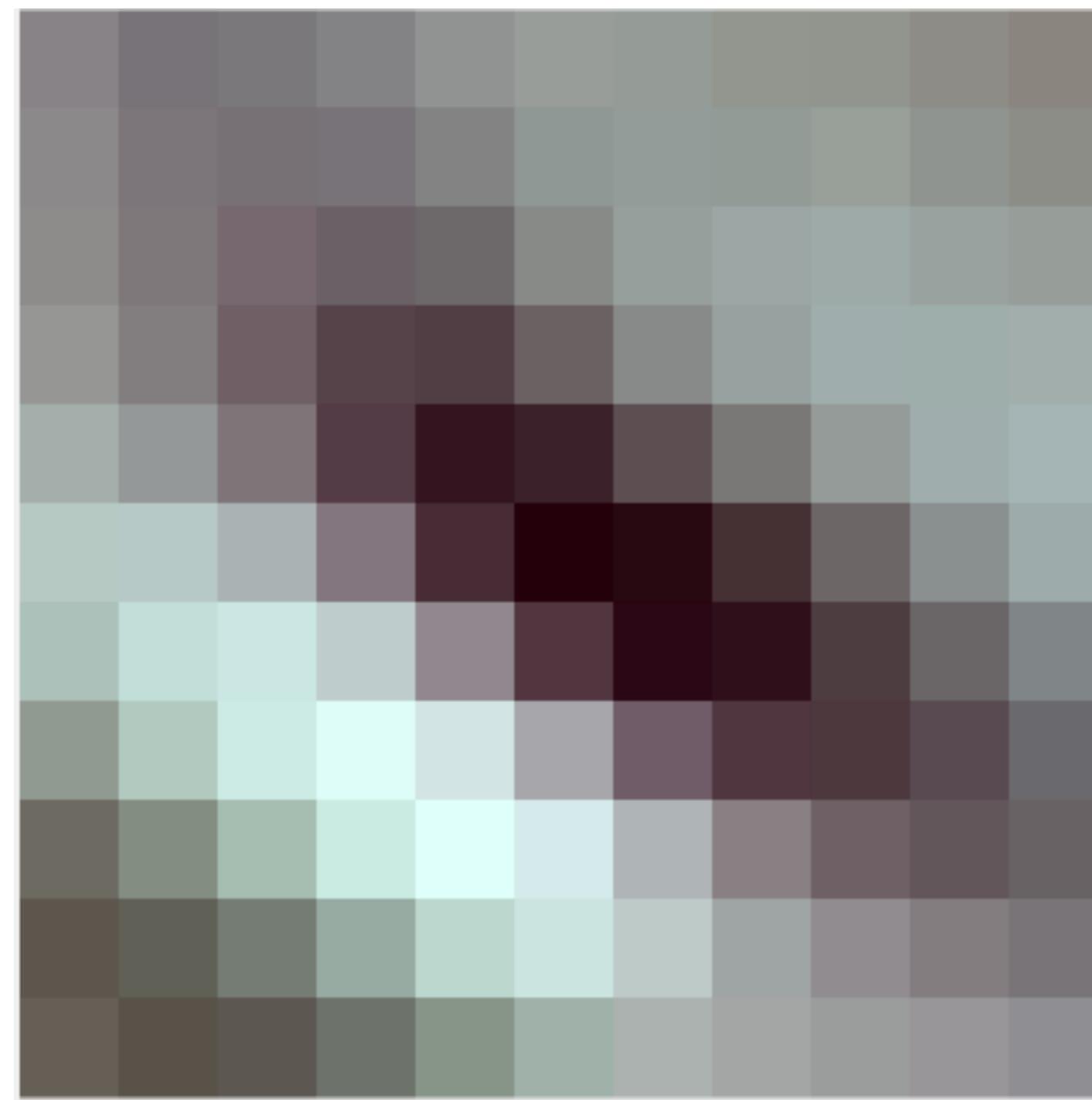
Filters in first layer



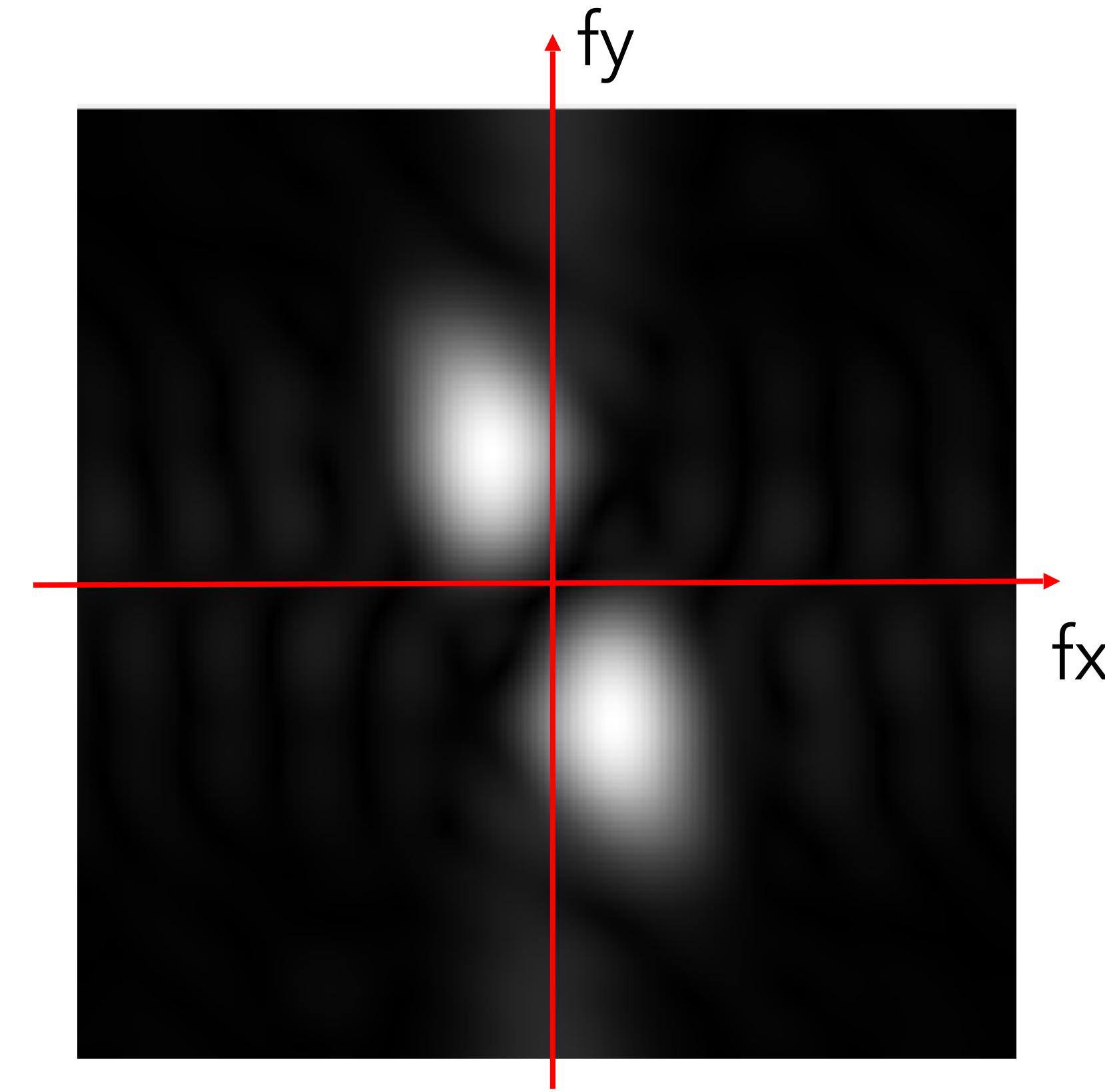
Filters in first layer



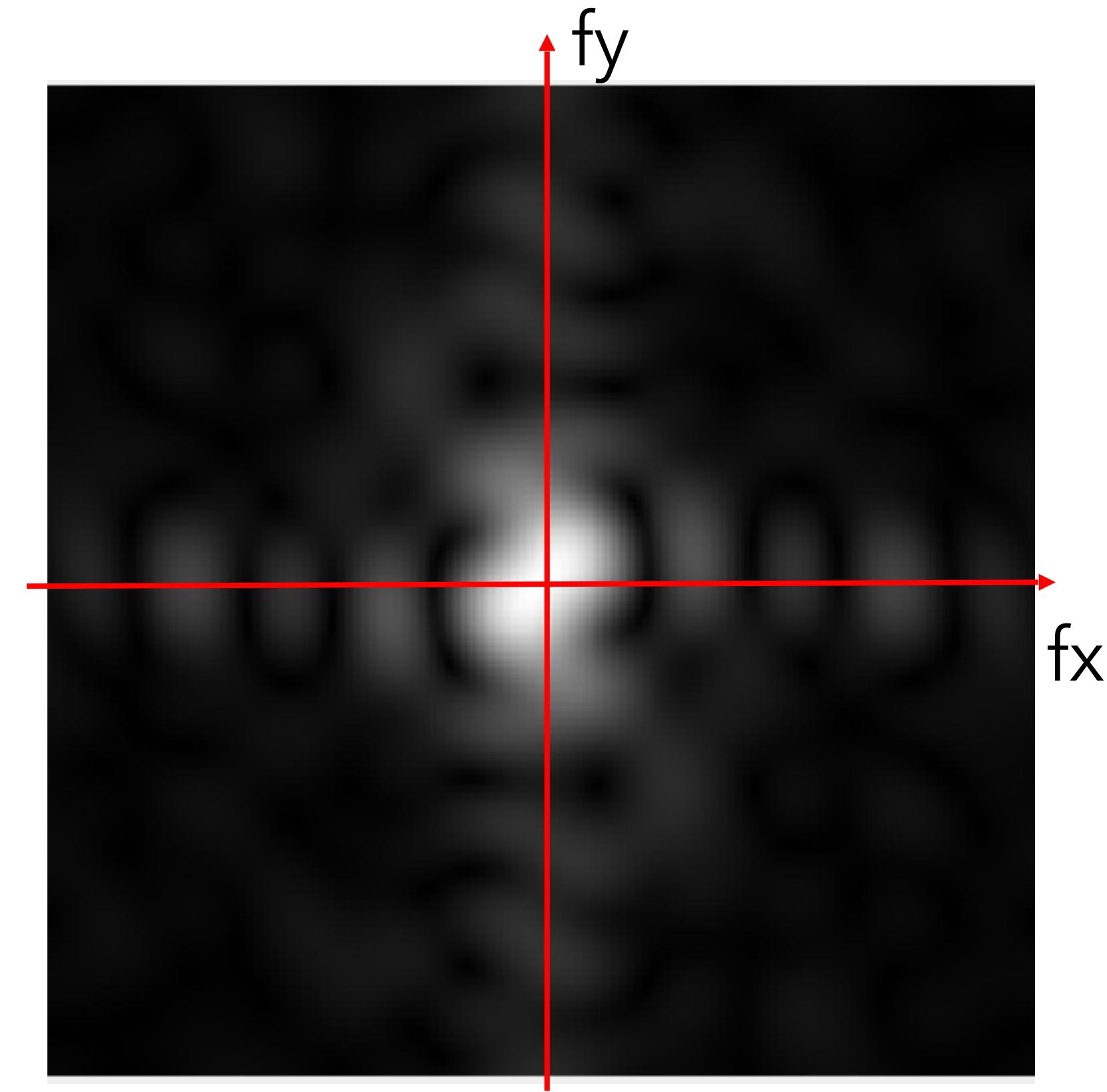
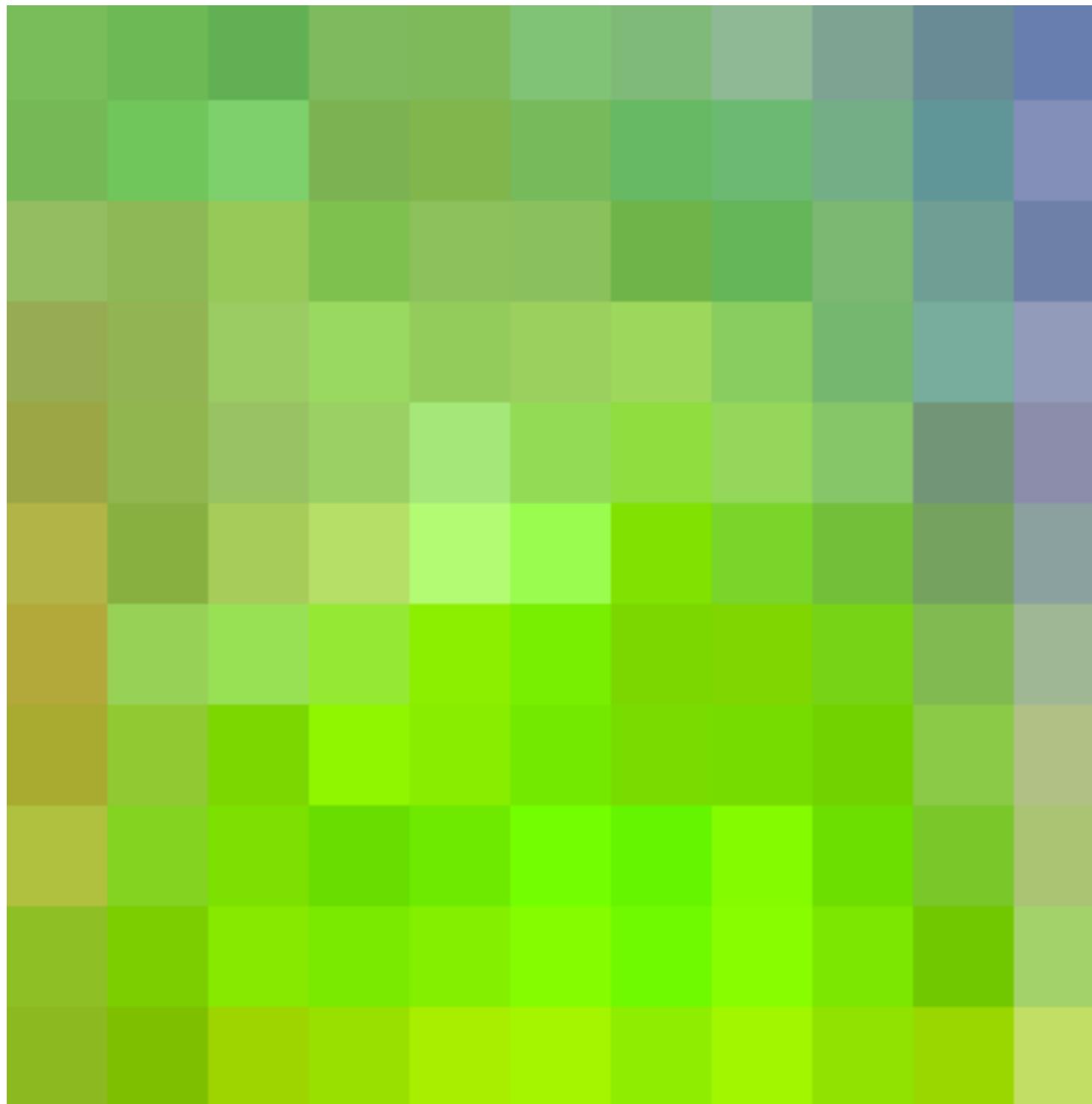
Filters in first layer



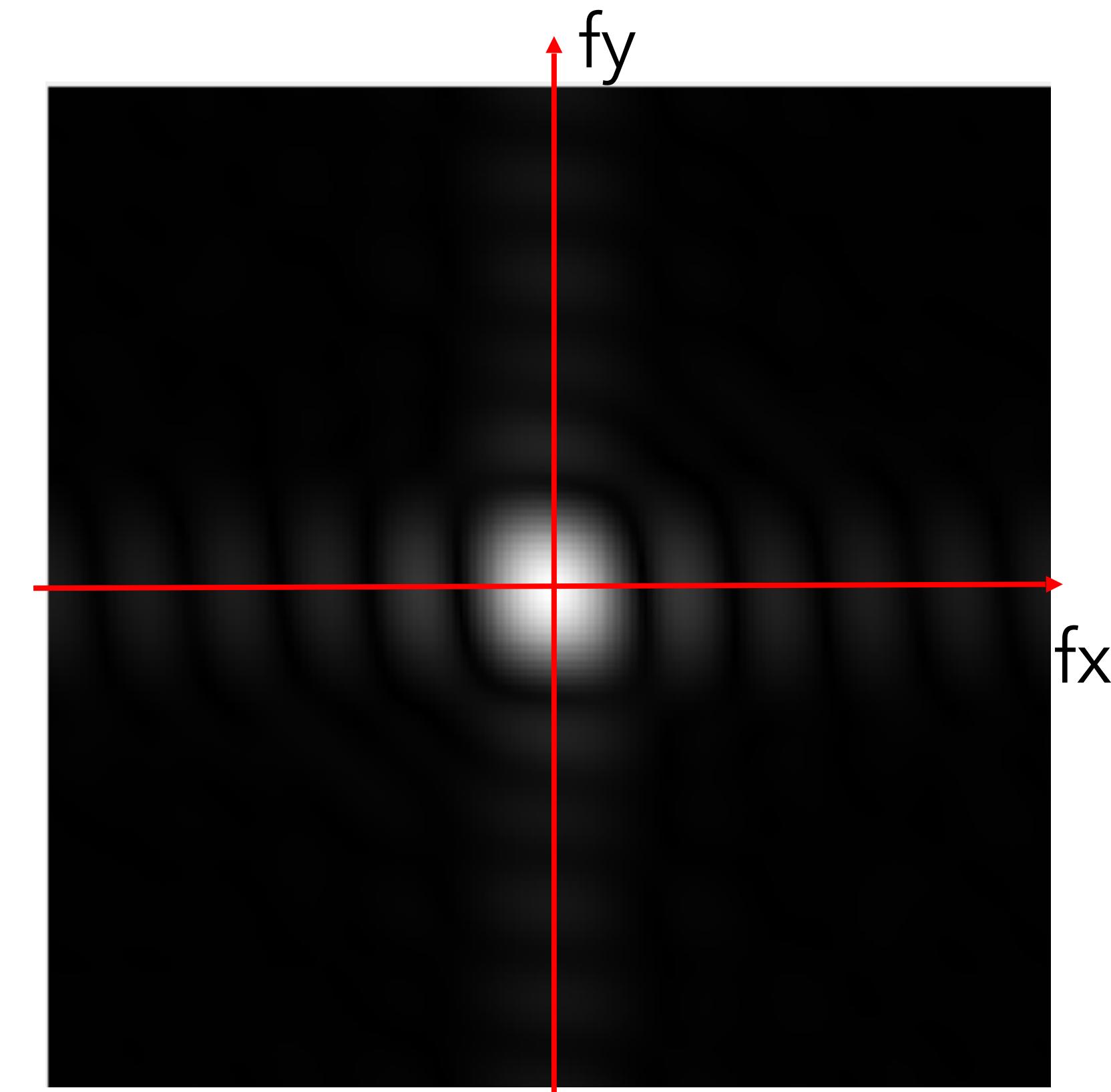
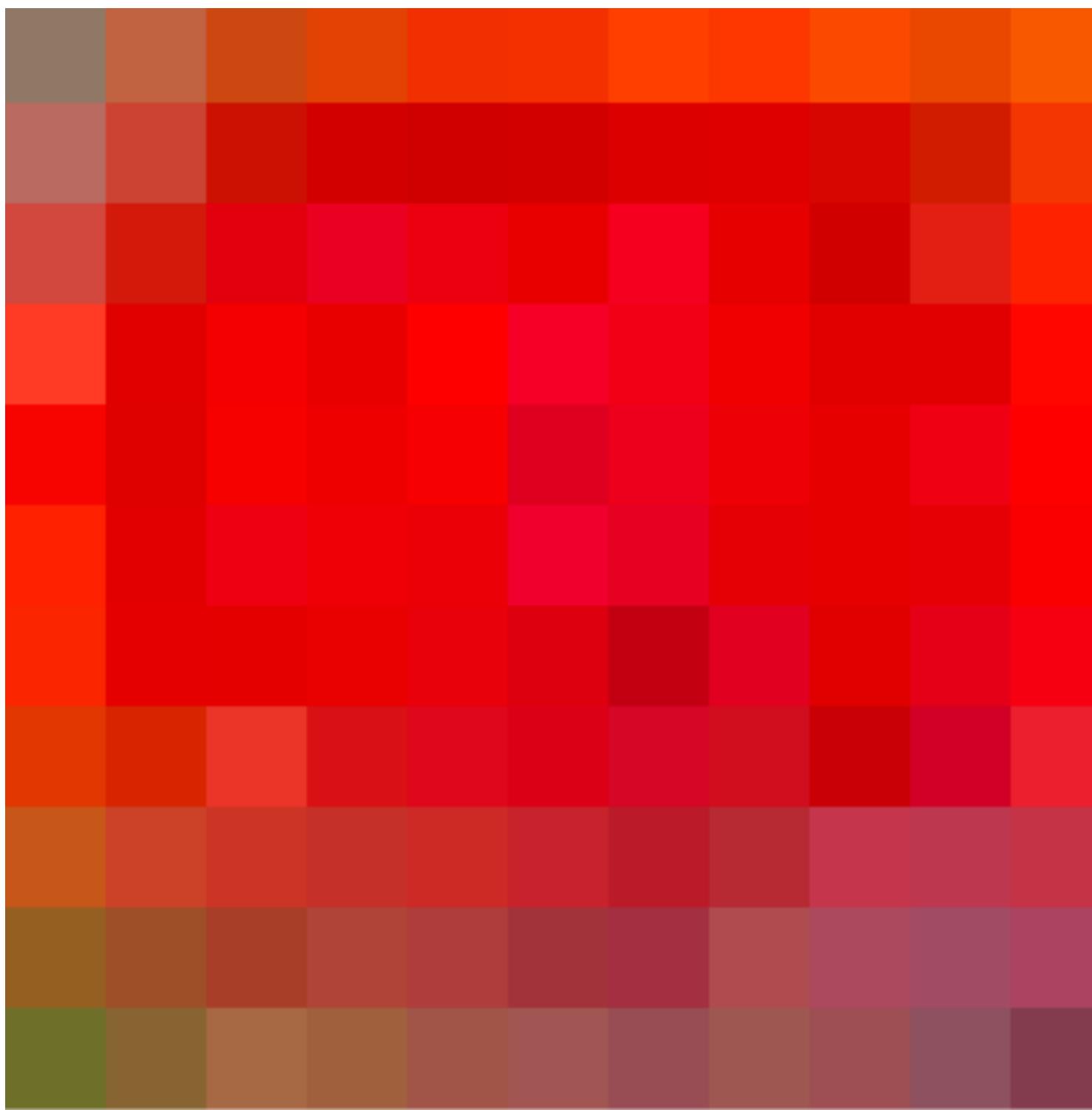
Filters in first layer



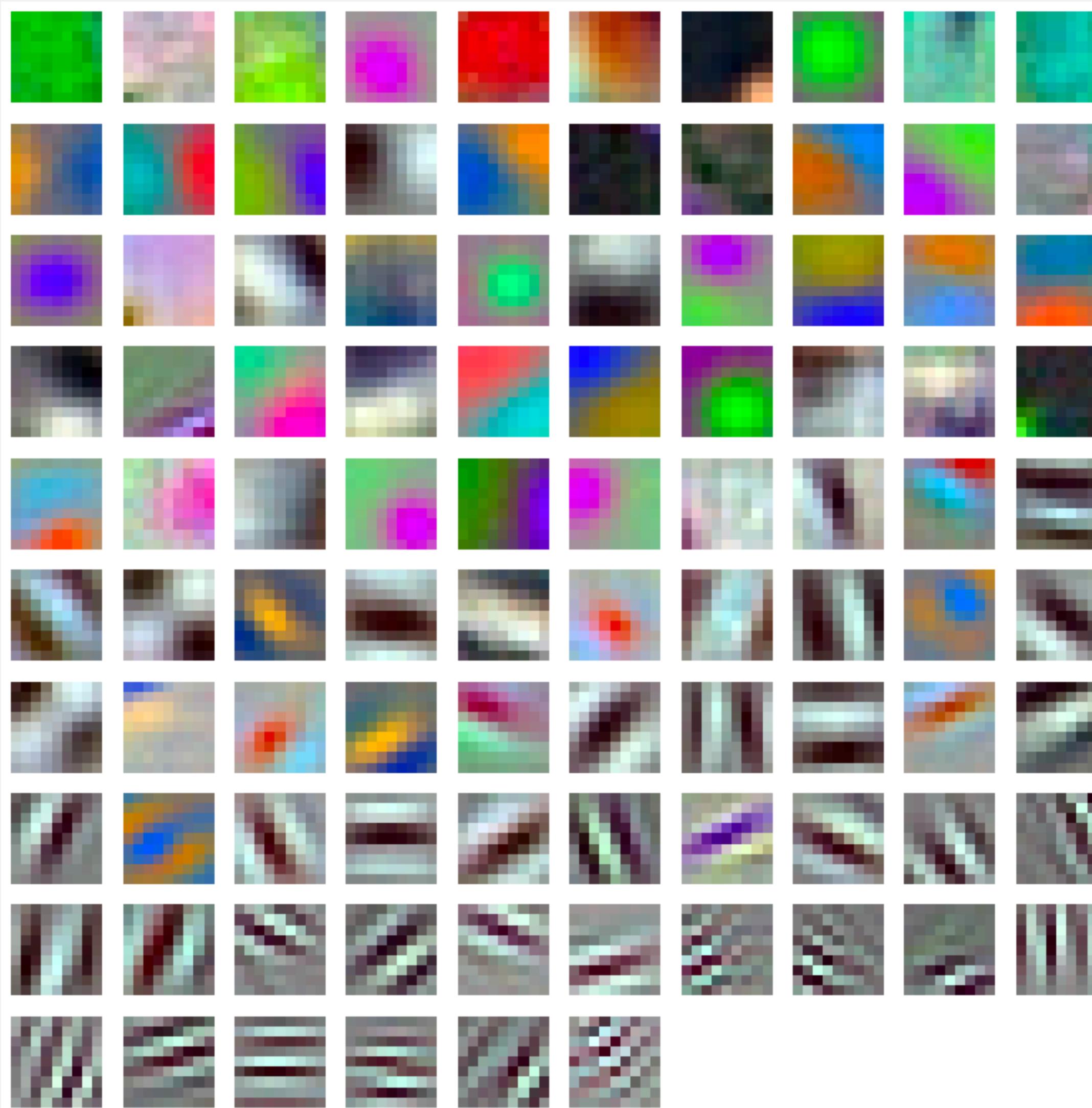
Filters in first layer



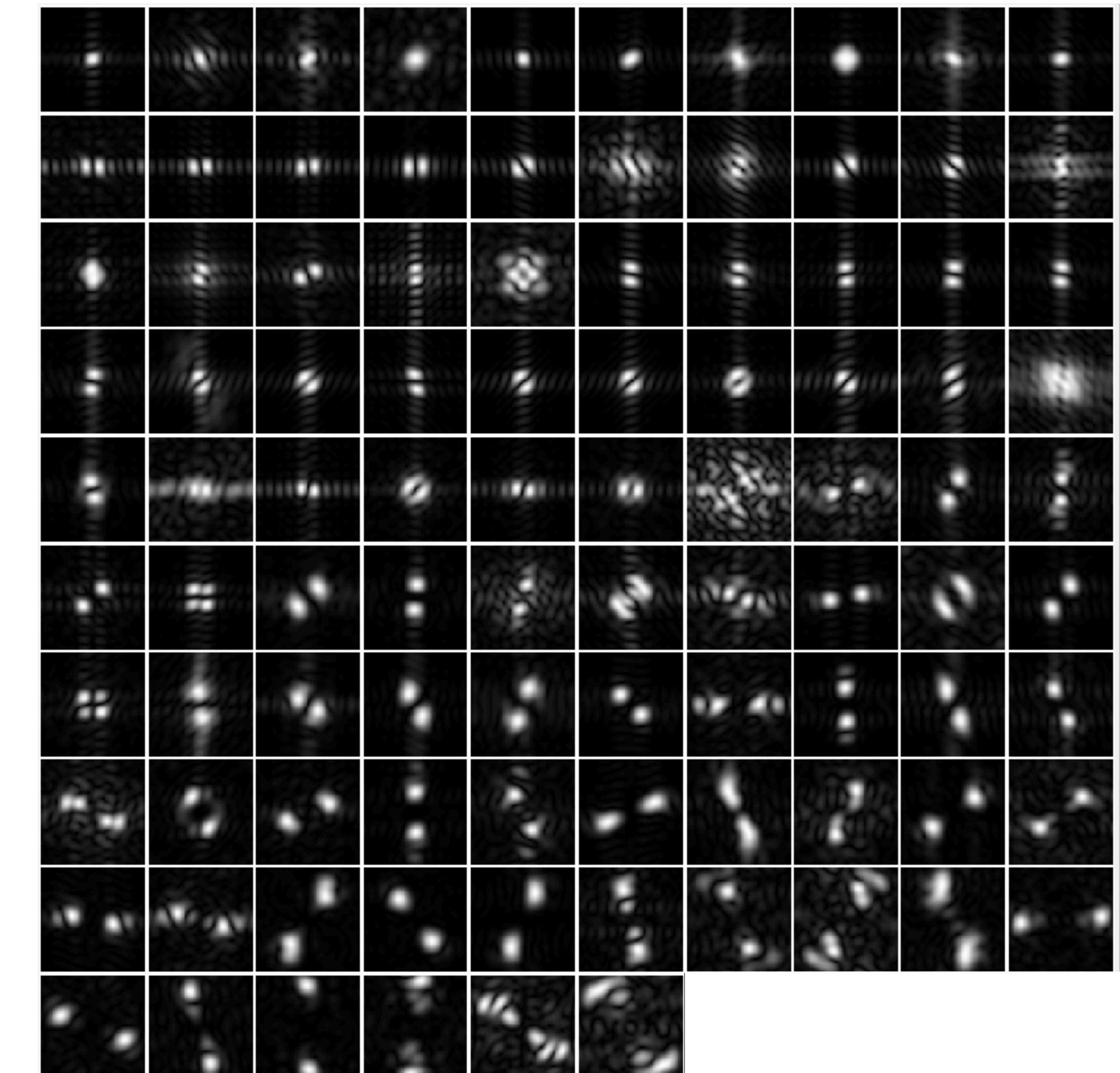
Filters in first layer



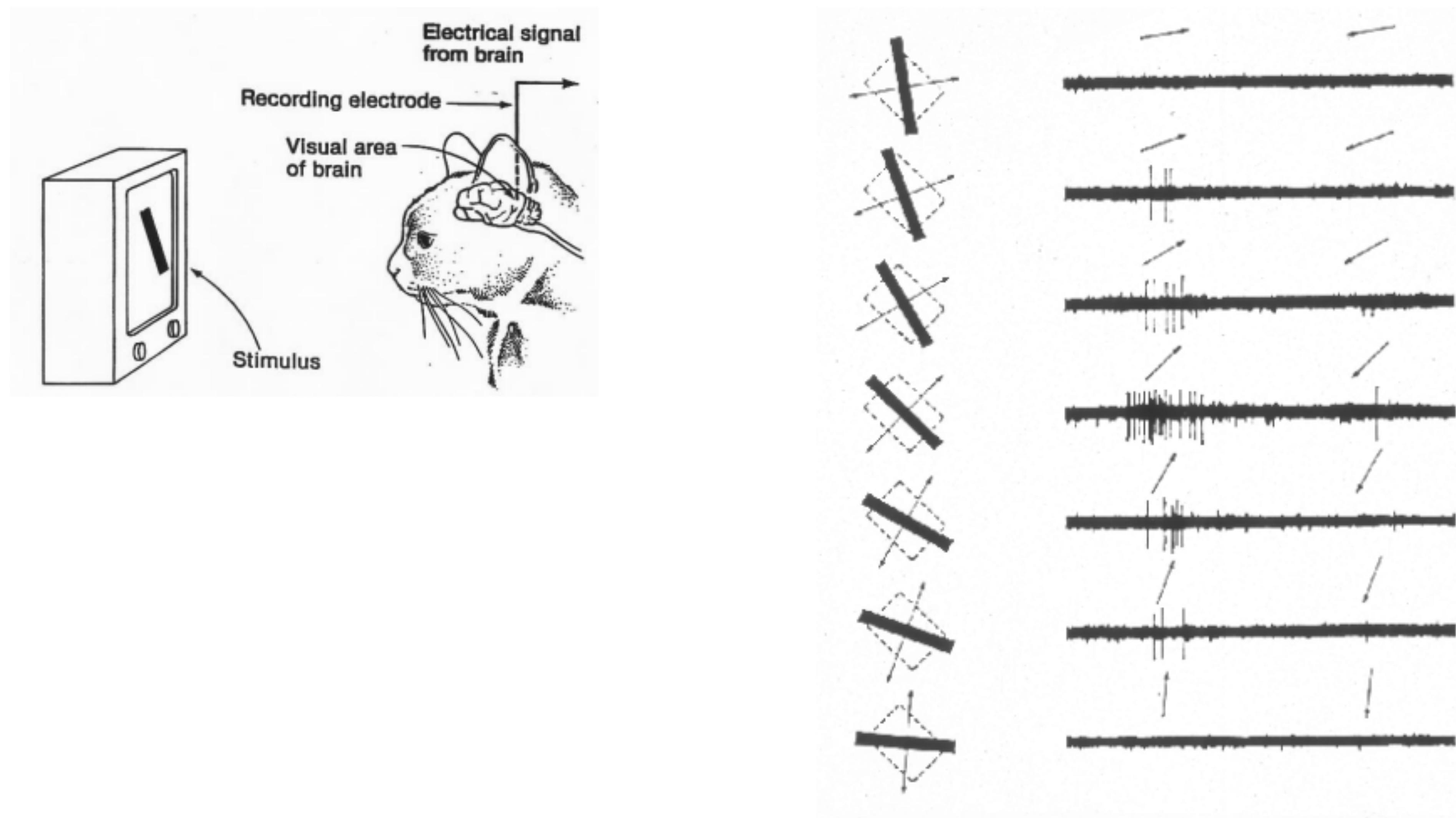
Filters in first layer



96 Units in conv1

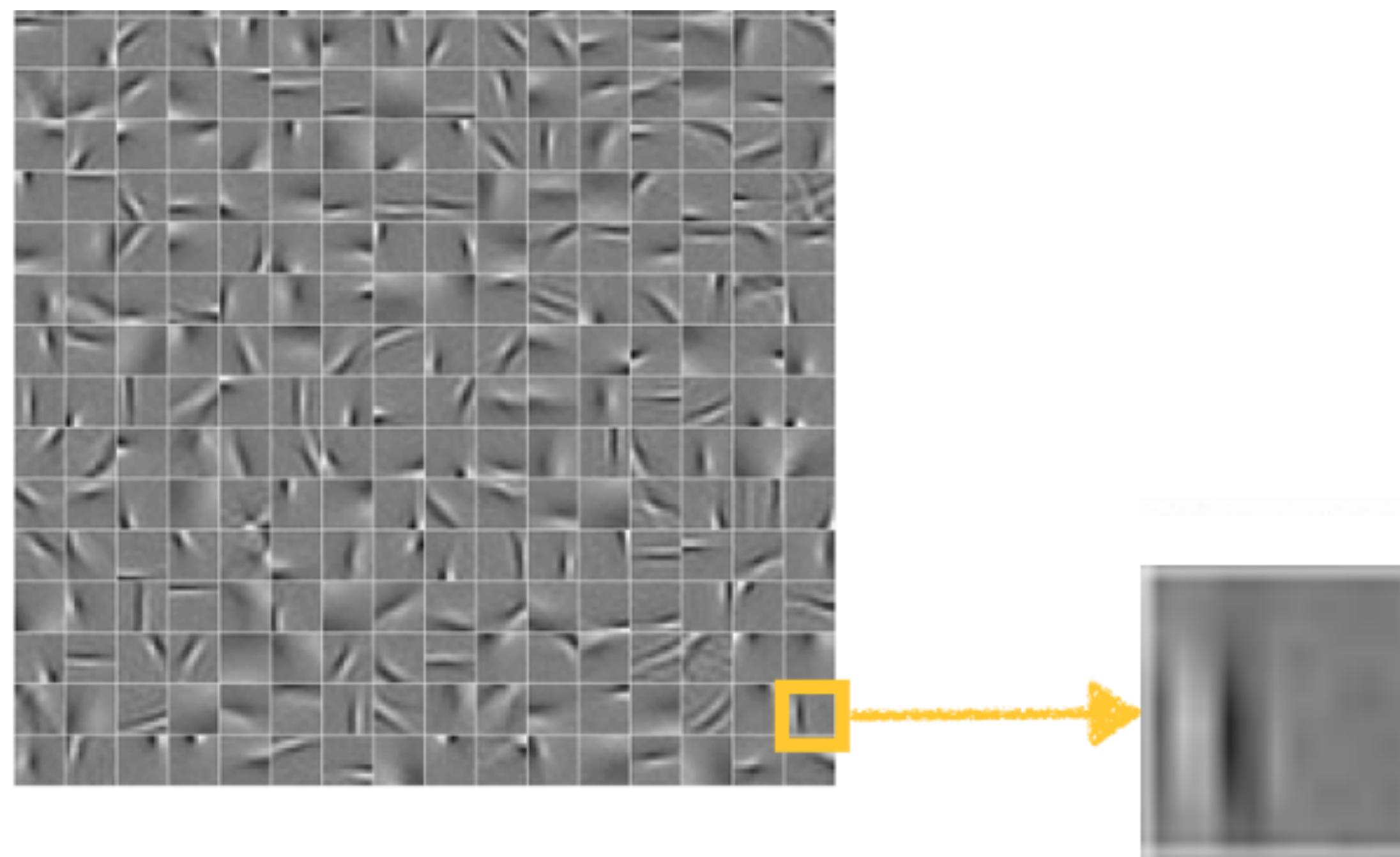


[Hubel and Wiesel 59]

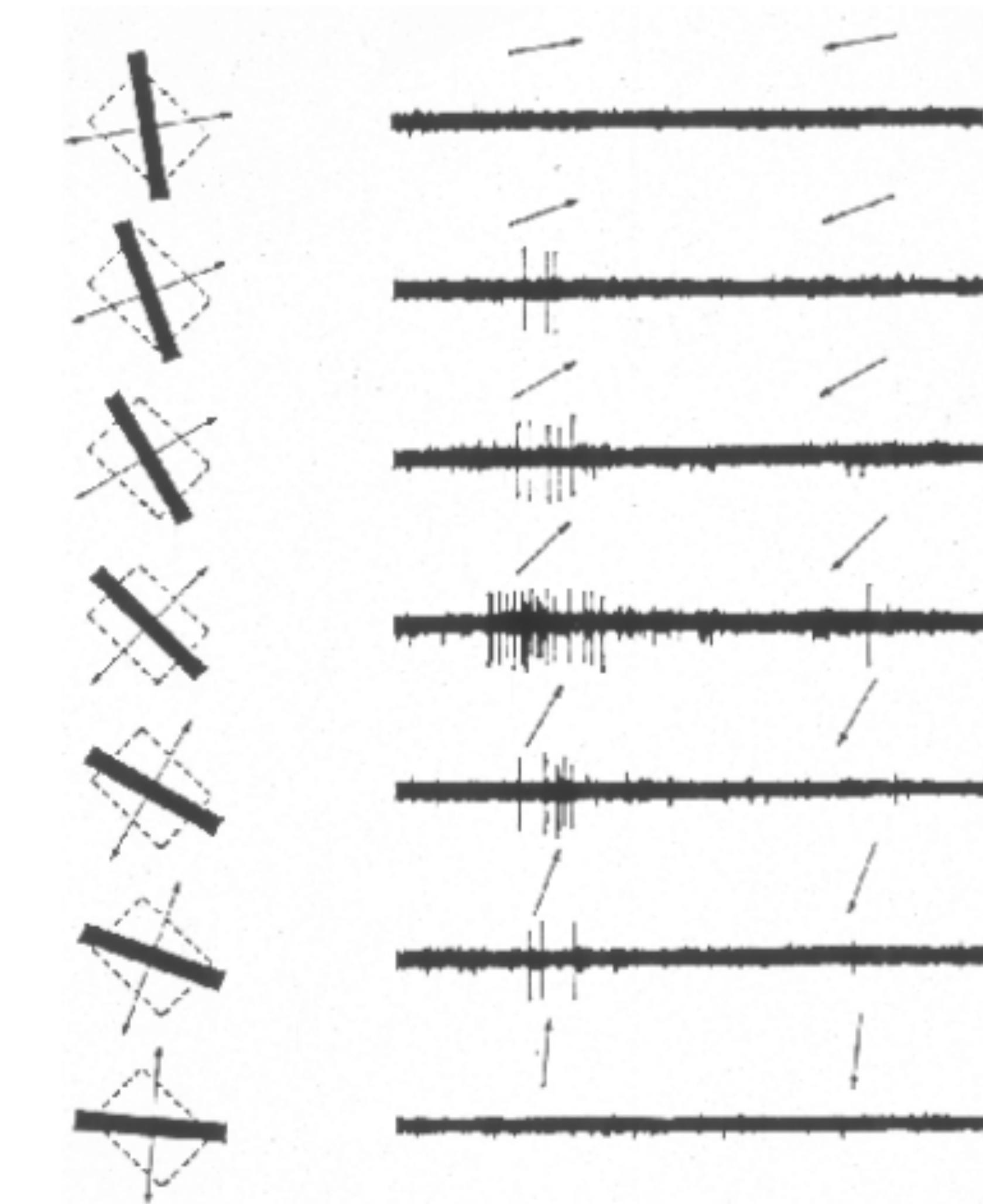


Source: Andrea Vedaldi

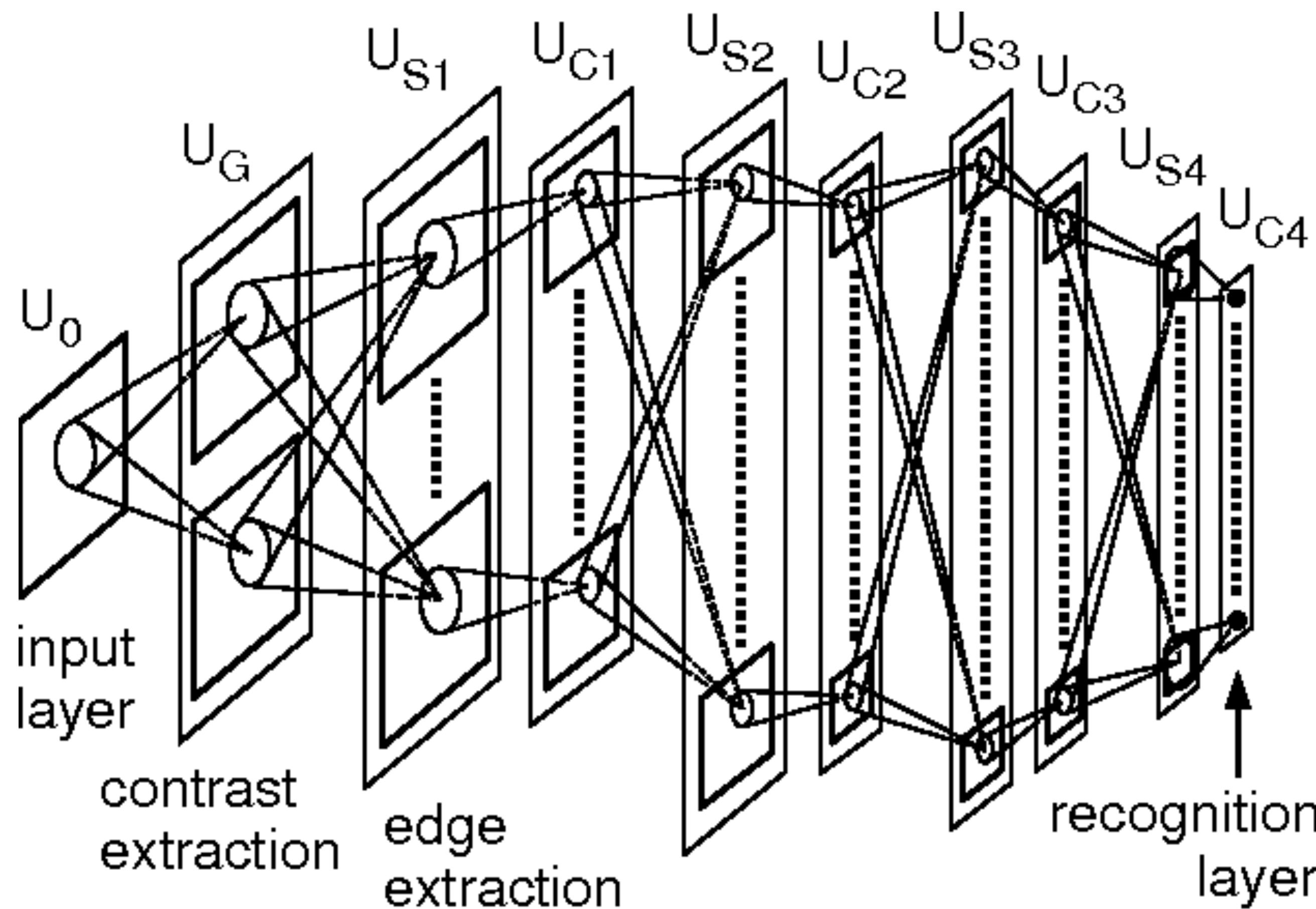
A good fit to the experimental results



Oriented filters

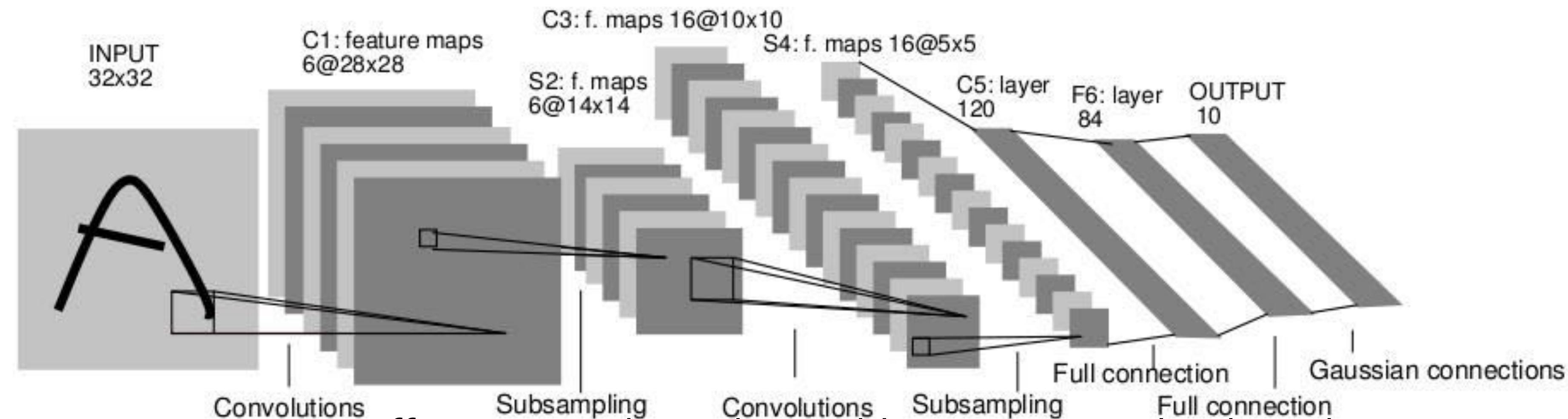


History: Neocognitron



K. Fukushima, 1980s

History: LeNet

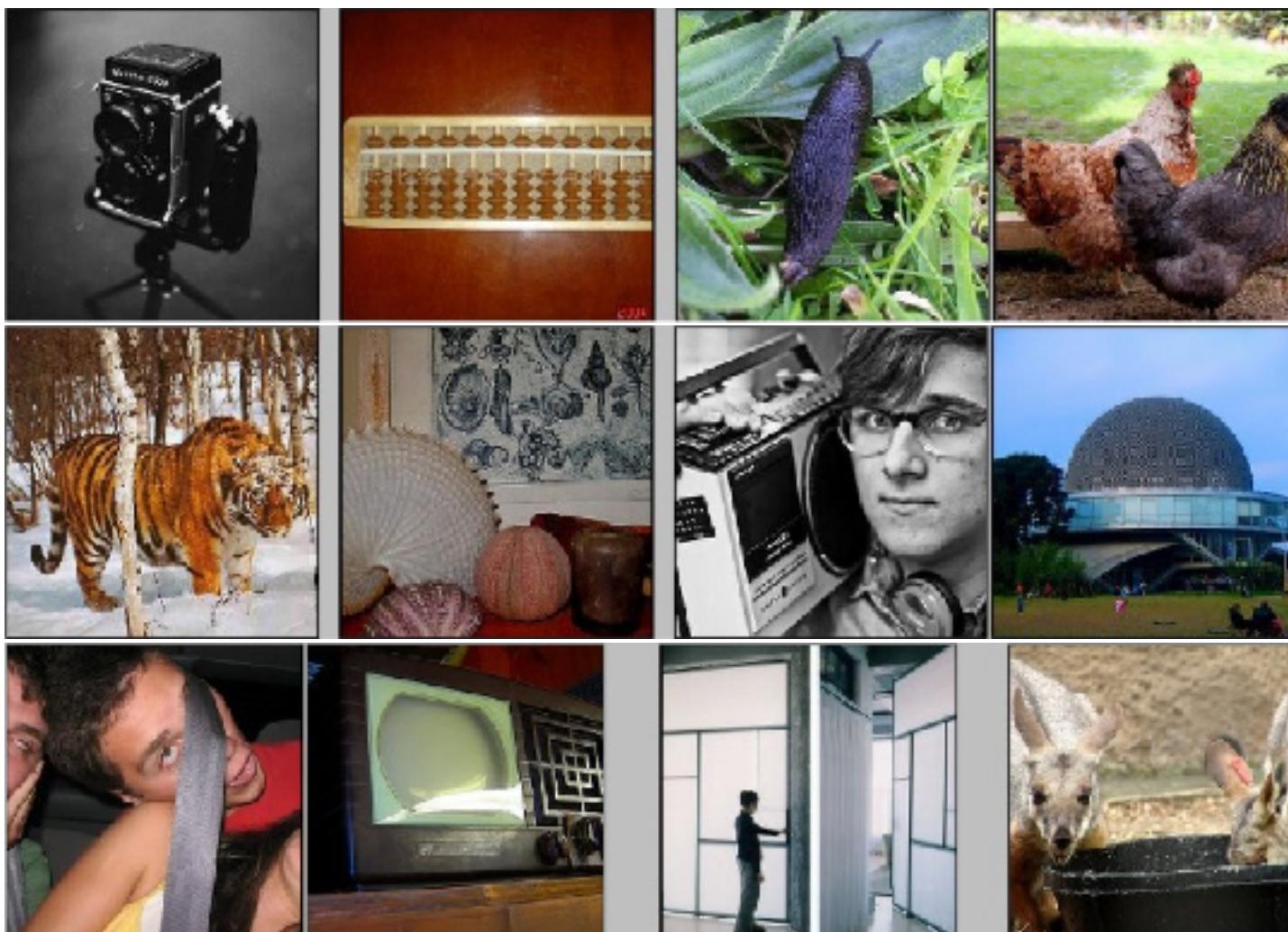


[LeCun, Bottou, Bengio, Haffner. "Gradient-based learning applied to document recognition", 1998]

- Neocognitron + backpropagation
- Average pooling
- Sigmoid or tanh nonlinearity
- Fully connected layers at the end
- Trained on MNIST digit dataset with 60K training examples



ImageNet Challenge



- ~14 million labeled images, 20k classes
- Images gathered from Internet
- Human labels via Amazon MTurk
- ImageNet Large-Scale Visual Recognition Challenge (ILSVRC):
1.2 million training images, 1000 classes

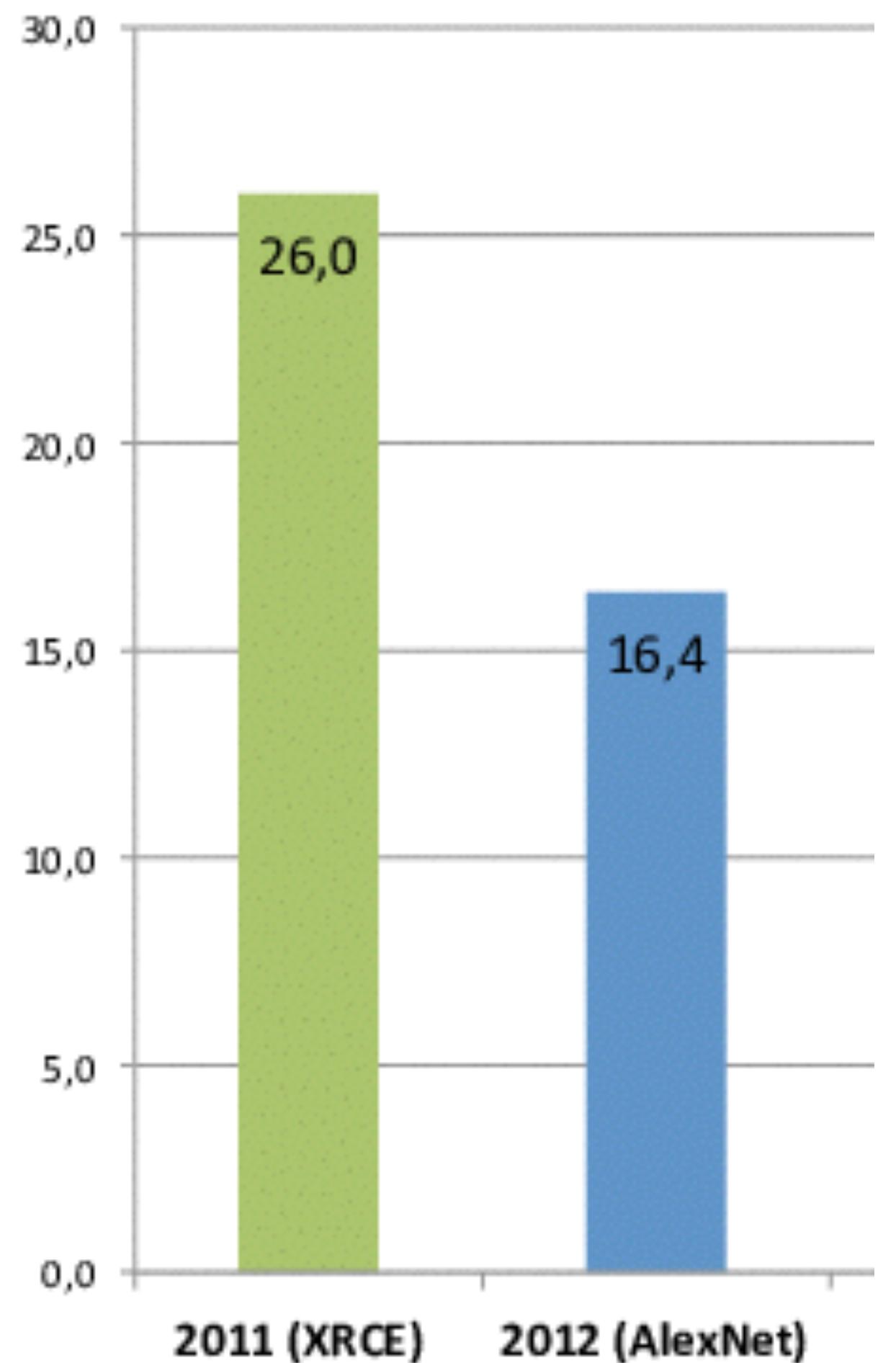
[Russakovsky, Deng, Su, Krause, Satheesh, Ma, Huang, Karpathy, Khosla, Bernstein, Berge, Fei-Fei, "ImageNet Large Scale Visual Recognition Challenge", 2015]

More recent networks

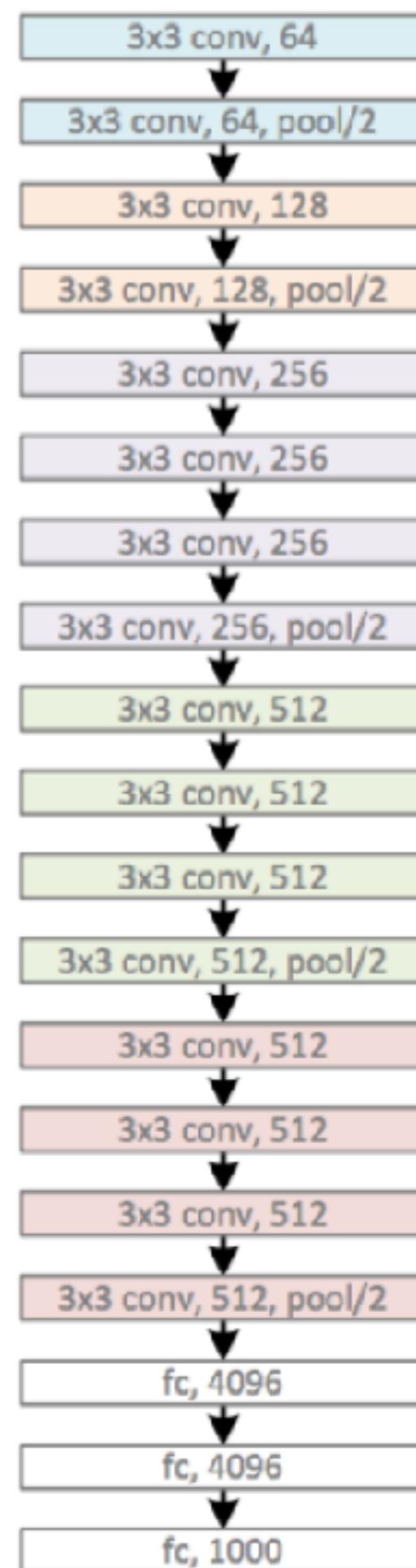
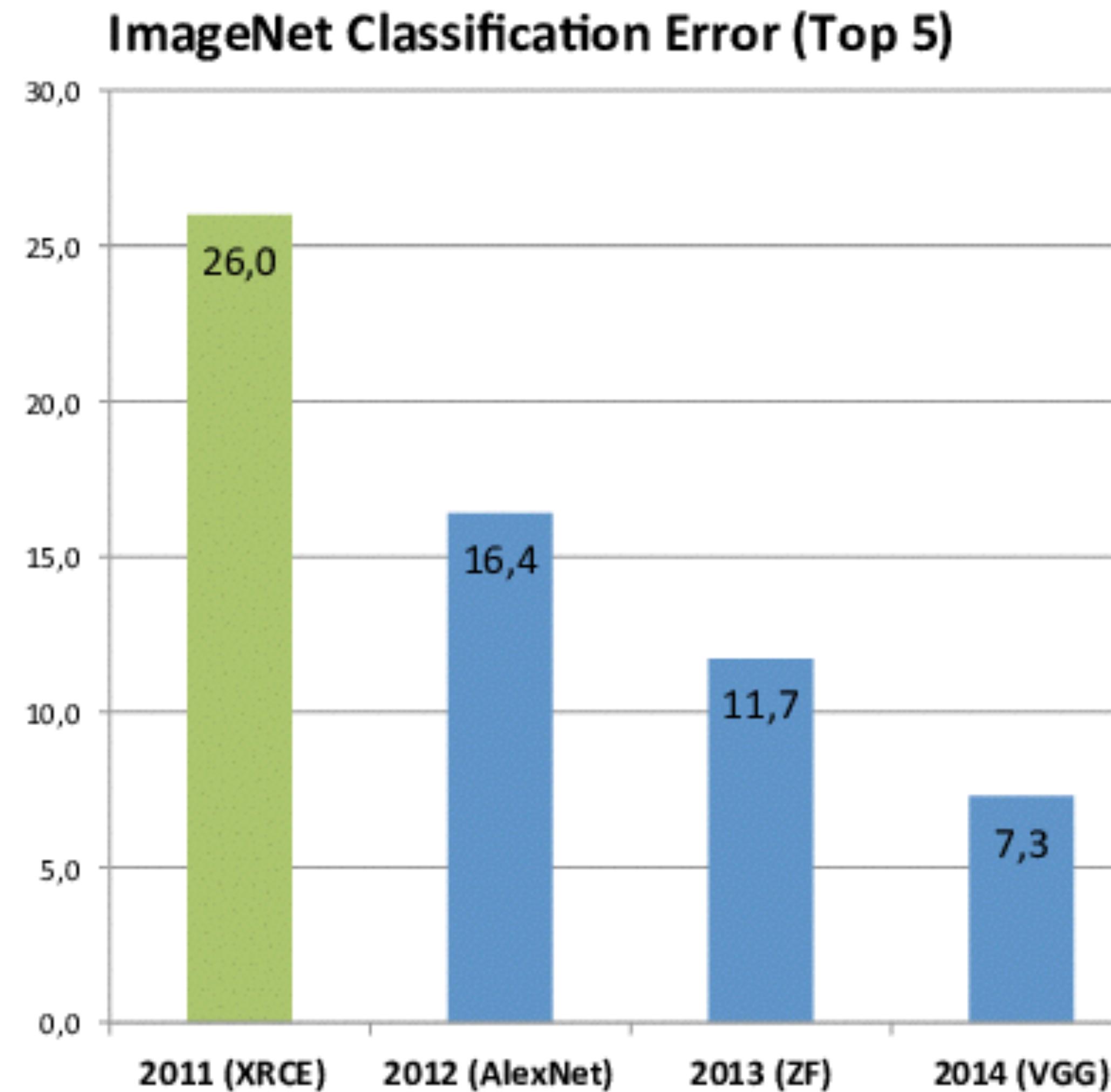
And advances that make them work:

- Chaining small filters
- Residual layers
- Normalization

ImageNet Classification Error (Top 5)



2014: VGG
16 conv. layers

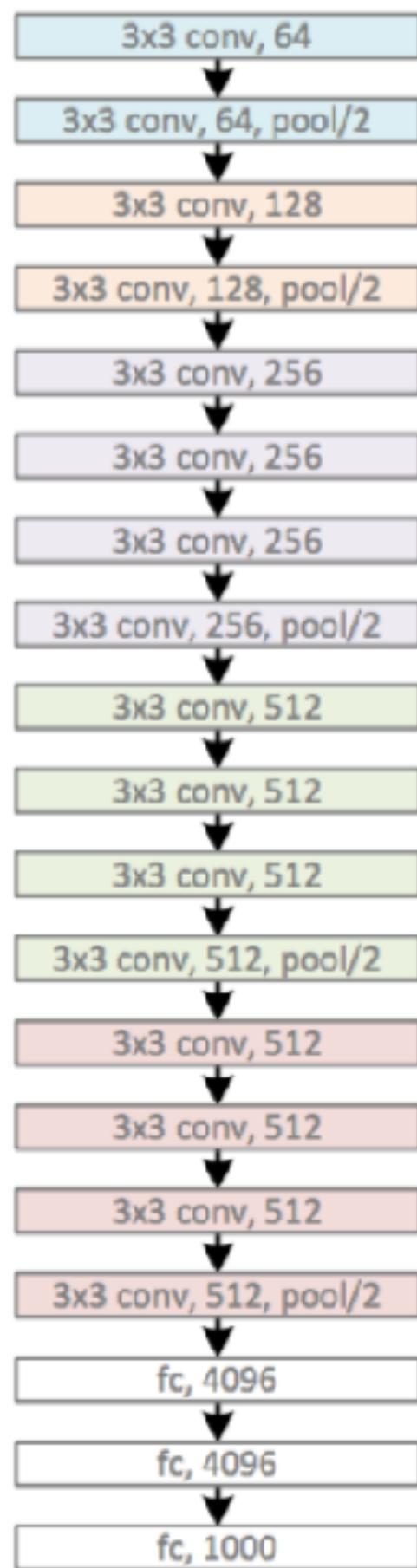


Error: 7.3%

[Simonyan & Zisserman: Very Deep Convolutional Networks
for Large-Scale Image Recognition, ICLR 2015]

VGG-Net [Simonyan & Zisserman, 2015]

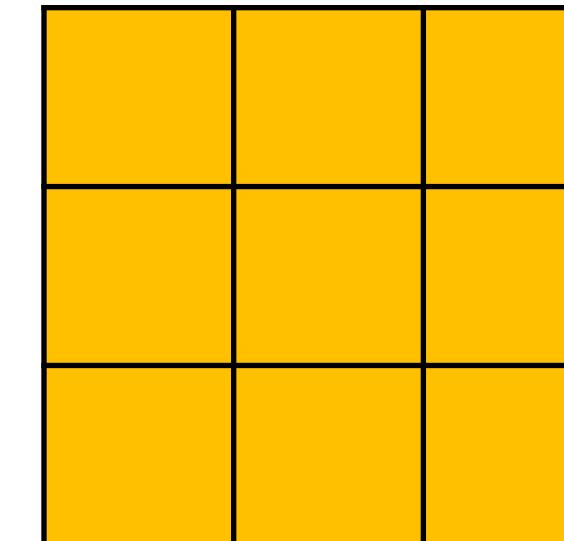
2014: VGG
16 conv. layers



Error: 7.3%

Main developments

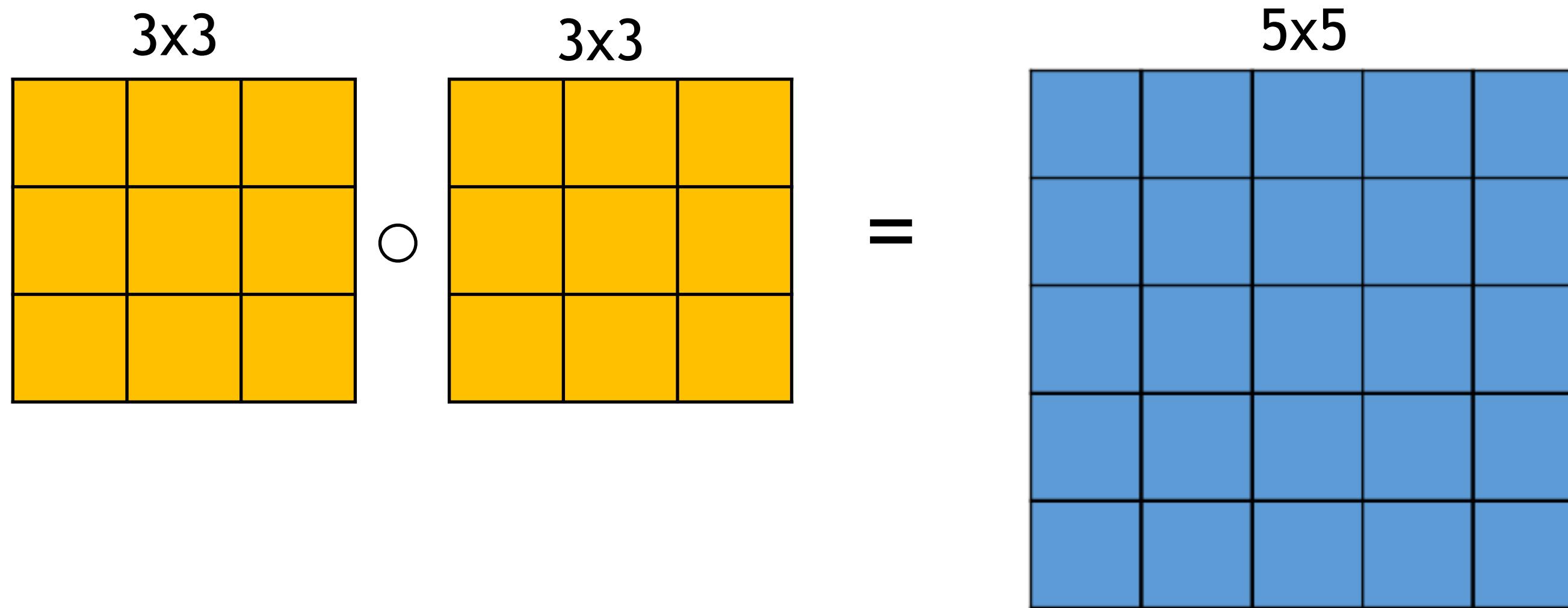
- Small convolutional kernels: only 3x3



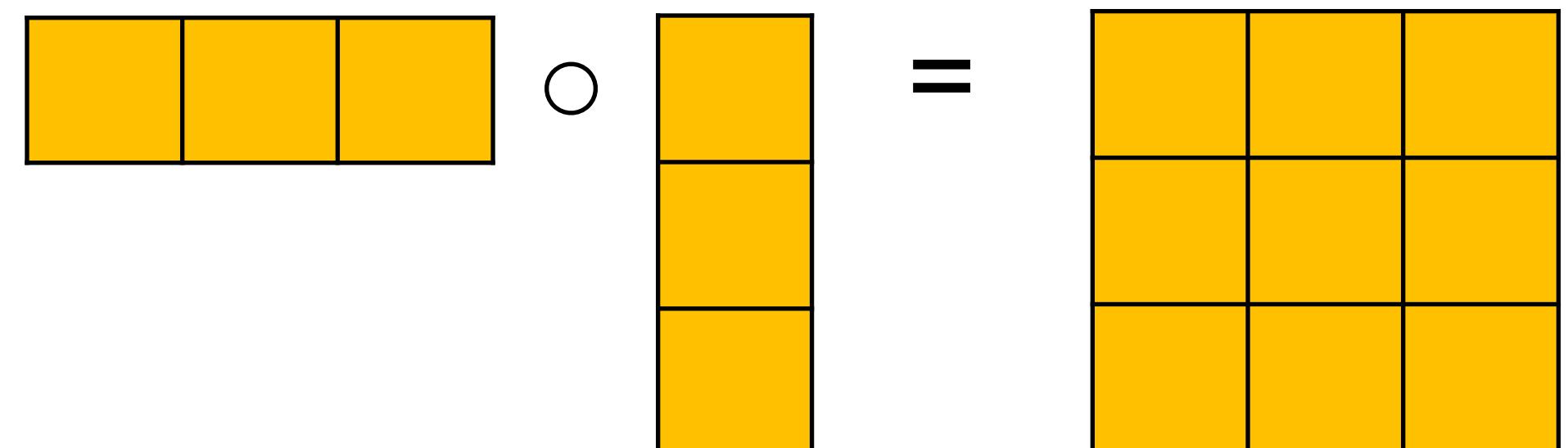
- Increased depth (5 -> 16/19 layers)

Other tricks for designing convolutional nets

Chaining convolutions

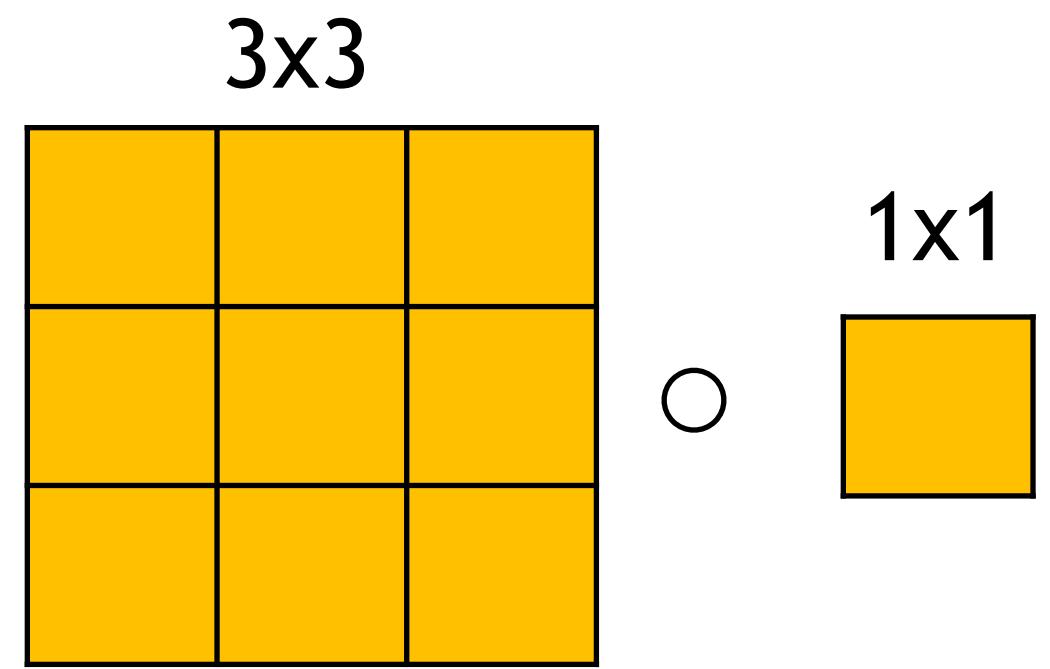


Nonlinearity in between.
25 coefficients, but only
18 degrees of freedom



9 coefficients, but only
6 degrees of freedom.
Less common.

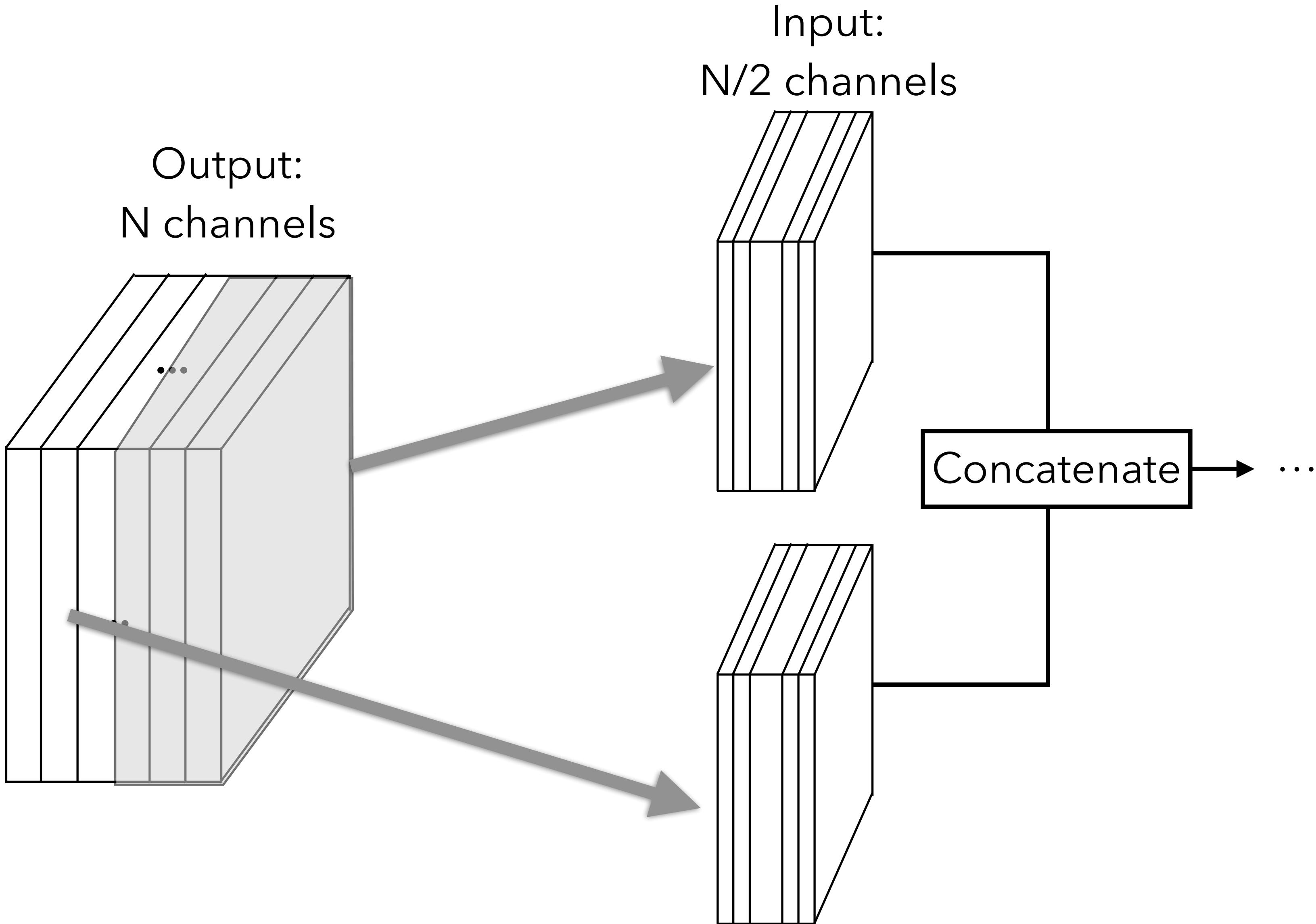
1×1 convolutions



Why do this?

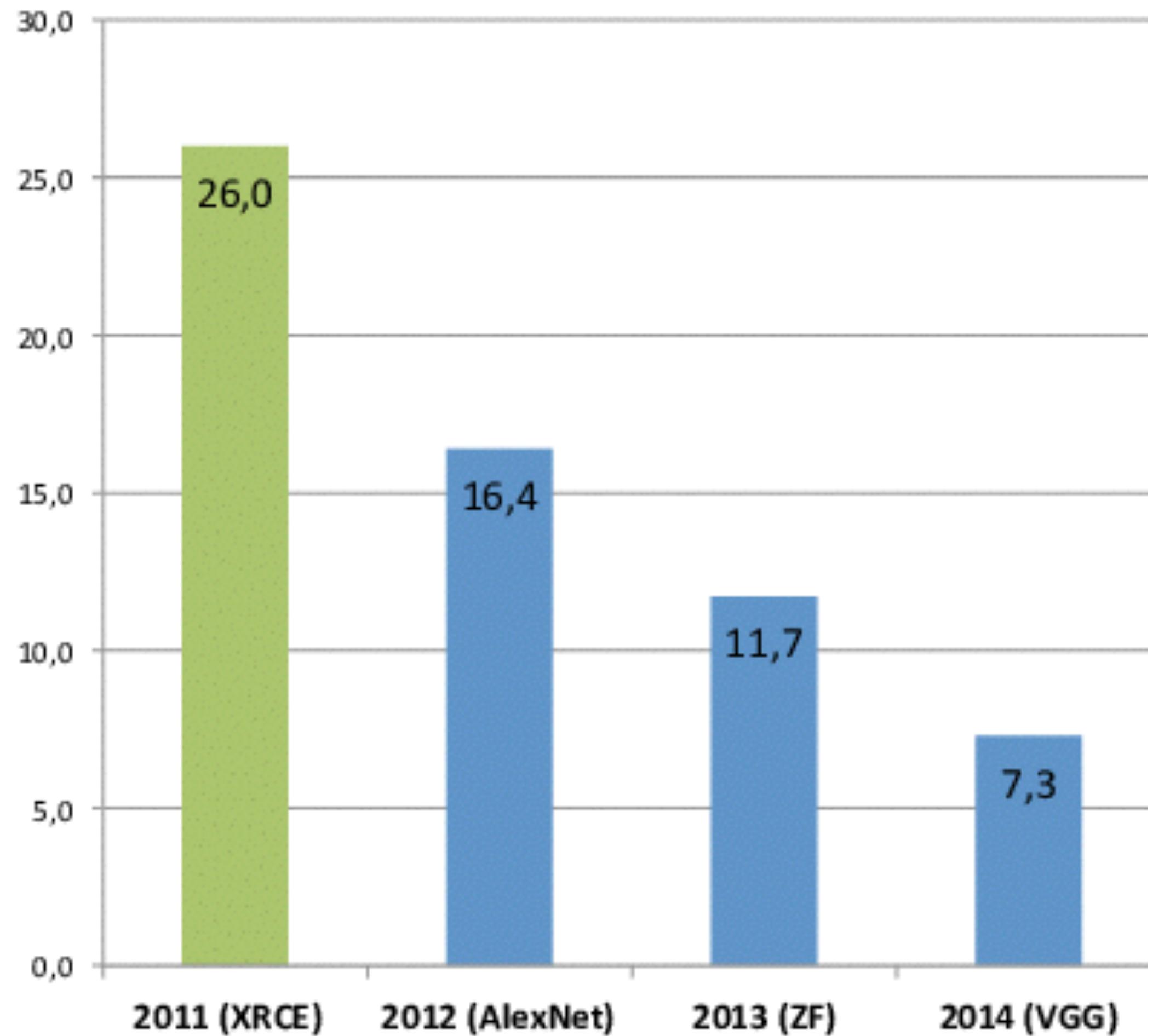
(nonlinearity in between)

Grouped Convolutions

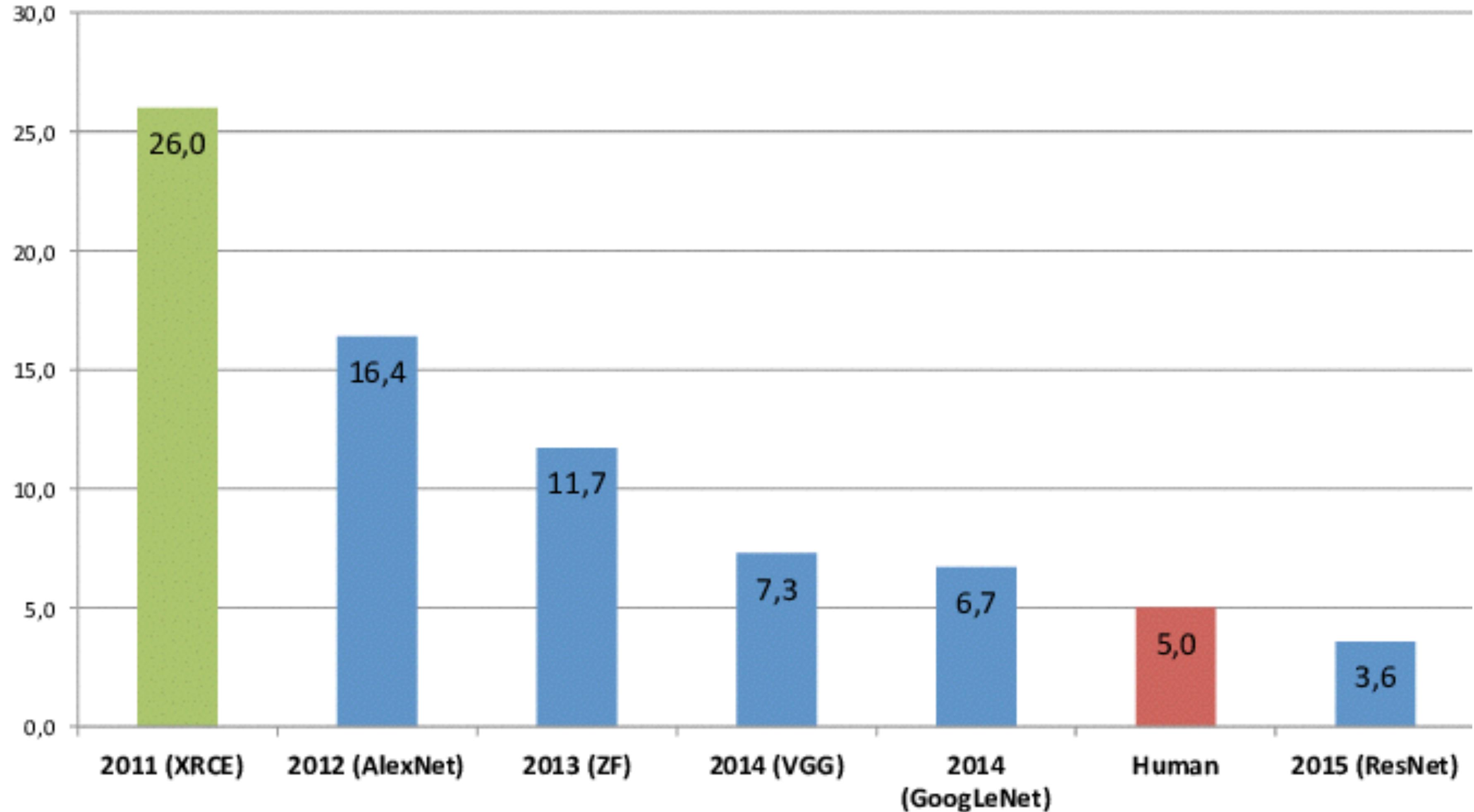


Split channels into N groups, and process separately with N convolution layers.

ImageNet Classification Error (Top 5)



ImageNet Classification Error (Top 5)



2016: ResNet
>100 conv. layers

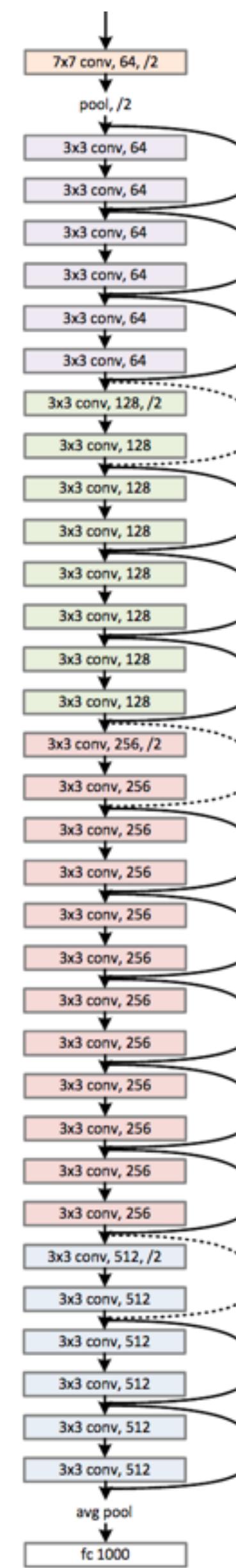
Error: 3.6%

2016: ResNet
>100 conv. layers



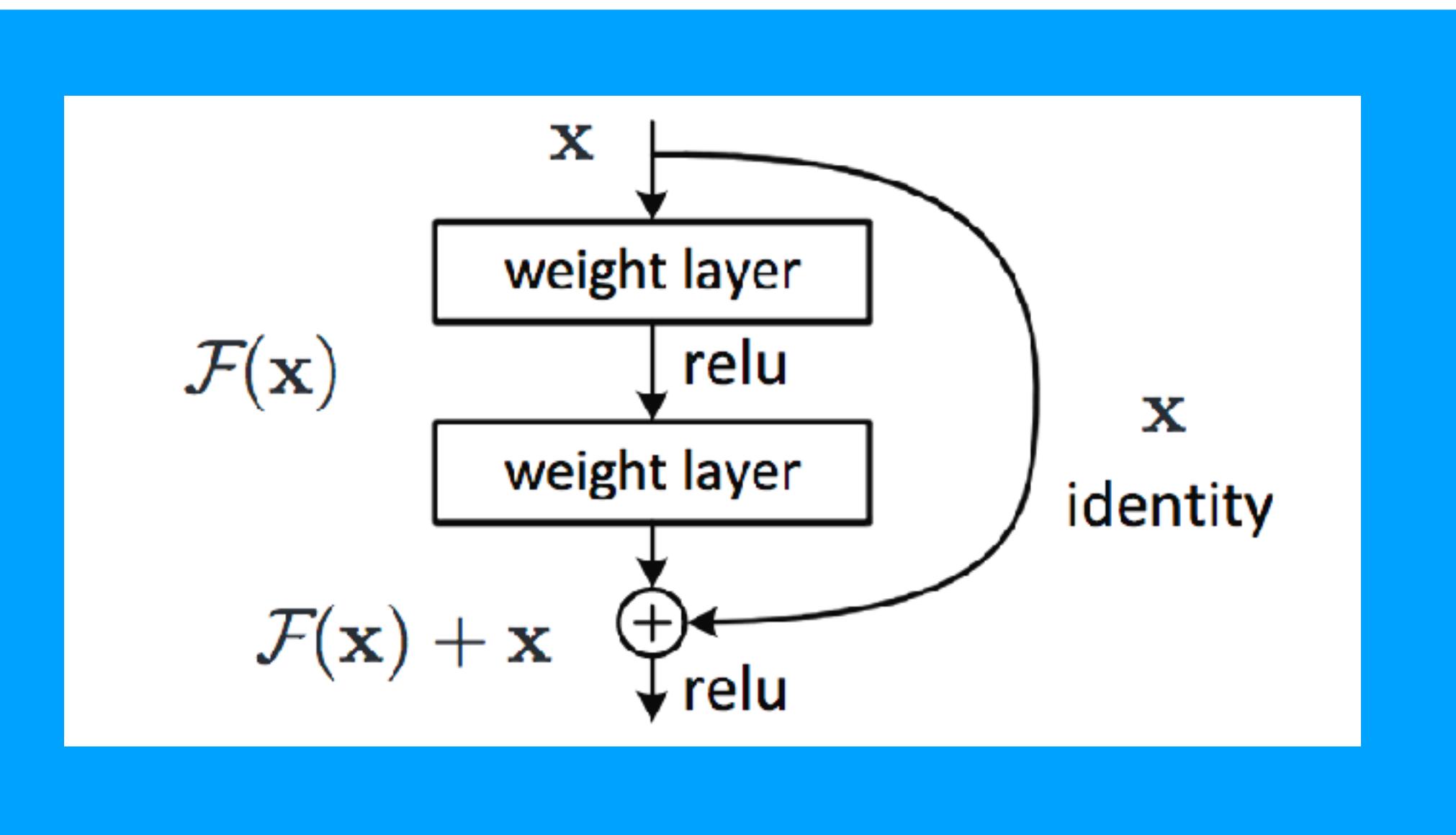
Error: 3.6%

ResNet [He et al, 2016]



Main developments

- Increased depth possible through residual blocks



Residual Blocks

Problem: Hard to train very deep nets (50+ layers). This is an optimization issue, not overfitting: shallow models often get higher *training* accuracy than deep ones!

Idea: Make it easy to represent for the network to implement the identity.

Normal convolution + relu:

$$x_{i+1} = \text{relu}(x_i \circ f)$$

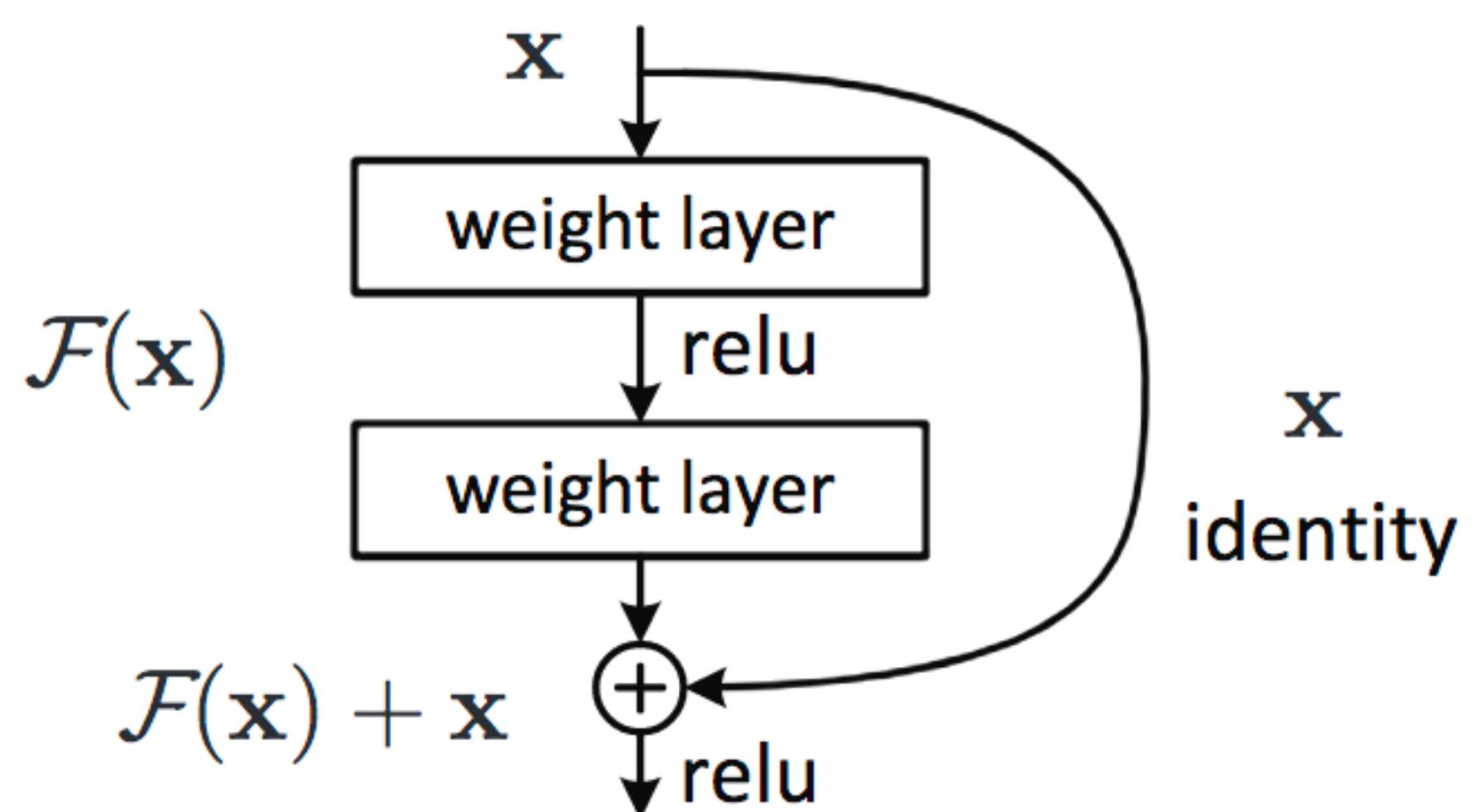
Residual connection

$$x_{i+1} = \text{relu}((x_i \circ f) + x_i)$$

More generally: do multiple convolutions (with nonlinearities) before summing:

$$x_{i+1} = \text{relu}(\mathcal{F}(x_i) + x_i)$$

Residual Blocks



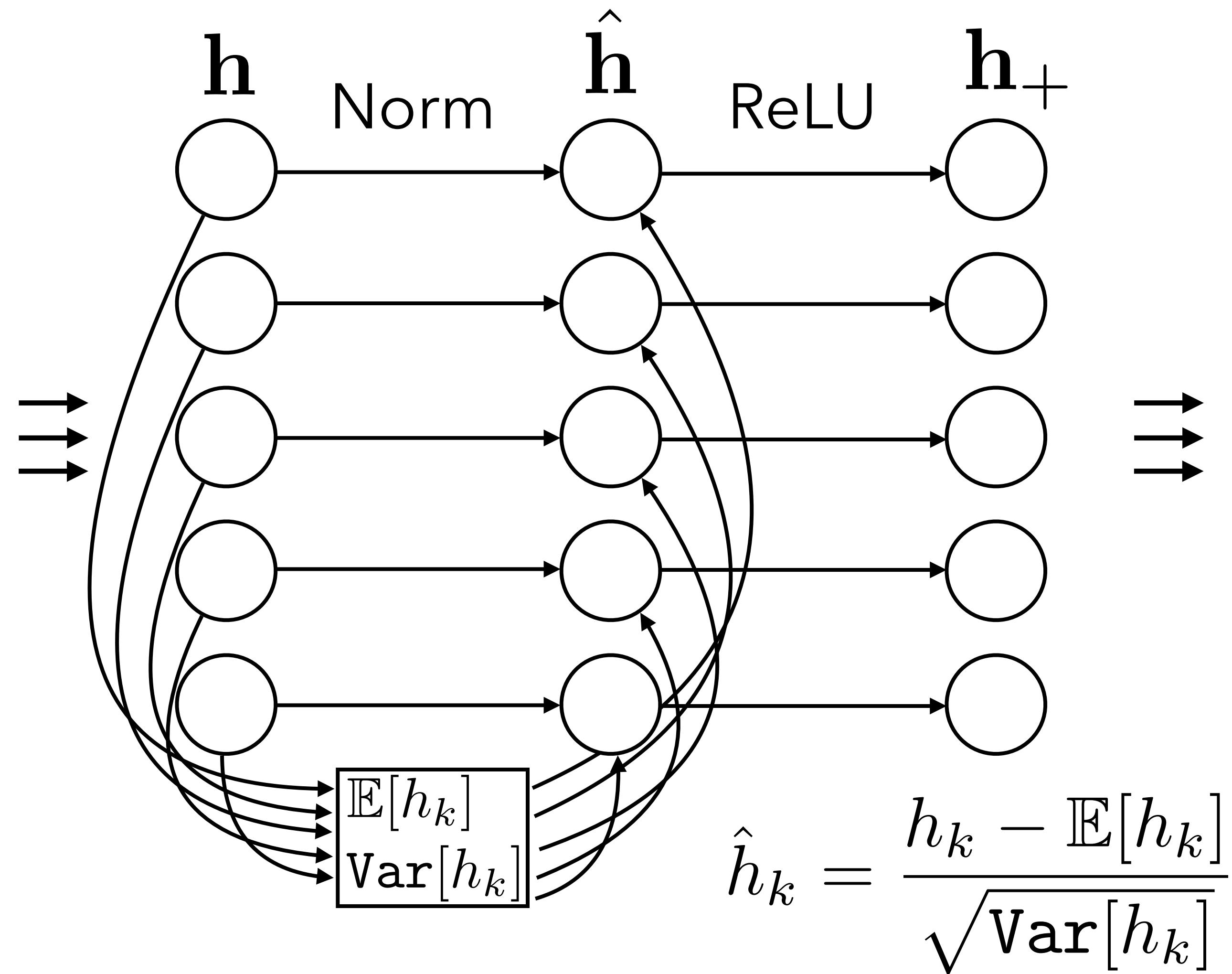
Why do they work?

- Gradients can propagate faster (via the identity mapping)
- Within each block, only small residuals have to be learned

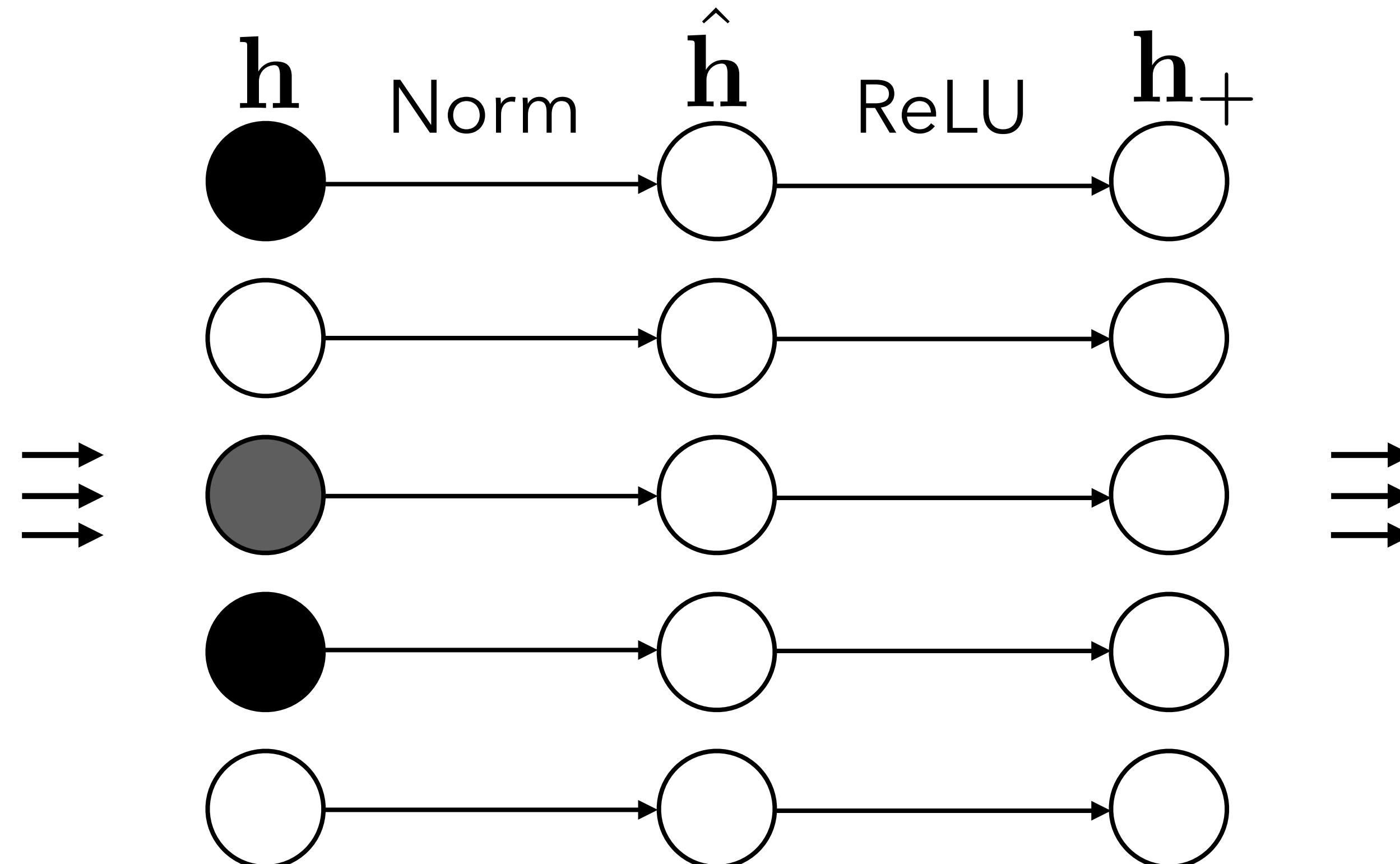
Normalization layers

- Standardize activations by subtracting mean and dividing by standard deviation (averaged over all spatial locations).
- This provides a constant “interface” for later layers of the networks. Ensures that the previous layer will have unit variance and zero mean.
- Obtains invariance to mean and variance.
- Can allow you to train with larger learning rate and significantly speed up training!

Normalization layers

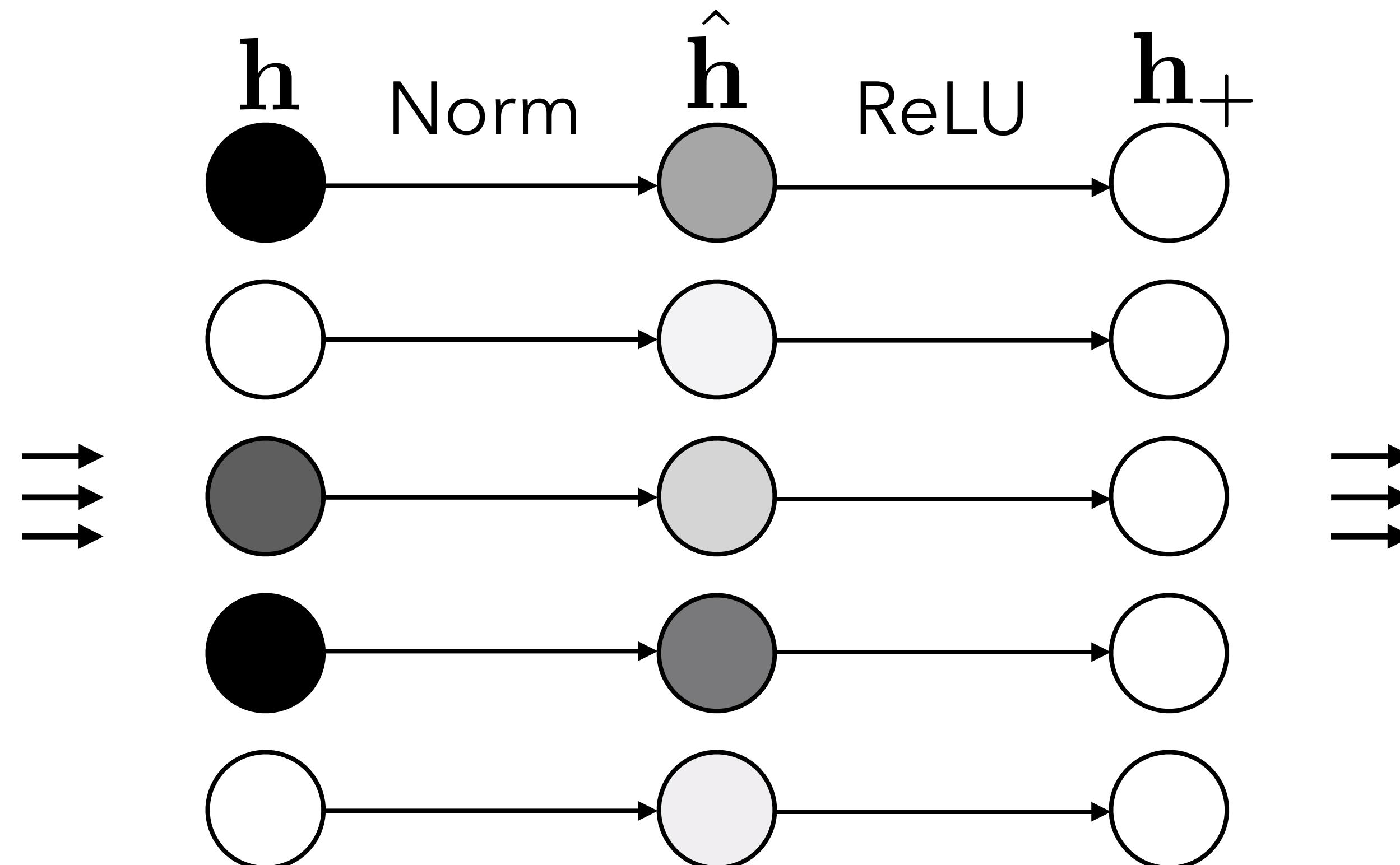


Normalization layers



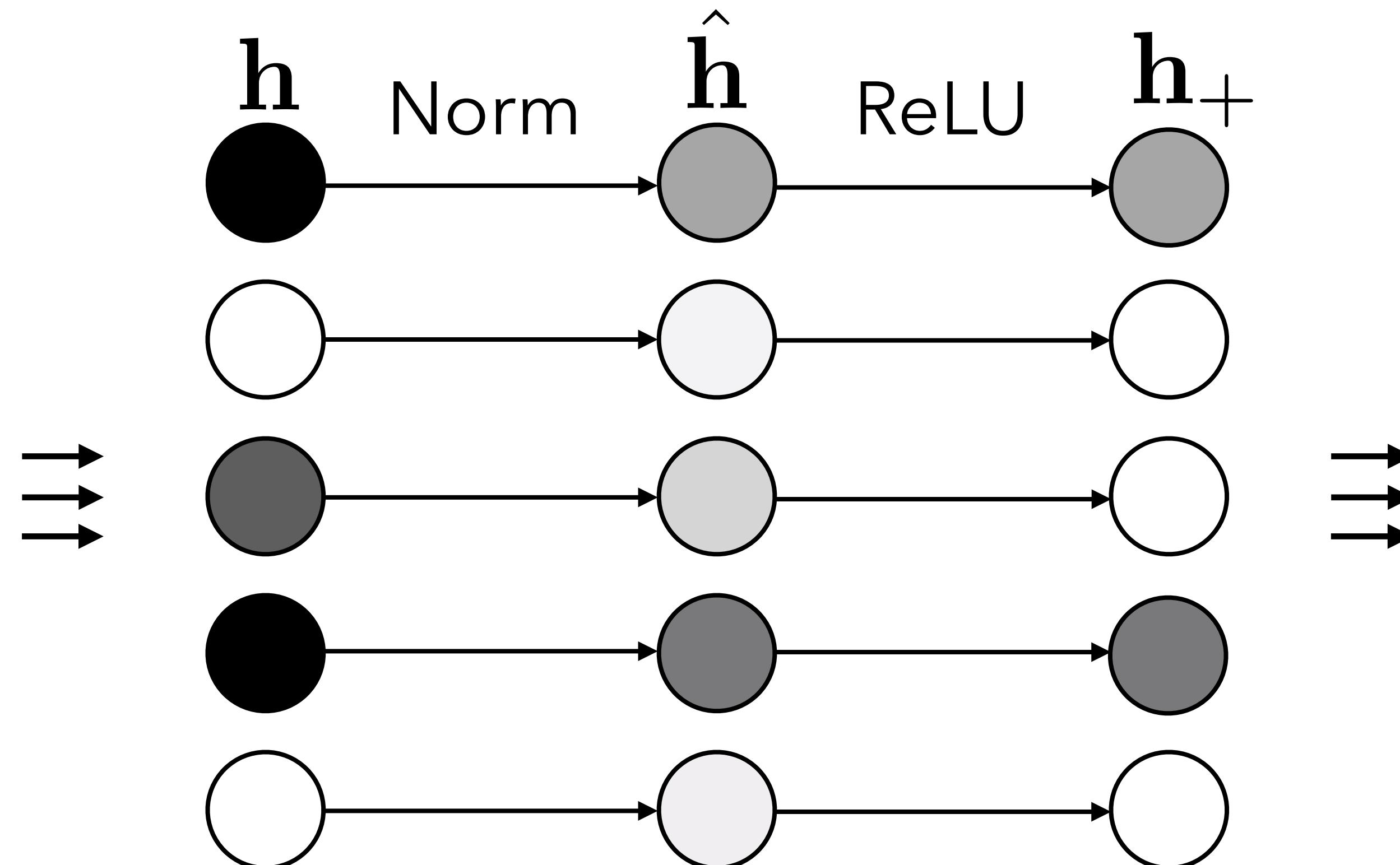
$$\hat{h}_k = \frac{h_k - \mathbb{E}[h_k]}{\sqrt{\text{Var}[h_k]}}$$

Normalization layers



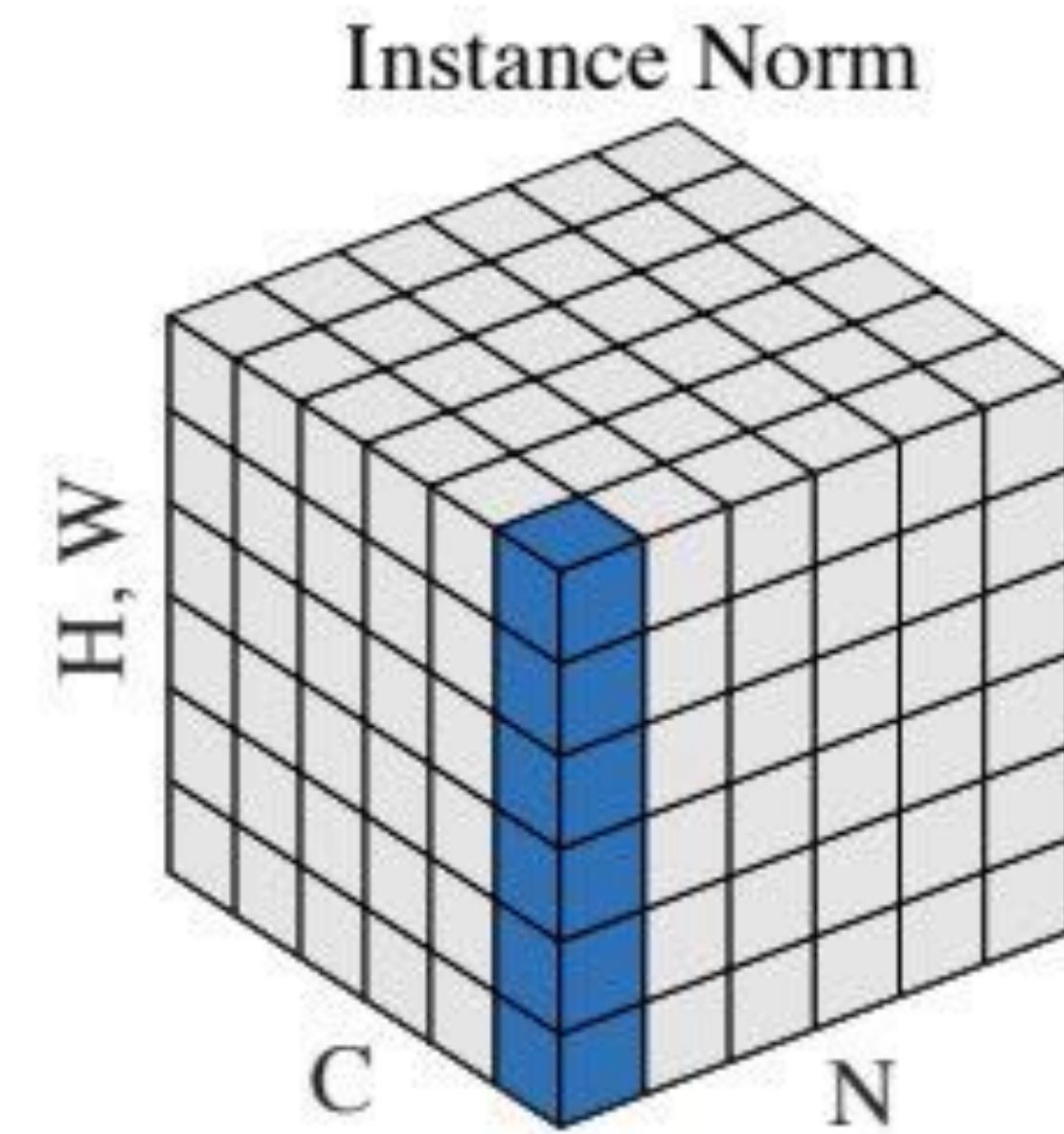
$$\hat{h}_k = \frac{h_k - \mathbb{E}[h_k]}{\sqrt{\text{Var}[h_k]}}$$

Normalization layers



$$\hat{h}_k = \frac{h_k - \mathbb{E}[h_k]}{\sqrt{\text{Var}[h_k]}}$$

Instance normalization



Filter value in a CNN layer

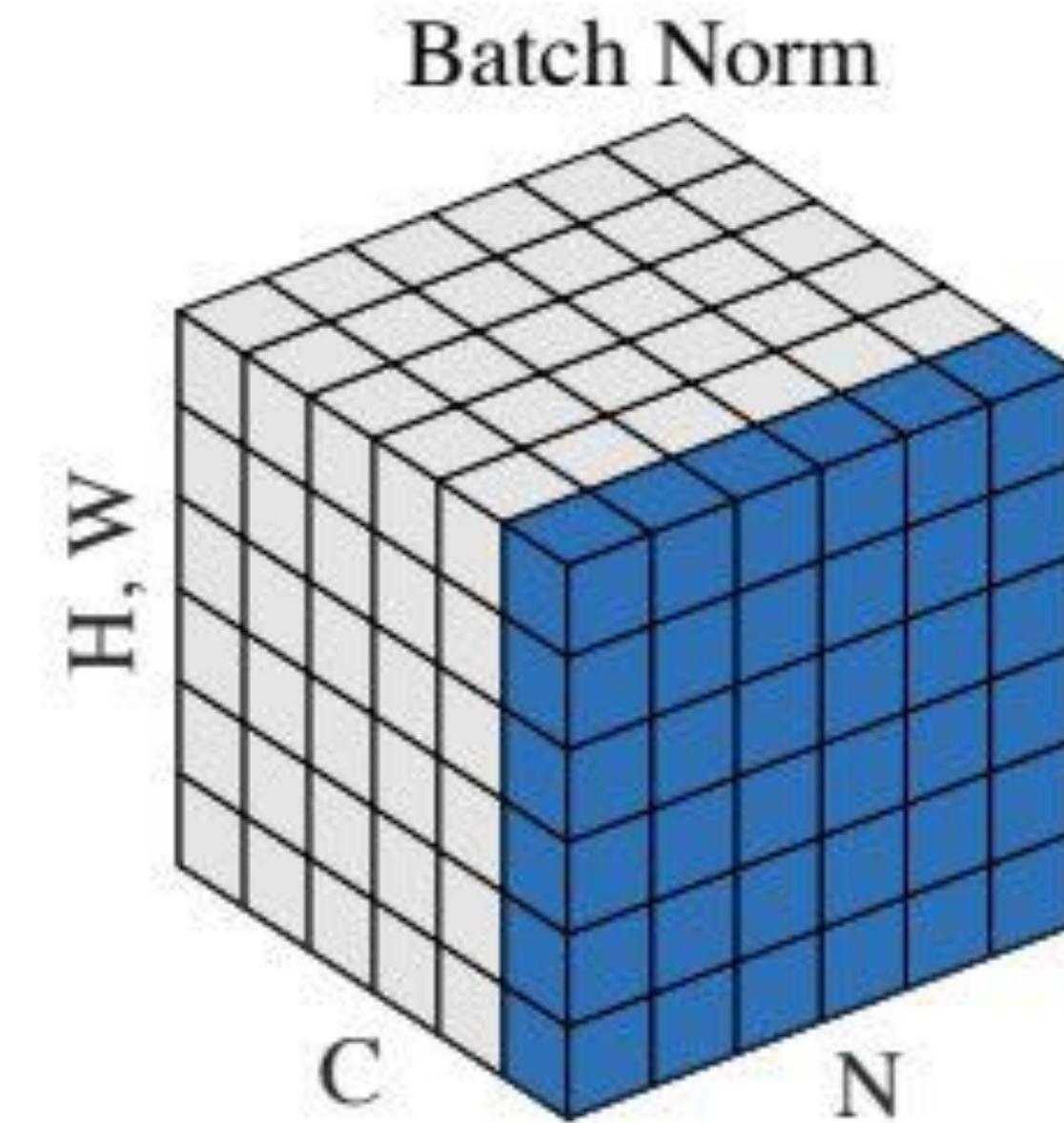
$$\hat{h}_k = \frac{h_k - \mathbb{E}[h_k]}{\sqrt{\text{Var}[h_k]}}$$

Average and standard deviation each filter response, estimated over **instance**.

Normalize a single hidden unit's activations to be mean 0, standard deviation 1.

[Figure from Wu & He, 2018][Ulyanov et al., 2015]

Batch normalization



Filter value in a CNN layer

$$\hat{h}_k = \frac{h_k - \mathbb{E}[h_k]}{\sqrt{\text{Var}[h_k]}}$$

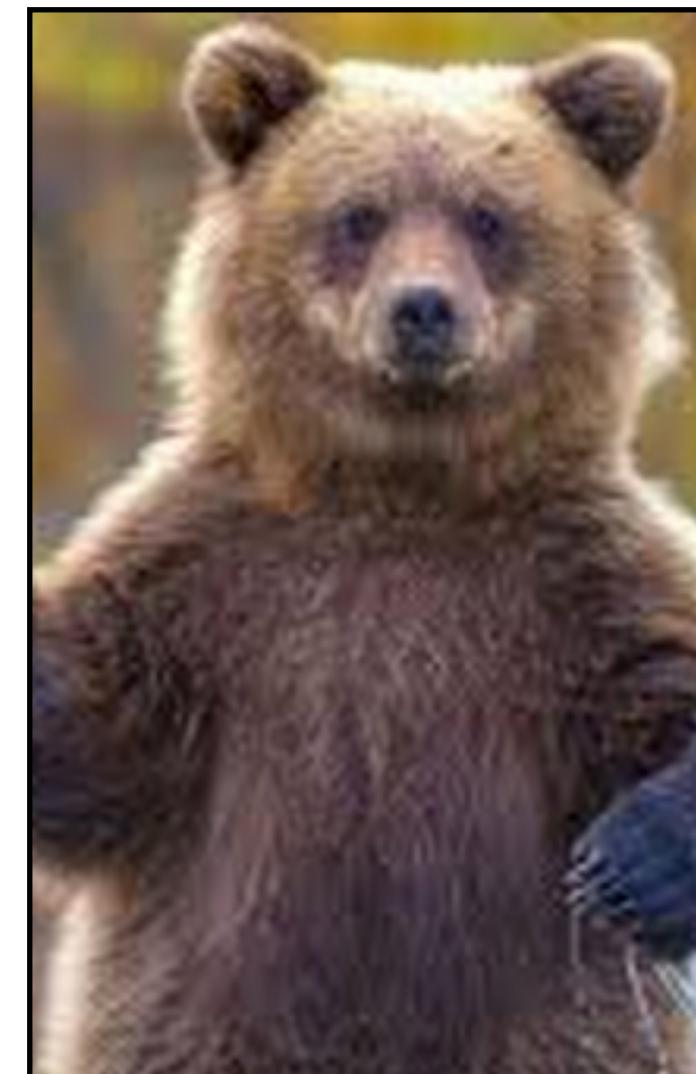
Average and standard deviation each filter response, estimated over **whole batch**.

Normalize a single hidden unit's activations to be mean 0, standard deviation 1. At test time, remember the mean and standard deviation seen during training.

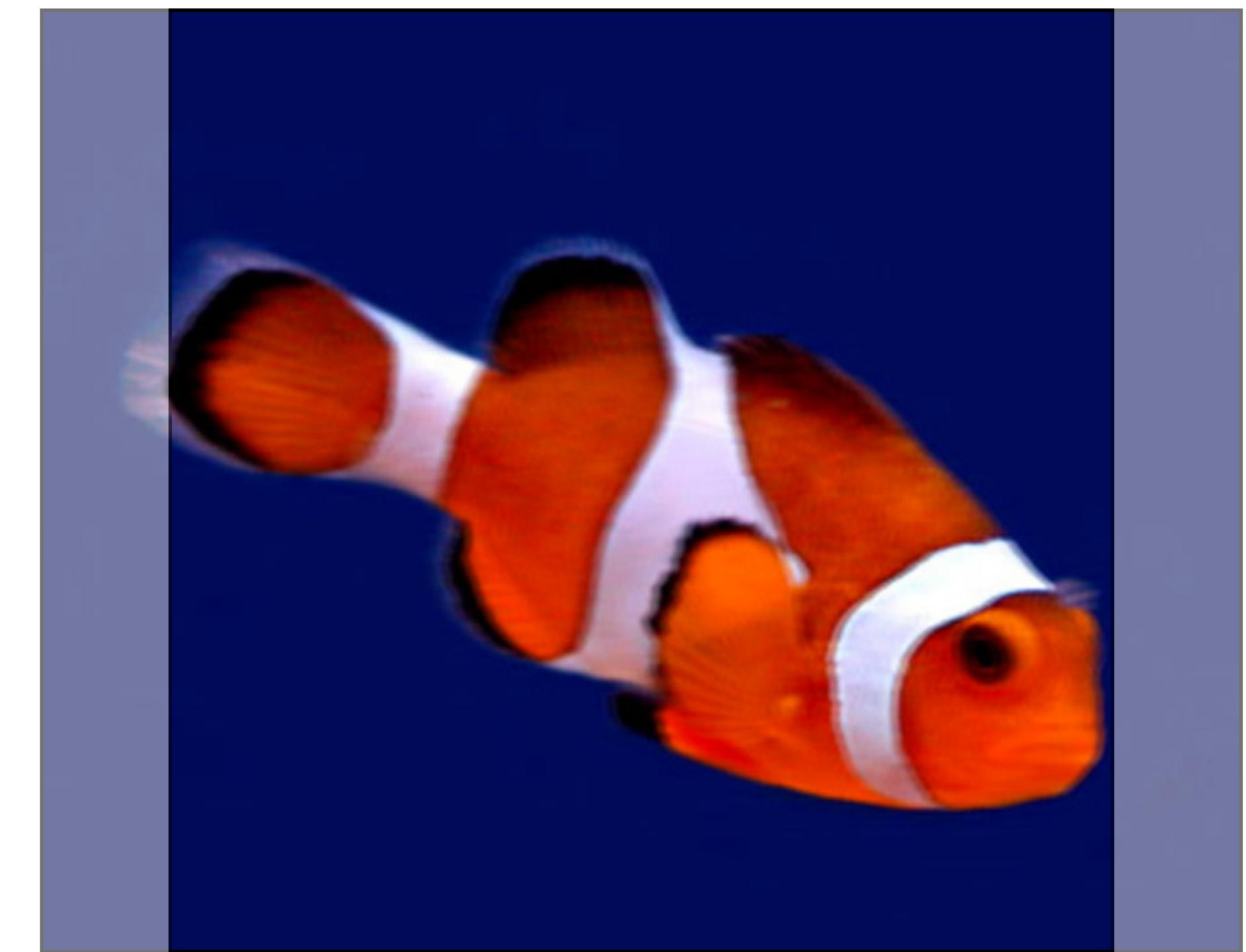
[Figure from Wu & He, 2018] [Ioffe & Szegedy, 2015]

A few practical issues

Dealing with rectangular images



Rectangular images

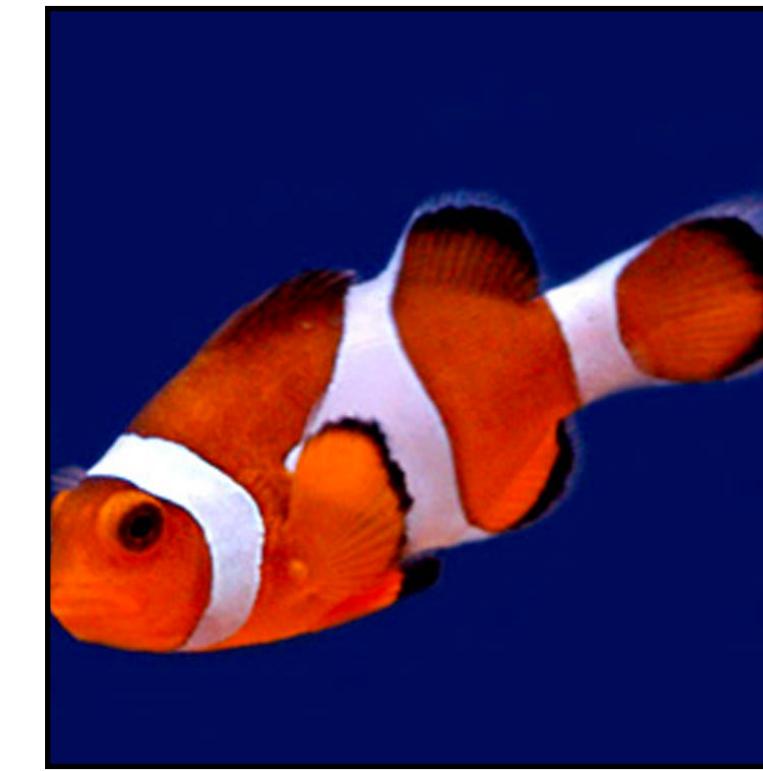


Resize, then take square
crop from center

Training with data augmentation



Original image



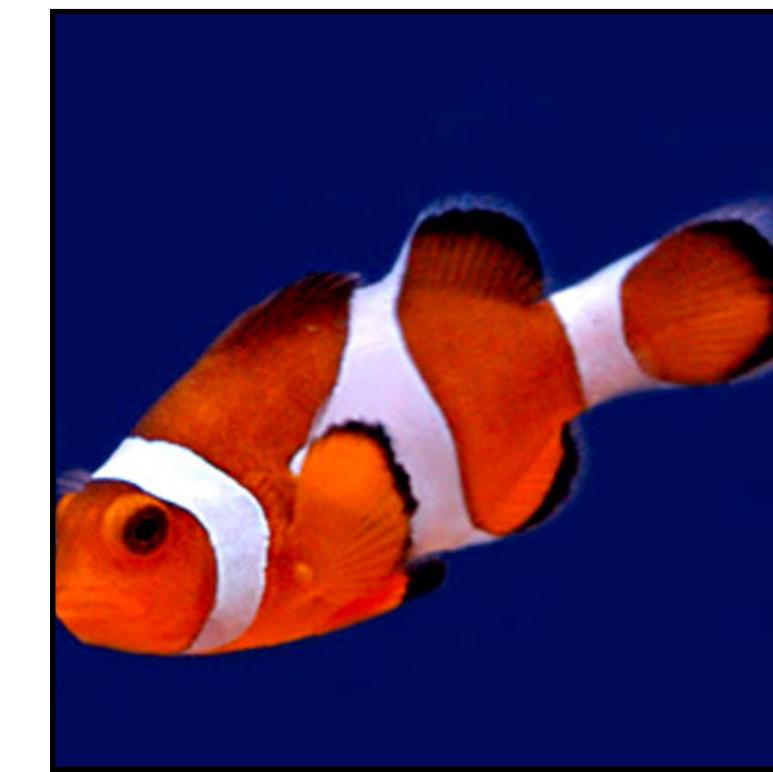
Flipping

- Less susceptible to overfitting.
- Improves performance by simulating examples

Training with data augmentation



Original image



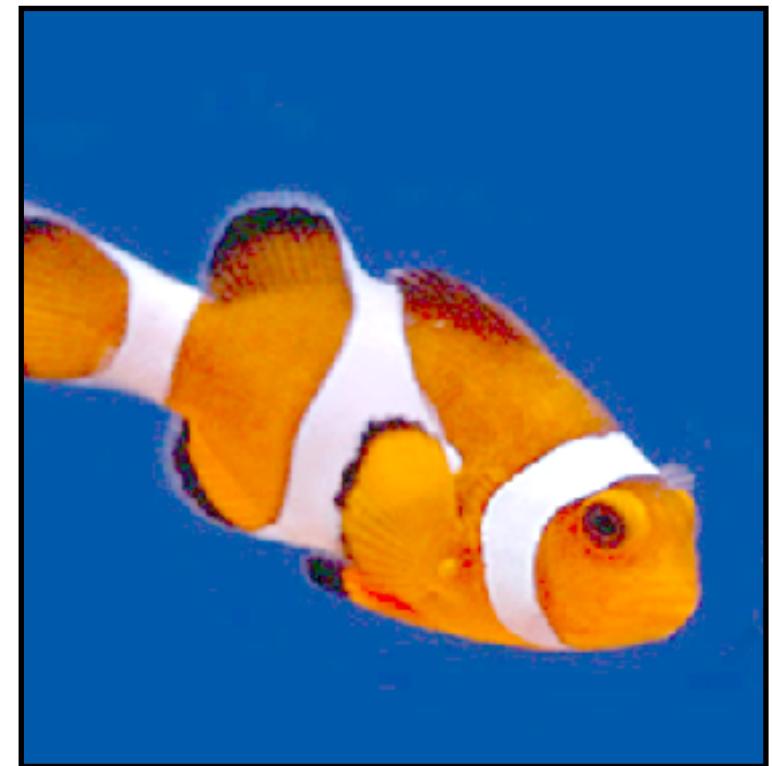
Flipping



Cropping



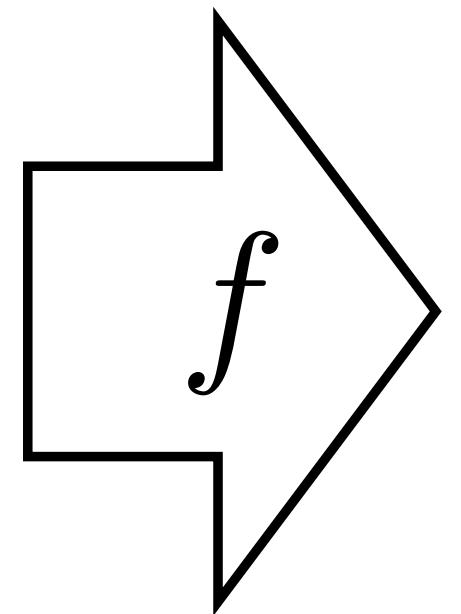
Scaling



Color jittering

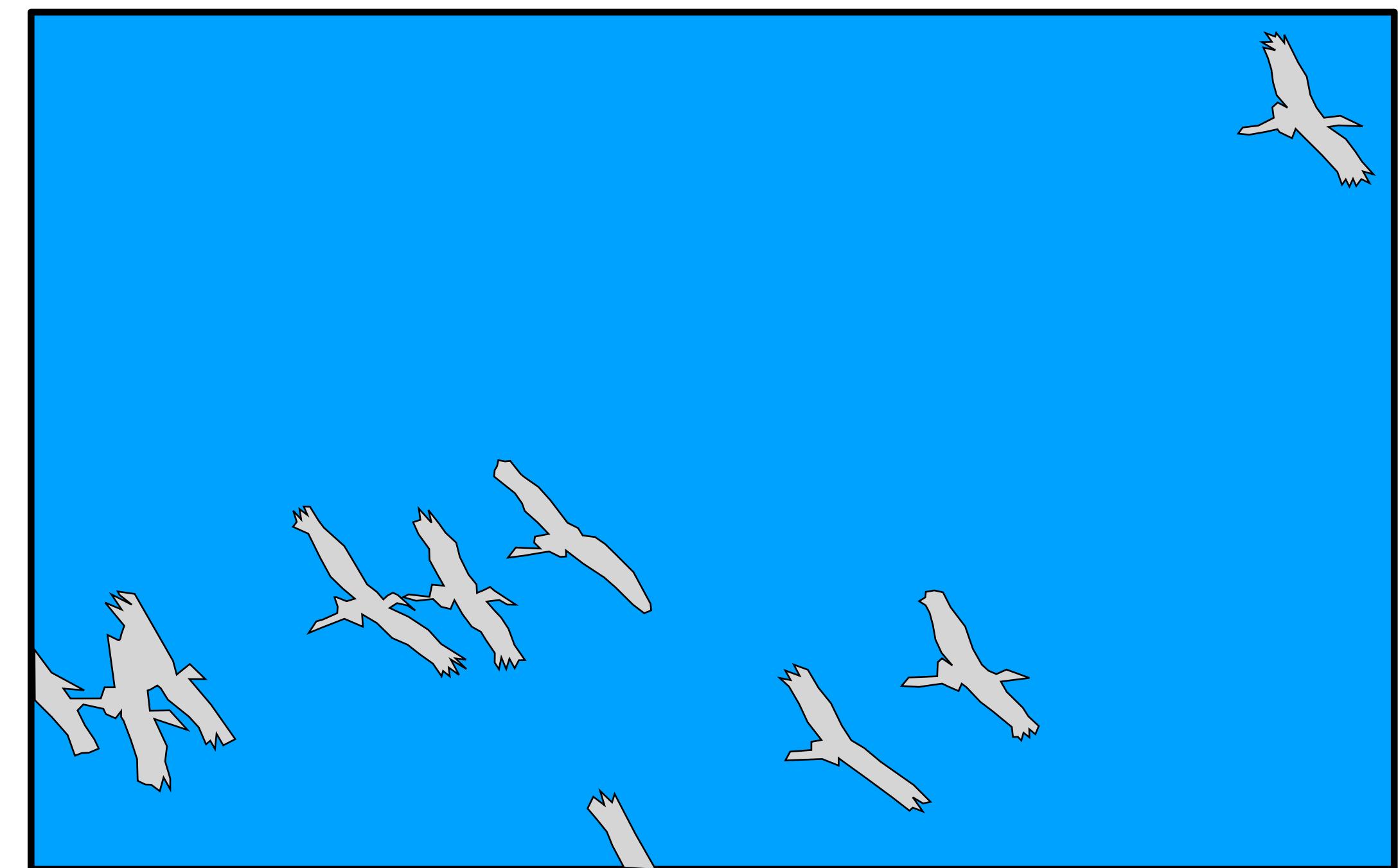
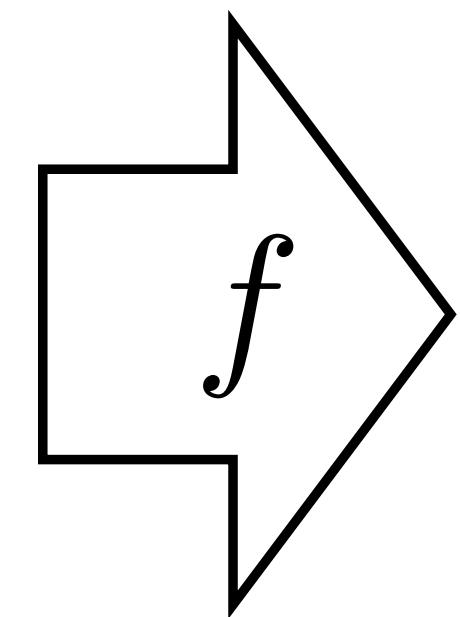
Beyond image labeling

Object recognition: what objects are in the image?



“Birds”

Semantic segmentation



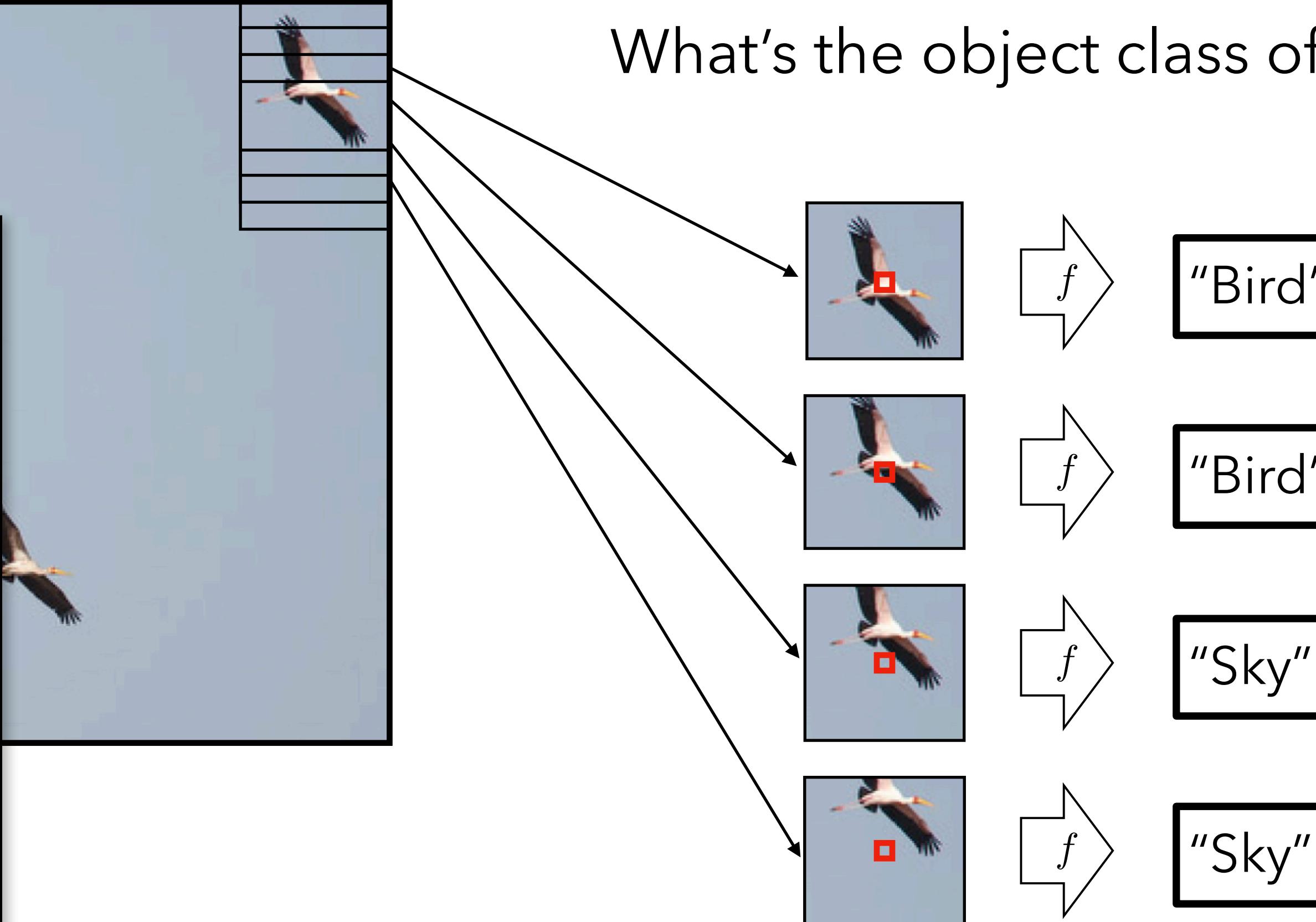
(Colors represent categories)

General technique: predict something at every pixel!⁷⁸

Idea #1: Independently classify windows

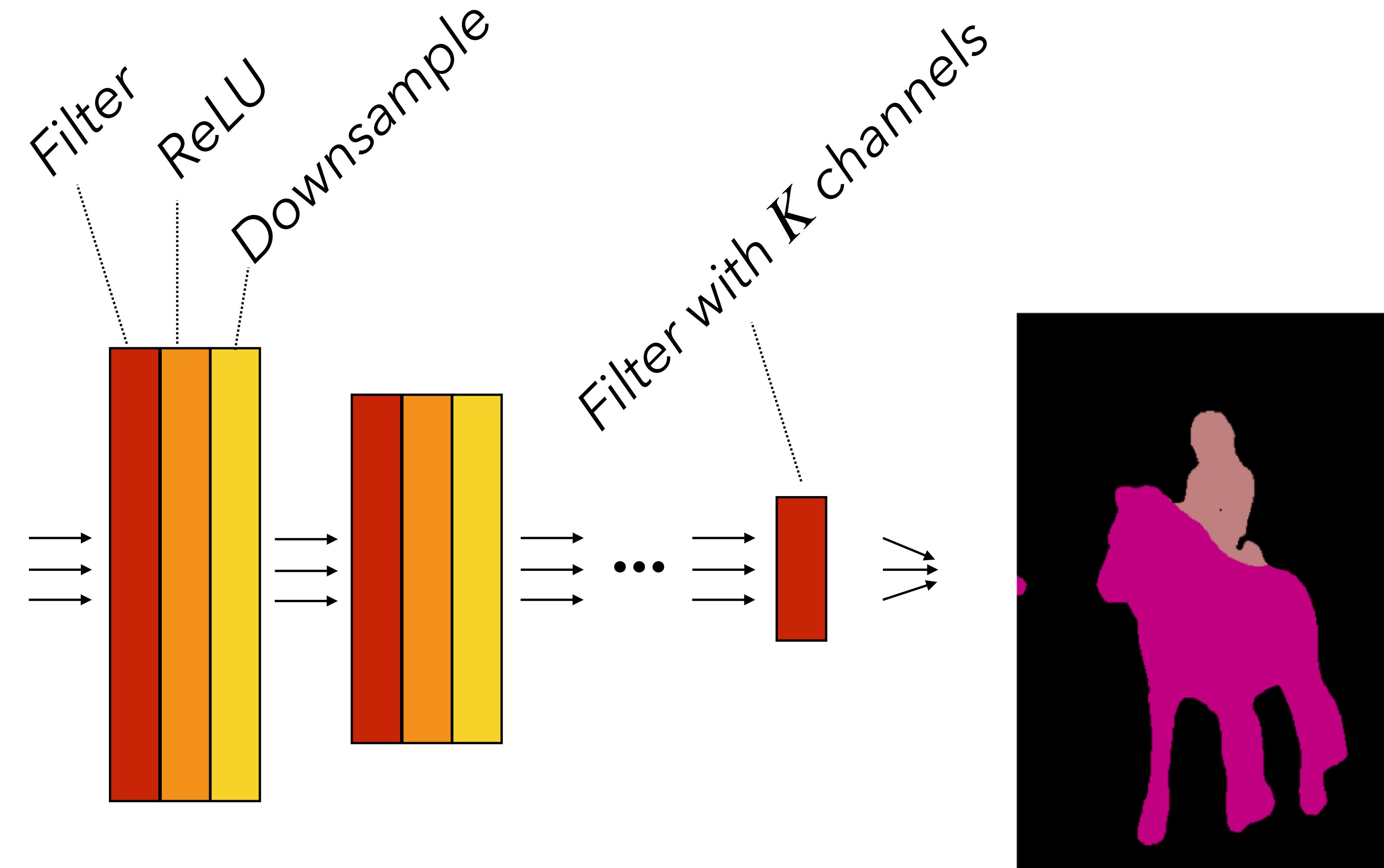
What's the object class of the center pixel?

| <i>Training data</i> | |
|----------------------|--------|
| x | y |
| { | "Bird" |
| { | "Bird" |
| { | "Sky" |
| : | |



Idea #2: Fully convolutional networks

Fully convolutional network



Classification problem with K classes

Idea #3: Dilated convolutions

Dilated convolutions

3x3

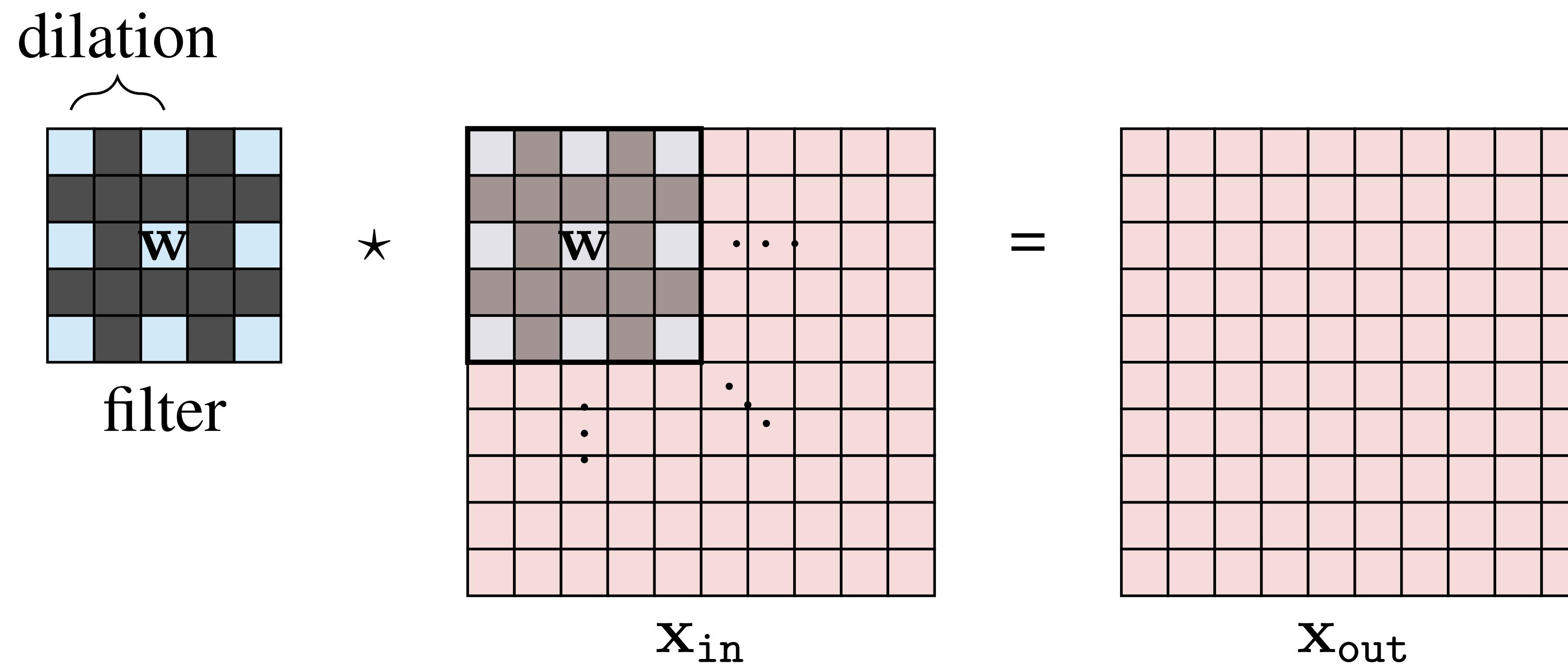
| | | |
|----------|----------|----------|
| a | b | c |
| d | e | f |
| g | h | i |

5x5

| | | | | |
|----------|---|----------|---|----------|
| a | 0 | b | 0 | c |
| 0 | 0 | 0 | 0 | 0 |
| d | 0 | e | 0 | f |
| 0 | 0 | 0 | 0 | 0 |
| g | 0 | h | 0 | i |

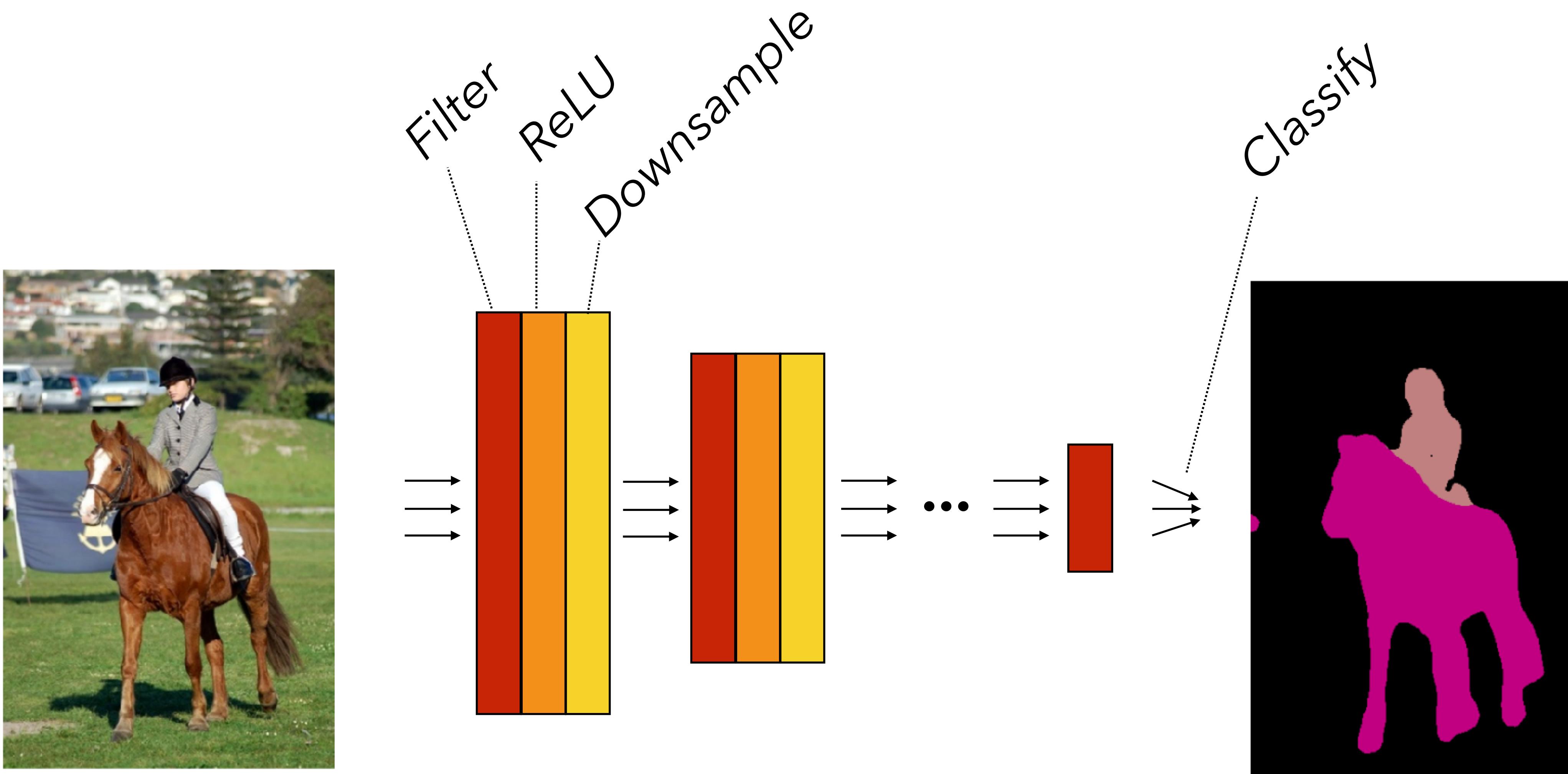
- Alternative to pooling that preserves input size
- 9 degrees of freedom
- 5x5 receptive field

Dilated filter

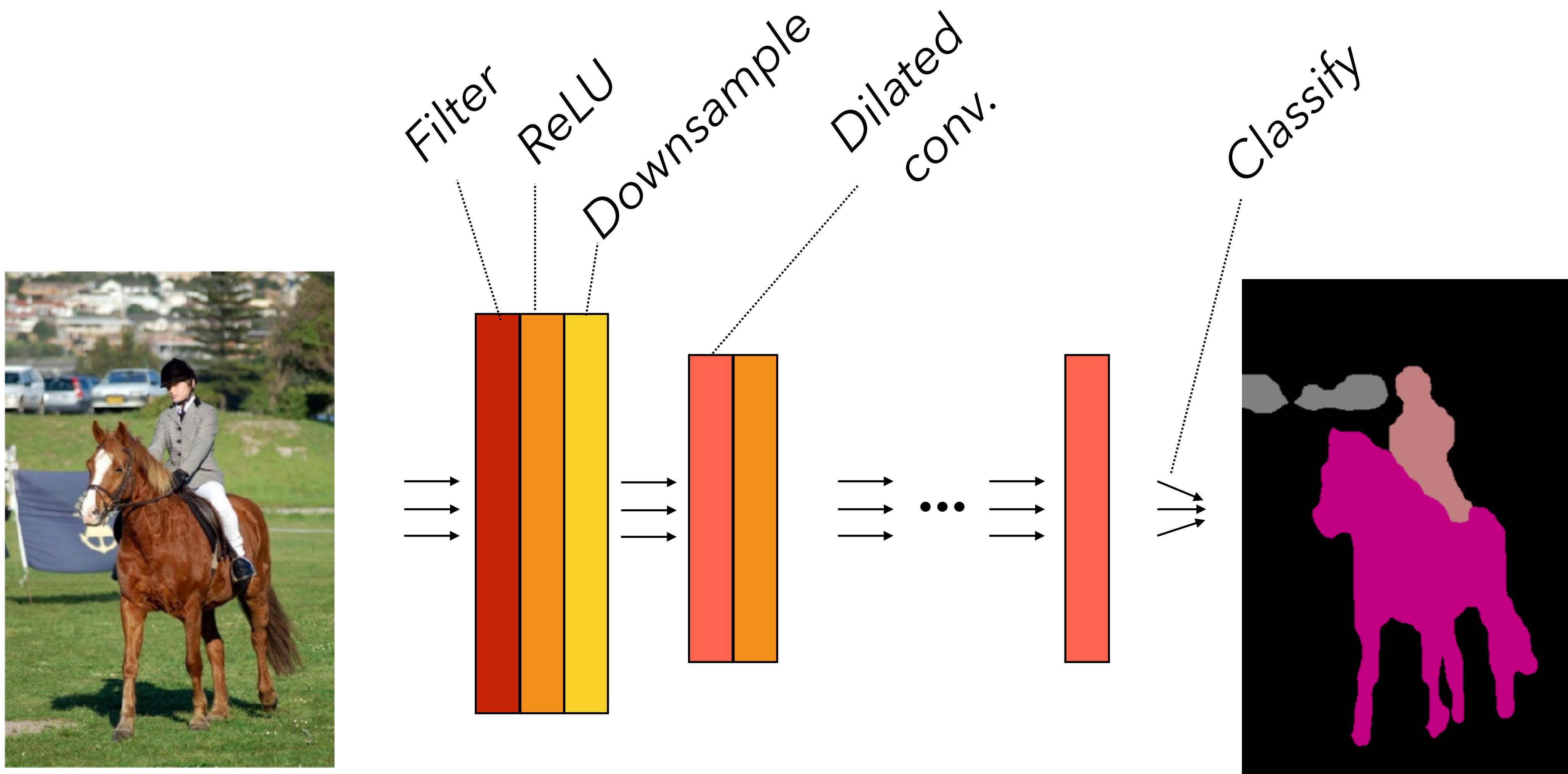


Covers a large receptive field with fewer parameters.

CNN without dilated convolutions

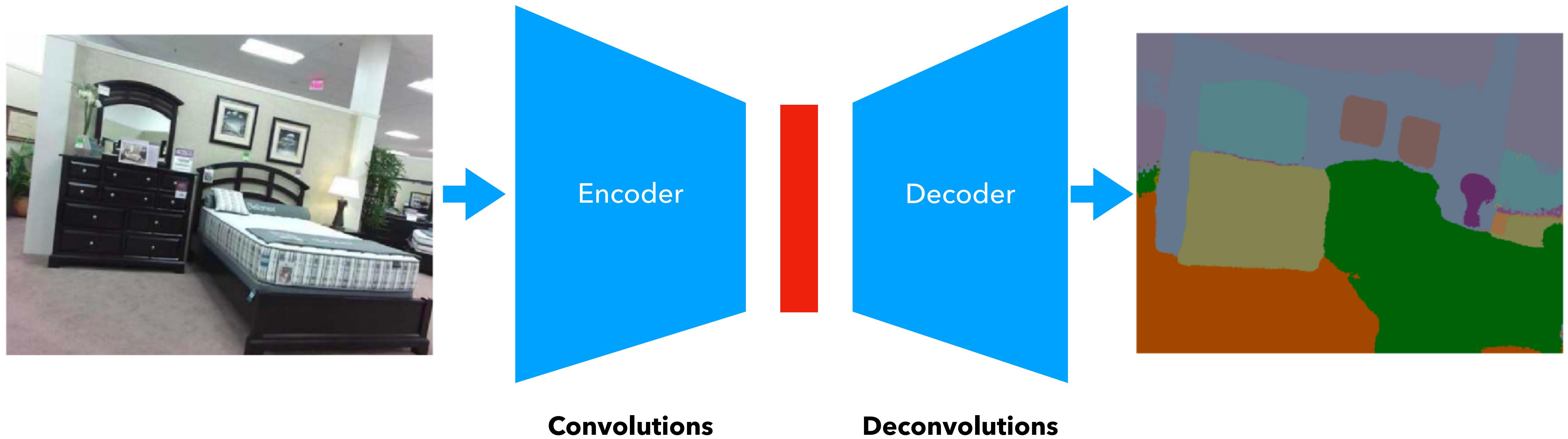


CNN with dilated convolutions

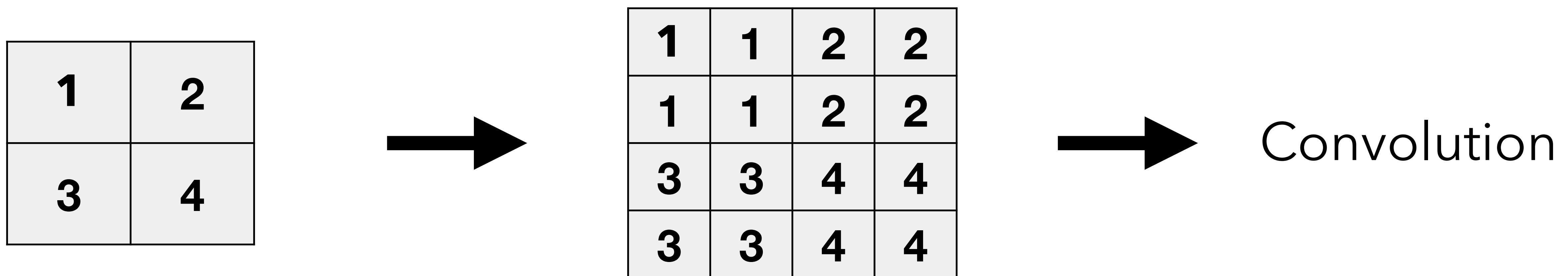


Idea #4: Encoder-decoder models

Encoder-decoder architectures

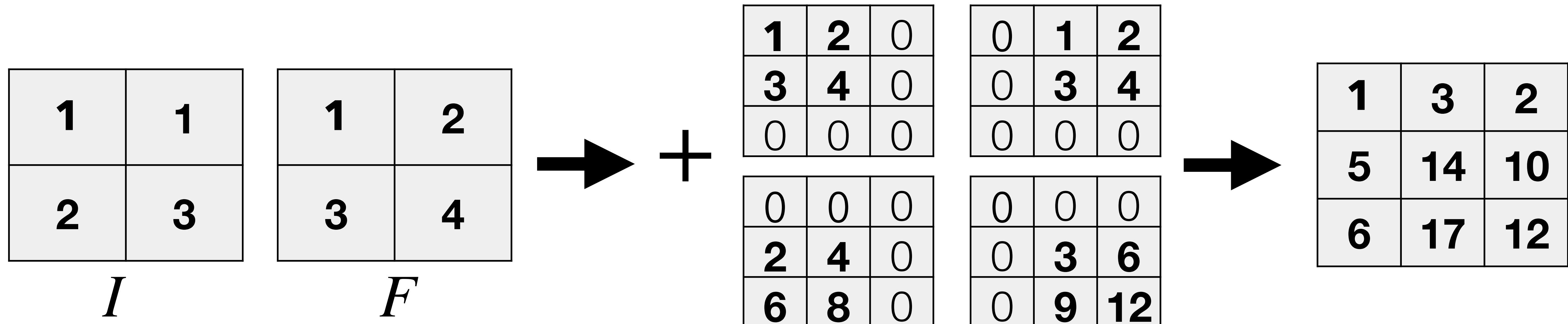


Upsampling



- Often using nearest-neighbor upsampling
- Can also use interpolation.
- Produces fewer “checkerboard” artifacts

Transposed convolution



- Weight the filter by the image coefficient and sum.
- Also sometimes called “upconvolution” or “deconvolution”.

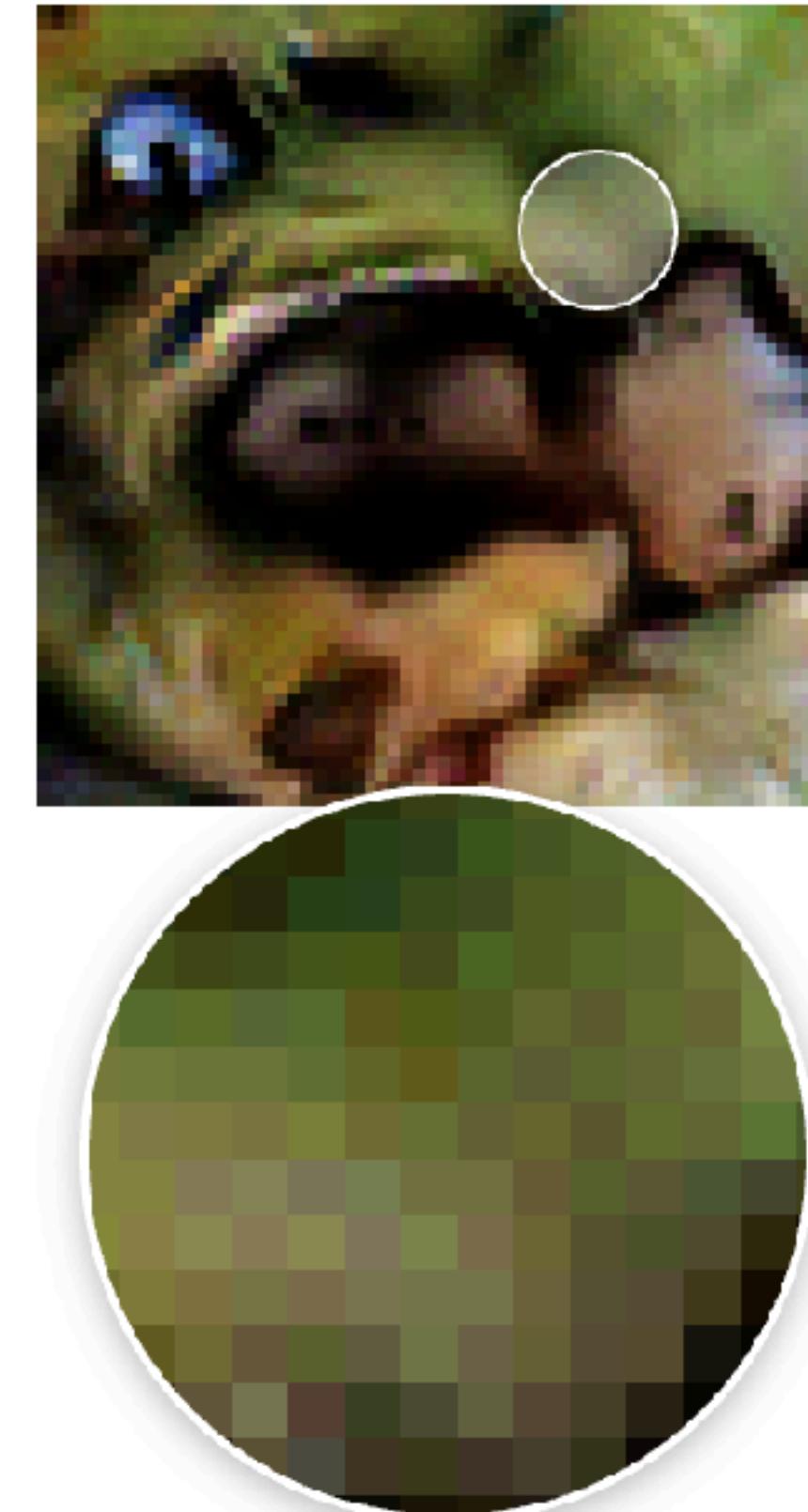
Transposed convolution

$$\begin{array}{|c|c|} \hline 1 & 1 \\ \hline \end{array} \quad \begin{array}{|c|c|} \hline 1 & 2 \\ \hline \end{array}$$

I

$$\begin{array}{|c|c|} \hline 2 & 3 \\ \hline 3 & 4 \\ \hline \end{array} \quad F$$

F

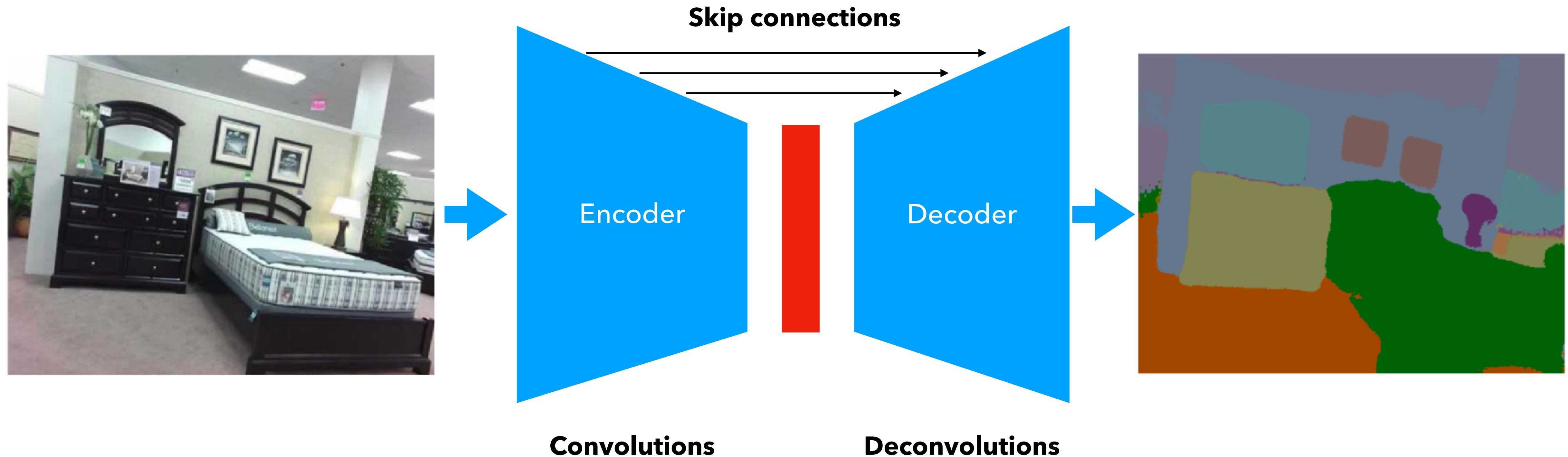


Can lead to “checkerboard” artifacts.

Donahue, et al., 2016 [3]

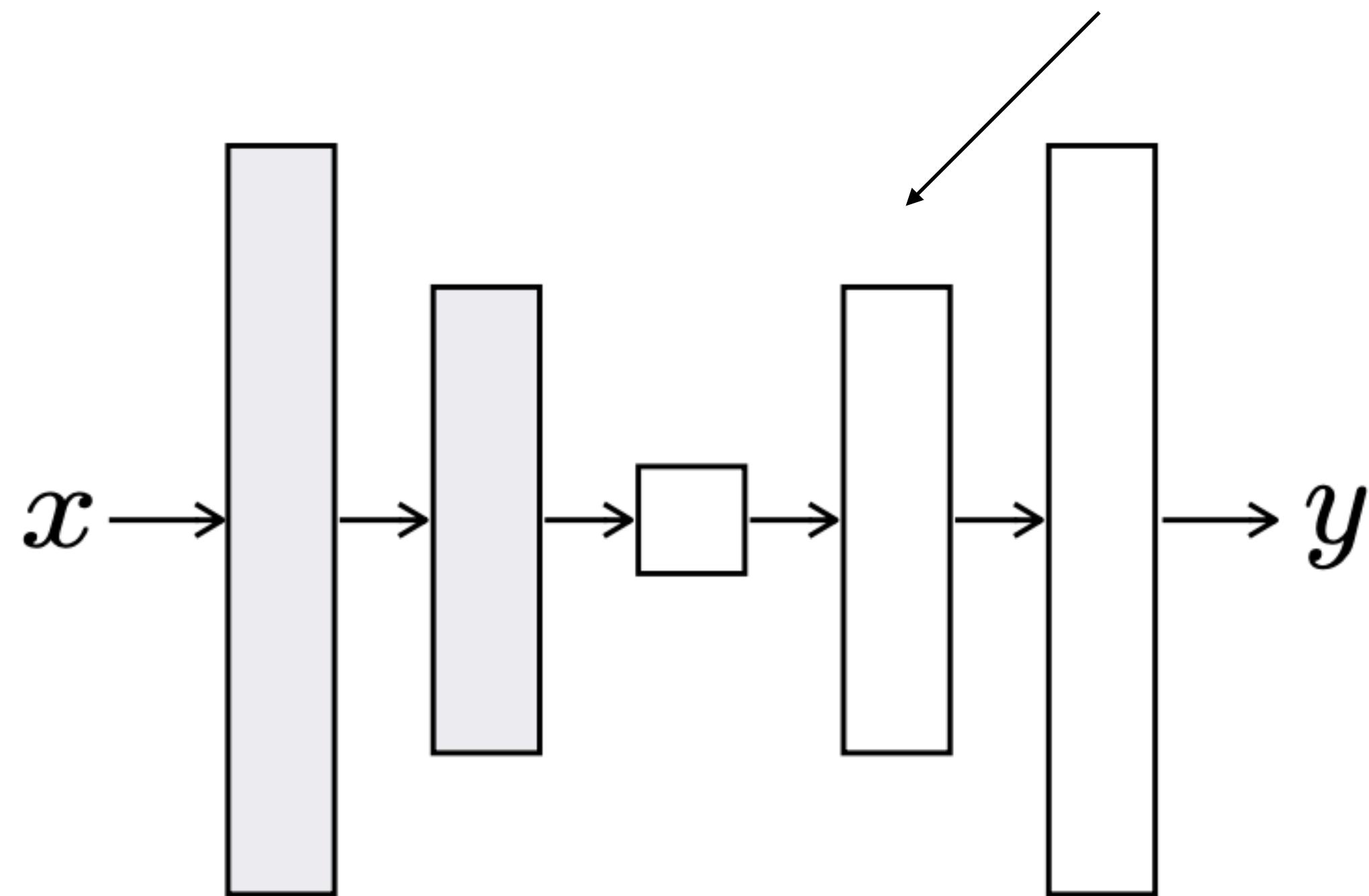
[Odena et al. Distill article]

Encoder-decoder architectures



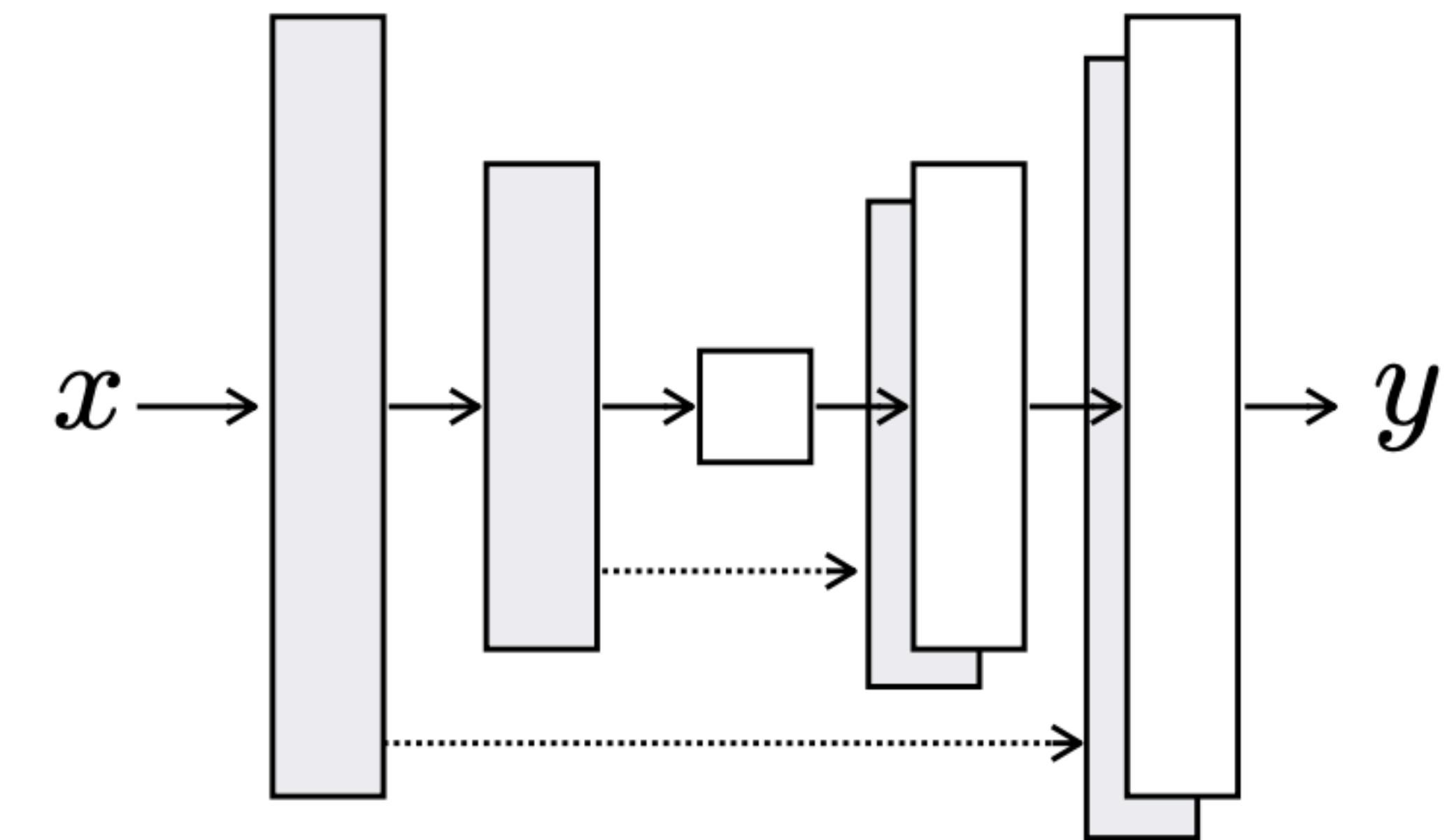
Encoder-decoder architectures

Transposed convolution



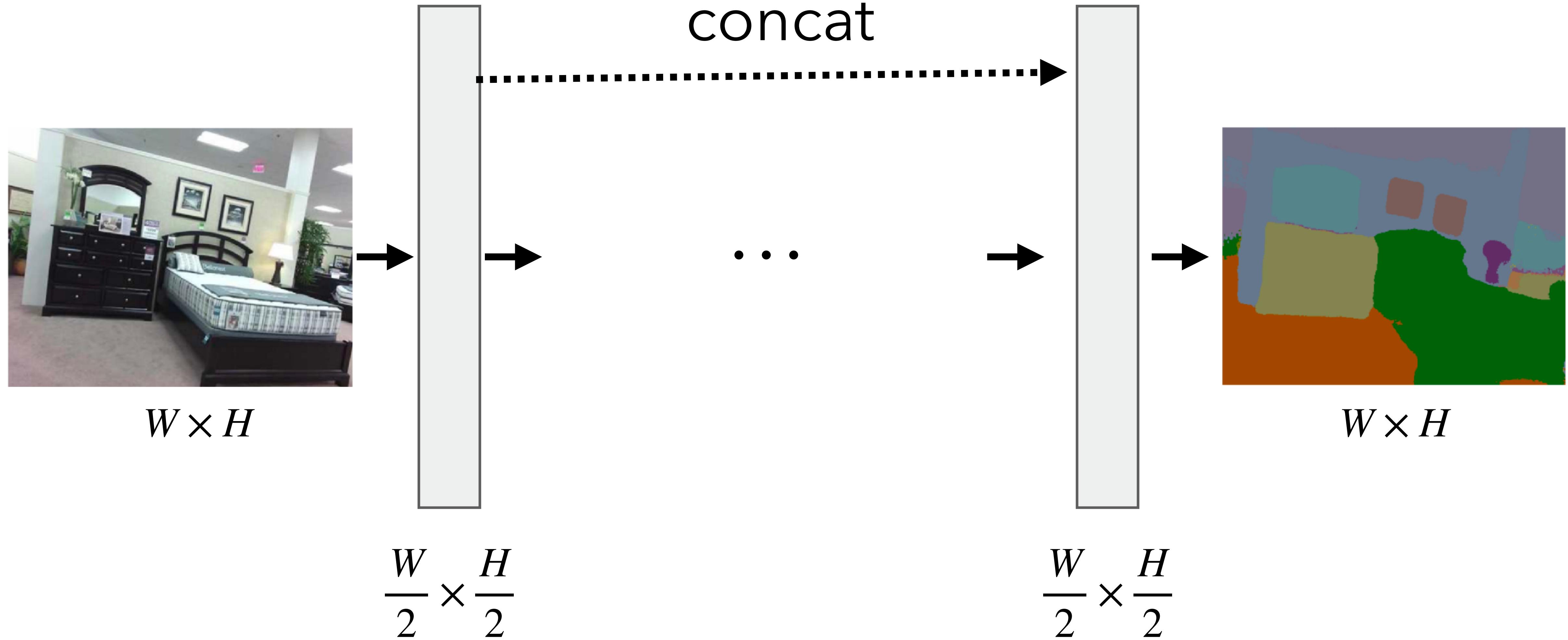
“Vanilla” encoder-decoder architecture

Early layers and late layers have
same shape. Concatenate channel-wise!

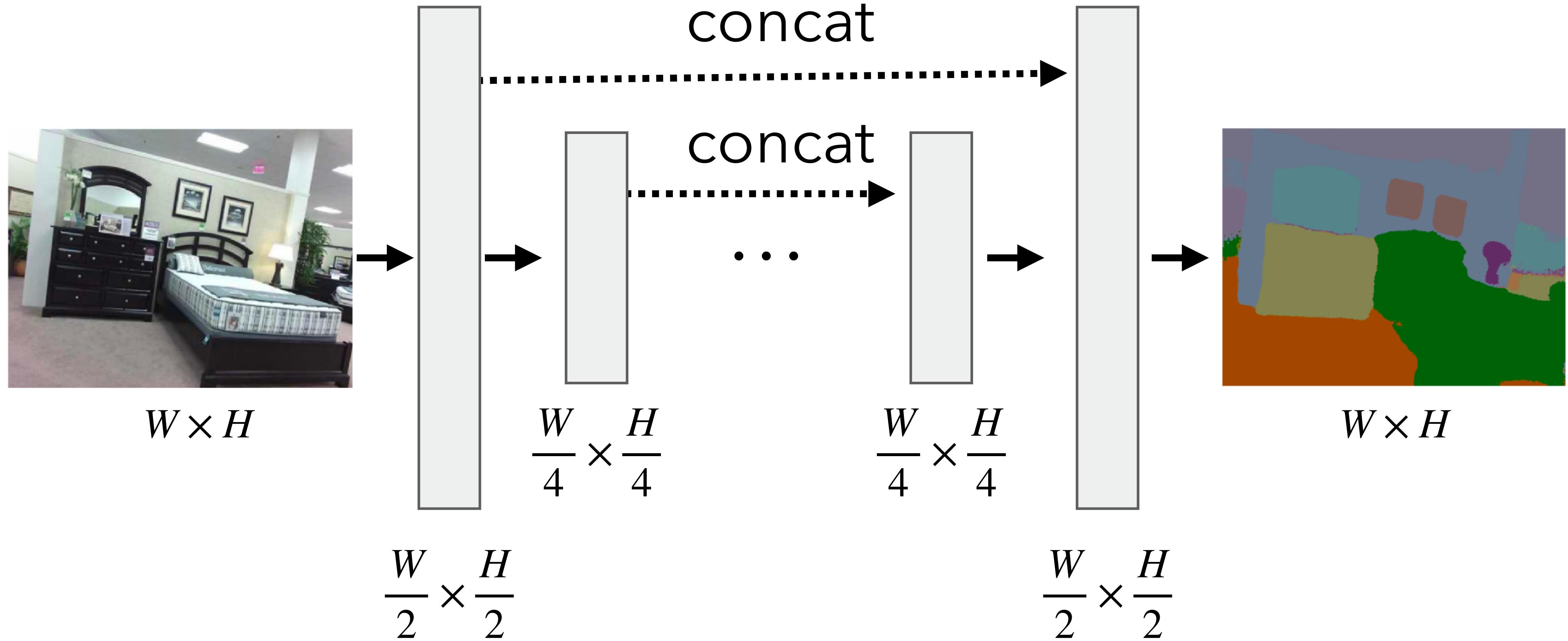


U-Net

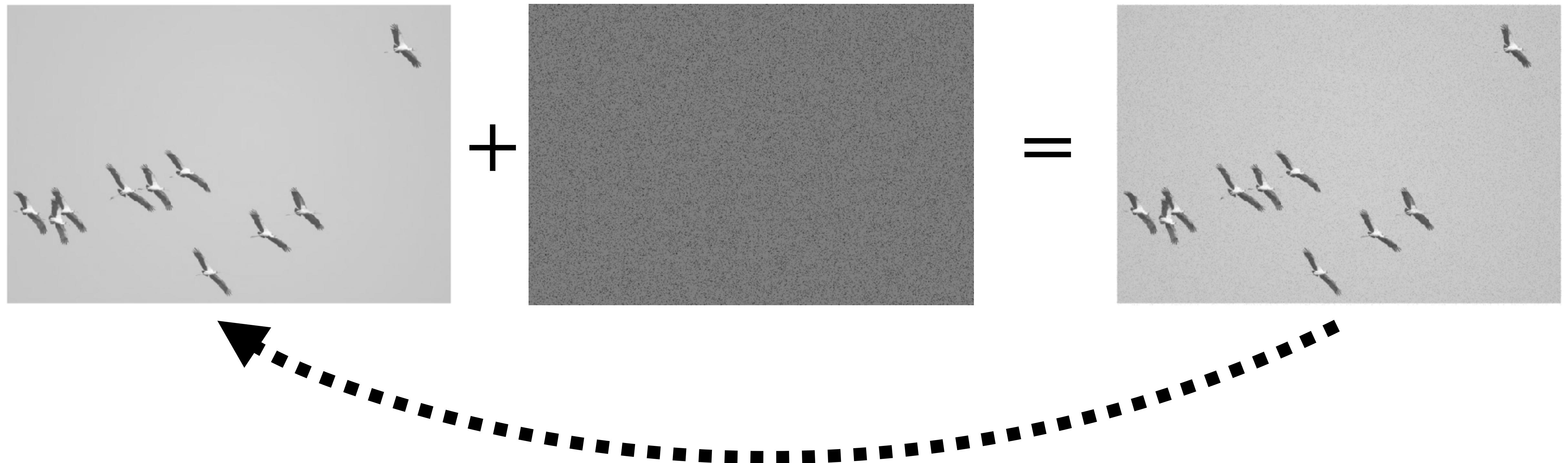
U-net



U-net

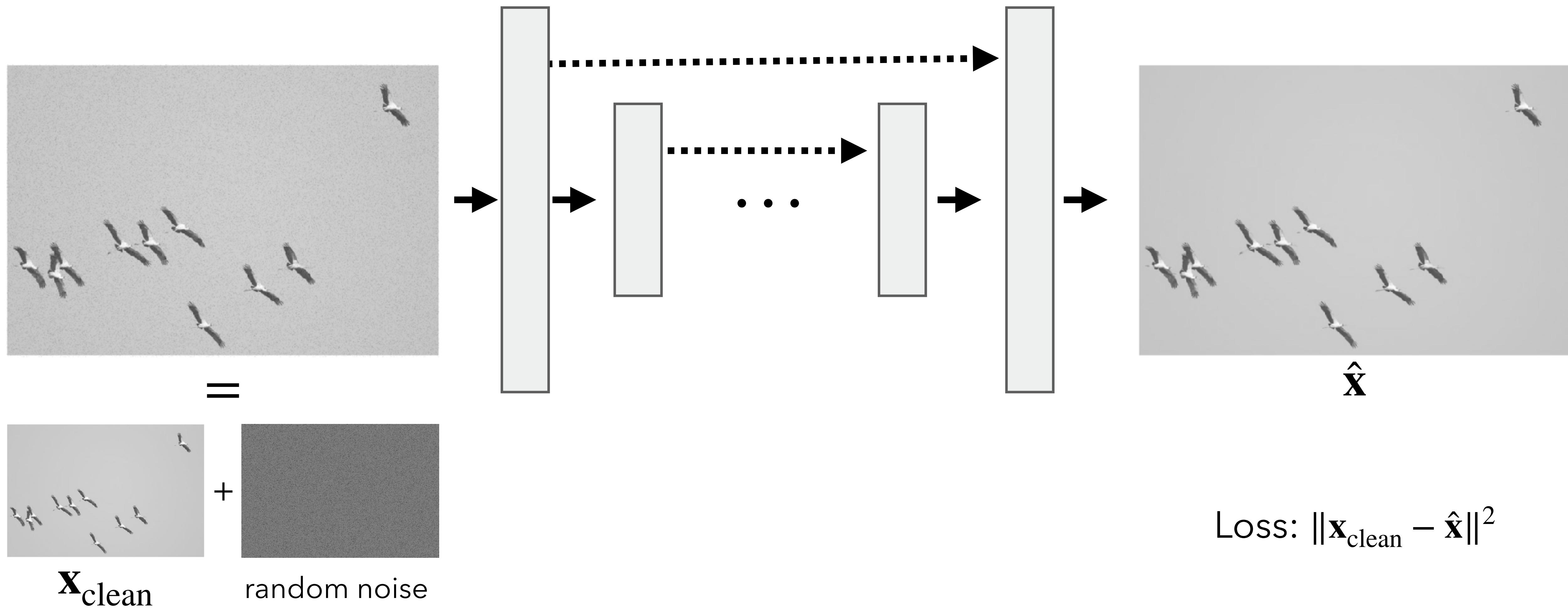


Other uses for U-nets



Goal: recover the original image
Recall: denoising problem

Denoising



Next class: vision and language