

# Lecture 9: Convolutional networks

# Announcements

- PS3 due tonight, PS4 out tonight
- **Next week:** guest lectures on image generation by Daniel Geng, Sarah Jabbour, and Yiming Dou

# Image classification

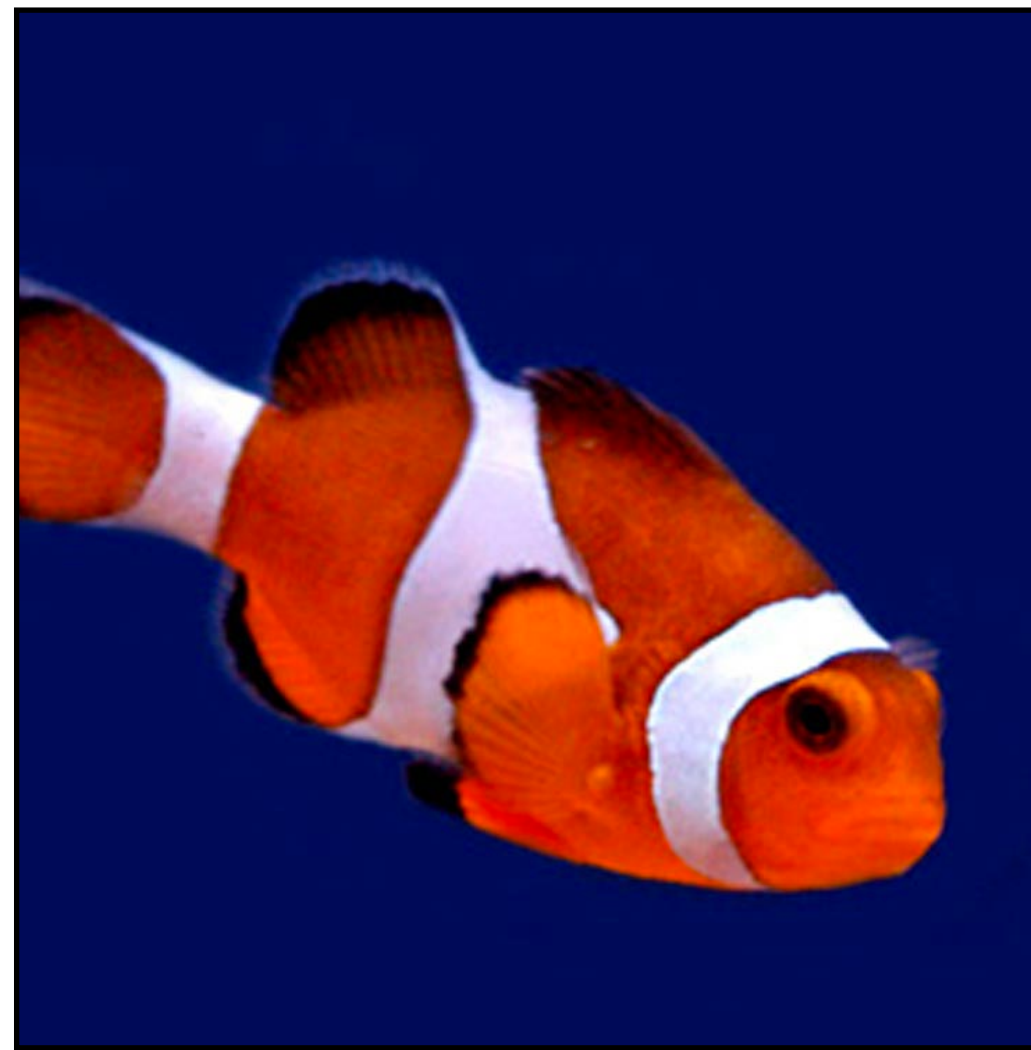


image  $x$



"Fish"

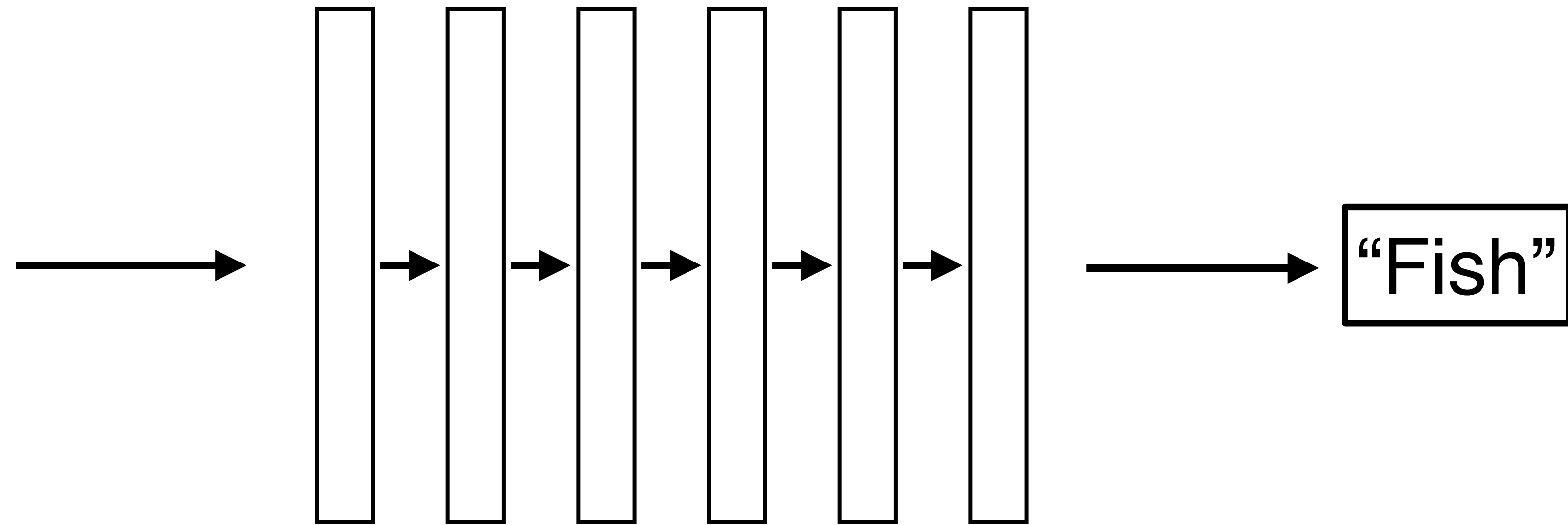
label  $y$

# Image classification

What should these be?



image  $x$

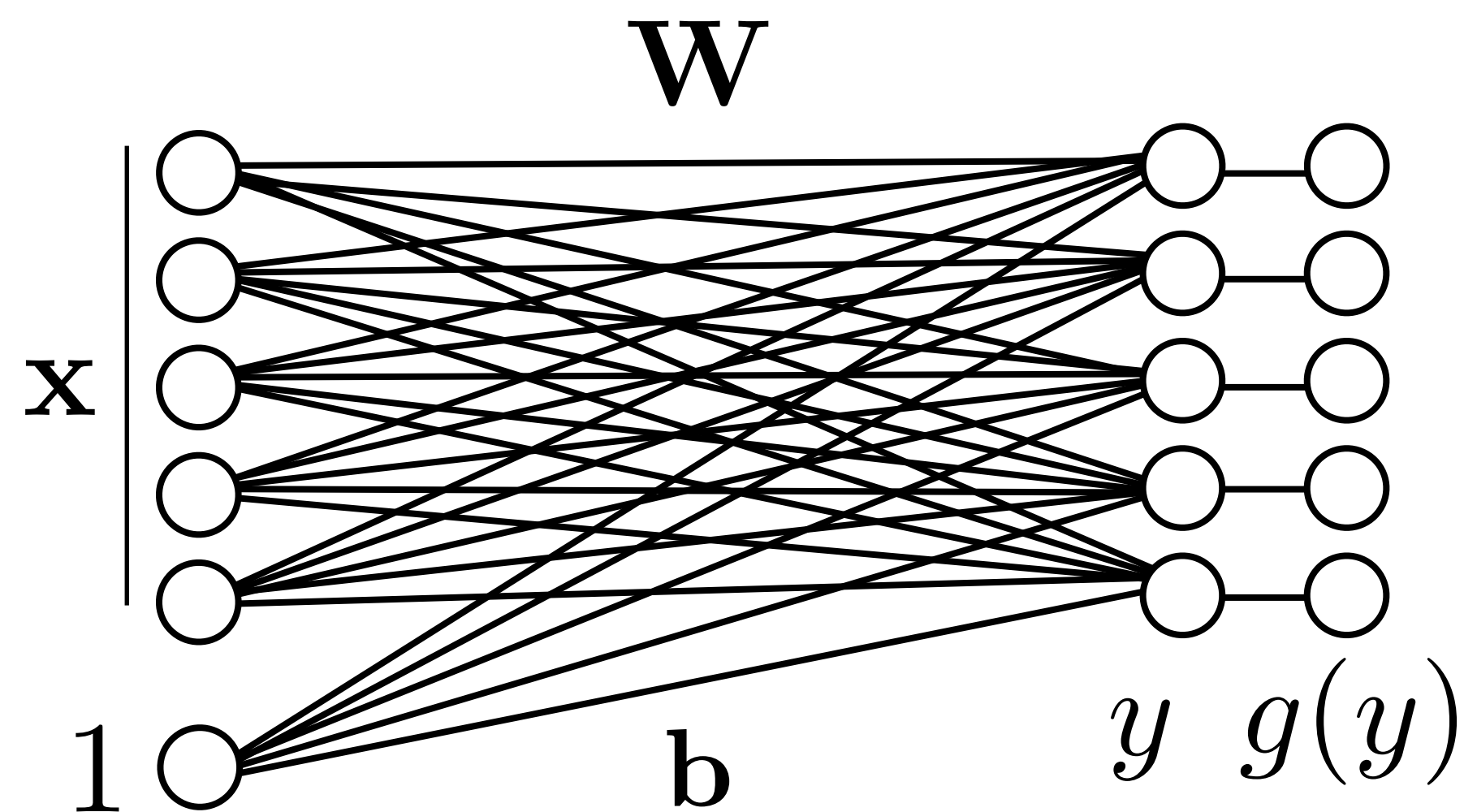


label  $y$

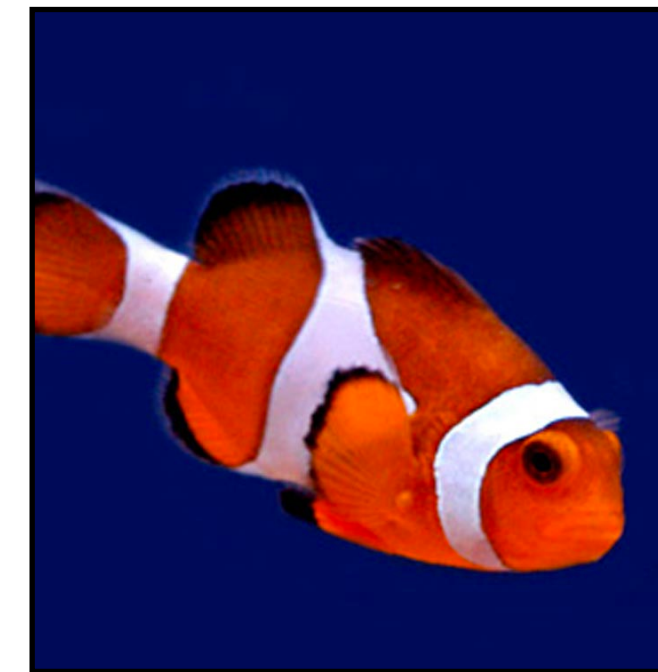


# Idea #1: Fully-connected network

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



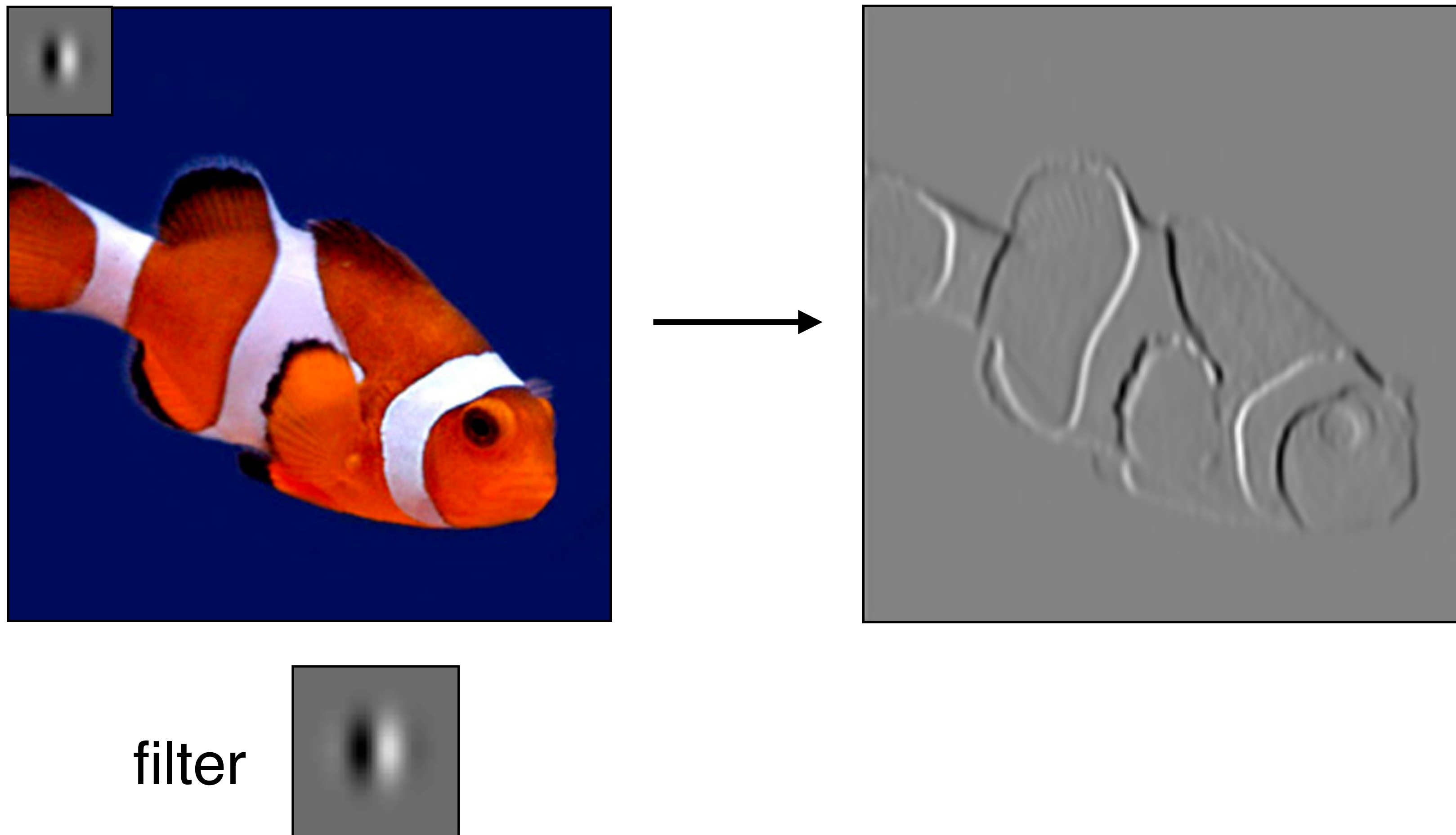
$\mathbf{X} =$



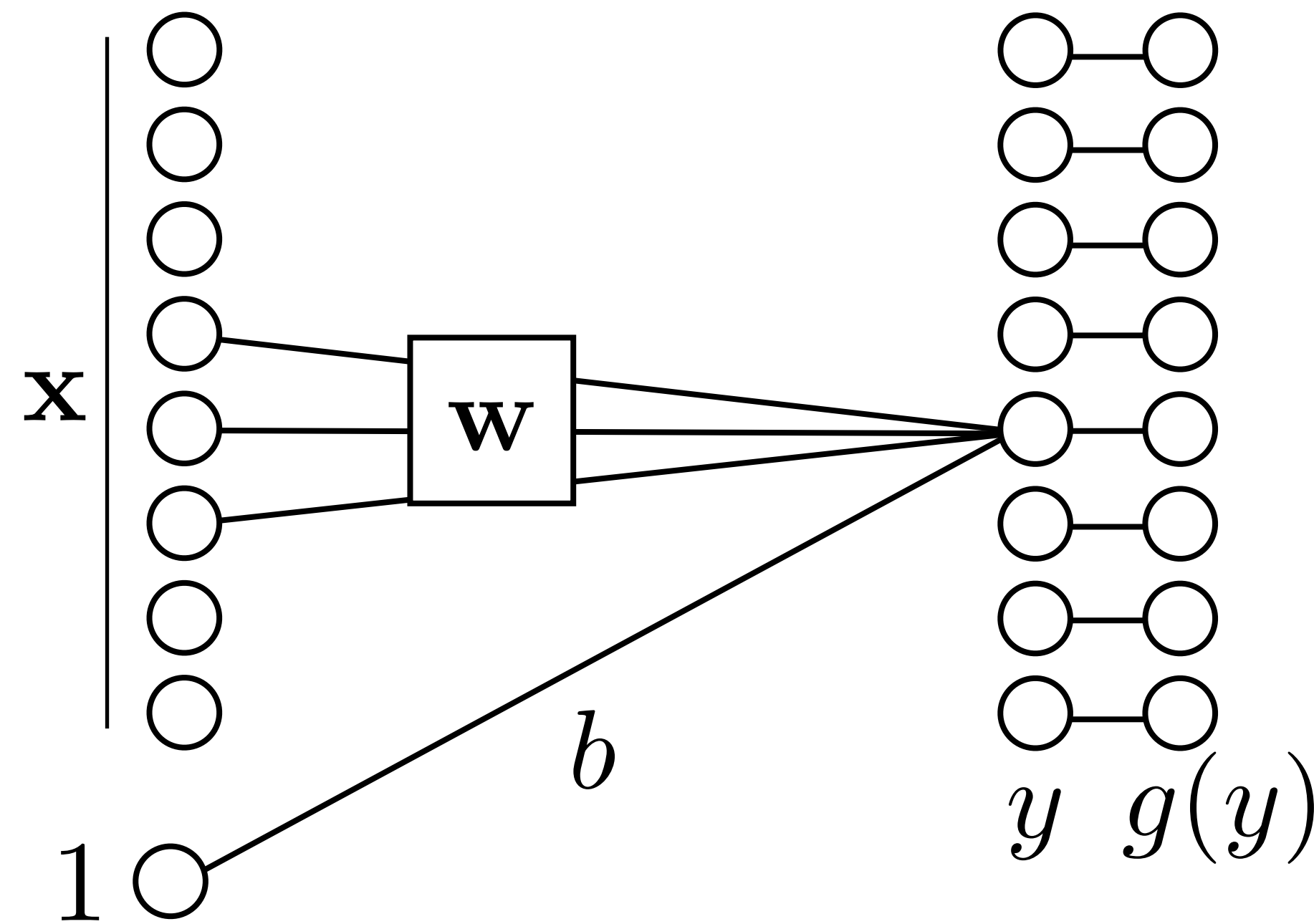
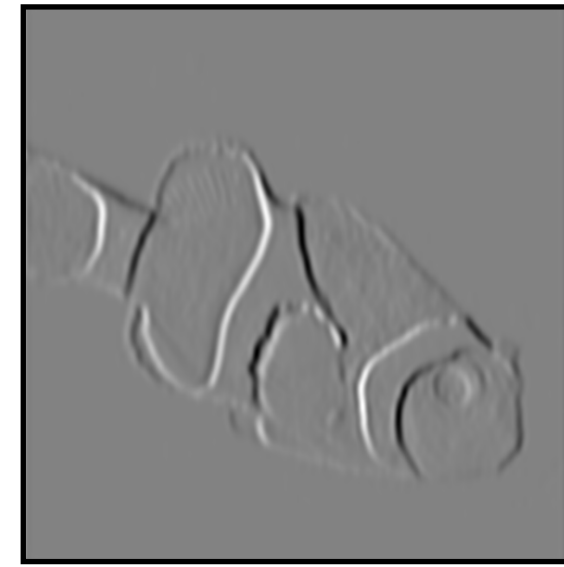
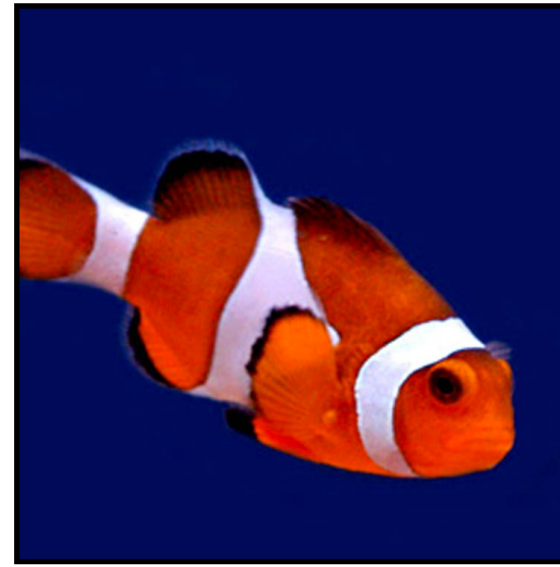
But  $\mathbf{X}$  is really big!

Say,  $256 \times 256 \times 3 = 197\text{k}$

# Can we use convolution in a neural network?



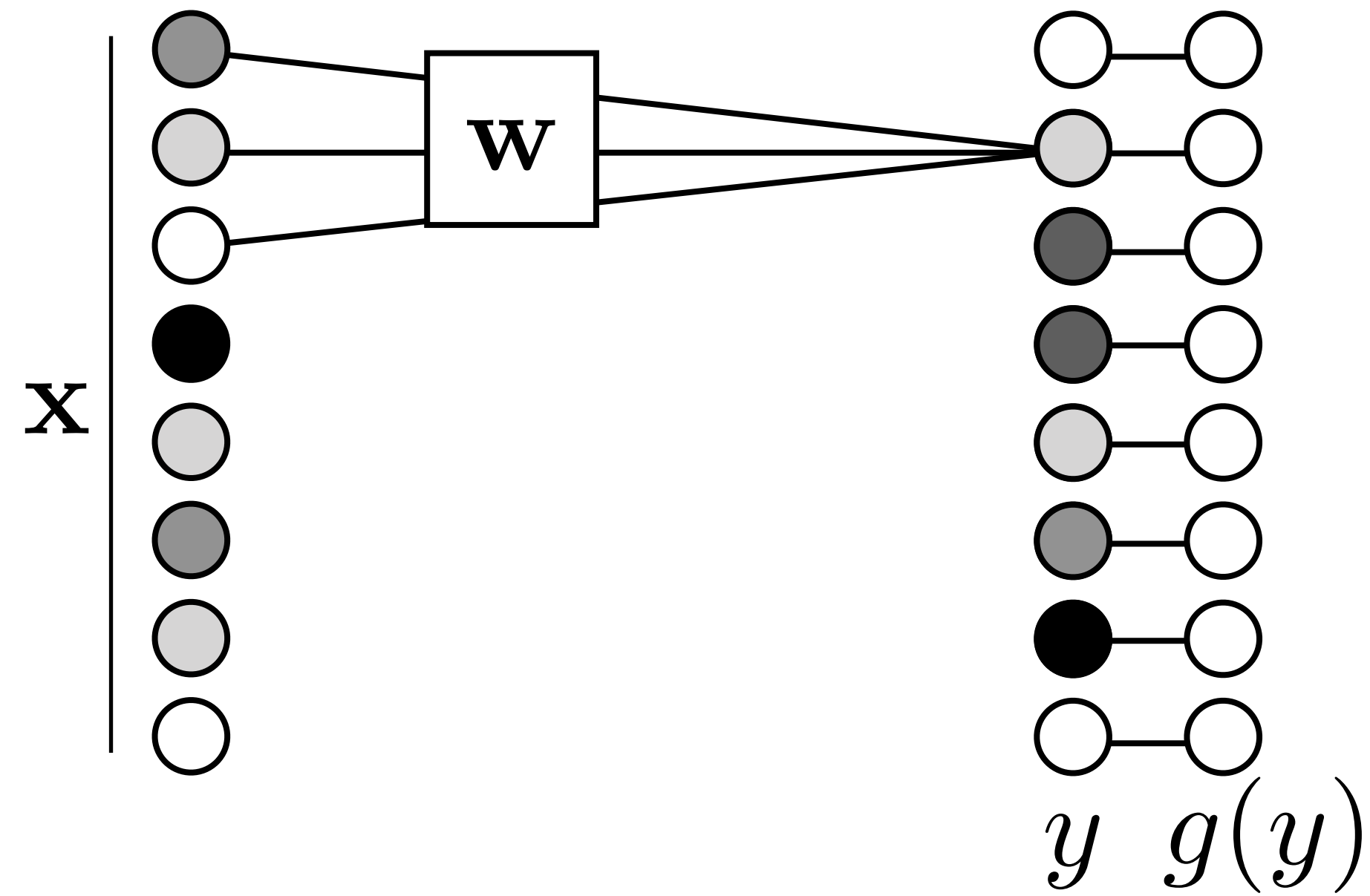
# Recall: Sparsely connected network



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

# Convolutional neural network

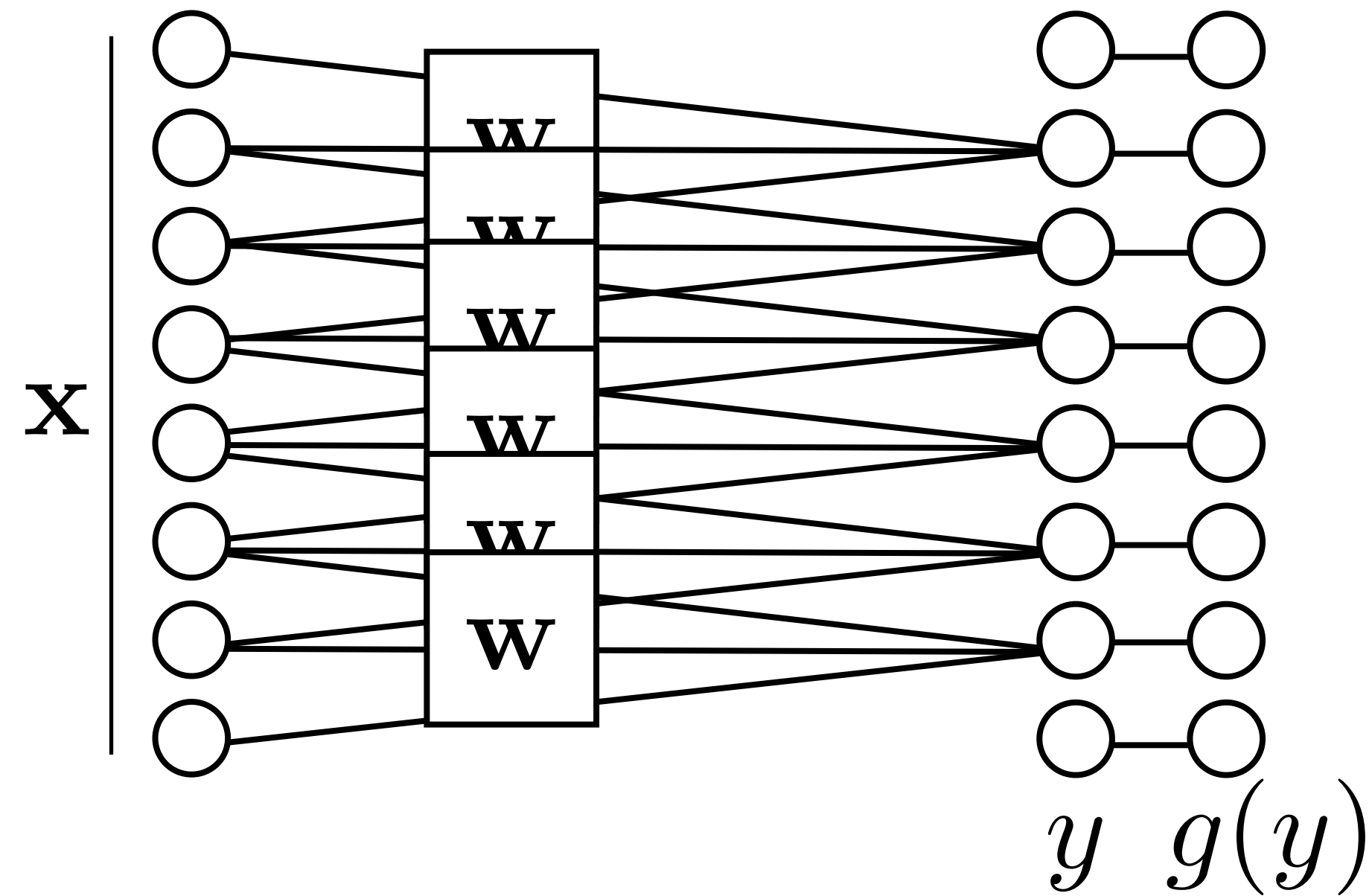
## Conv layer



Each output unit is computed from an image patch.

# Weight sharing

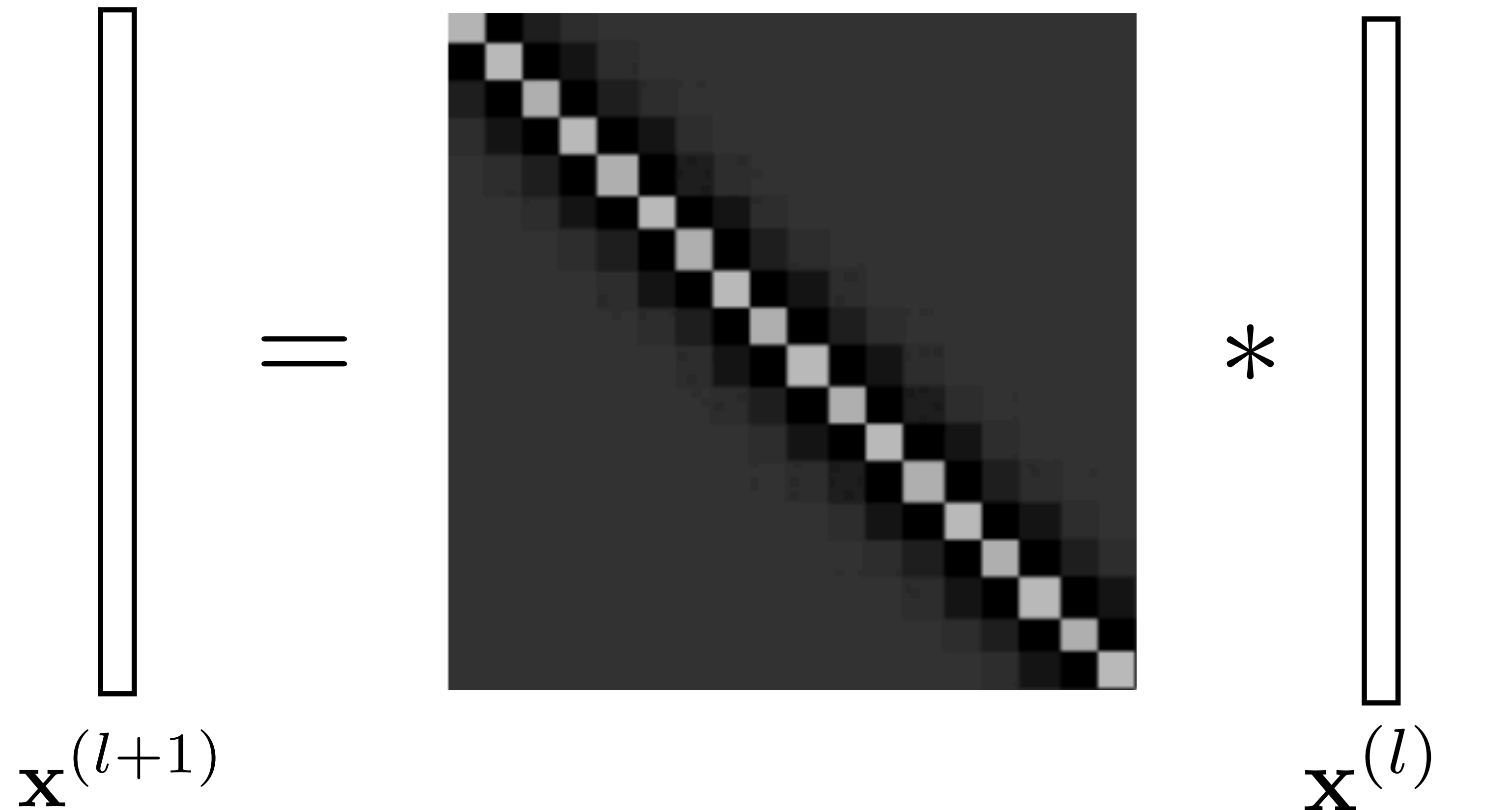
## Conv layer



We “share” weights for each patch.

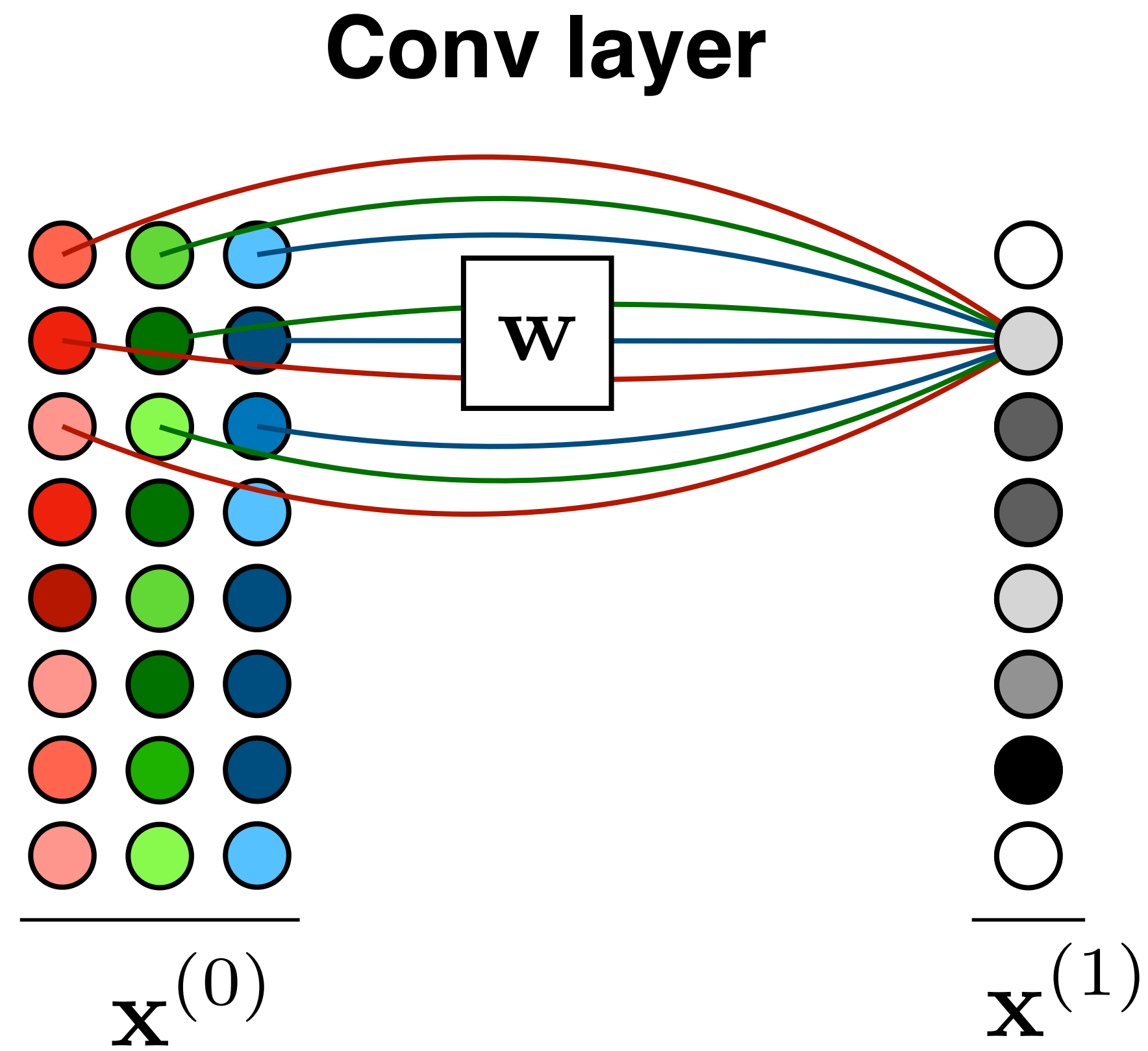
If a feature is useful in one position, it should be useful in others, too.

# Convolution is a linear function



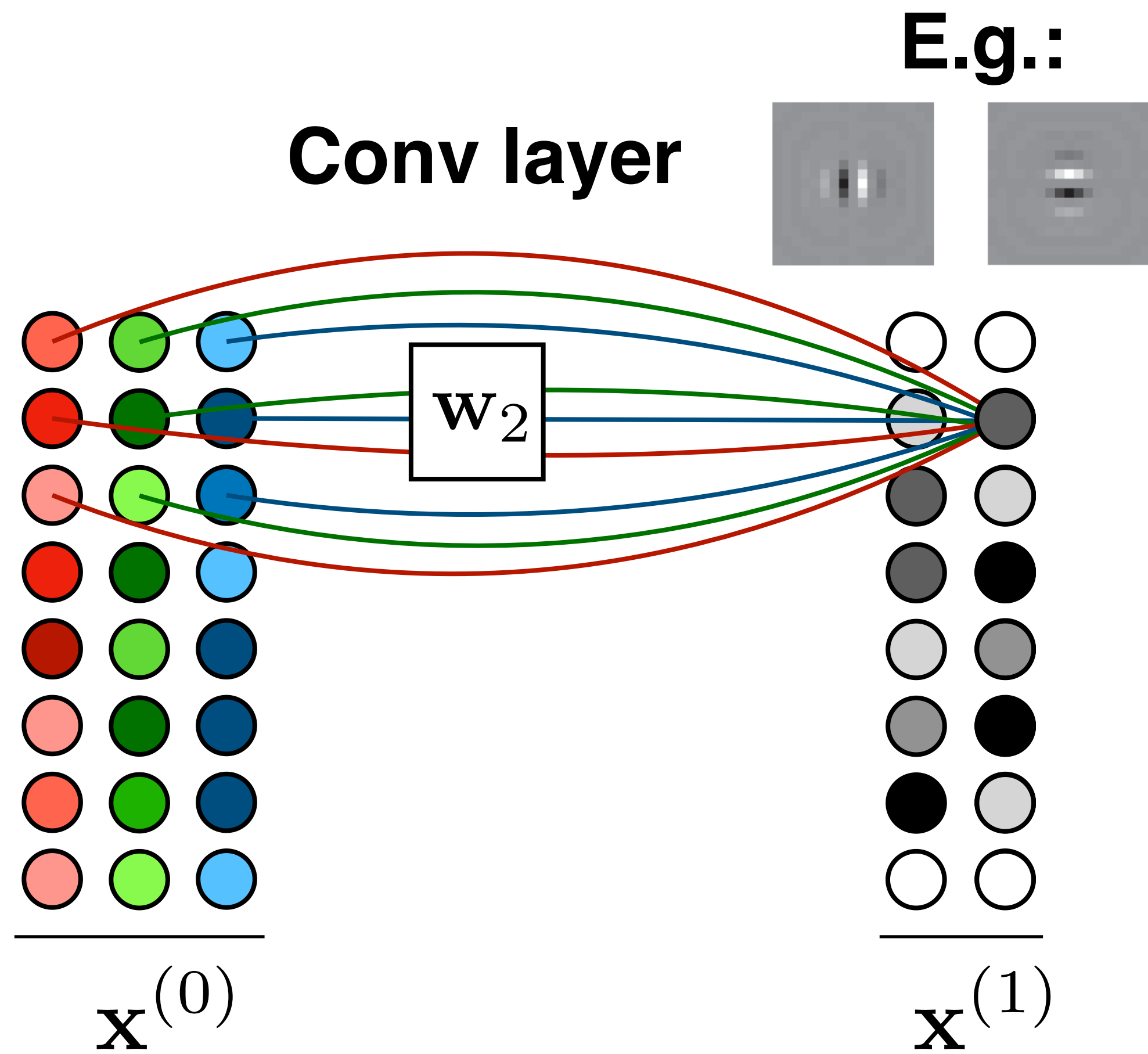
- Constrained linear layer
  - Fewer parameters: easier to learn, less overfitting
  - Usually use zero padding
- e.g., image

# Multiple channels



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

# Multiple channels

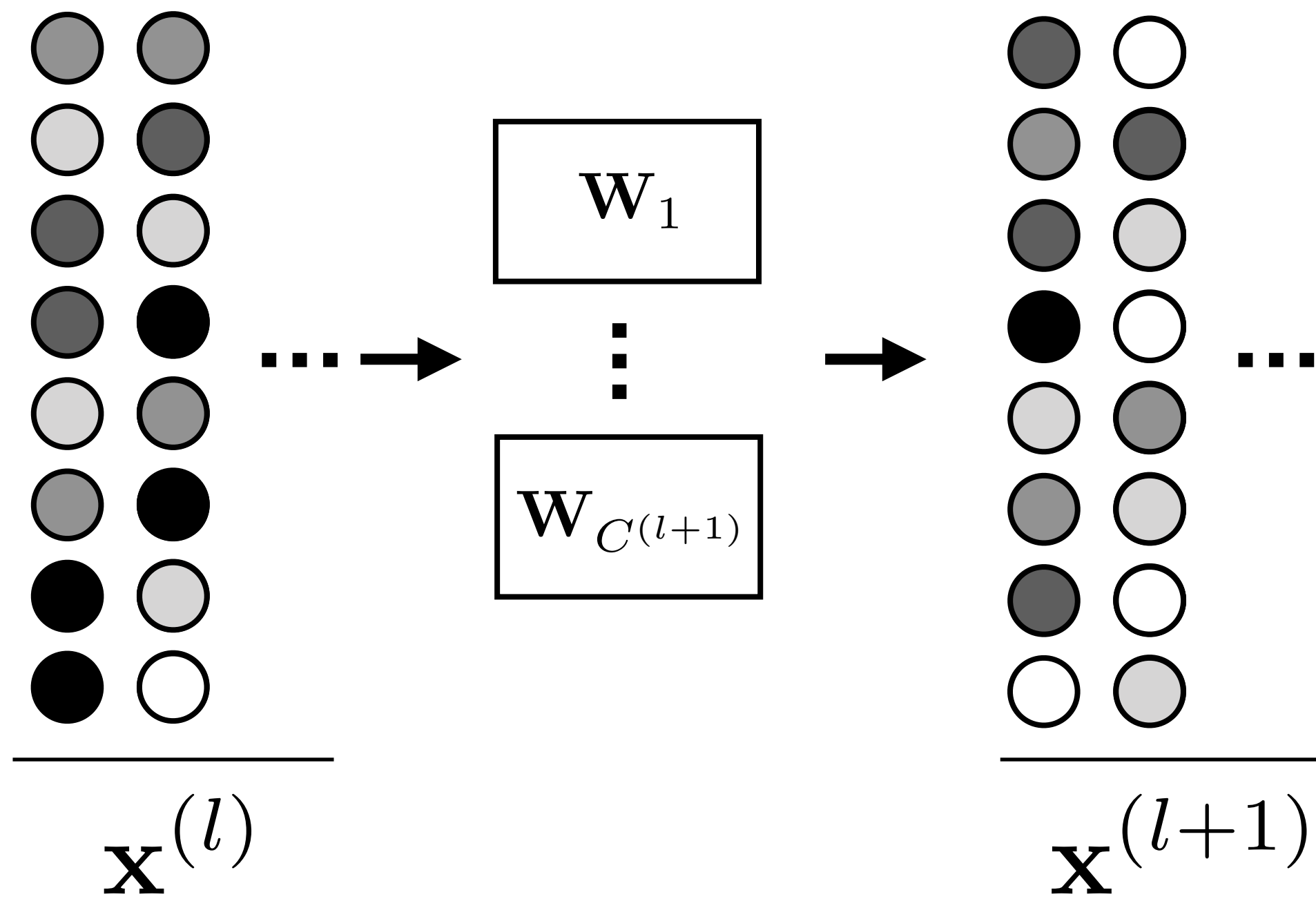


$$\mathbb{R}^{N \times C^{(0)}} \rightarrow \mathbb{R}^{N \times C^{(1)}}$$



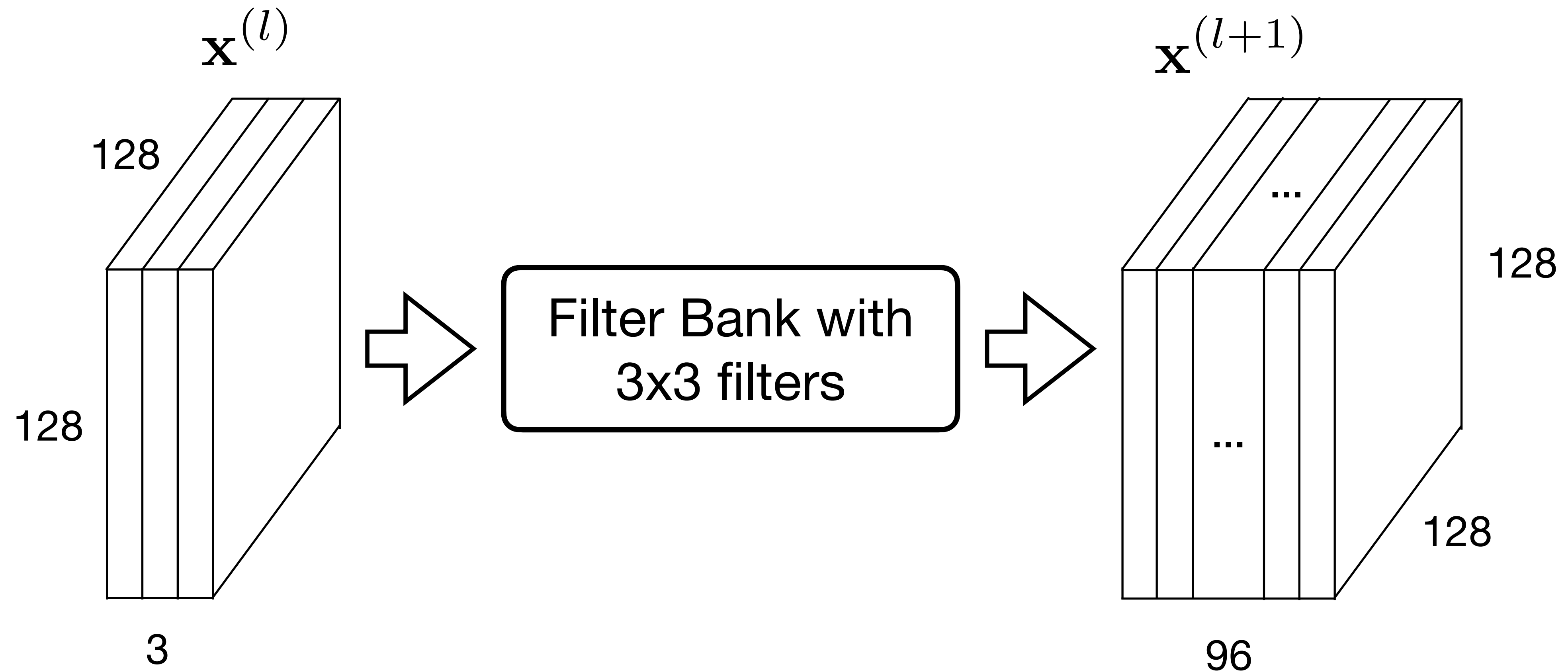
# Multiple channels

## Conv layer



$$\mathbb{R}^{N \times C^{(l)}} \rightarrow \mathbb{R}^{N \times C^{(l+1)}}$$

# Multiple channels: Example



How many parameters does *each filter* have?

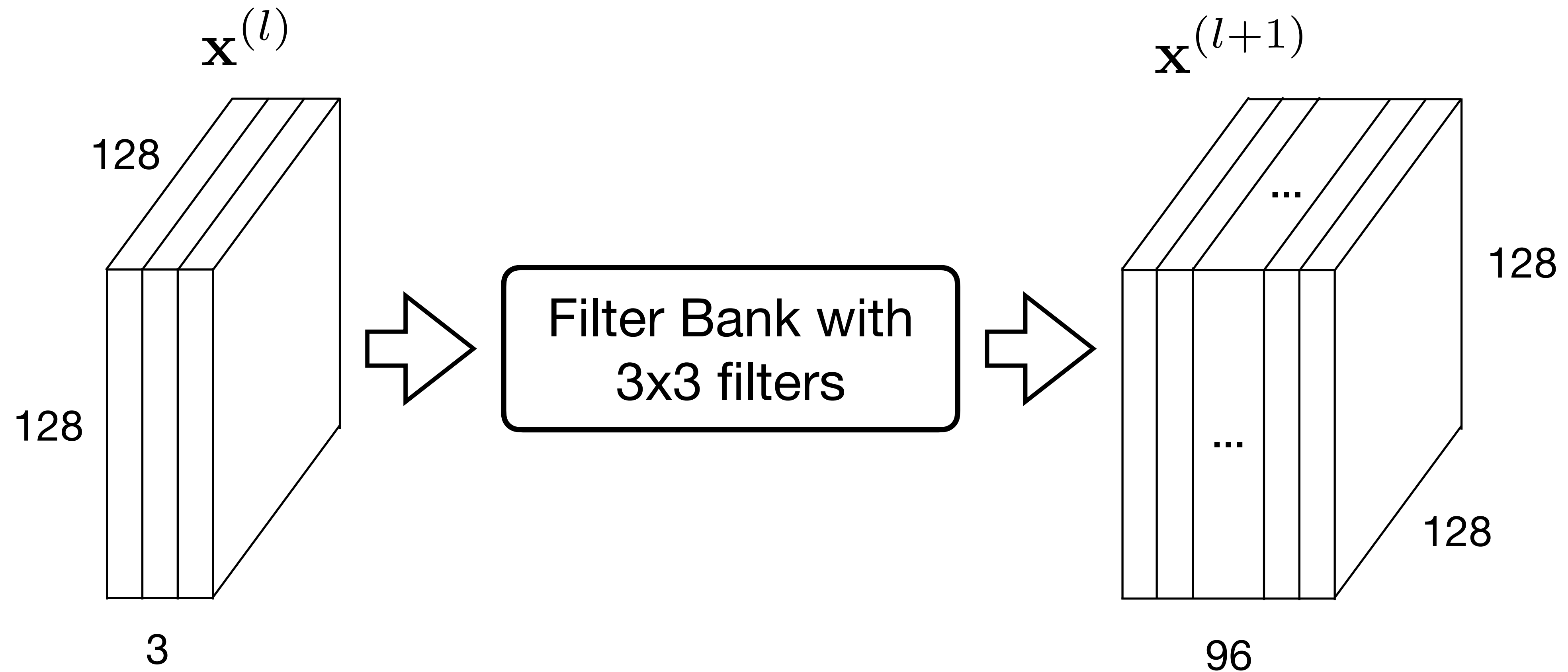
(a) 9

(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

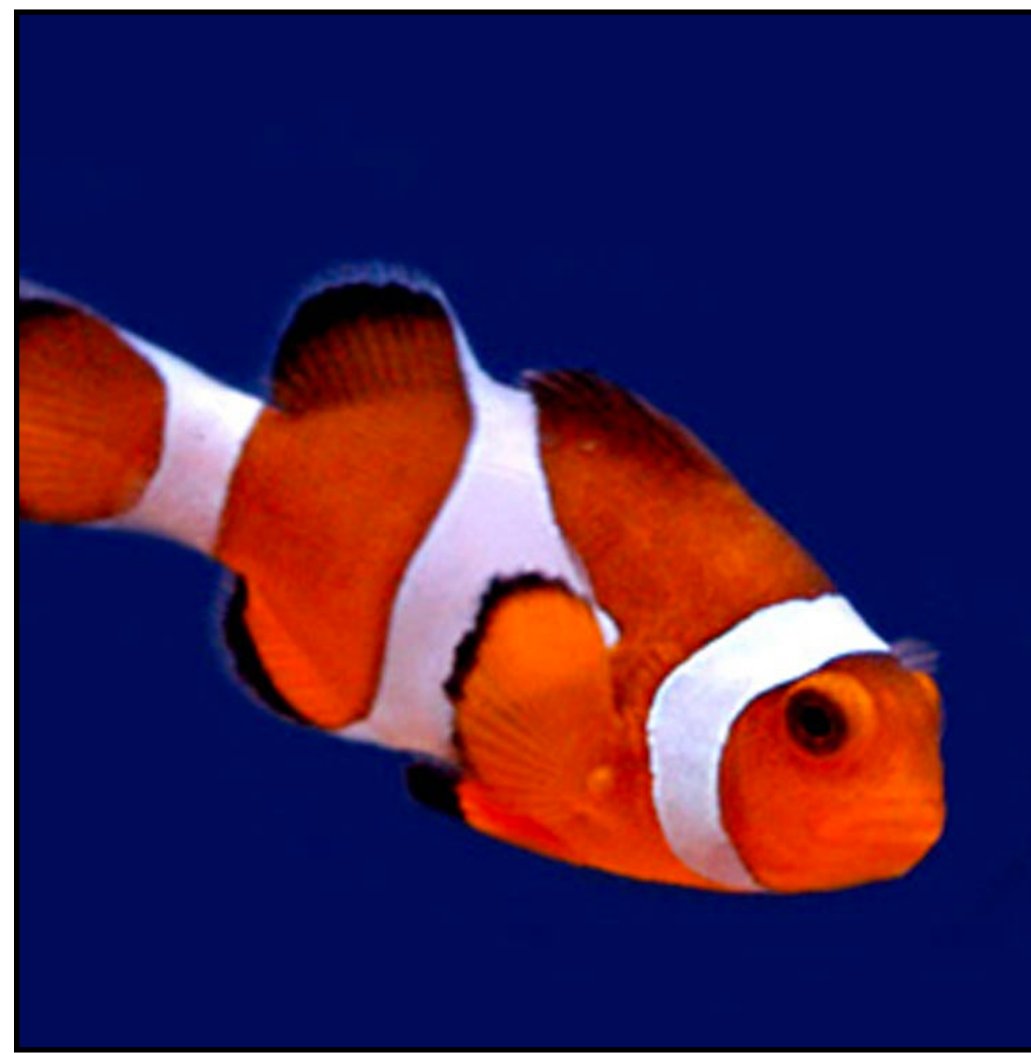
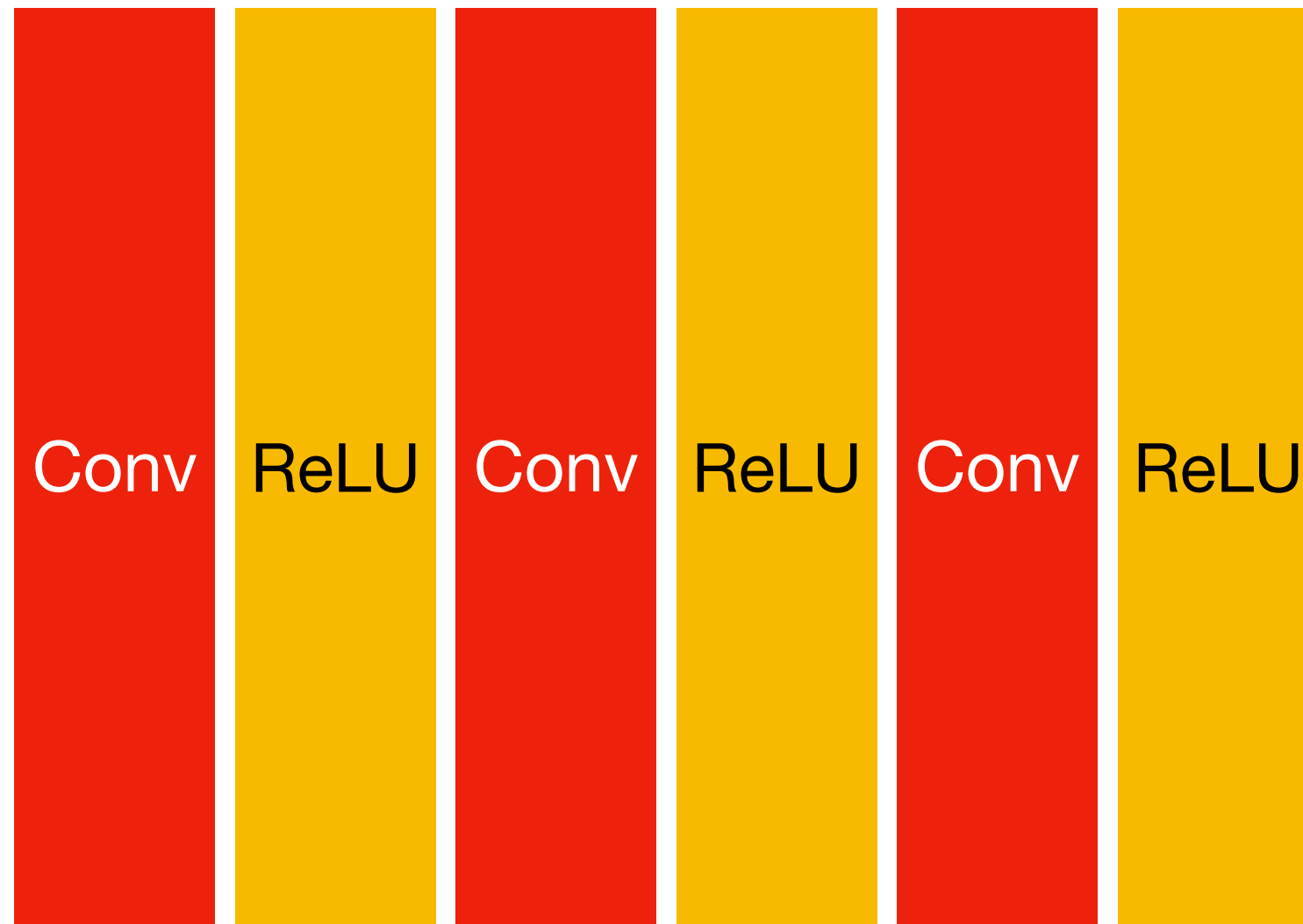


image  $x$



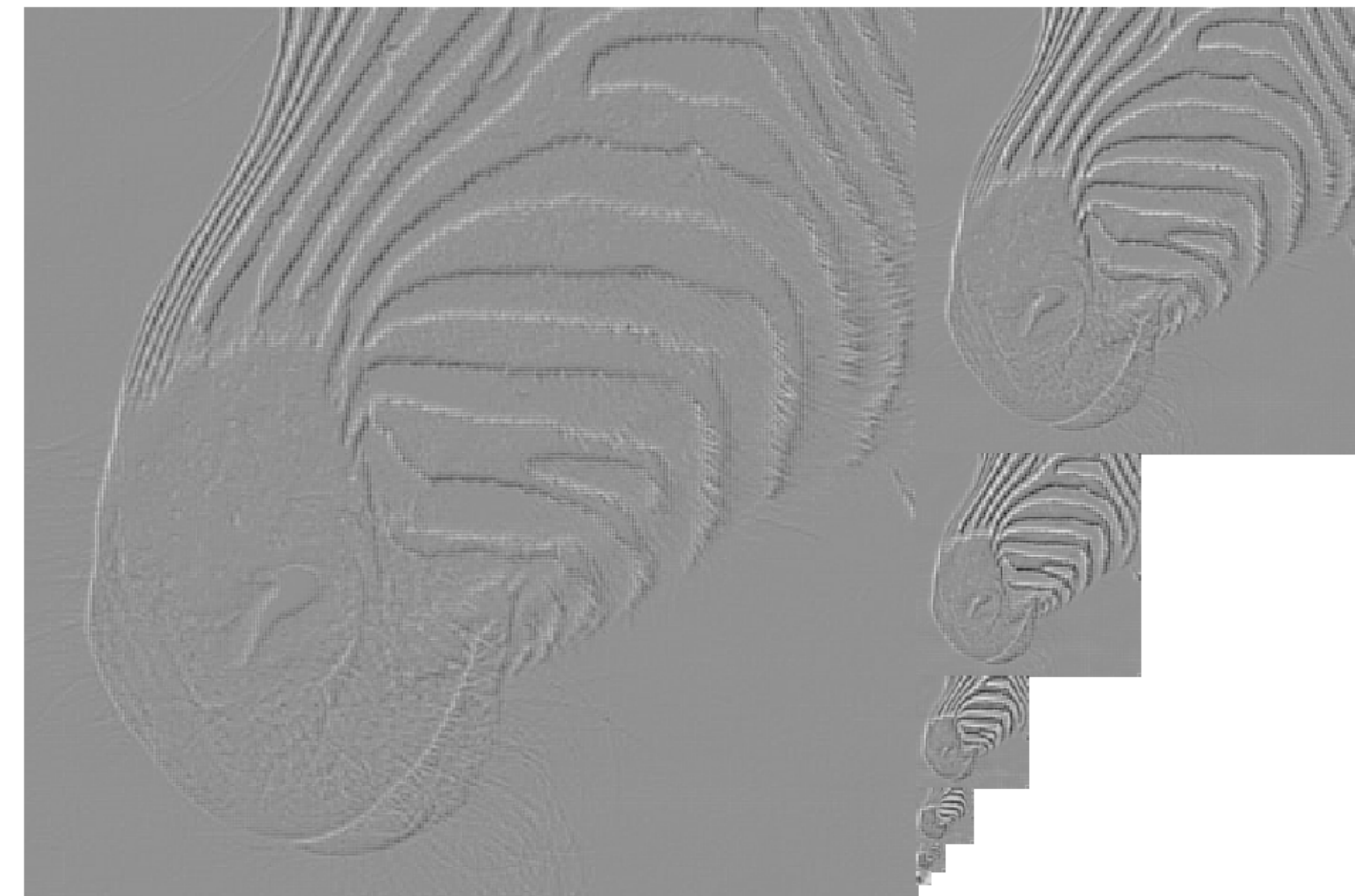
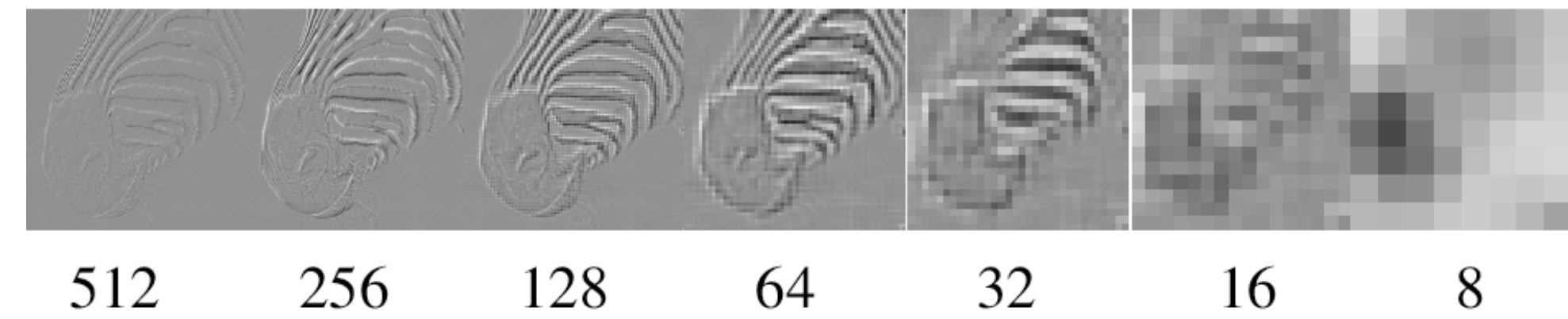
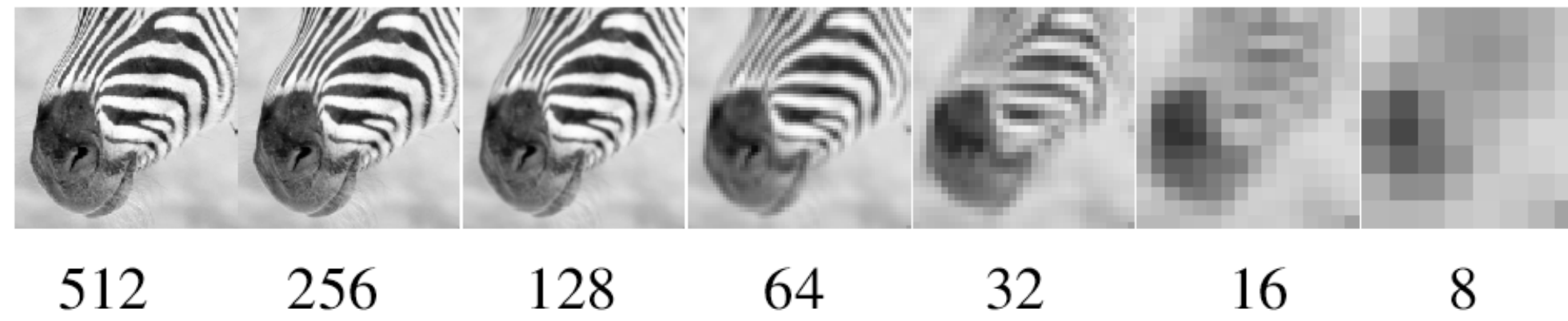
“Fish”

label  $y$

Problems with this idea:

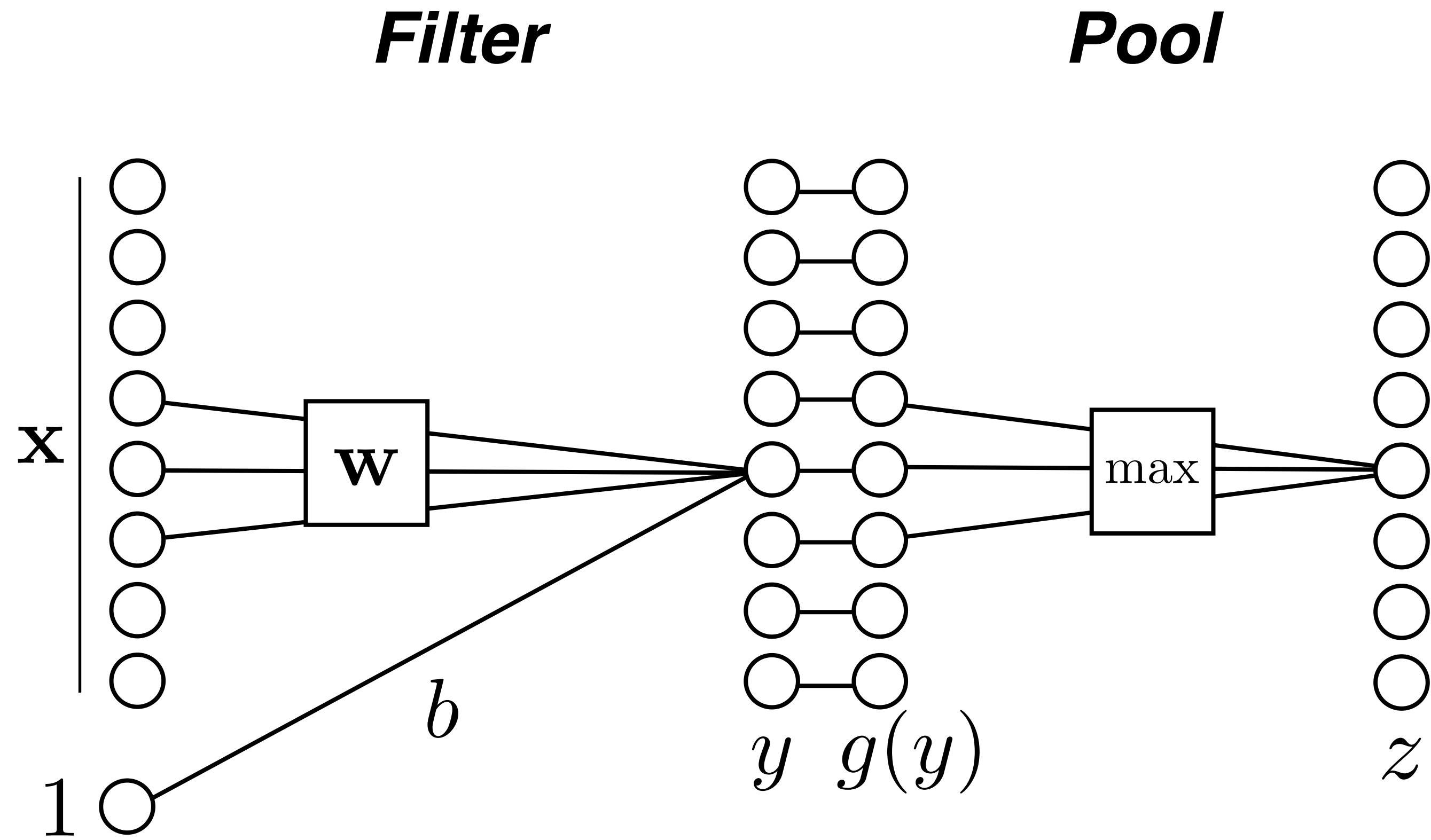
1. No “global” processing.
2. How do you get the final label?

# Recall: pyramid representations



Can we use a similar idea in CNNs?

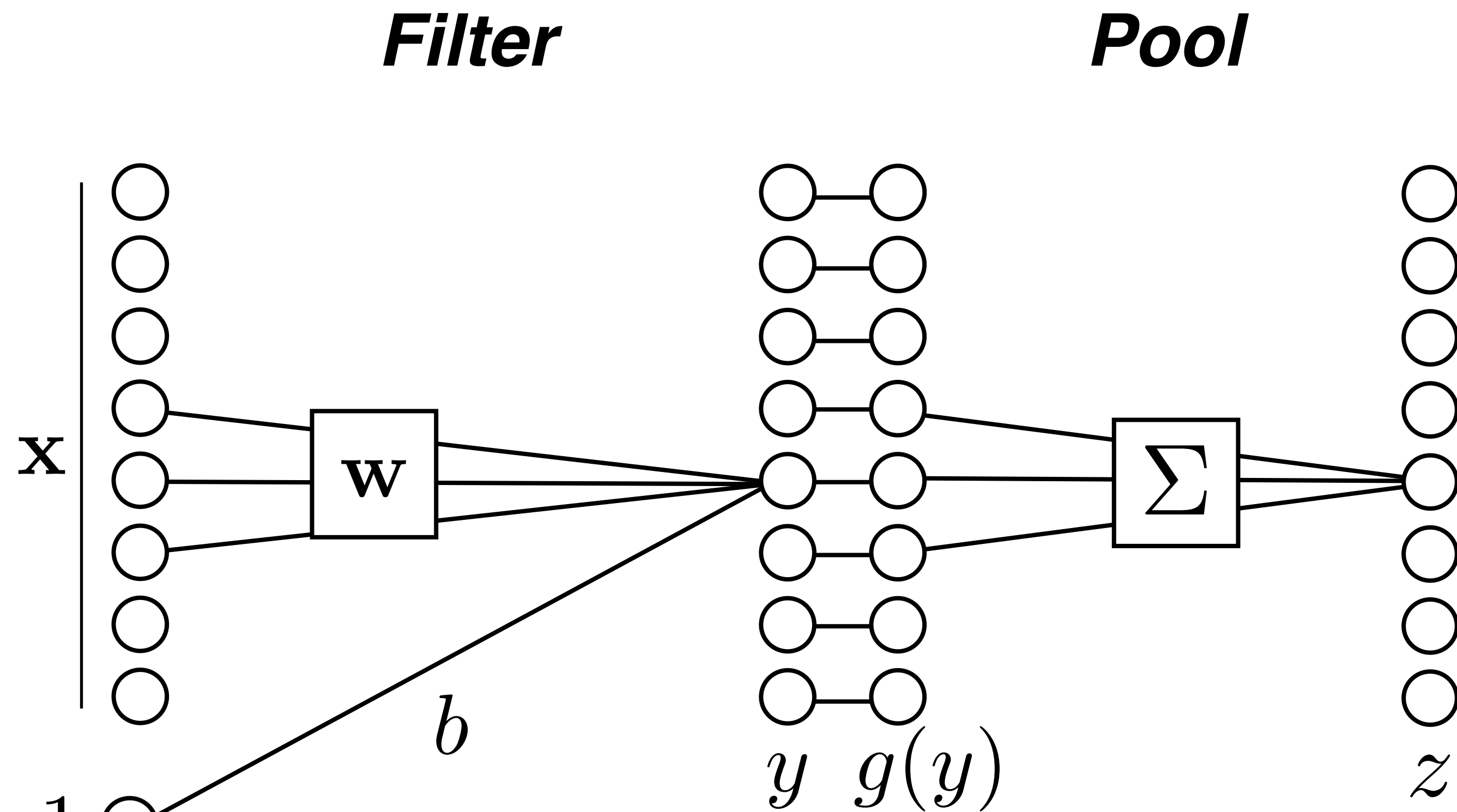
# Pooling



## Max pooling

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

# Pooling



## Max pooling

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

## Mean pooling

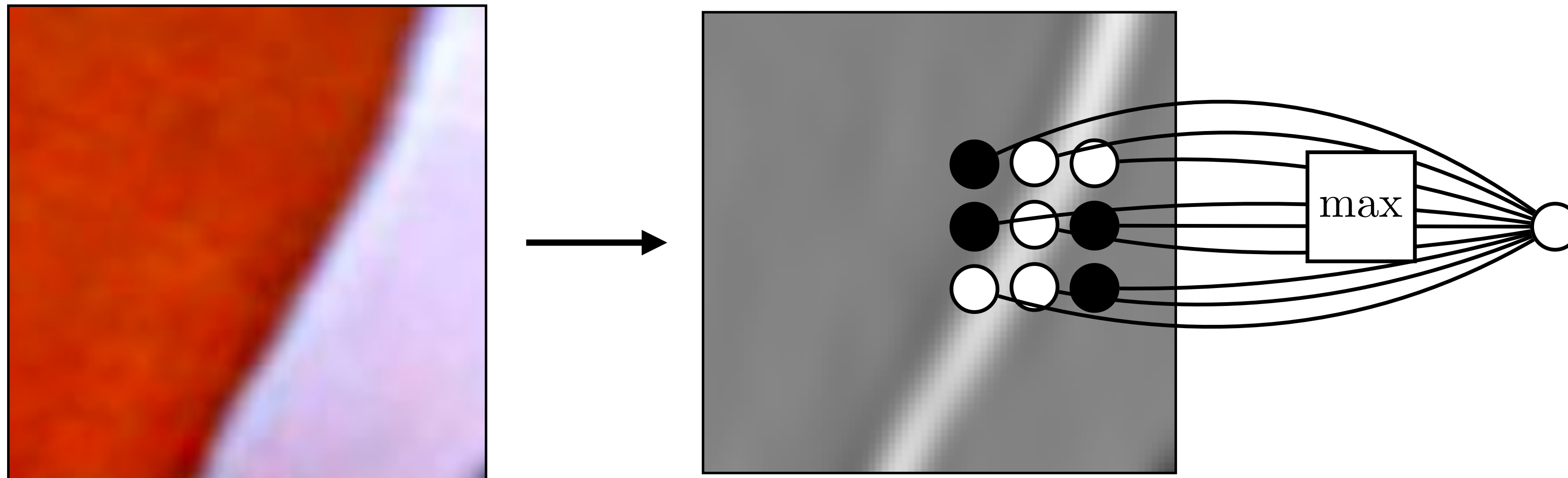
$$z_k = \frac{1}{|\mathcal{N}(k)|} \sum_{j \in \mathcal{N}(k)} g(y_j)$$

## 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?

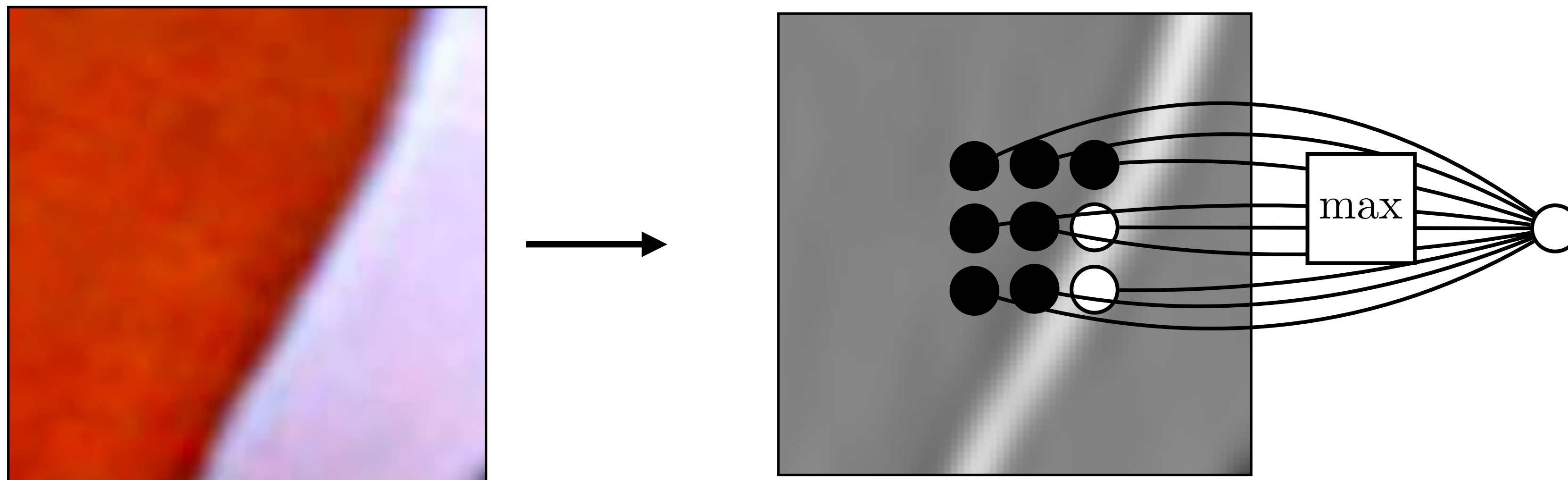
Pooling across spatial locations achieves stability w.r.t. small translations:





# Pooling — Why?

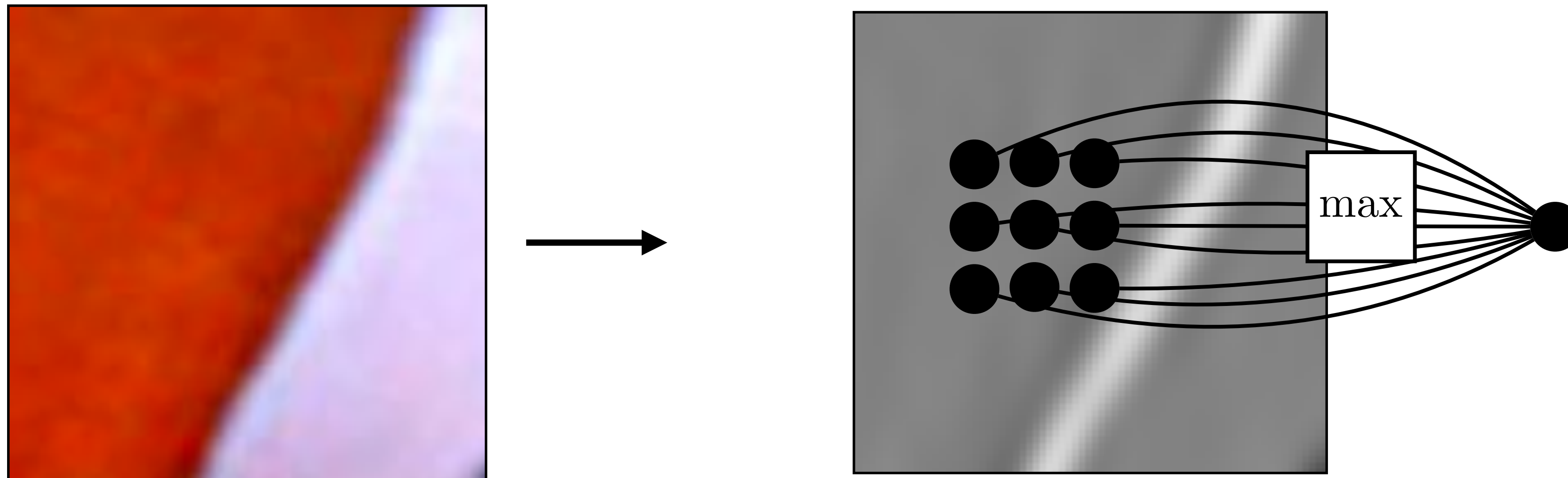
Pooling across spatial locations achieves stability w.r.t. small translations:



same large response  
regardless of exact  
position of edge

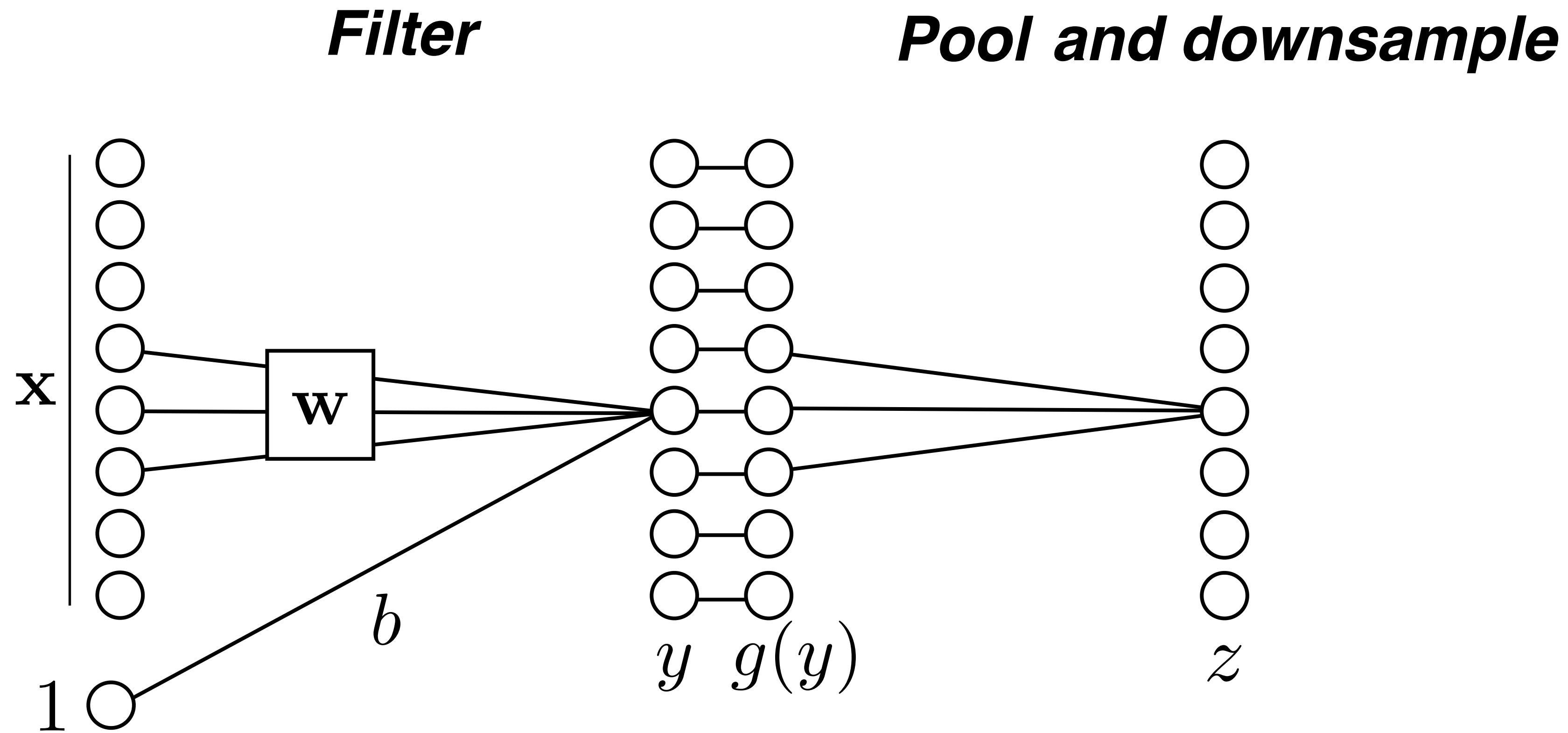
# Pooling — Why?

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

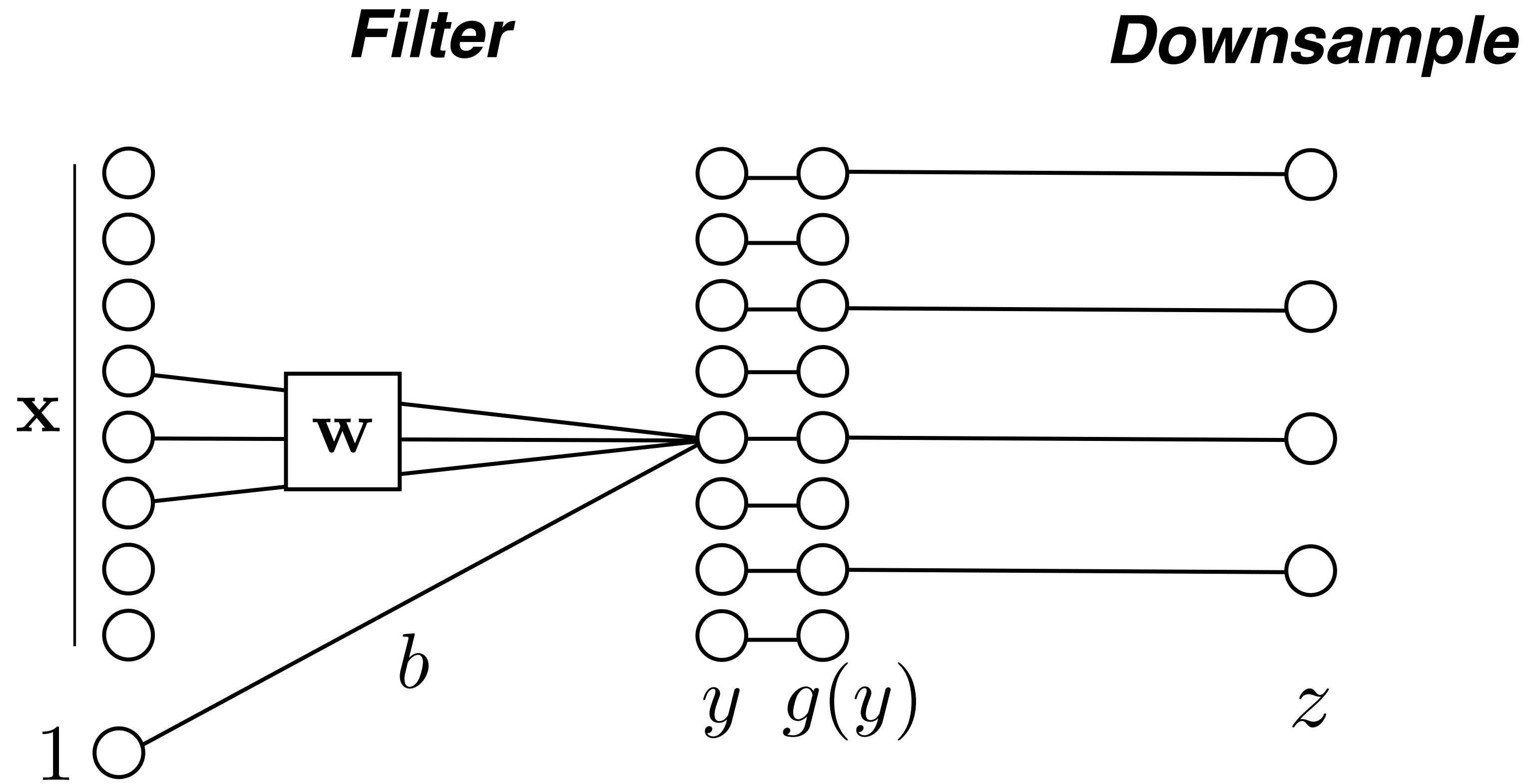


however, if the image is translated a lot...

# Downsampling

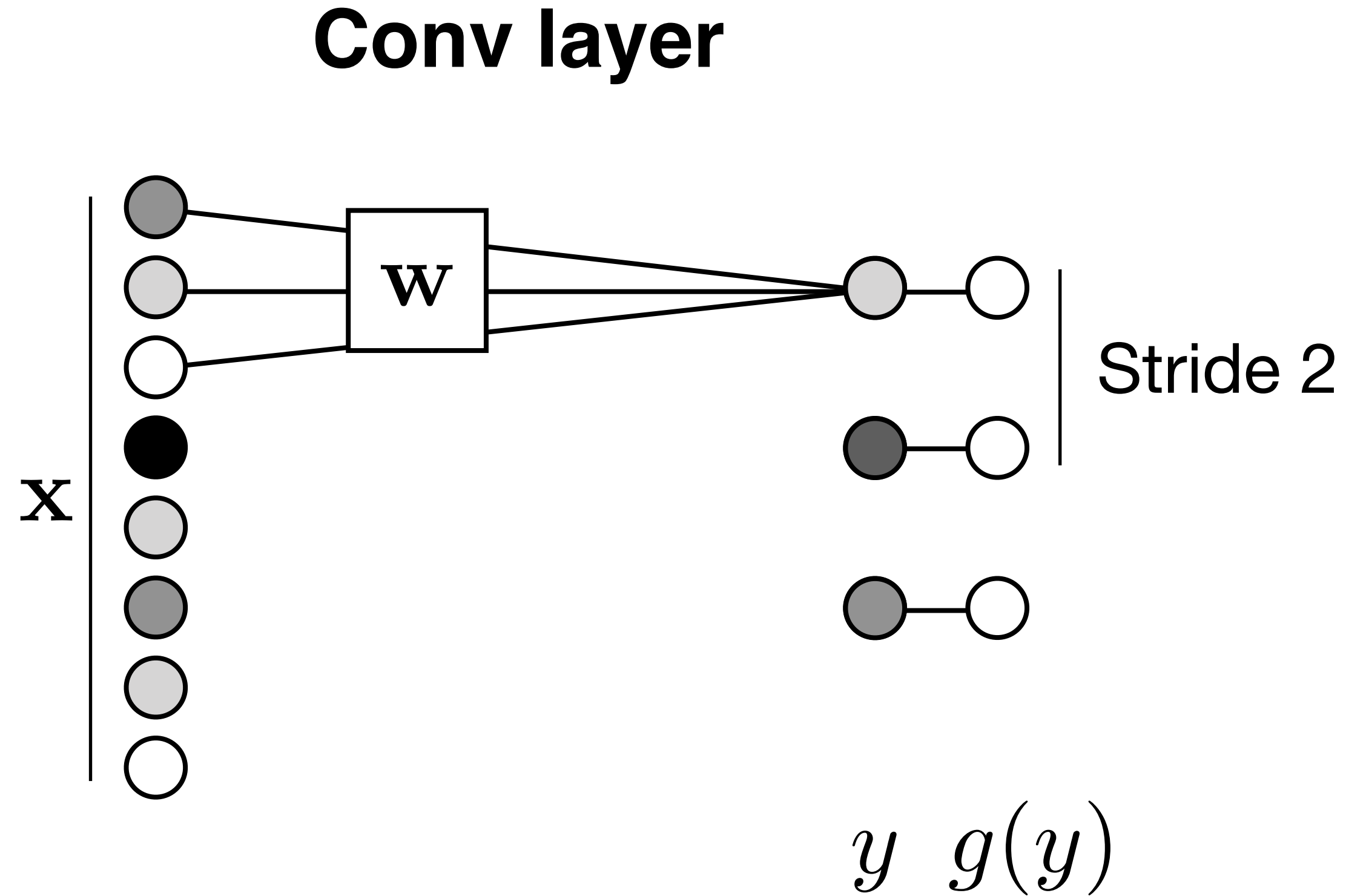


# Downsampling



$$\mathbb{R}^{H^{(l)} \times W^{(l)} \times C^{(l)}} \rightarrow \mathbb{R}^{H^{(l+1)} \times W^{(l+1)} \times C^{(l+1)}}$$

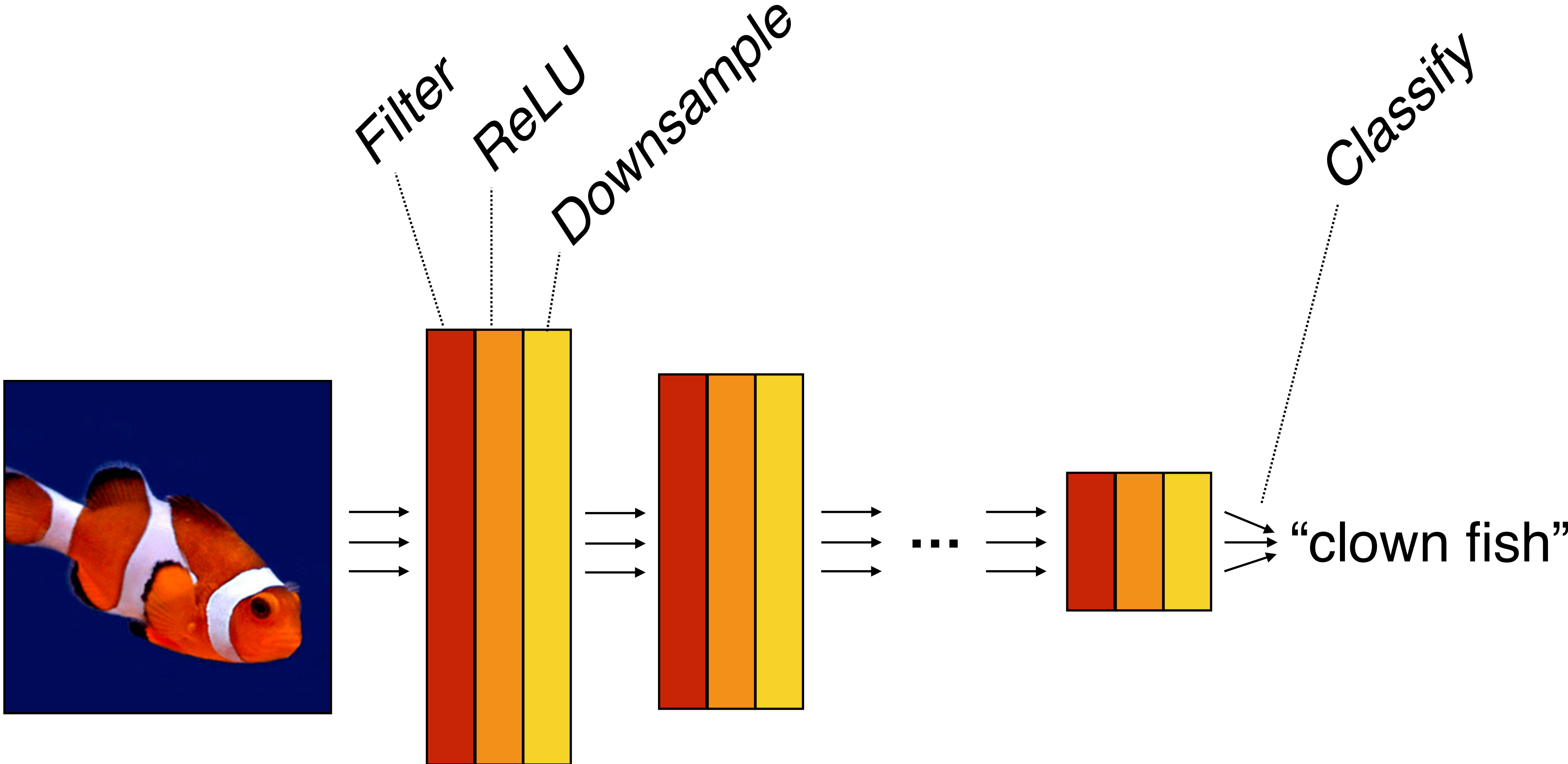
# Strided operations



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

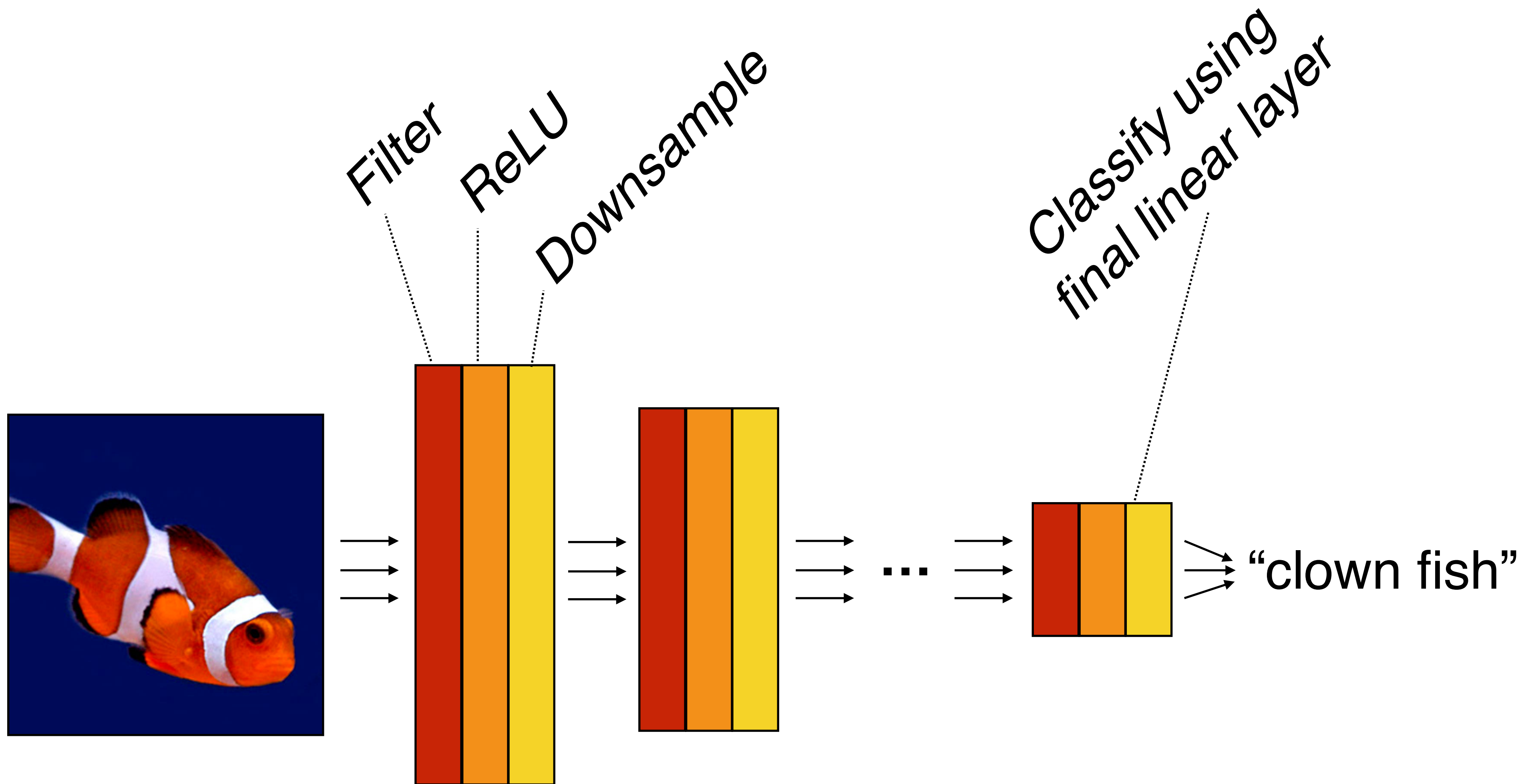
**Strided convolution is an alternative to pooling layers:**  
just do a strided convolution!

# Computation in a neural net



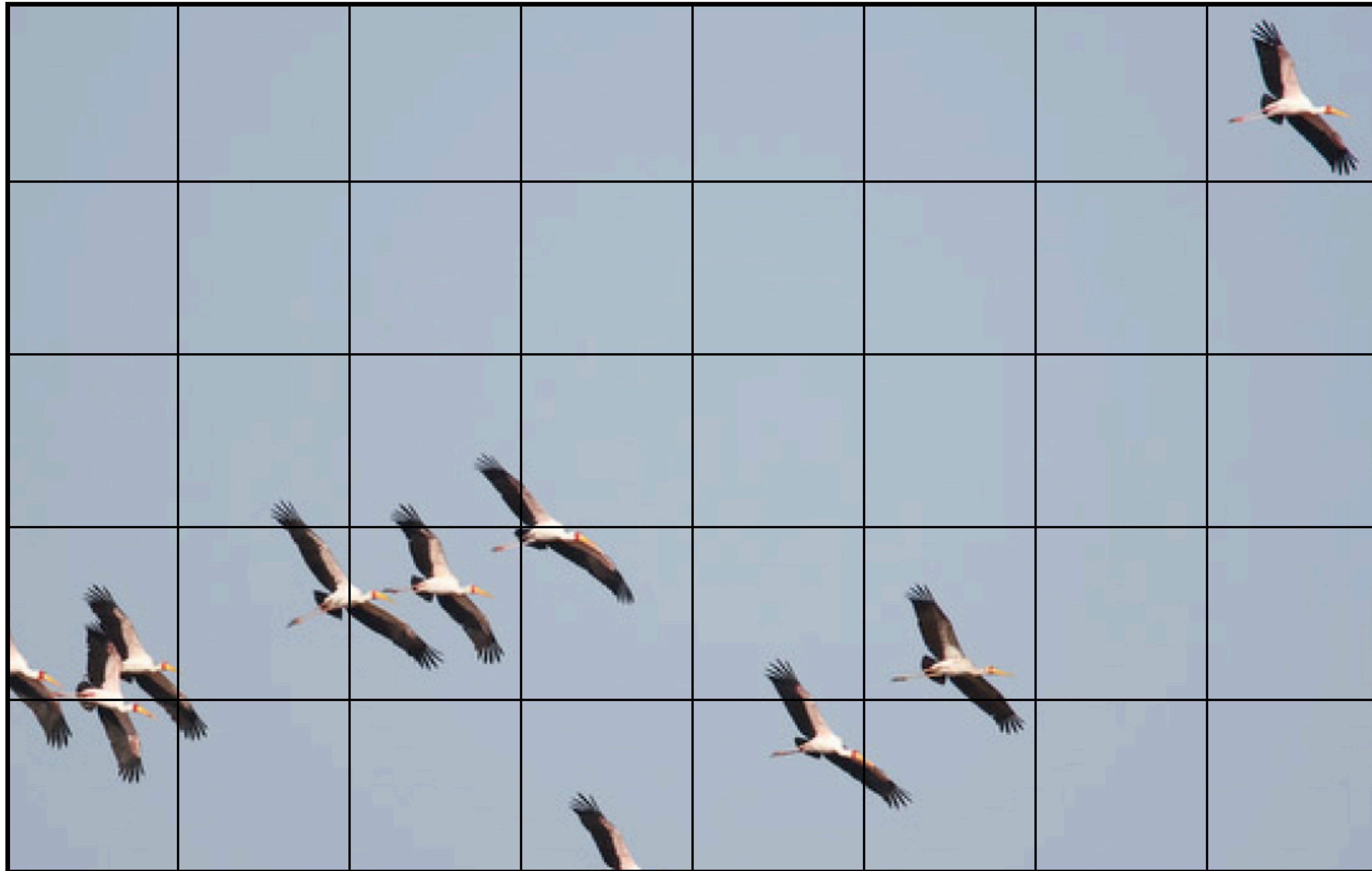
$$f(\mathbf{x}) = f_L(\dots f_2(f_1(\mathbf{x})))$$

# Computation in a neural net



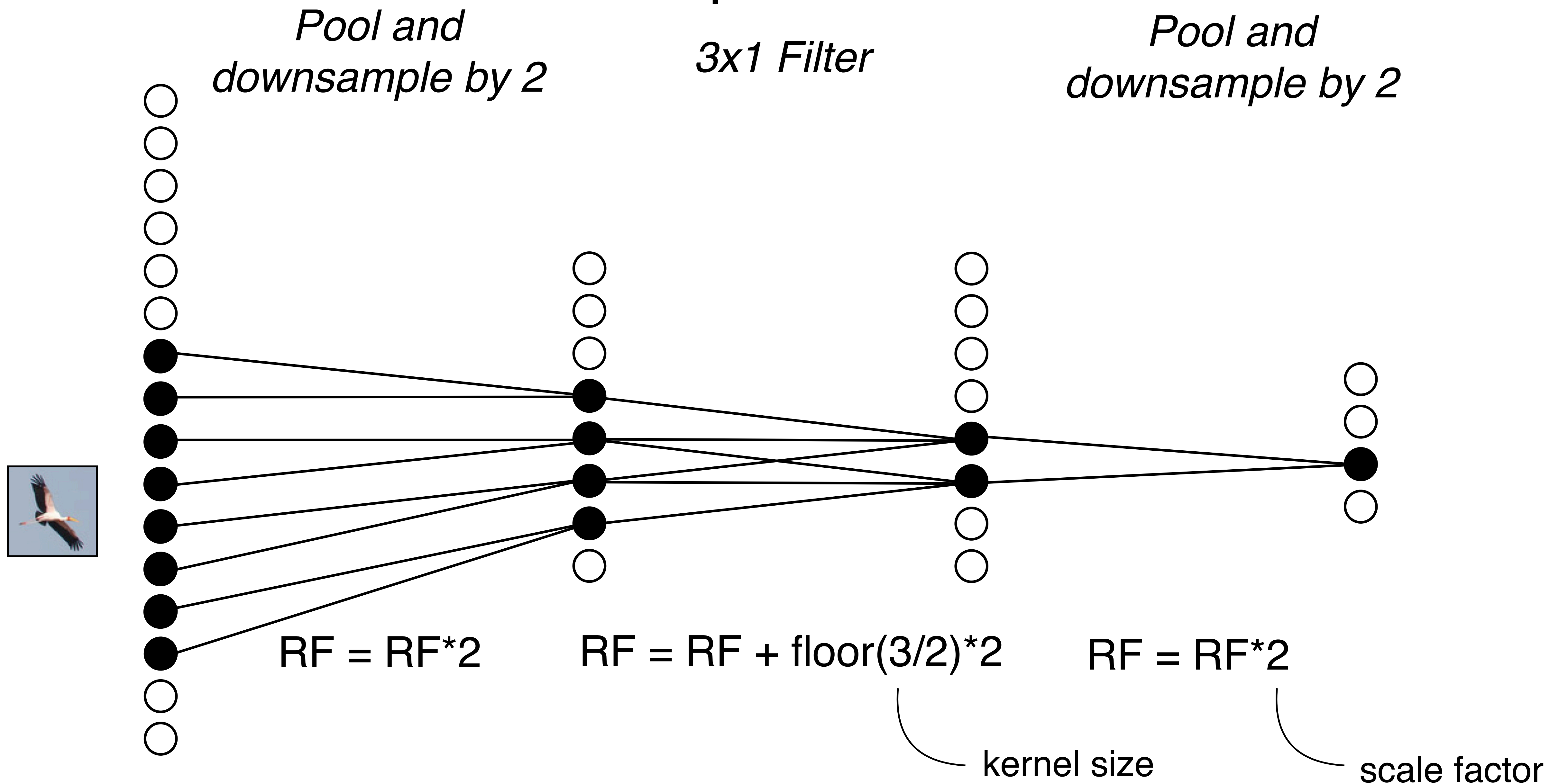
$$f(\mathbf{x}) = f_L(\dots f_2(f_1(\mathbf{x})))$$

# Receptive fields

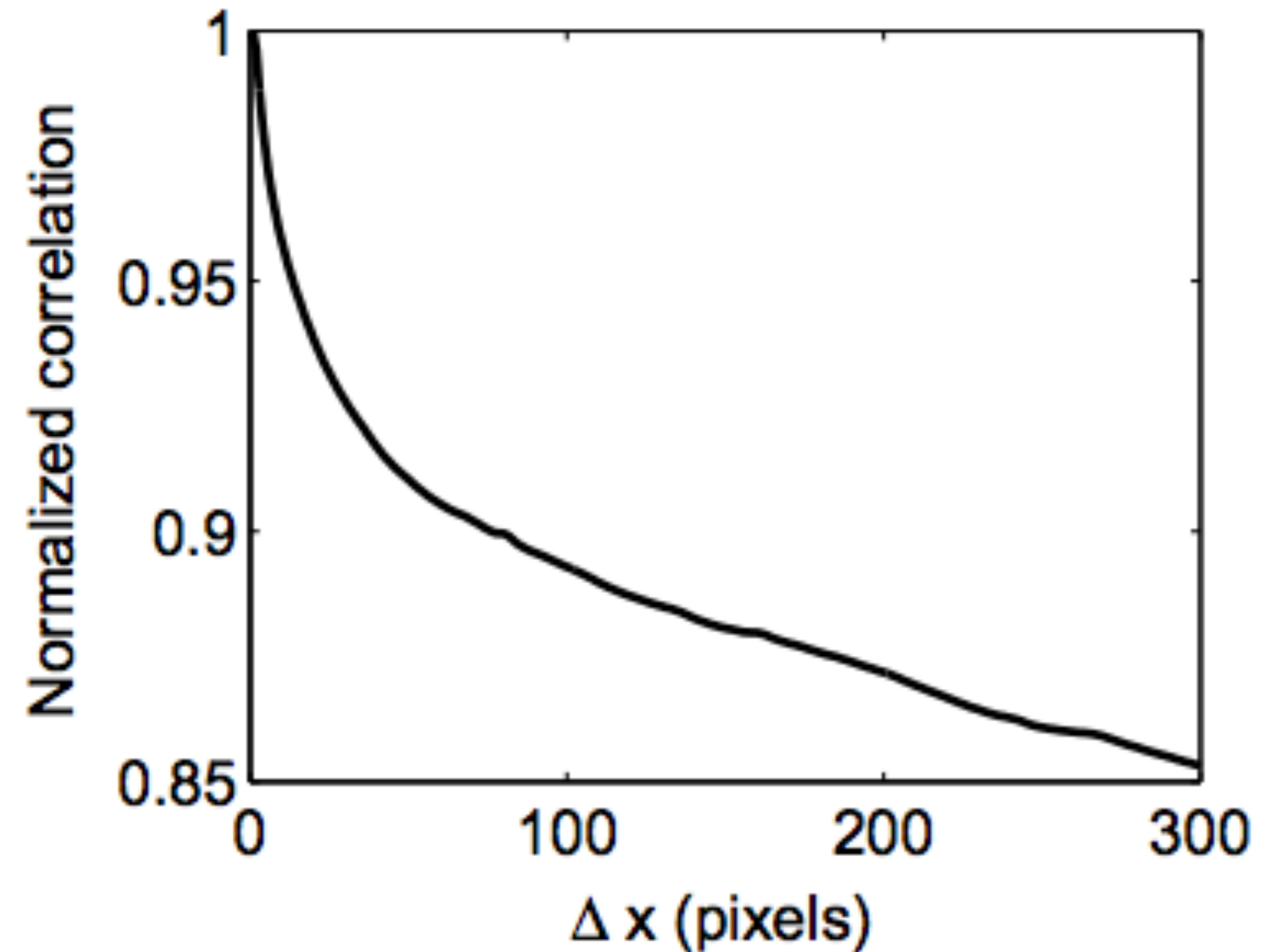
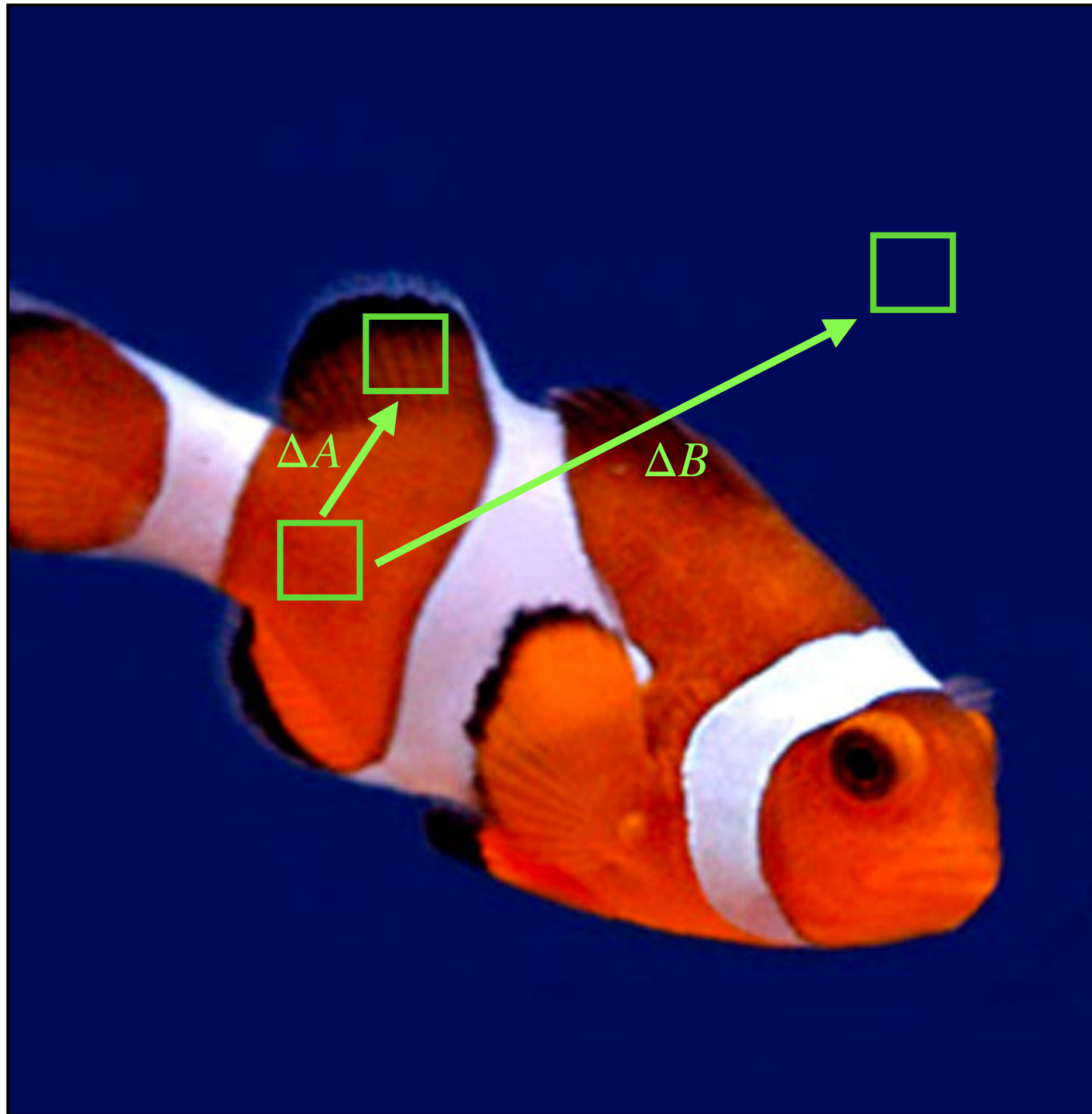




# Receptive fields



# Why have coefficients only for nearby pixels?



# Implementation details

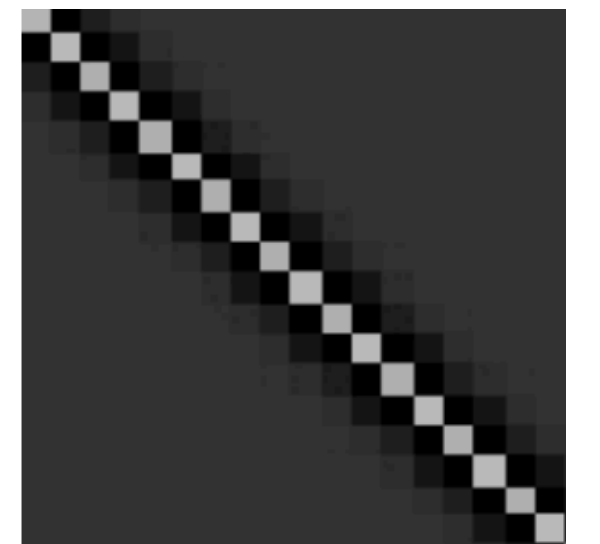
## Convolutional layers

**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.

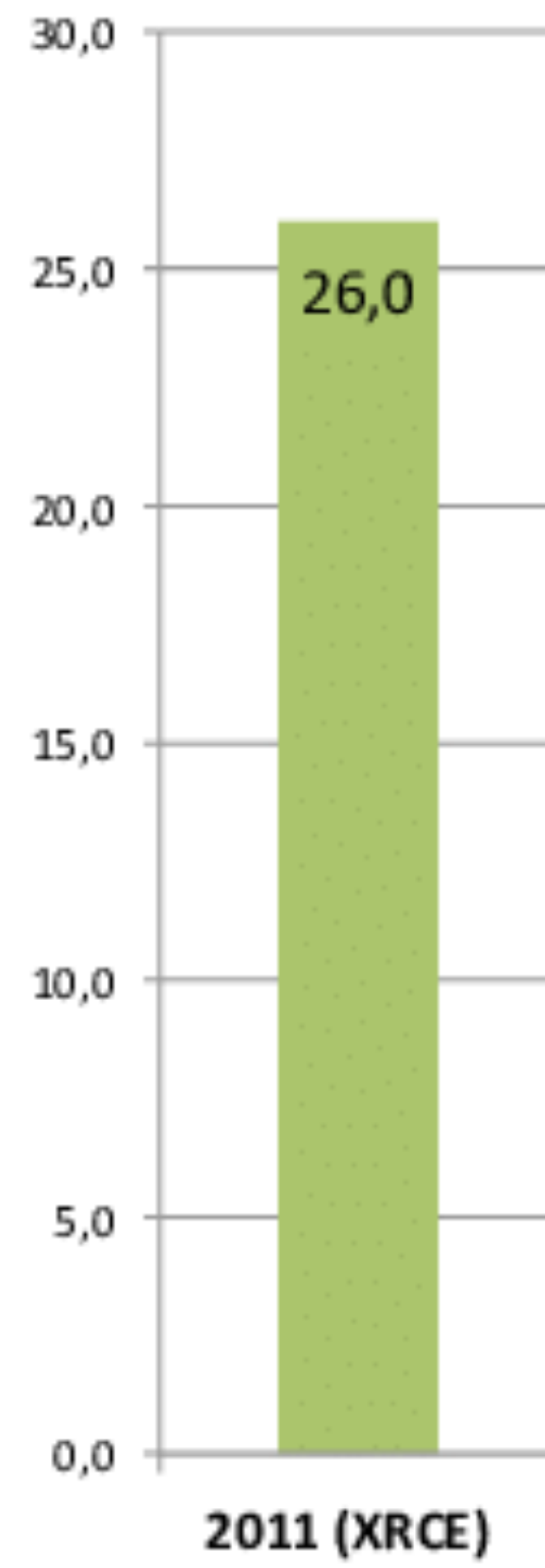
**Backwards pass:** These take the form of convolutions.

**Max pooling:** backwards pass similar to ReLU.

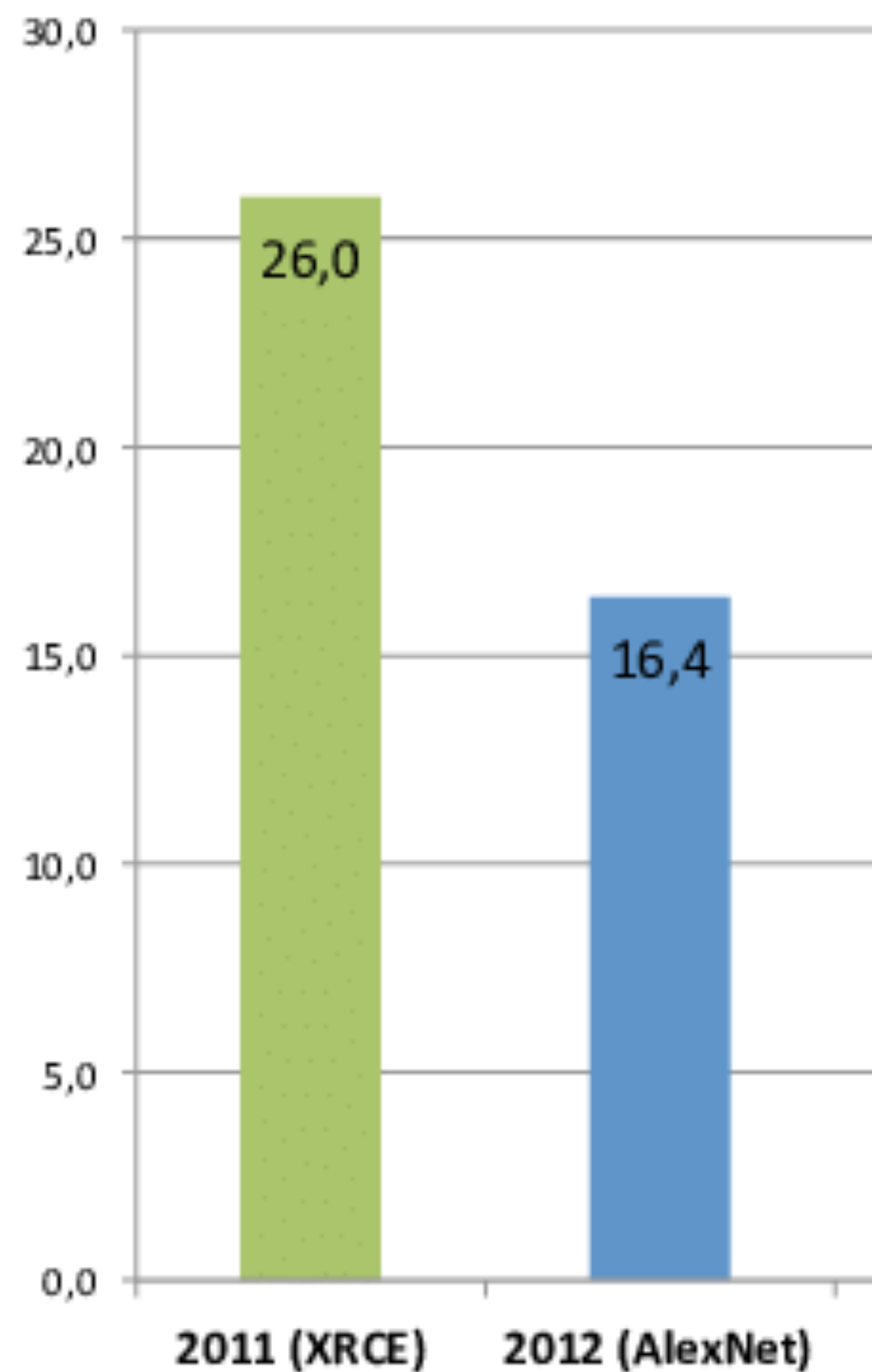


# Network designs

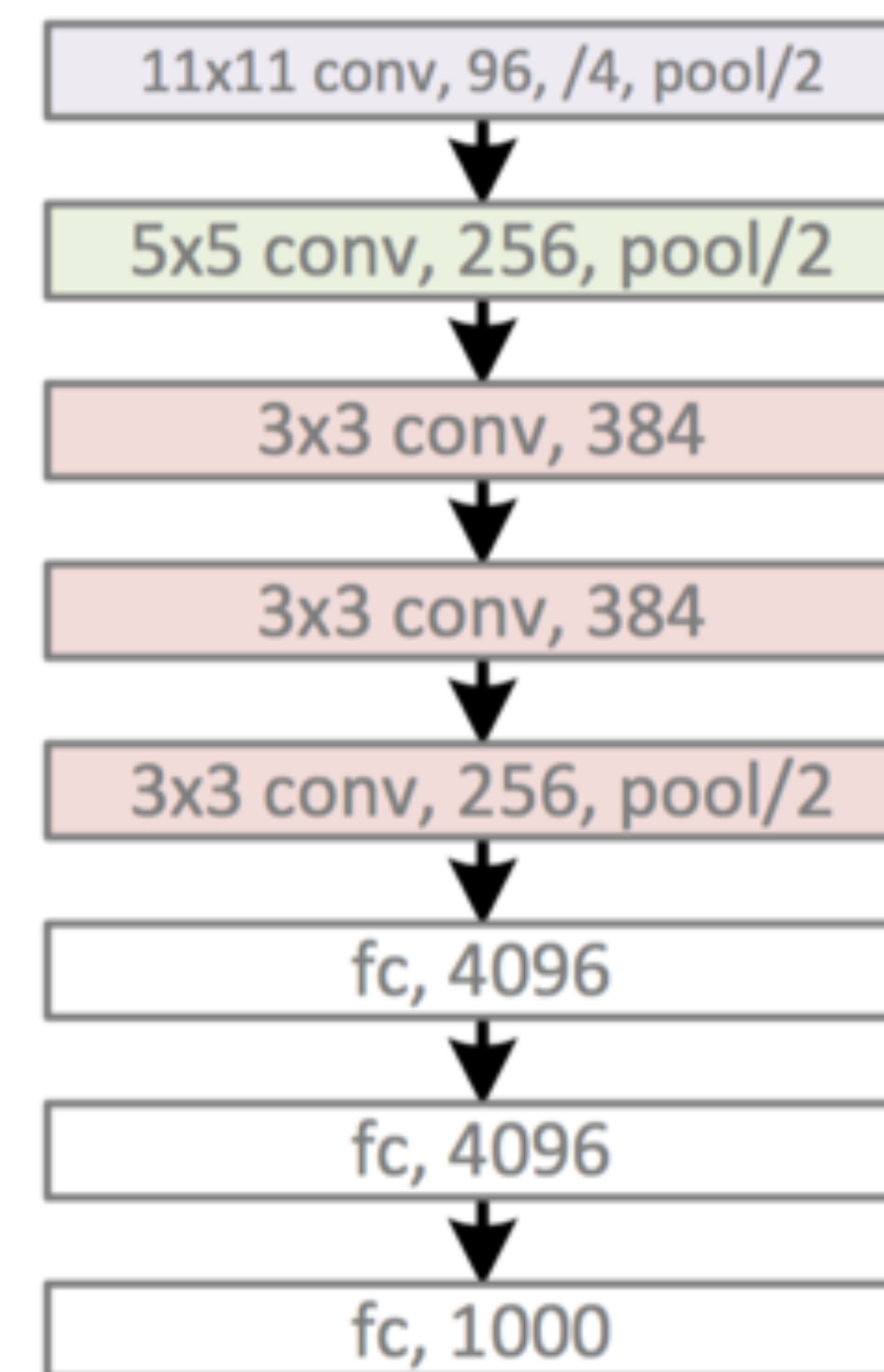
## ImageNet classification (top 5)



ImageNet classification (top 5)



2012: AlexNet  
5 conv. layers

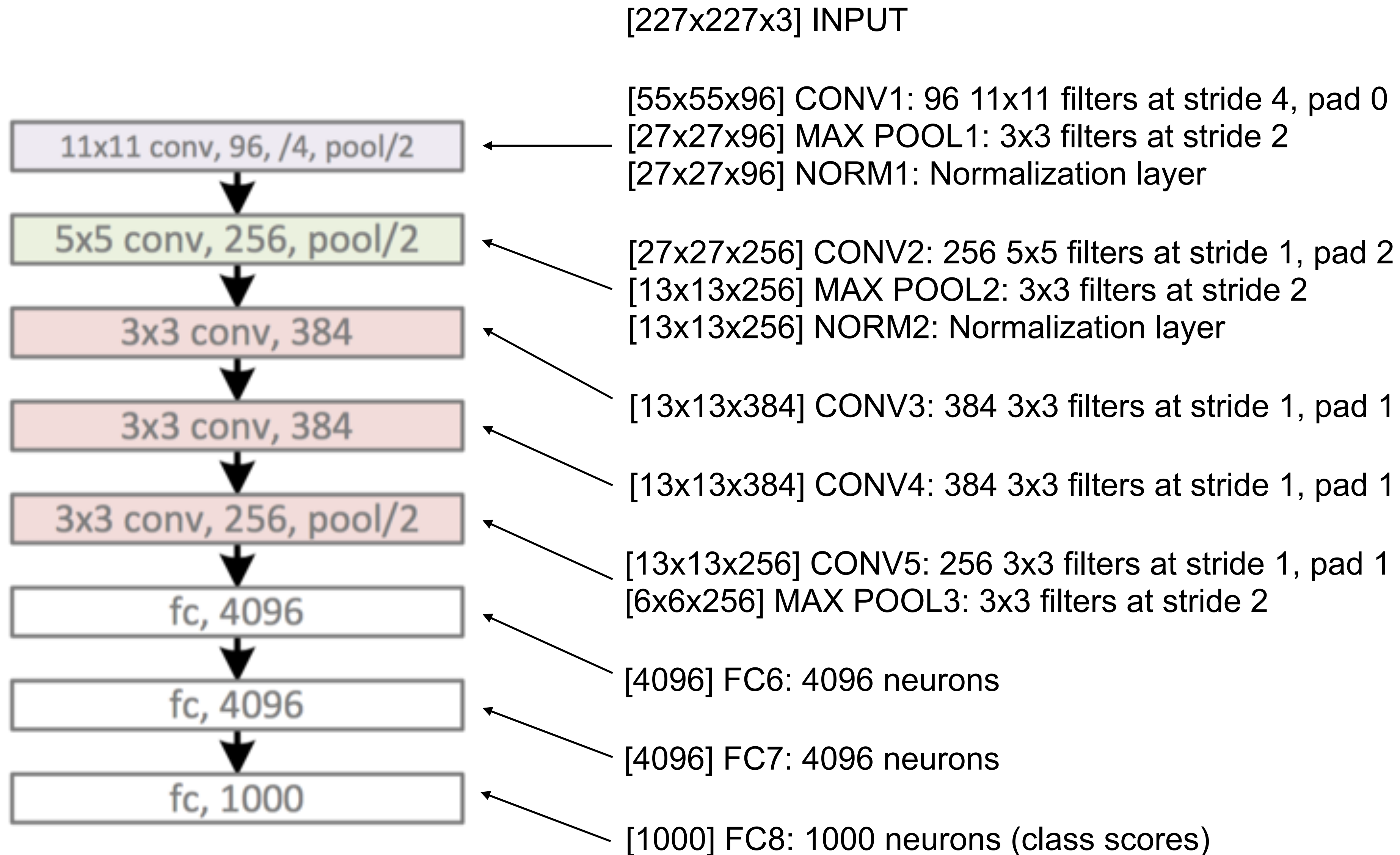


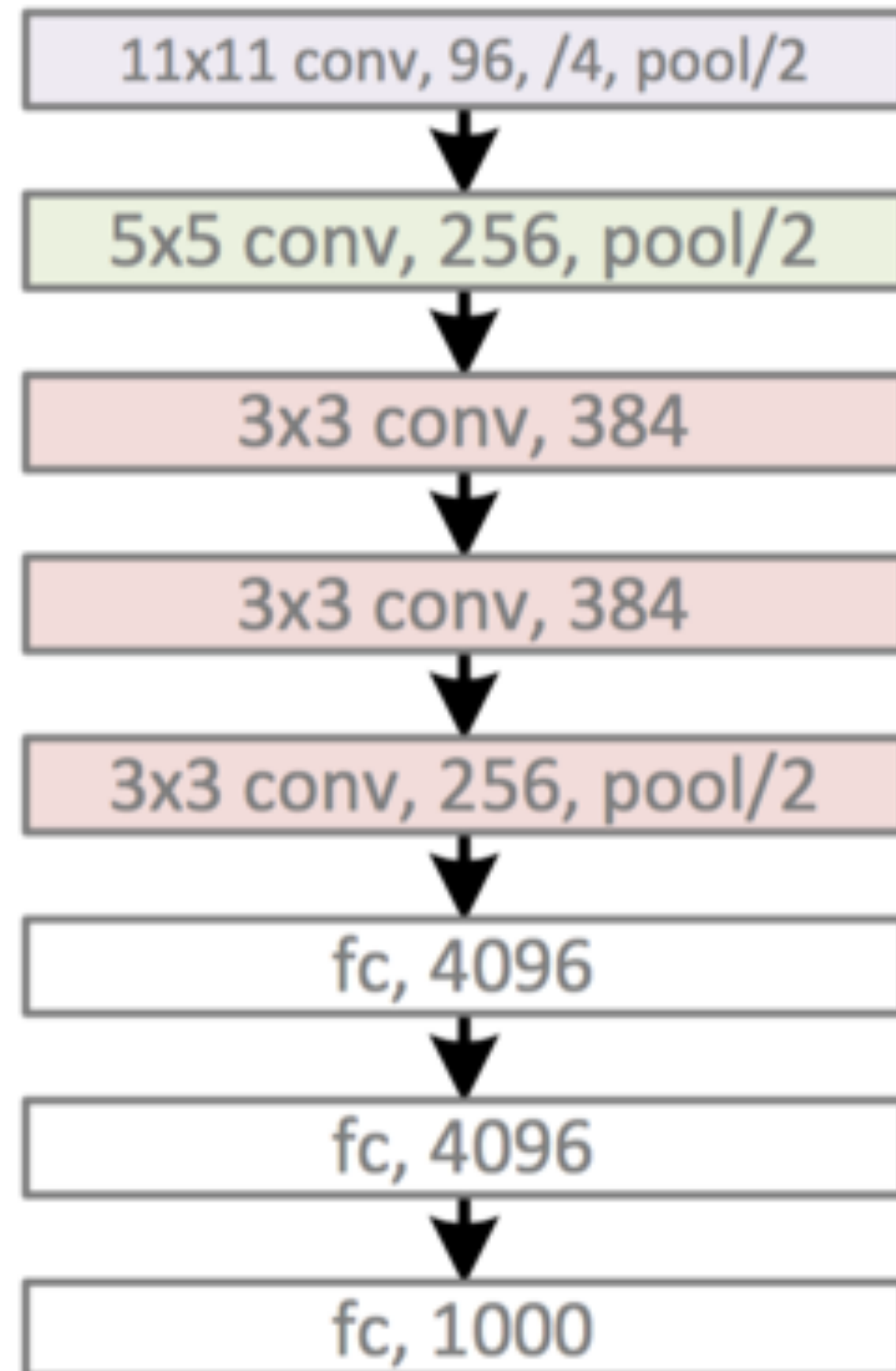
Error: 16.4%

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



# Alexnet — [Krizhevsky et al. NIPS 2012]



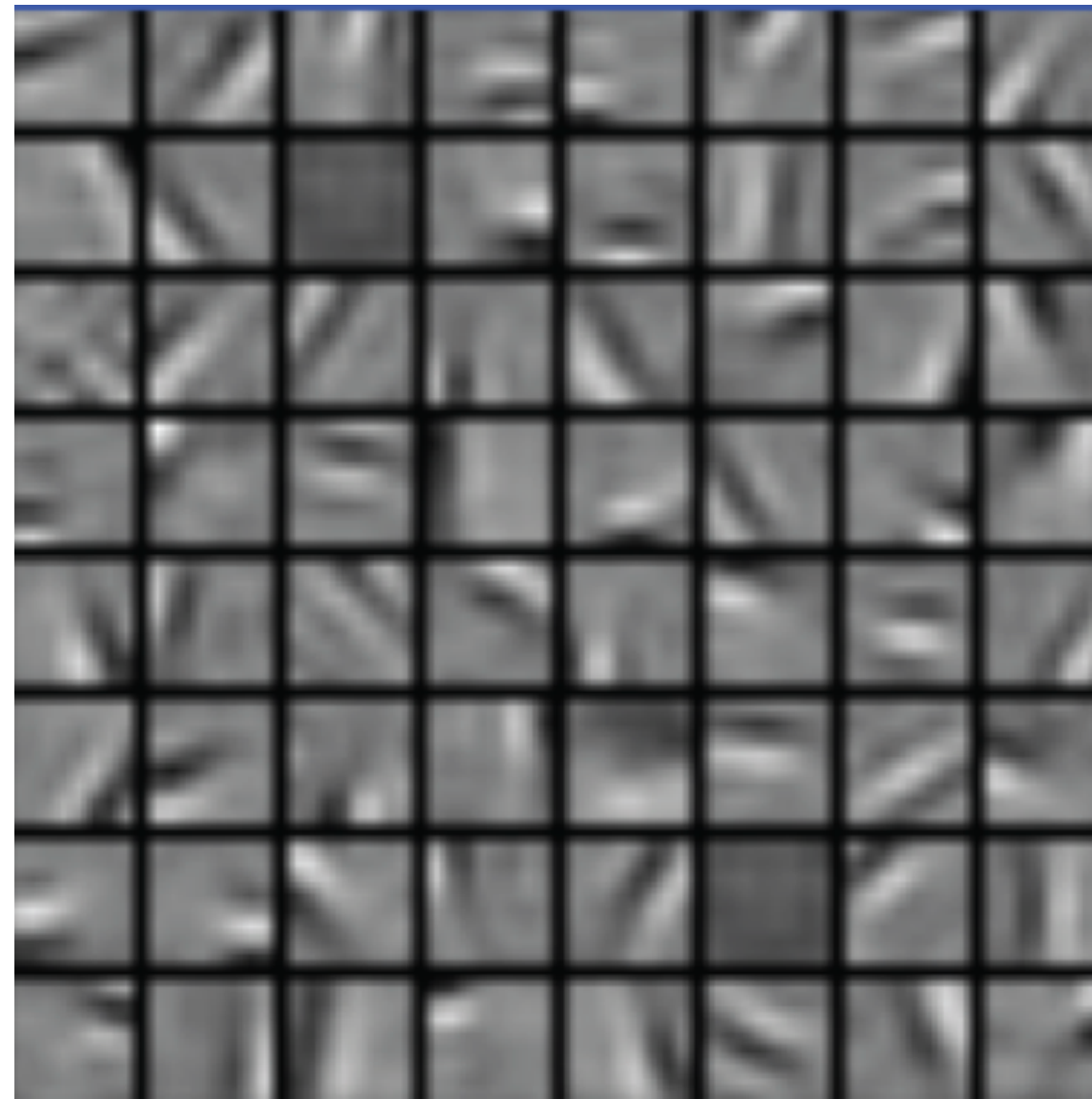


What filters are learned?

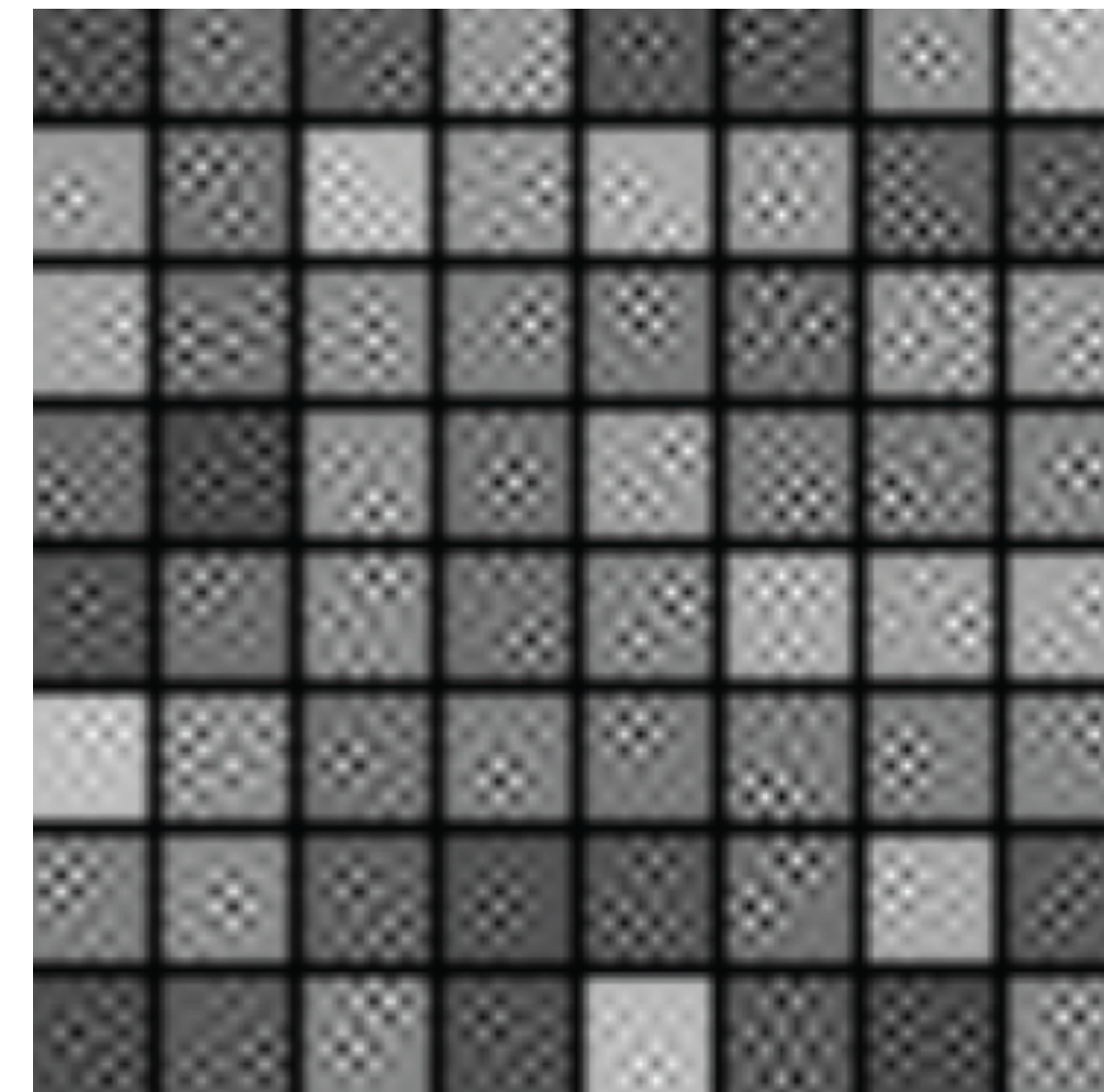


# What filters are learned?

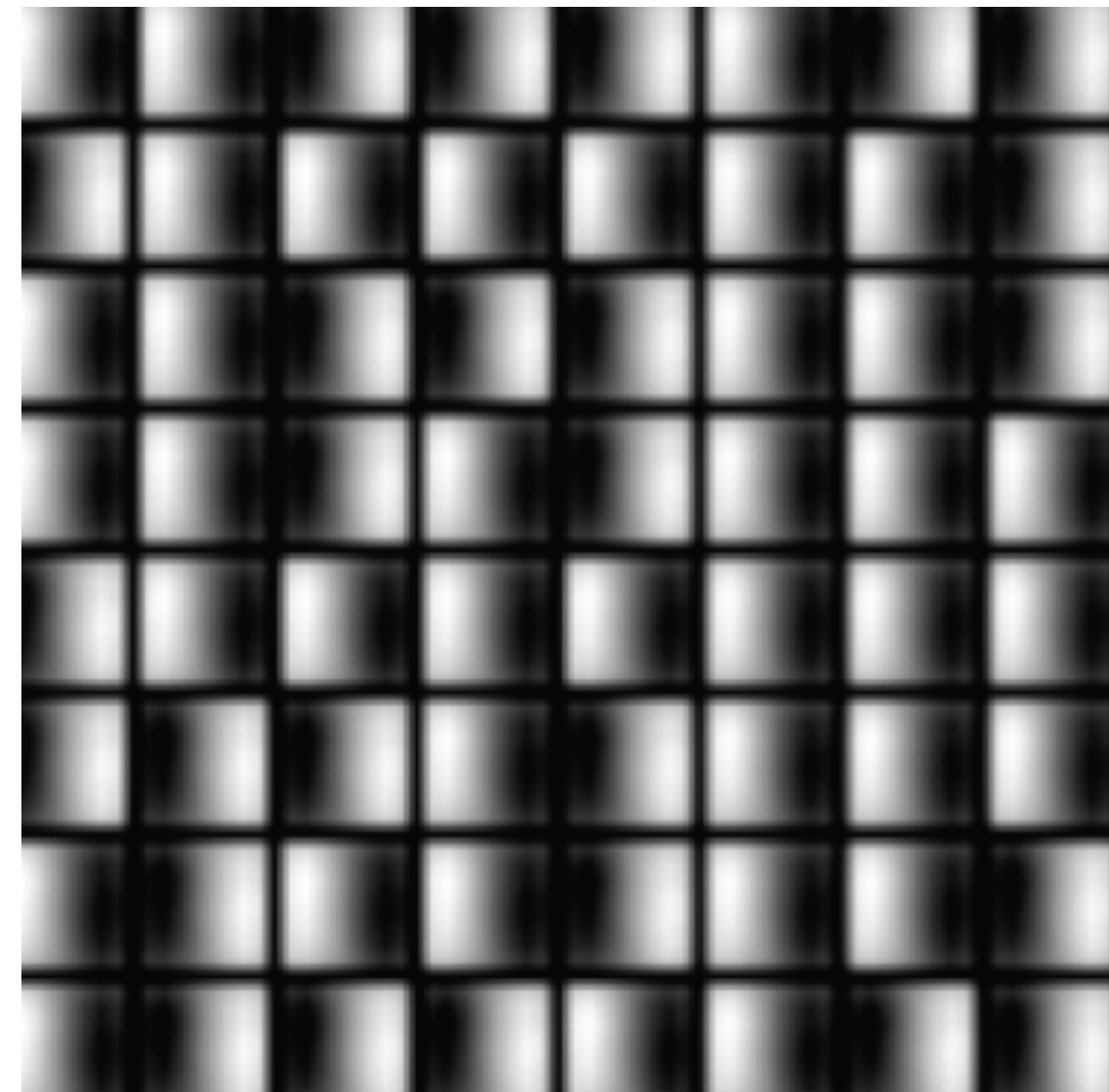
A



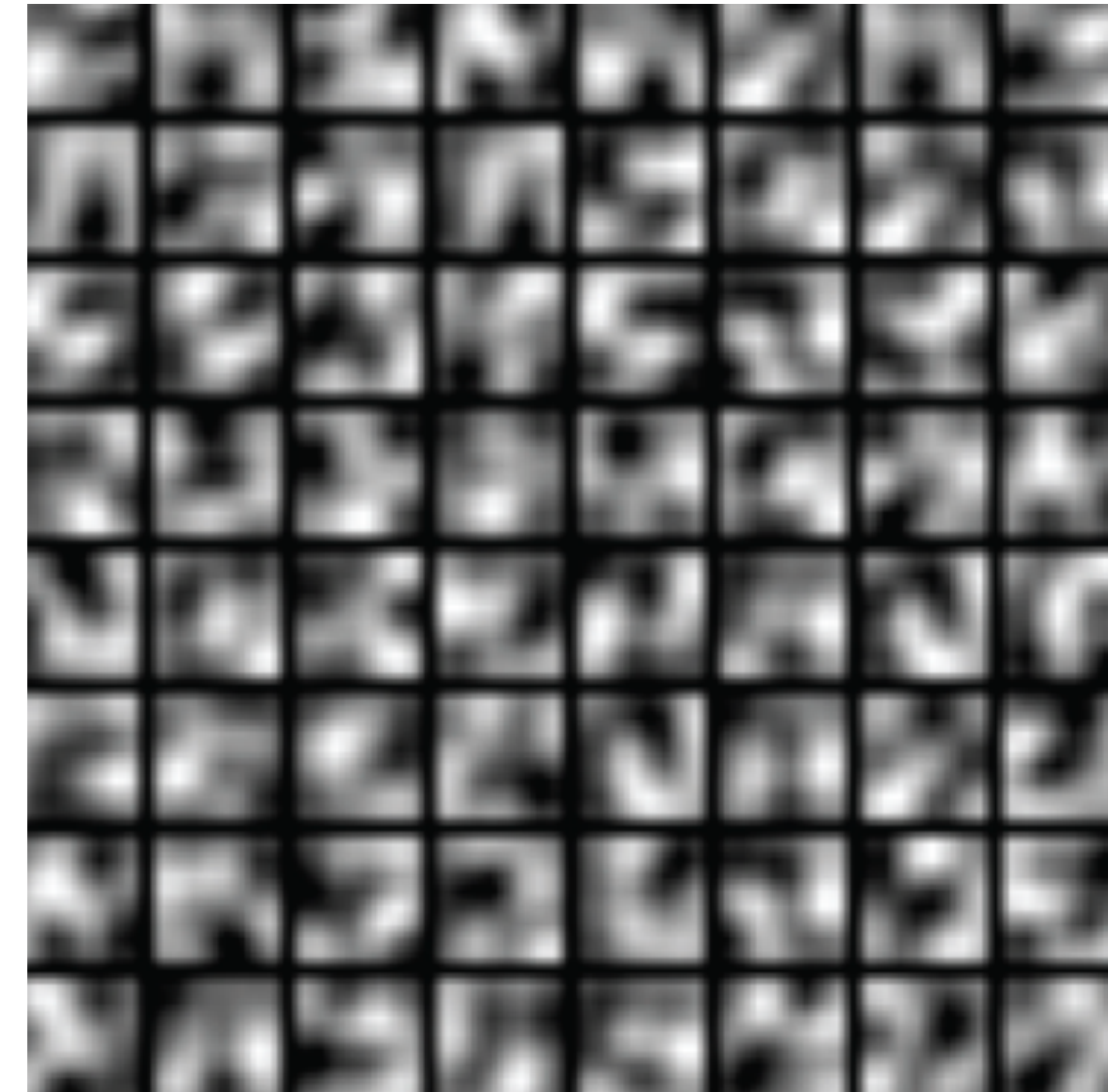
B



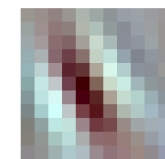
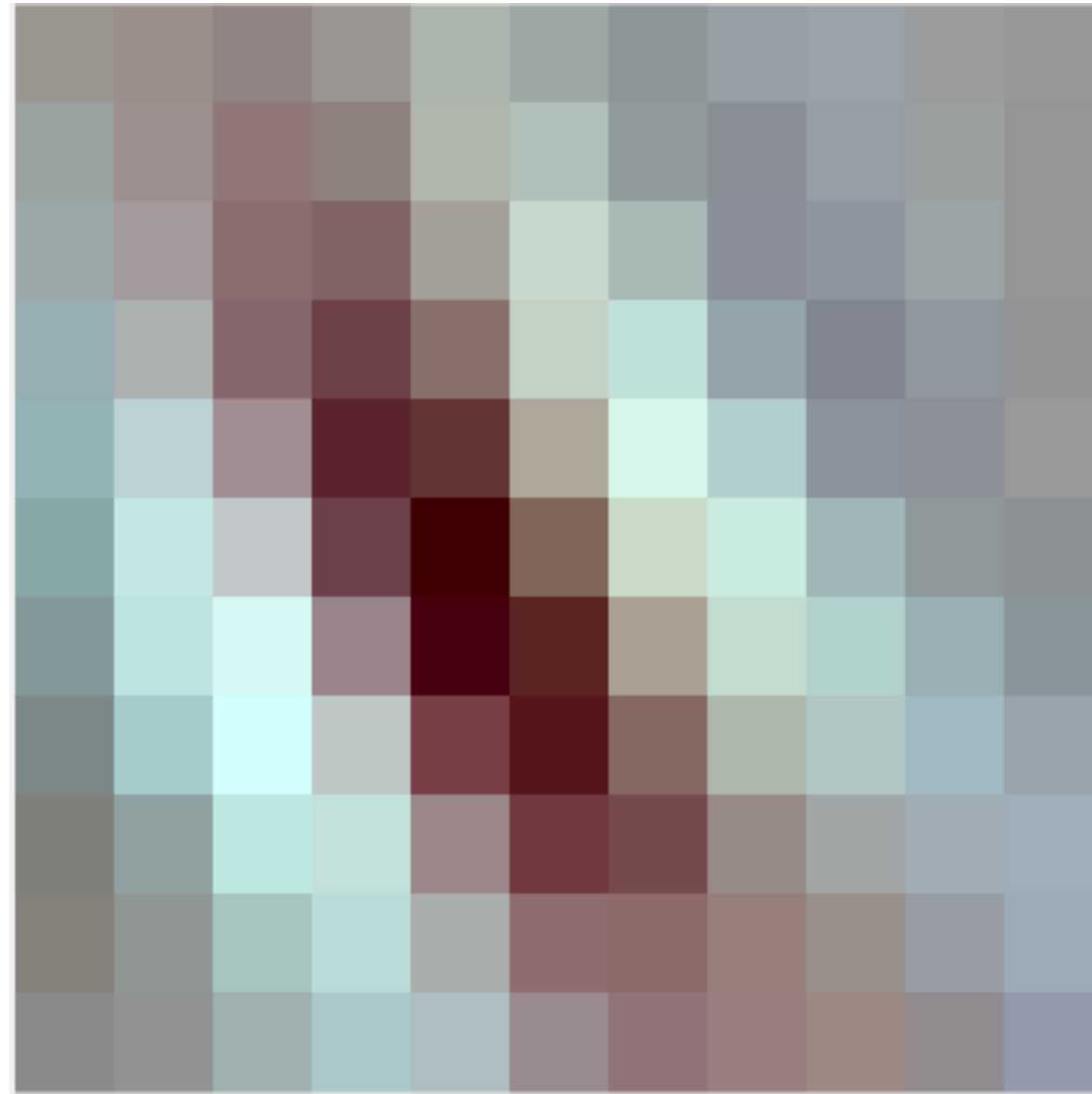
C



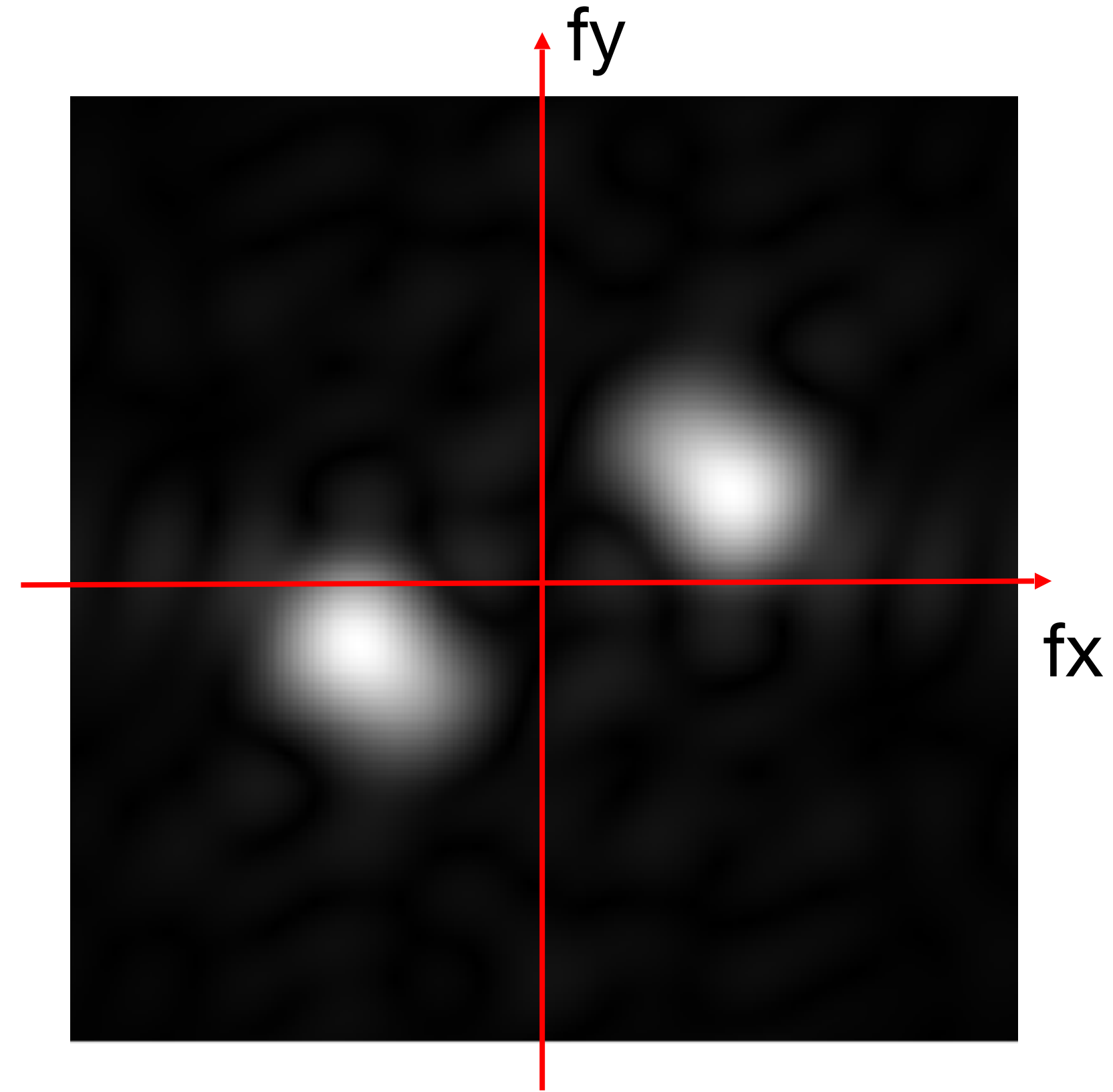
D



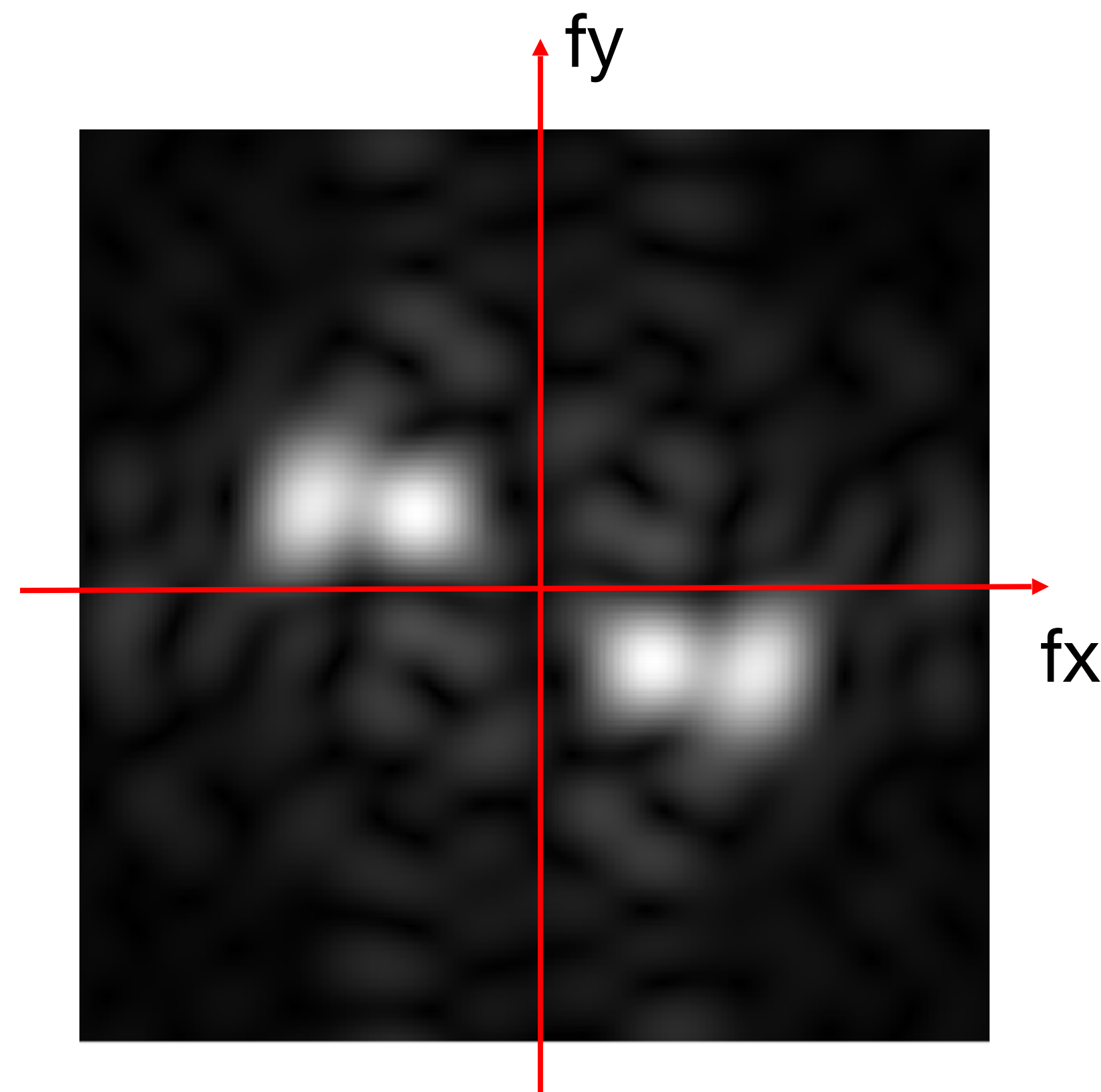
# Filters in first layer



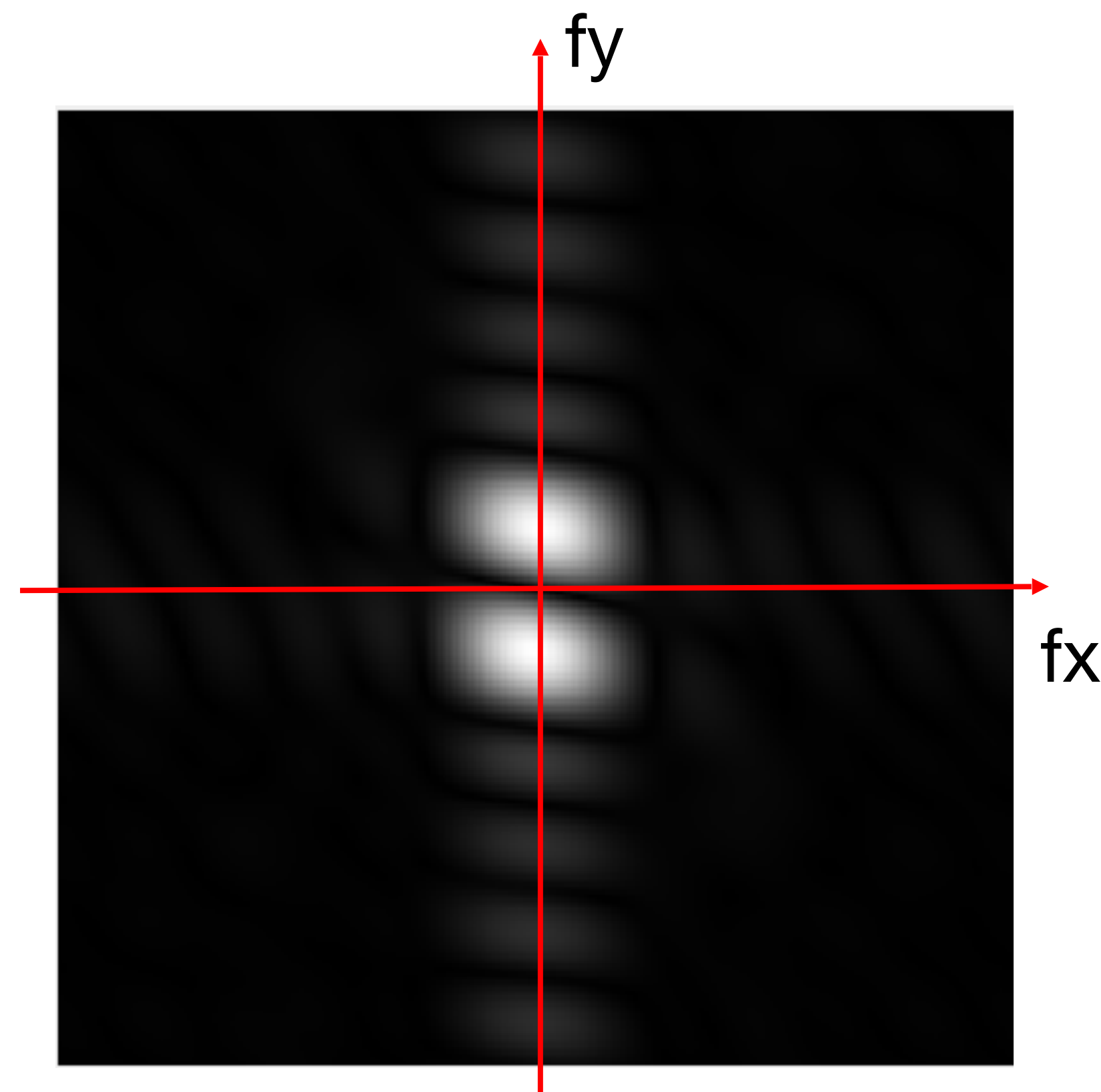
11x11 convolution kernel  
(3 color channels)



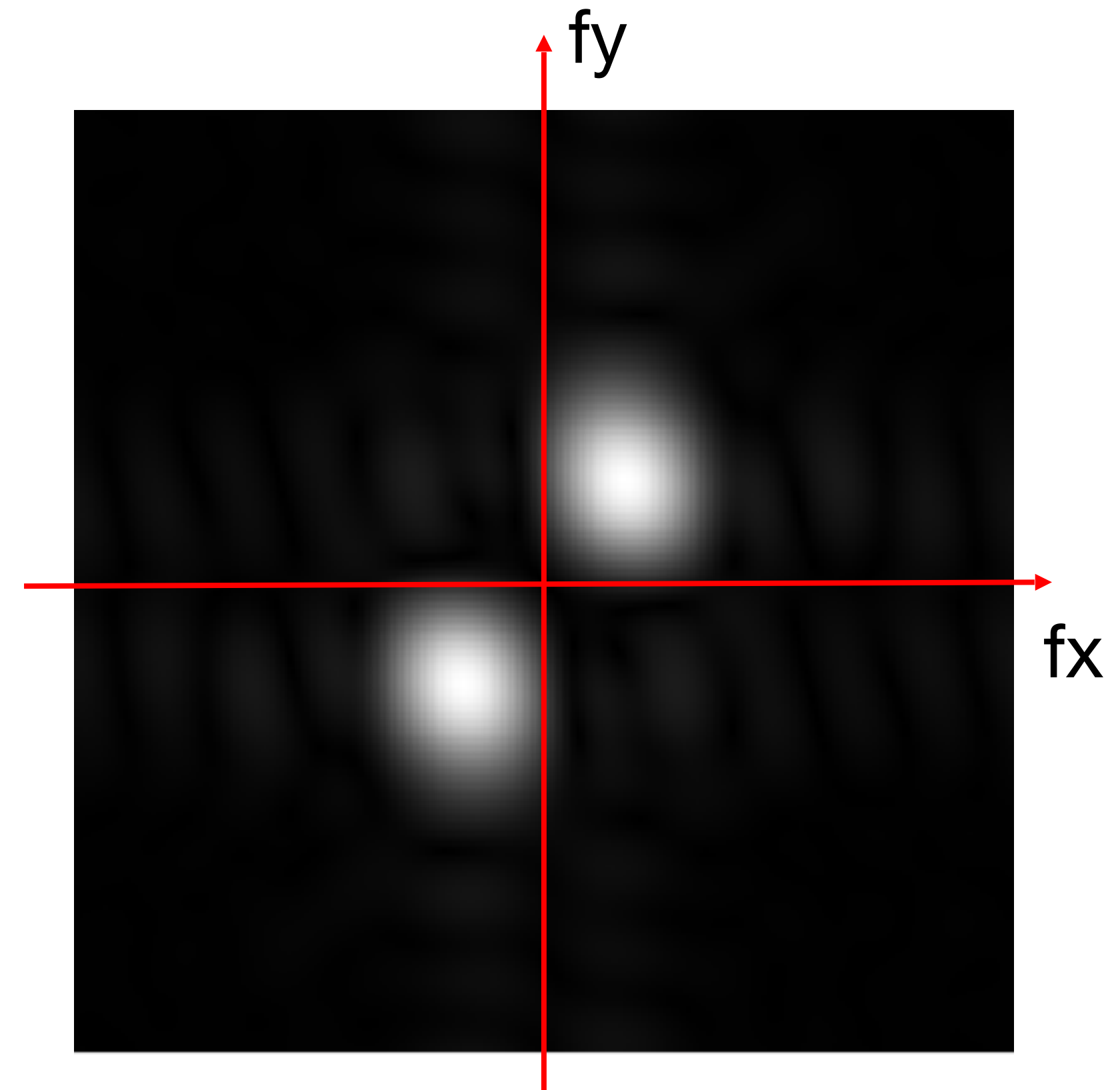
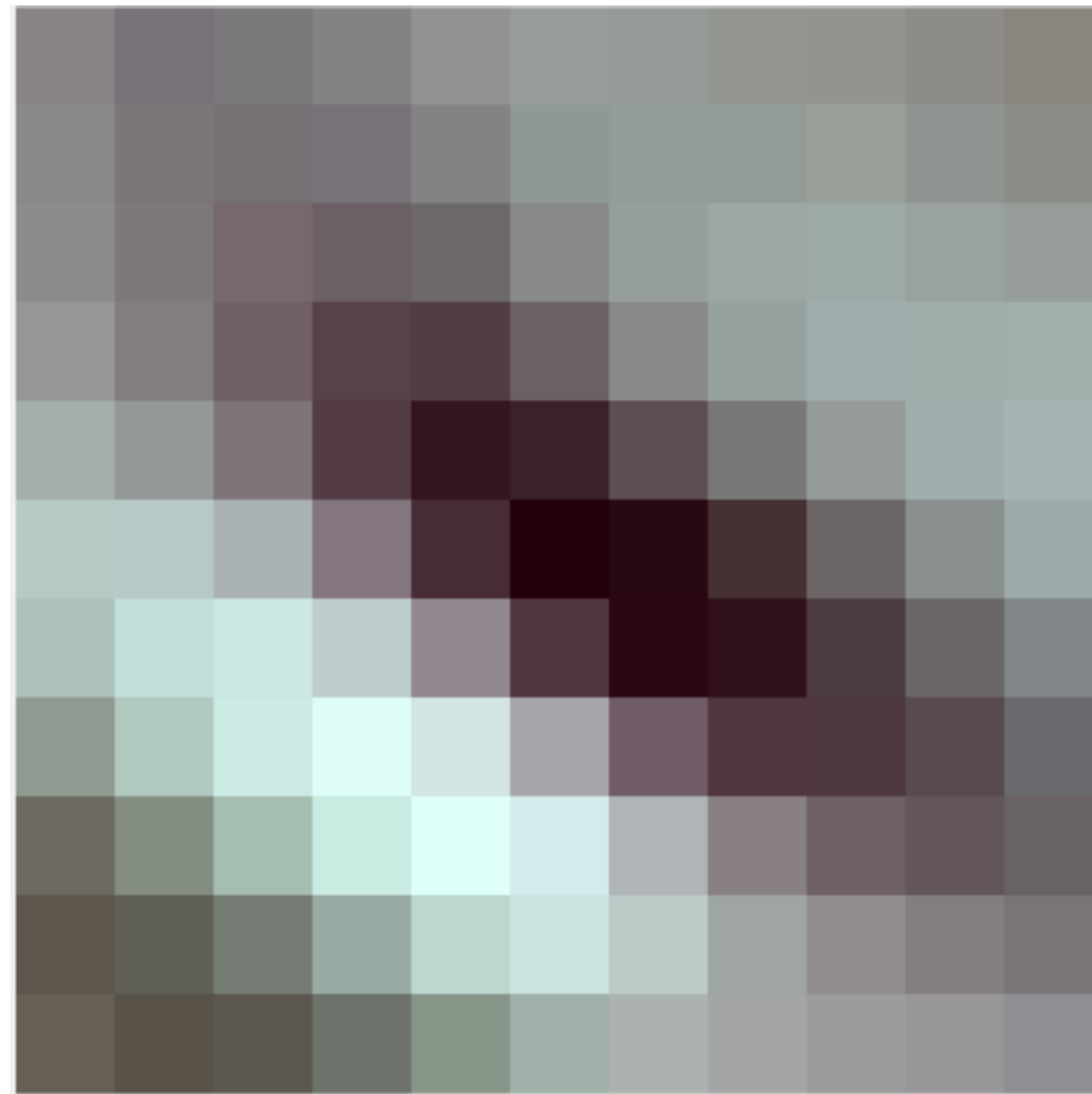
# Filters in first layer



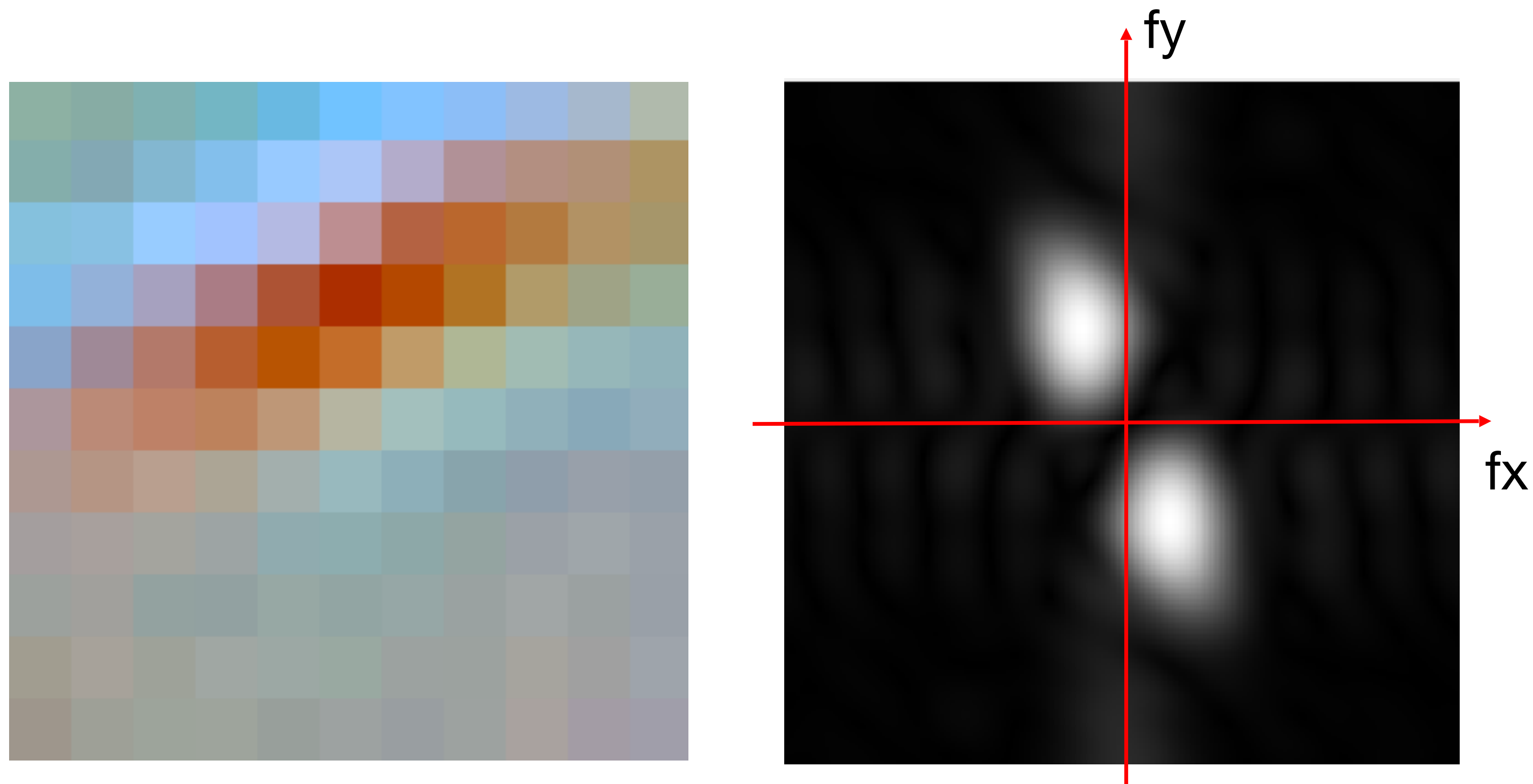
# Filters in first layer



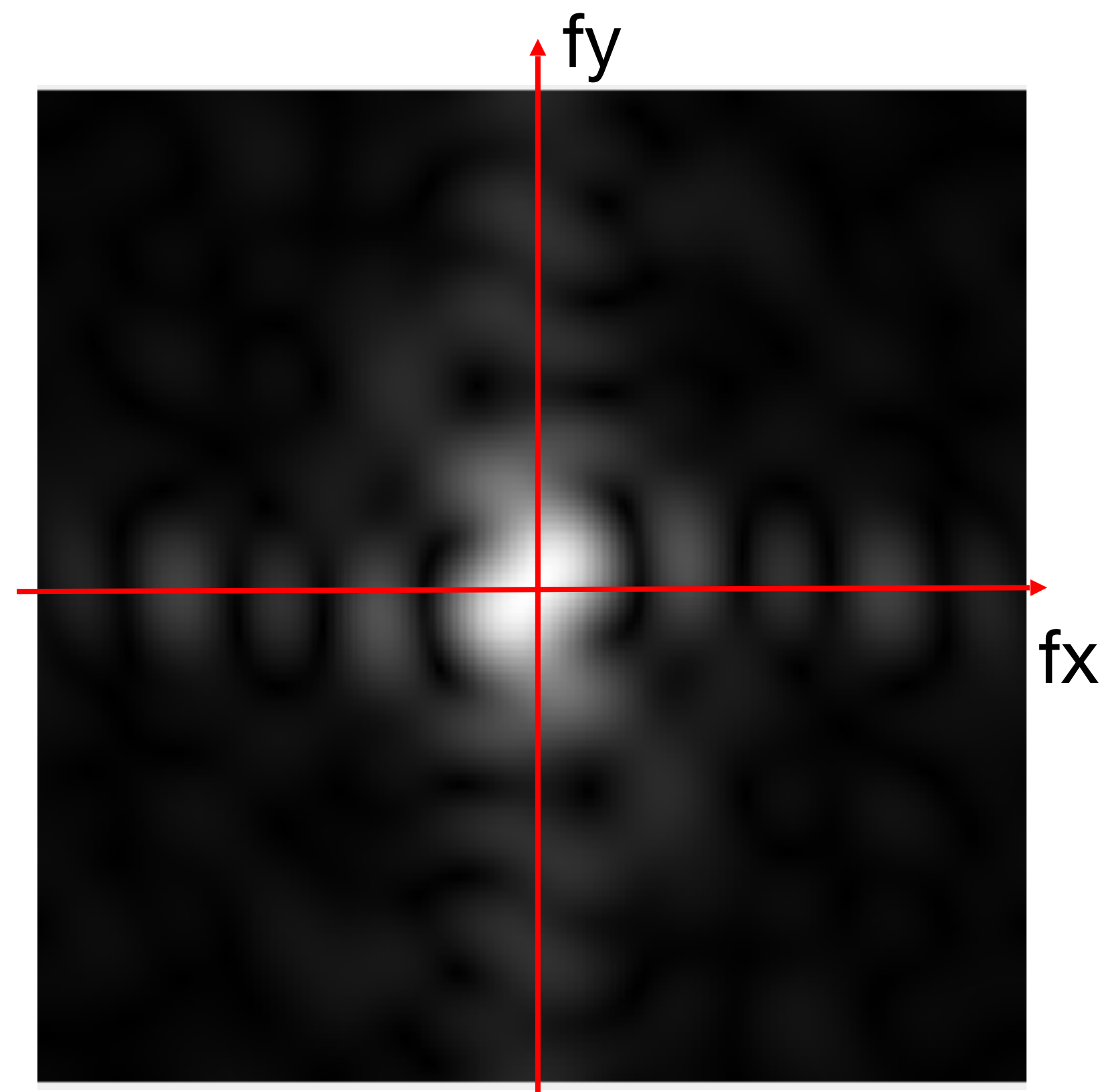
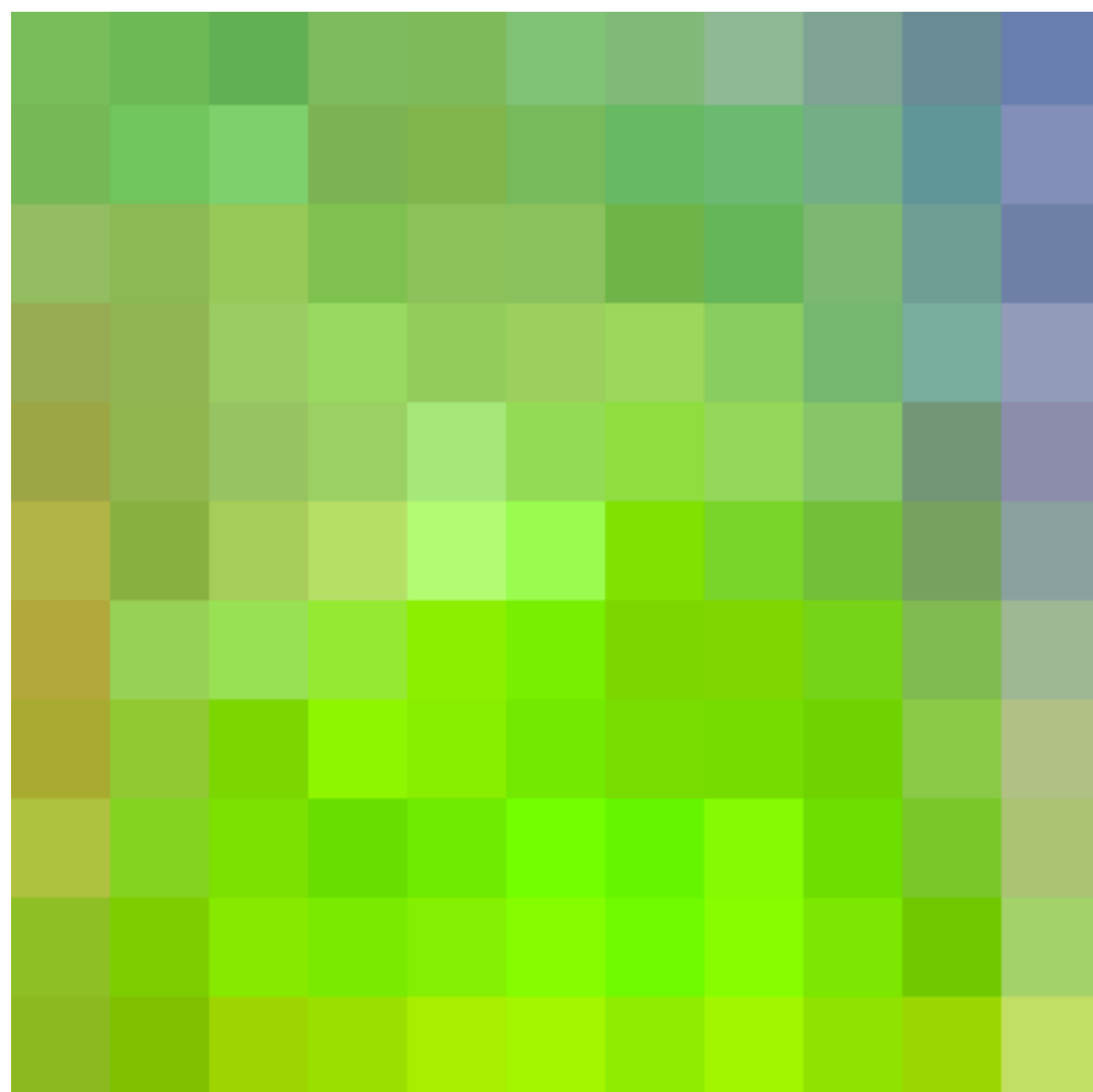
# Filters in first layer



# Filters in first layer

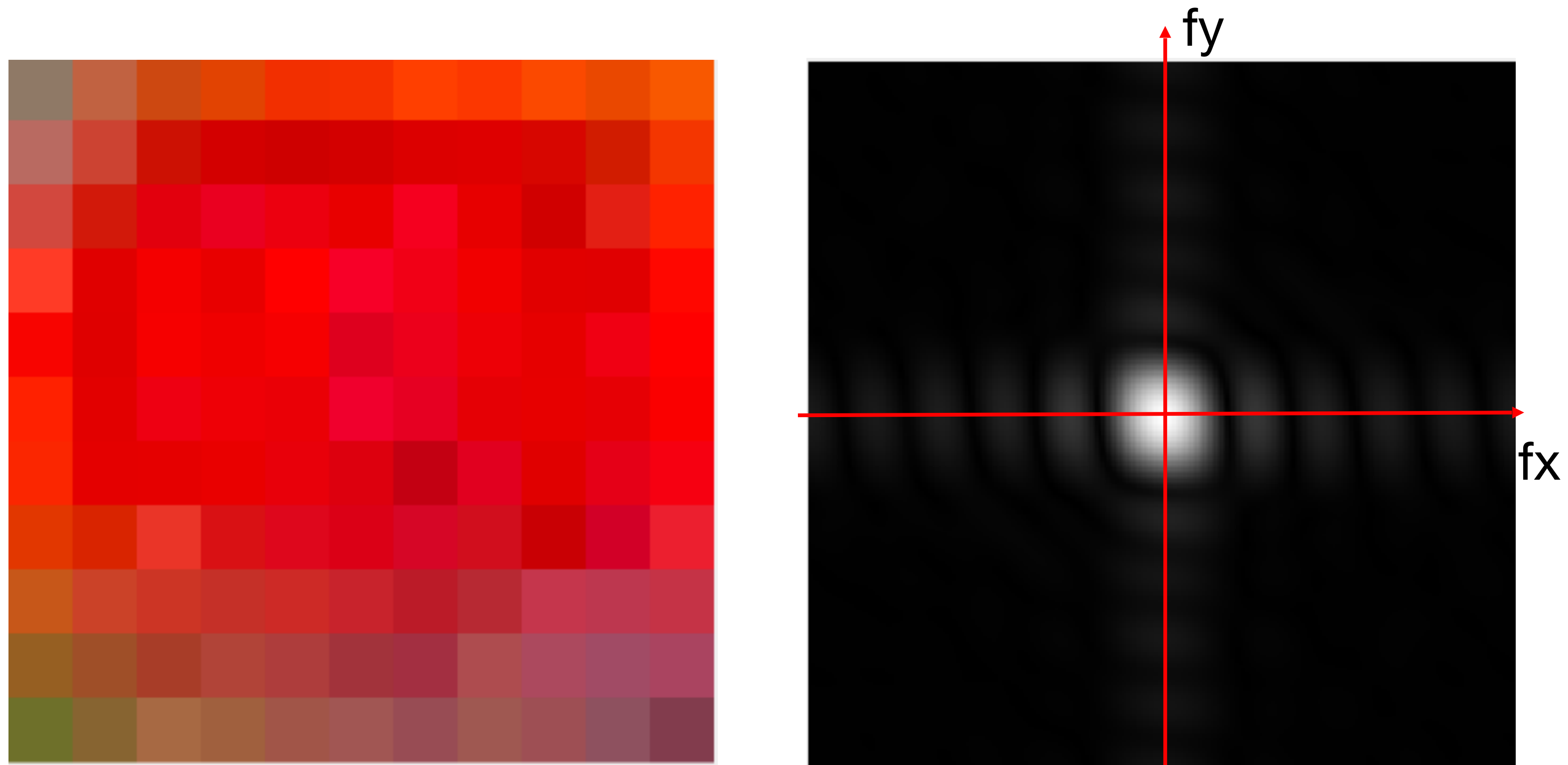


# Filters in first layer



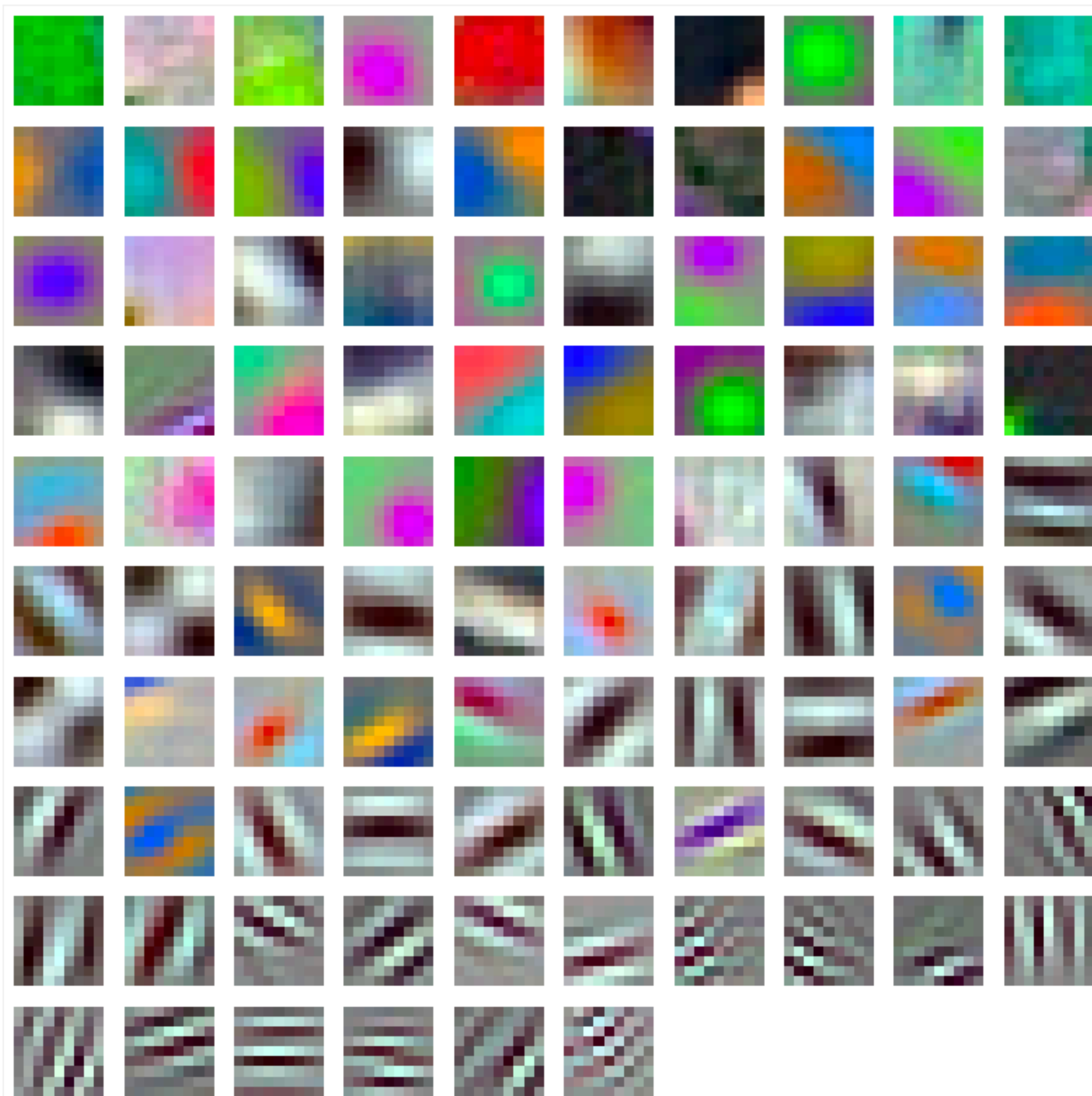


# Filters in first layer

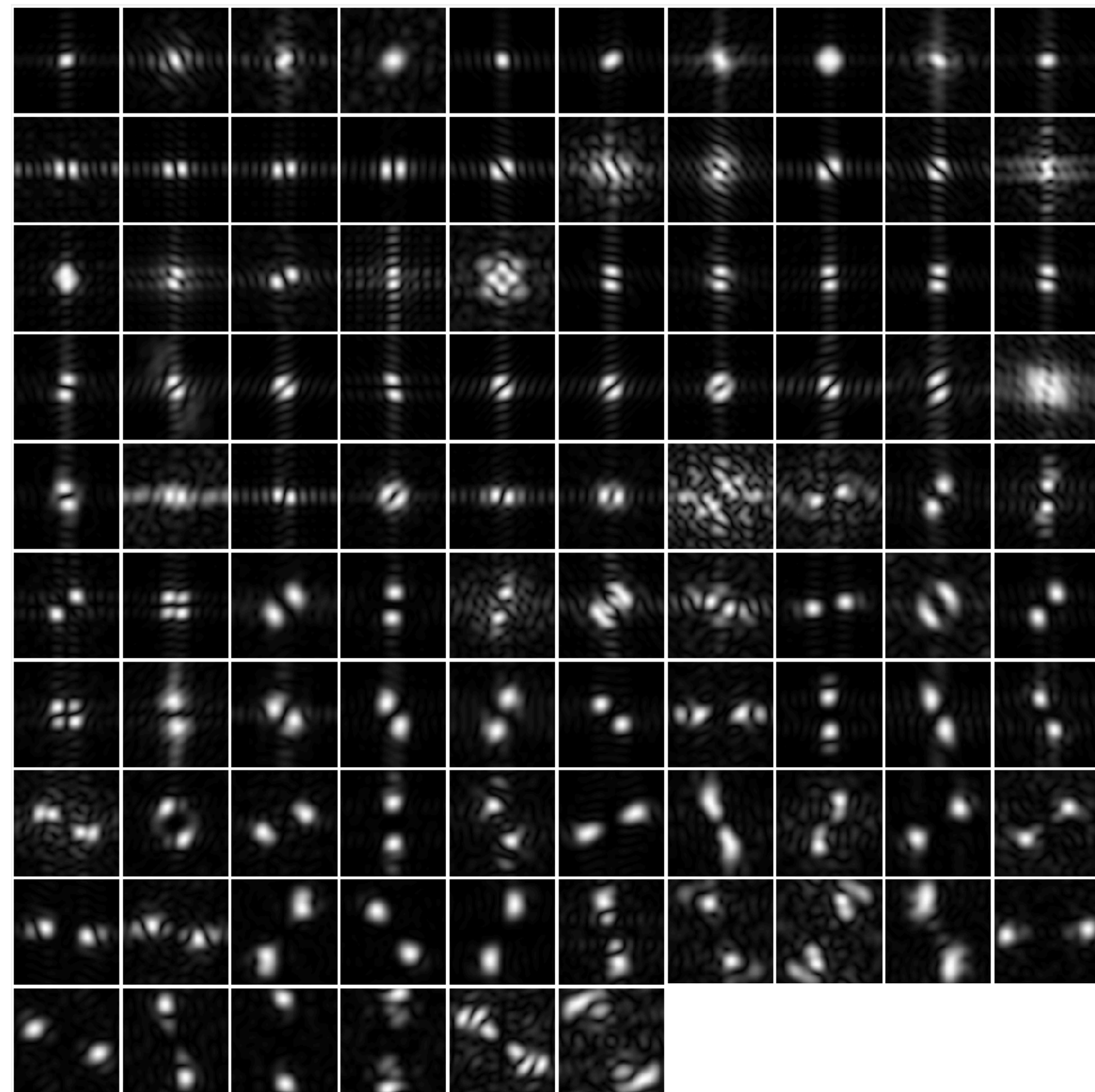




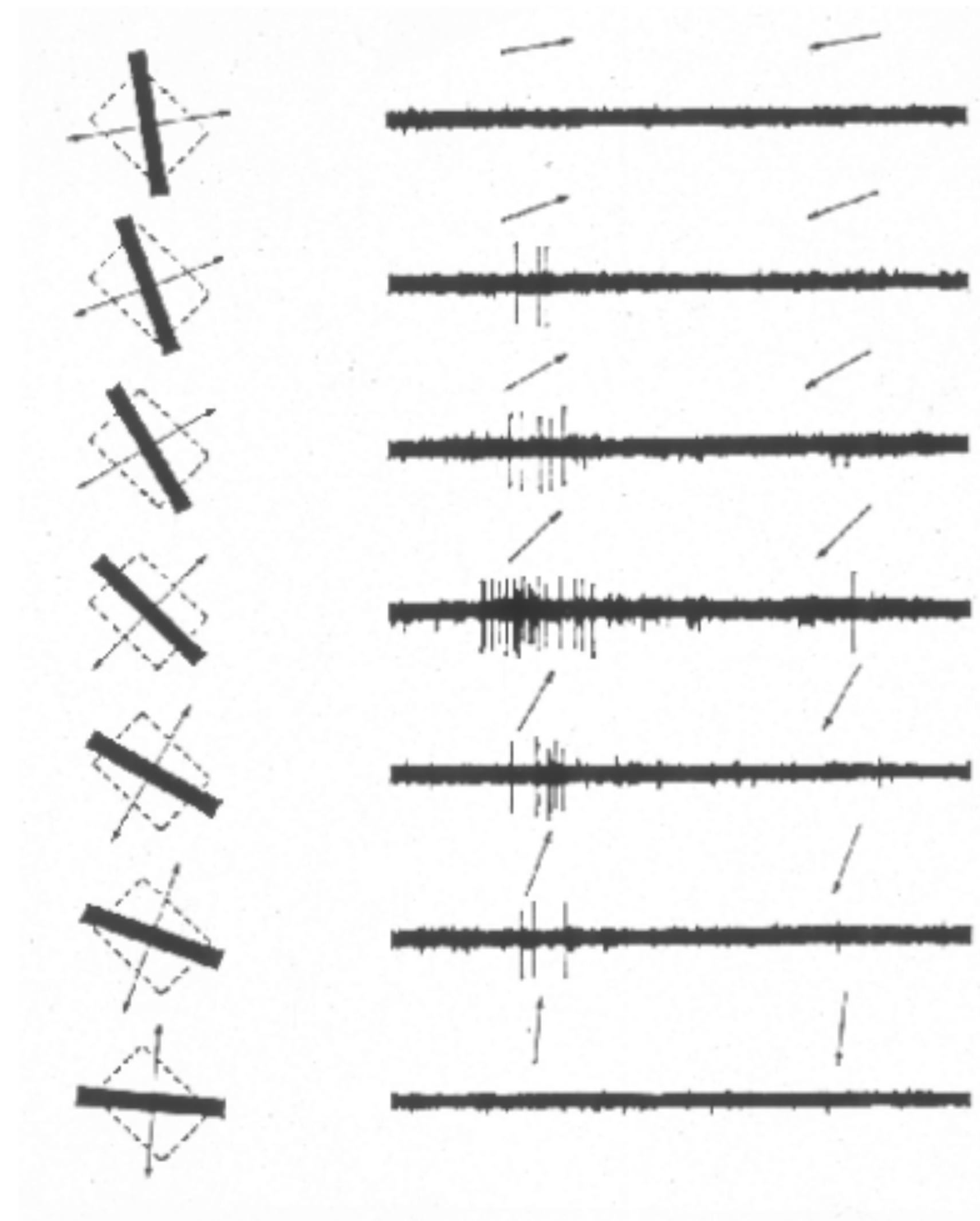
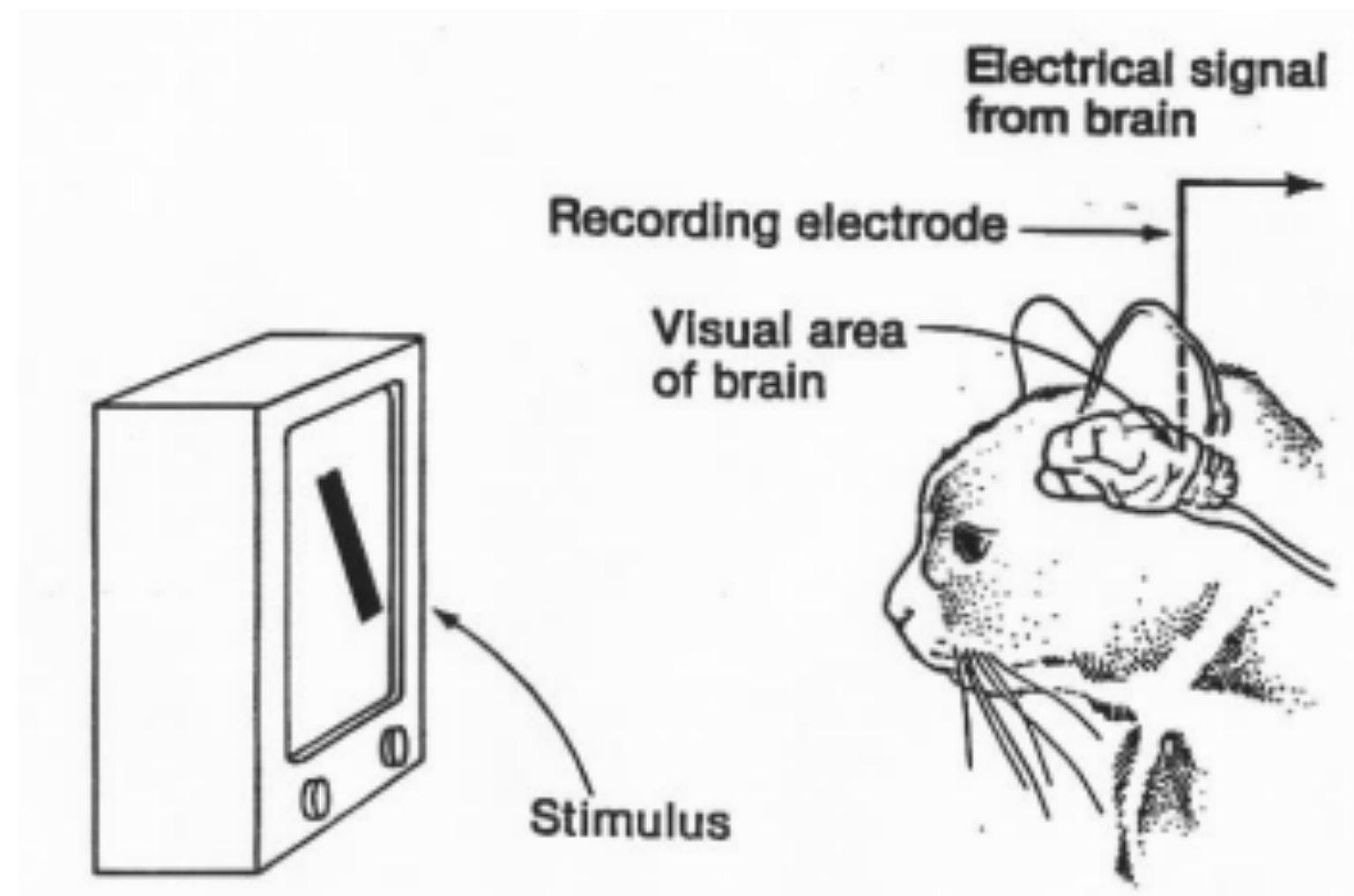
# Filters in first layer



96 Units in conv1

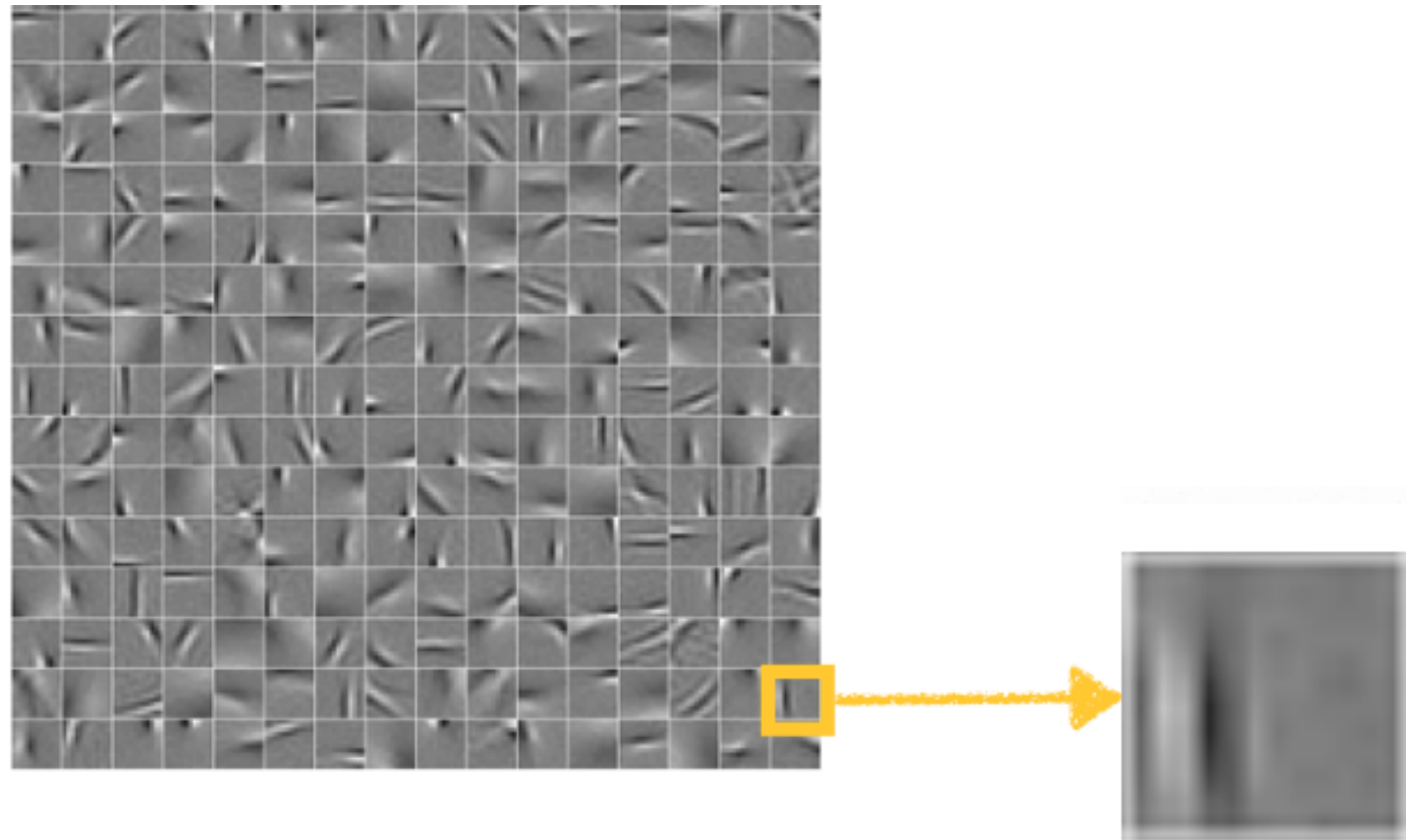


# [Hubel and Wiesel 59]

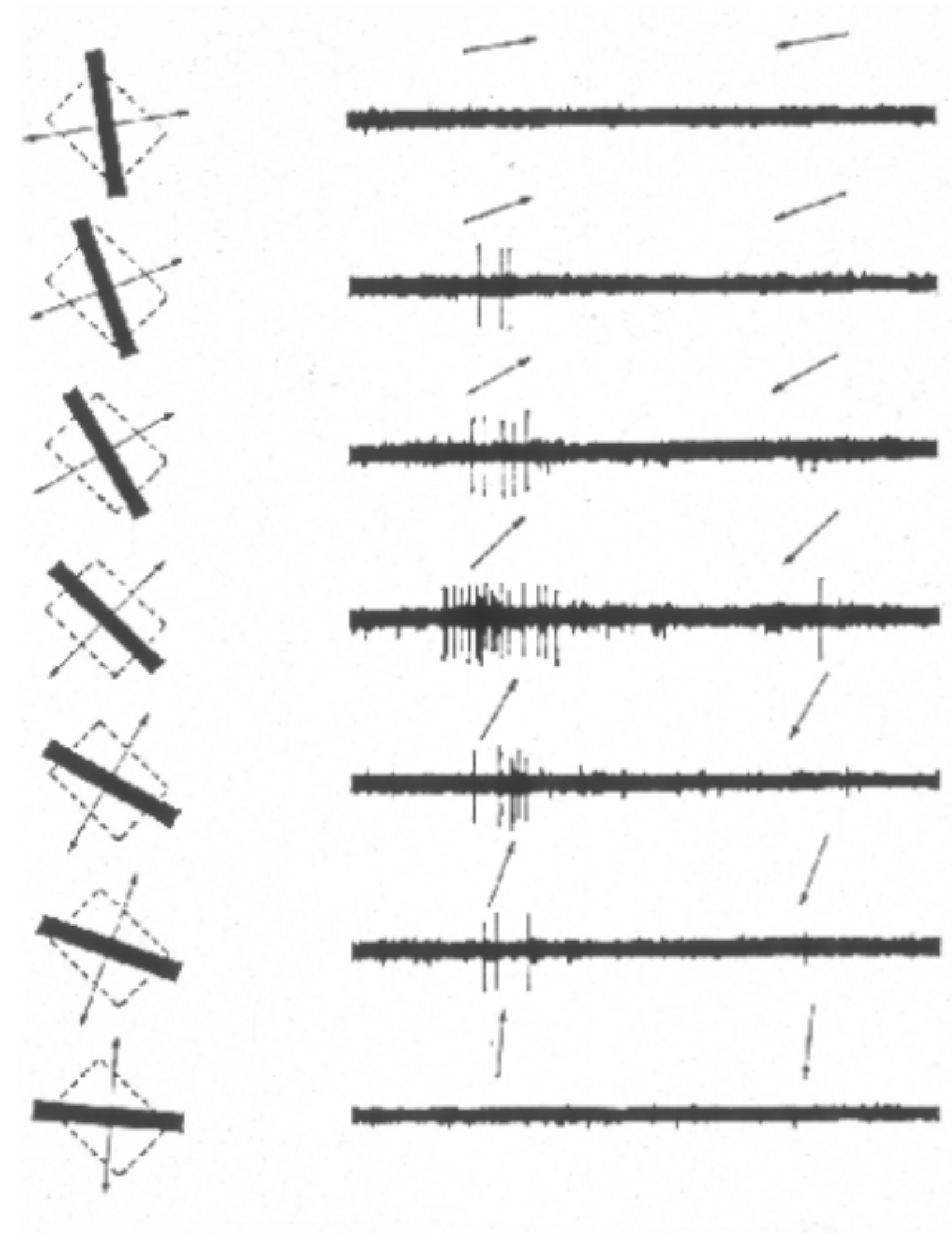




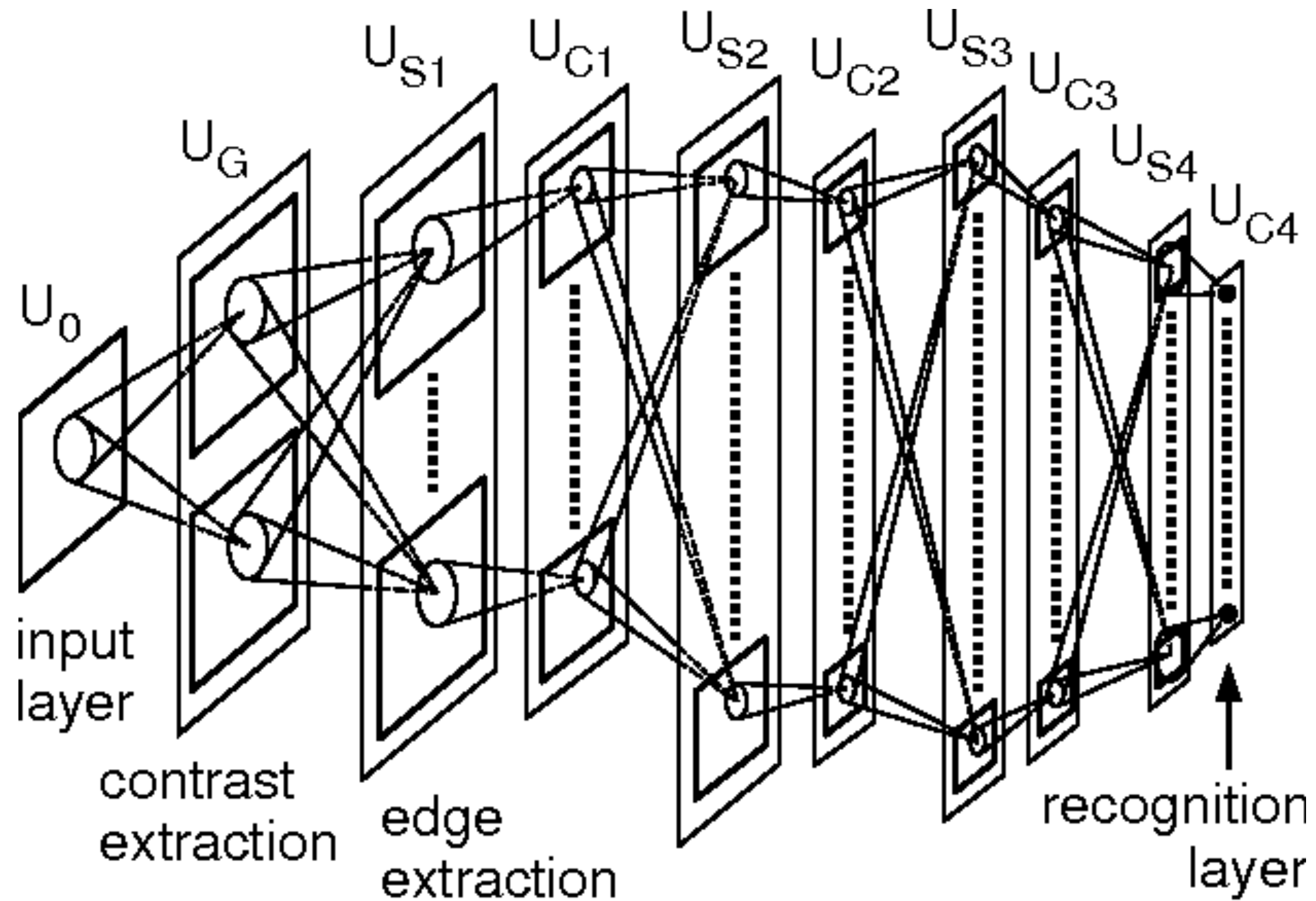
# A good fit to the experimental results



Oriented filters

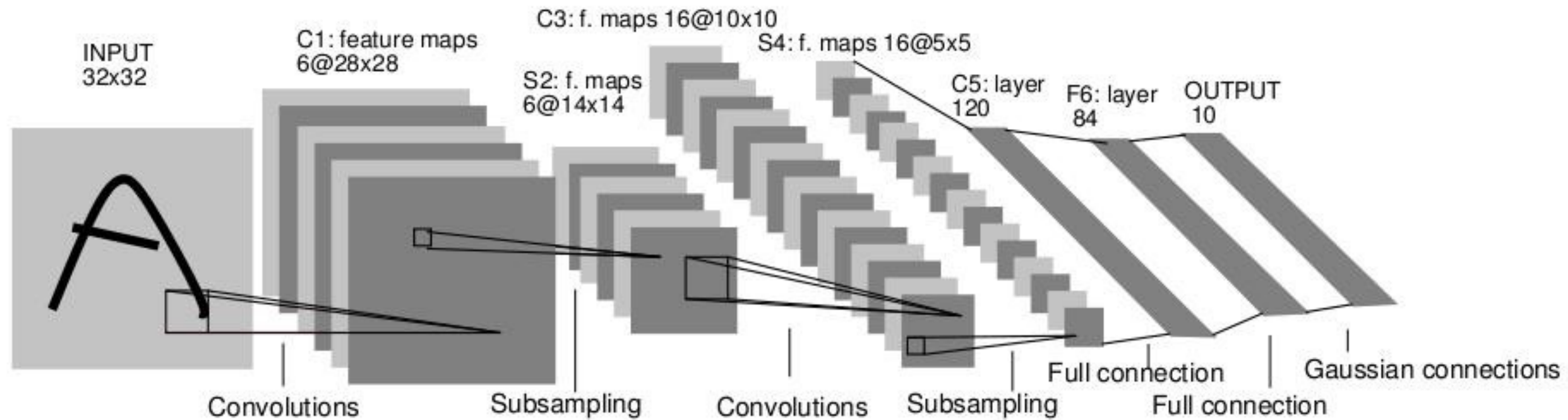


# History: Neocognitron



K. Fukushima, 1980s

# History: LeNet



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

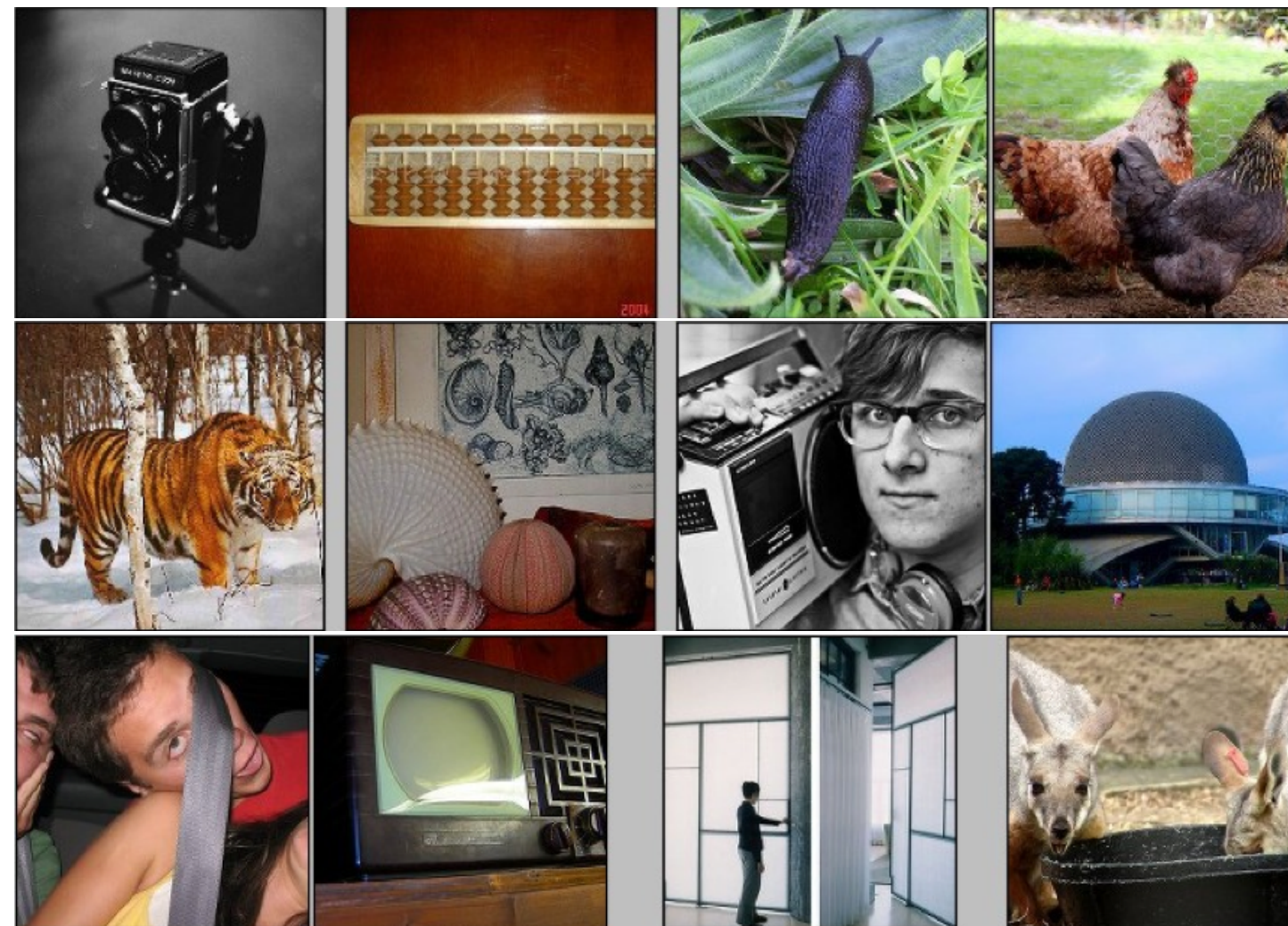
- Neocognitron + backpropagation
- Average pooling
- Sigmoid or tanh nonlinearity
- Fully connected layers at the end
- Trained on MNIST digit dataset with 60K training examples





# ImageNet Challenge

IMAGENET



- ~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

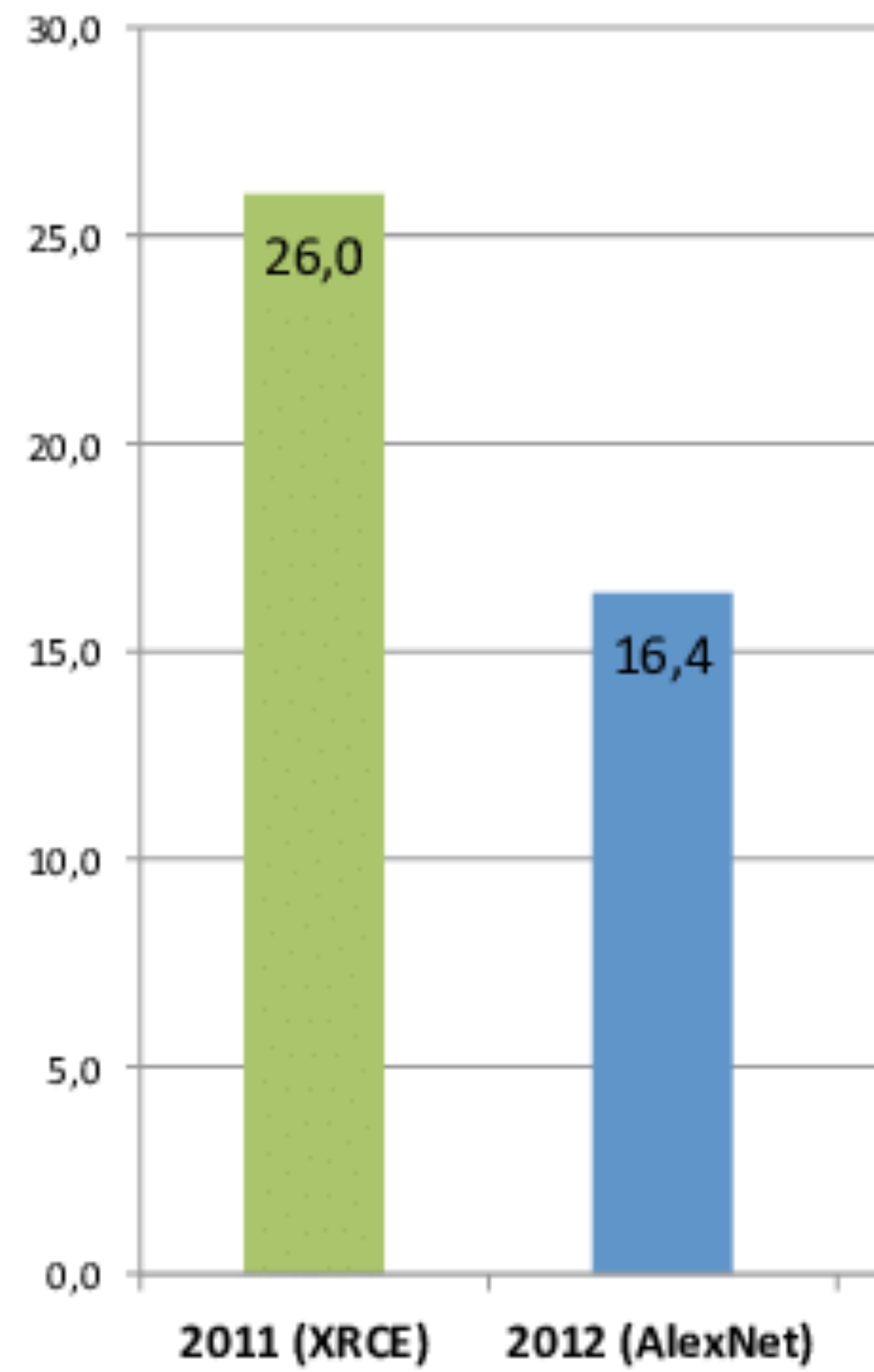
[Russakovsky, Deng, Su, Krause, Satheesh, Ma, Huang, Karpathy, Khosla, Bernstein, Berge, Fei-Fei, “ImageNet Large Scale Visual Recognition Challenge”, 2015]

# More recent networks

And advances that make them work:

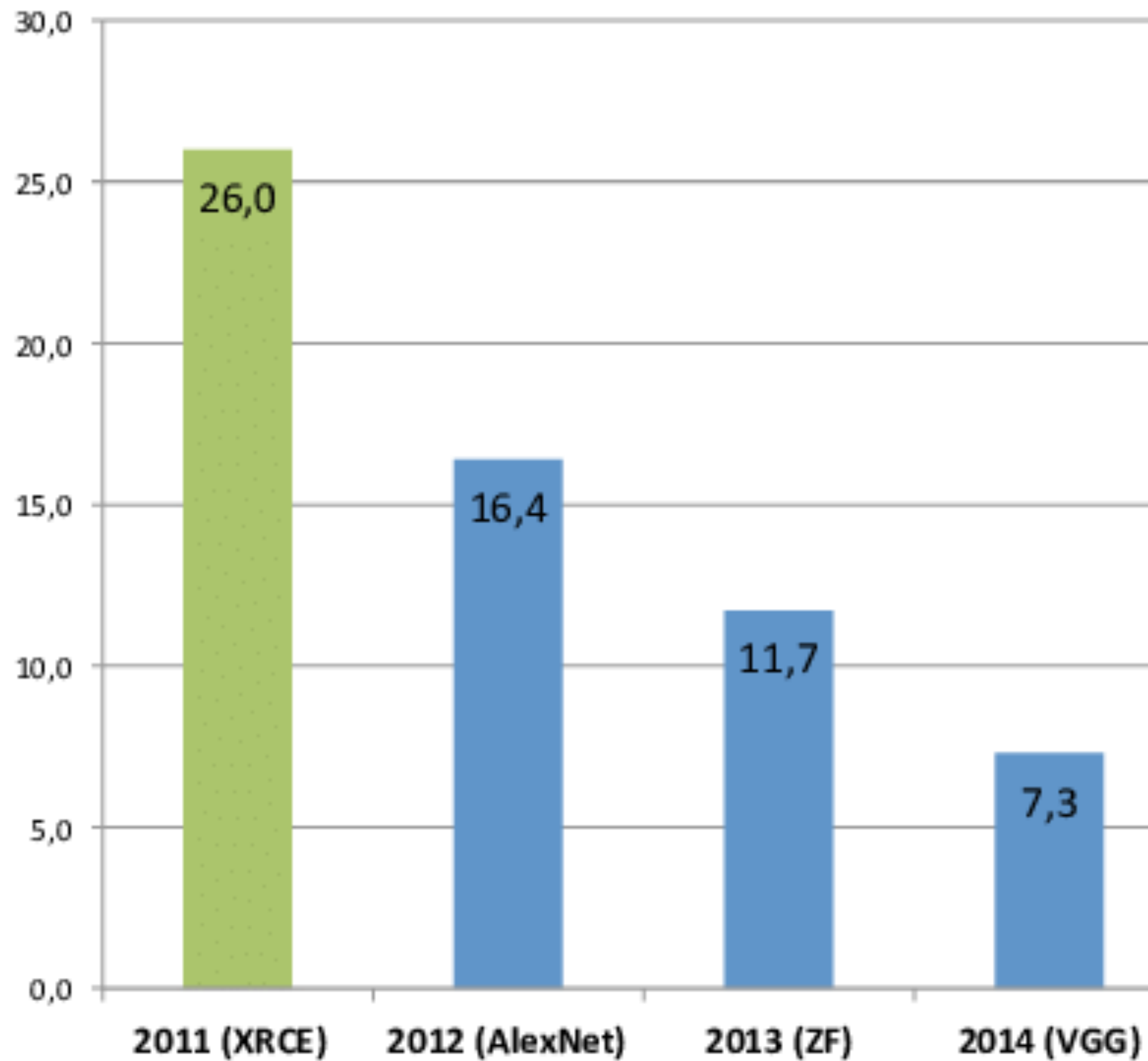
- Chaining small filters
- Residual layers
- Normalization

## ImageNet Classification Error (Top 5)

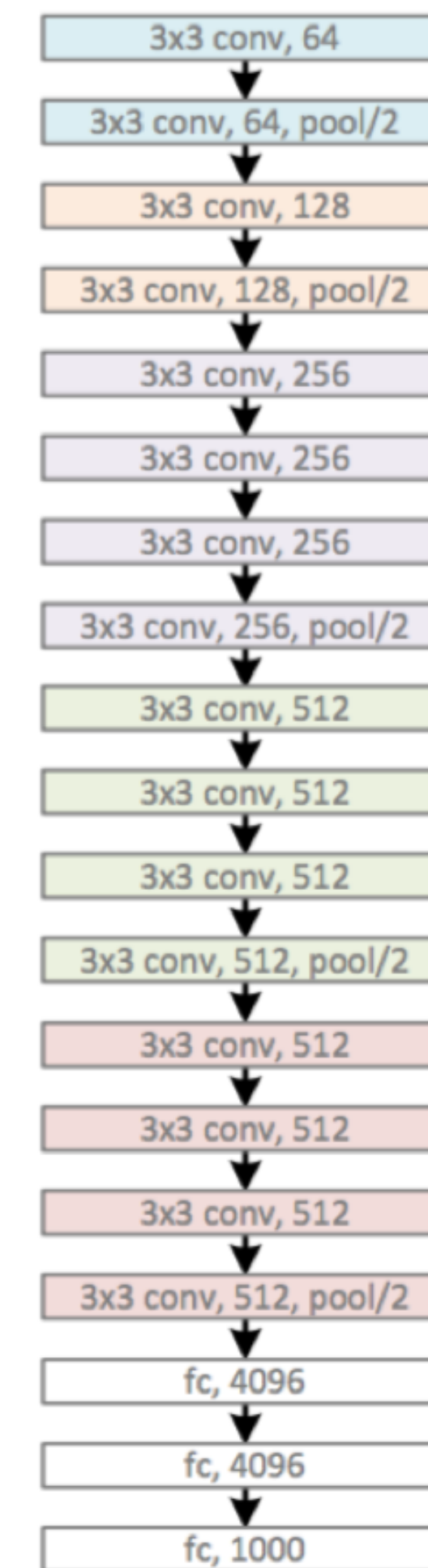




**ImageNet Classification Error (Top 5)**



**2014: VGG**  
16 conv. layers

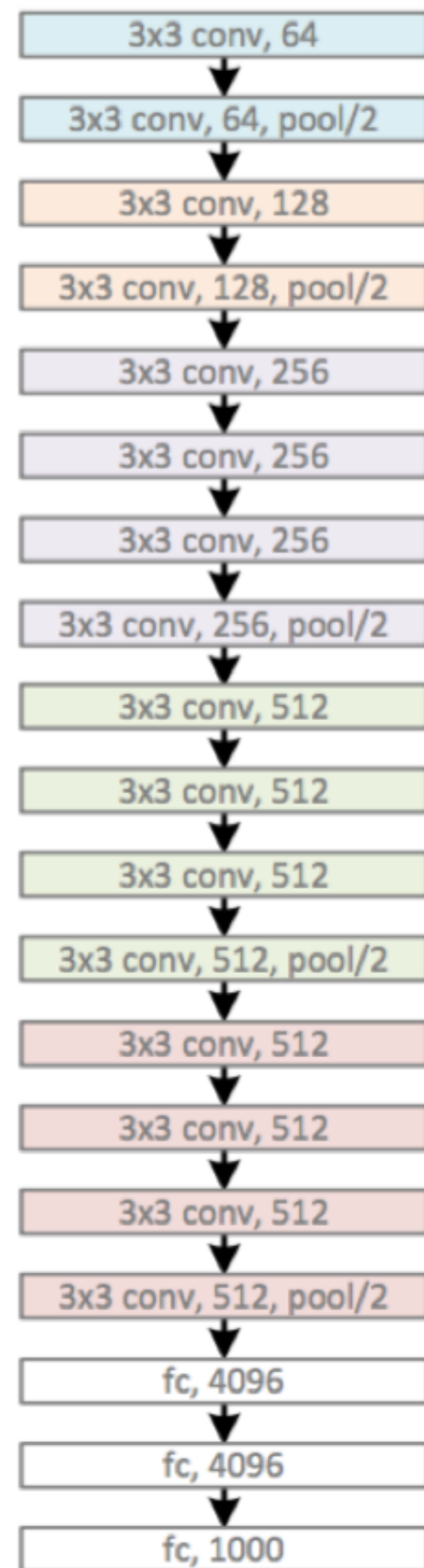


Error: 7.3%

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

# VGG-Net [Simonyan & Zisserman, 2015]

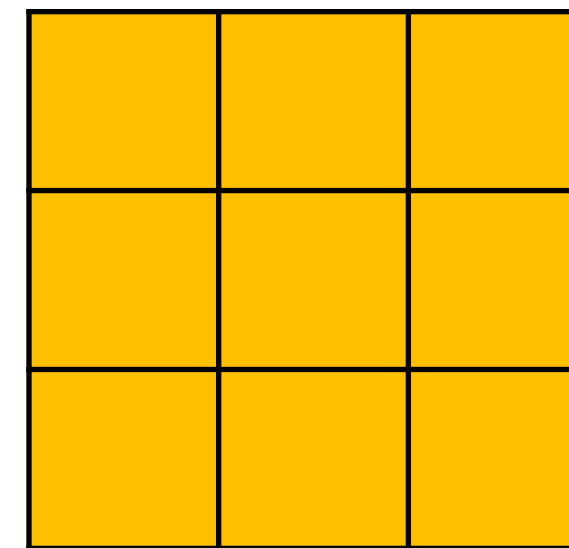
2014: VGG  
16 conv. layers



Error: 7.3%

## Main developments

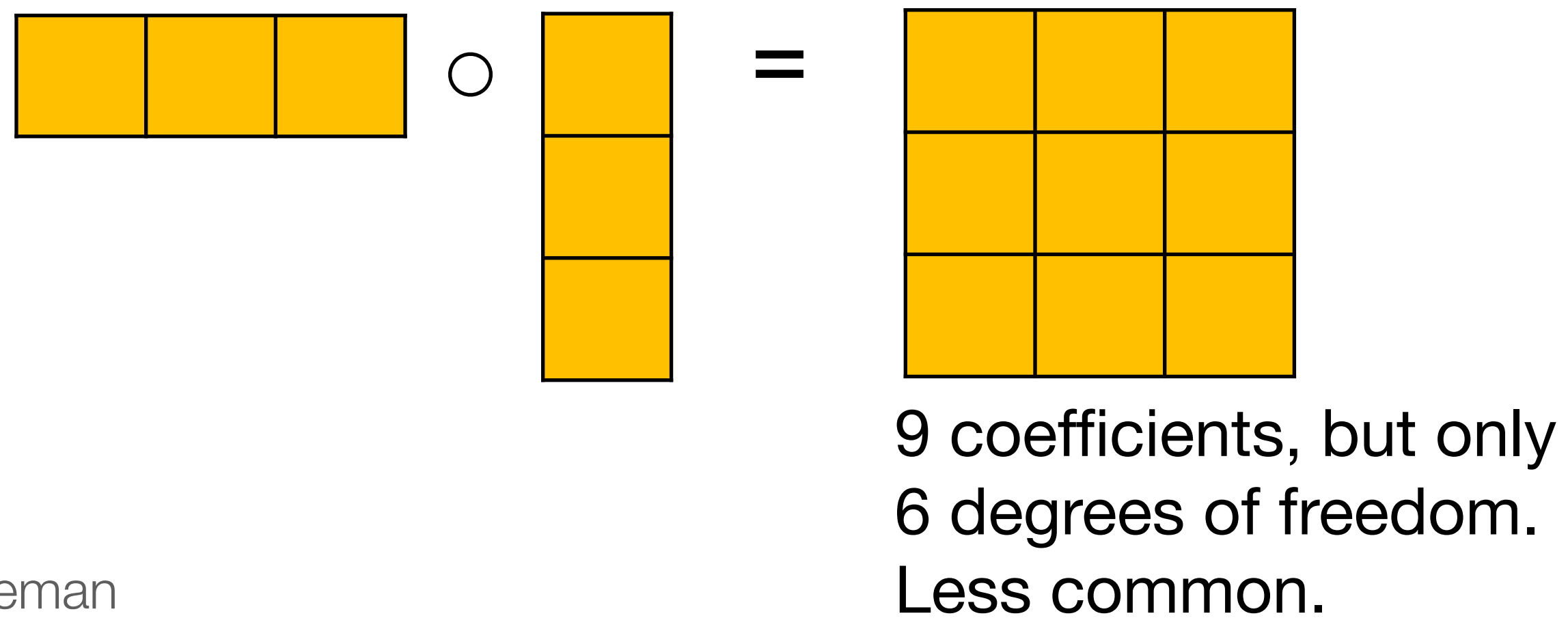
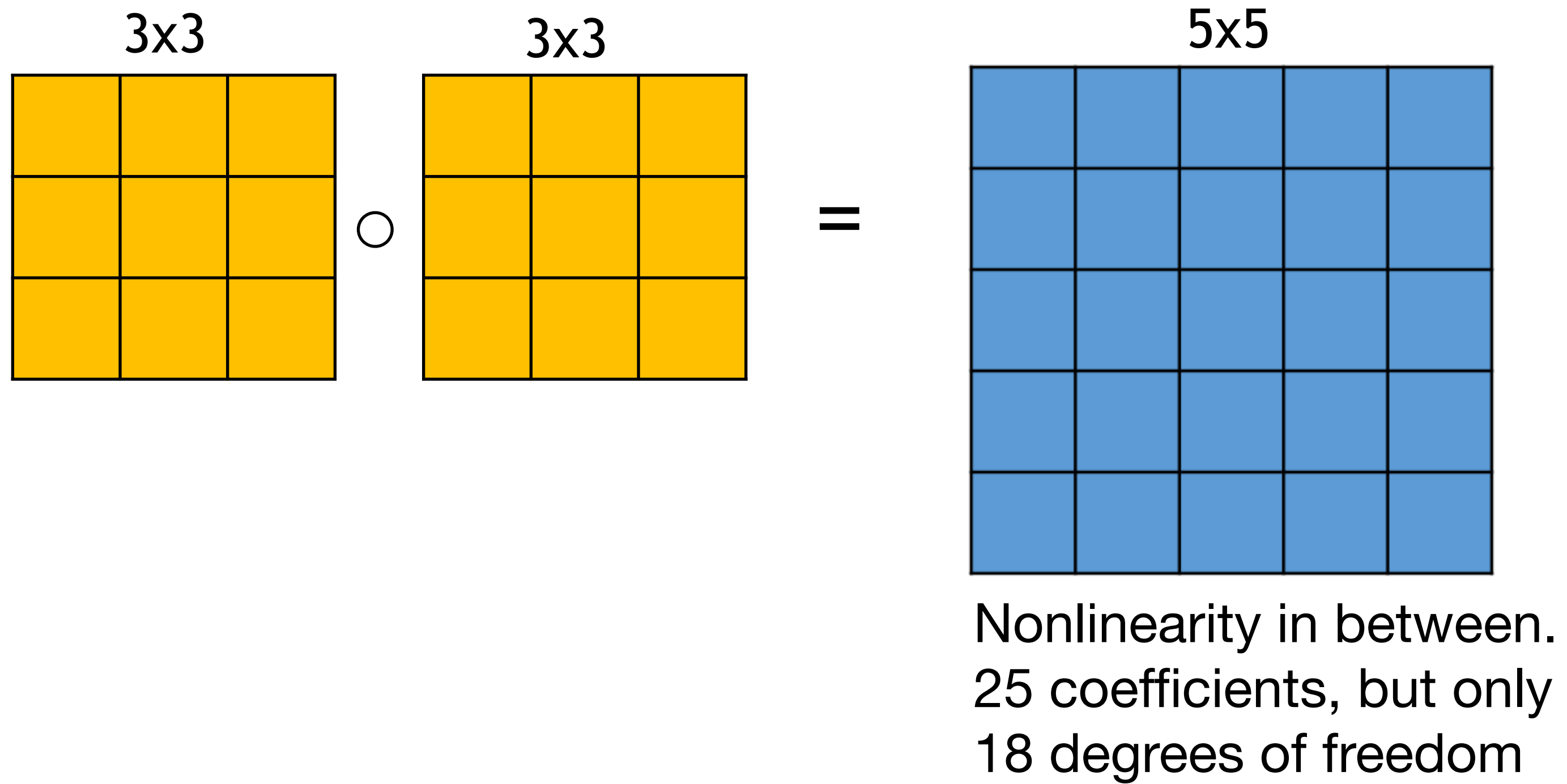
- Small convolutional kernels: only 3x3



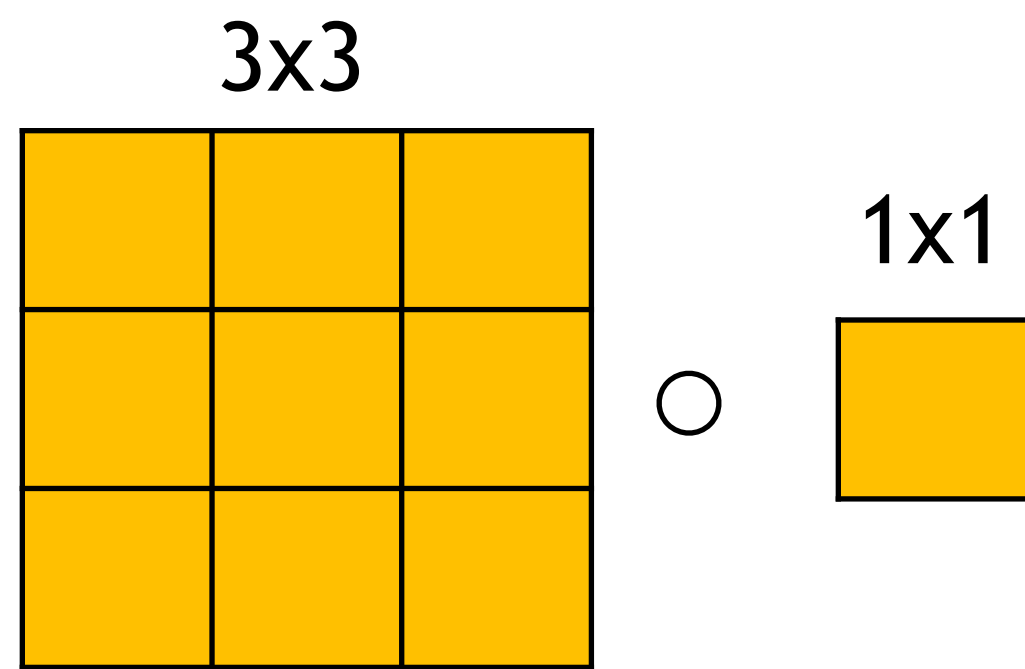
- Increased depth (5 -> 16/19 layers)

Other tricks for designing convolutional nets

# Chaining convolutions



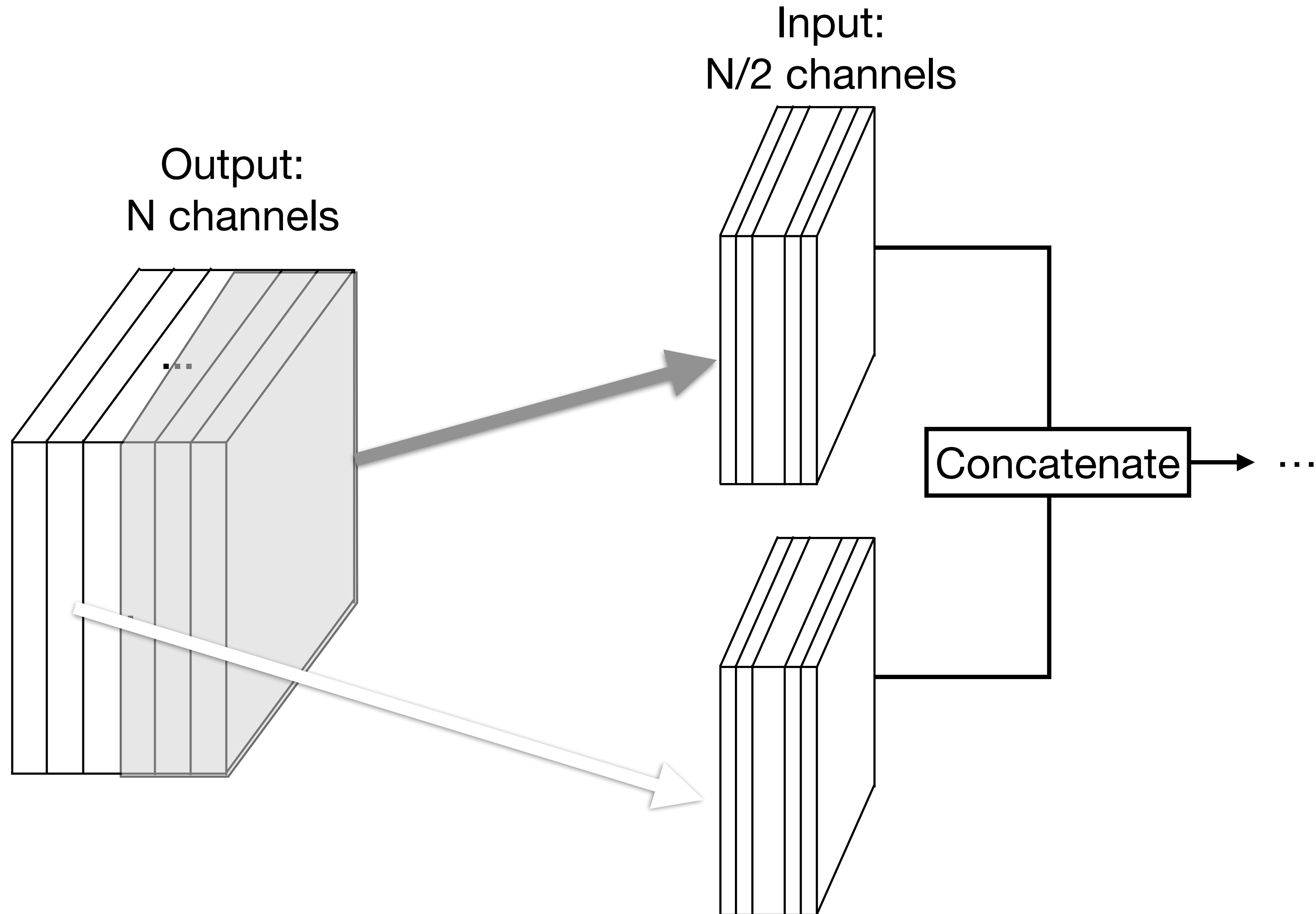
# 1x1 convolutions



Why do this?

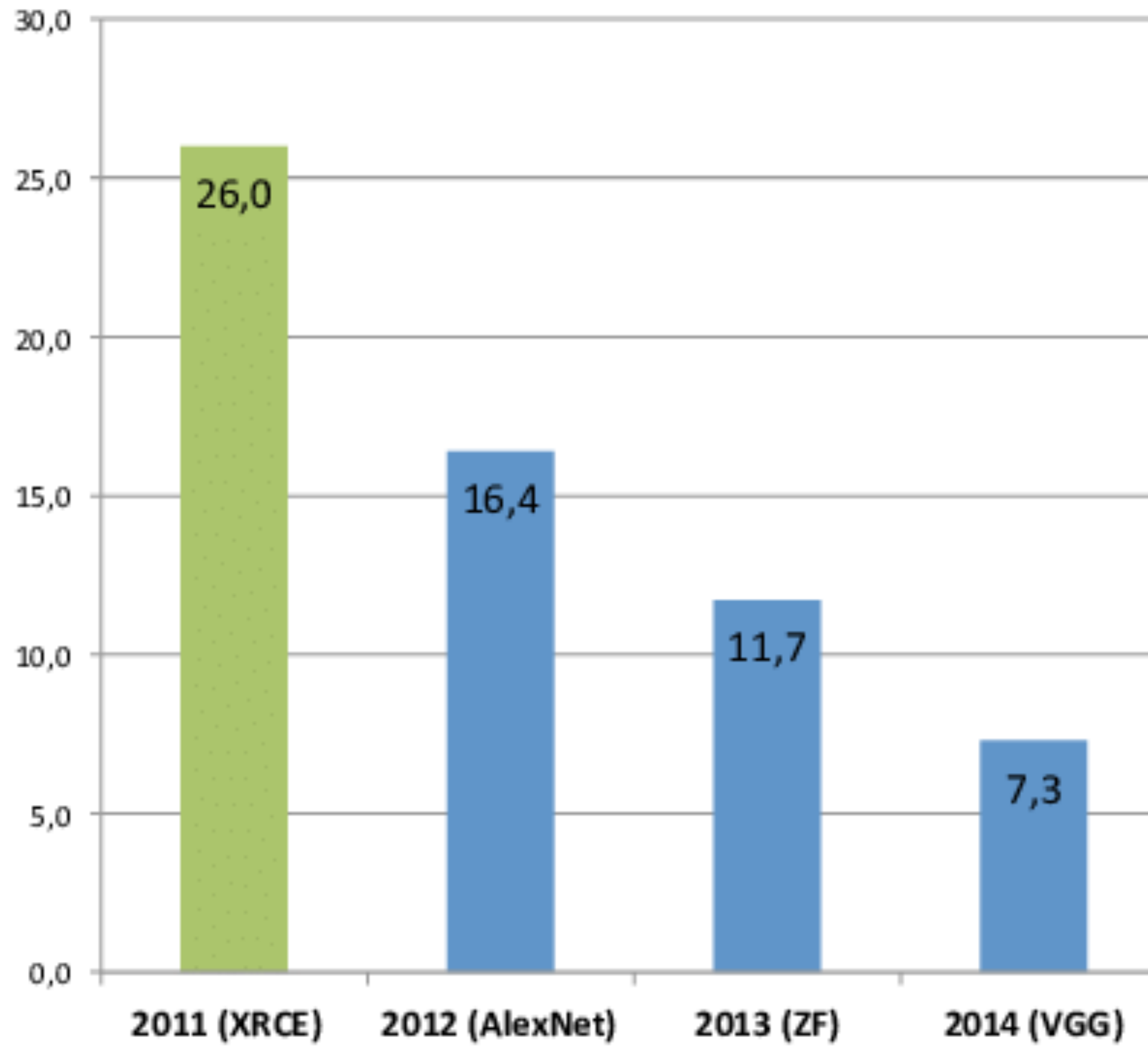
(nonlinearity in between)

# Grouped Convolutions

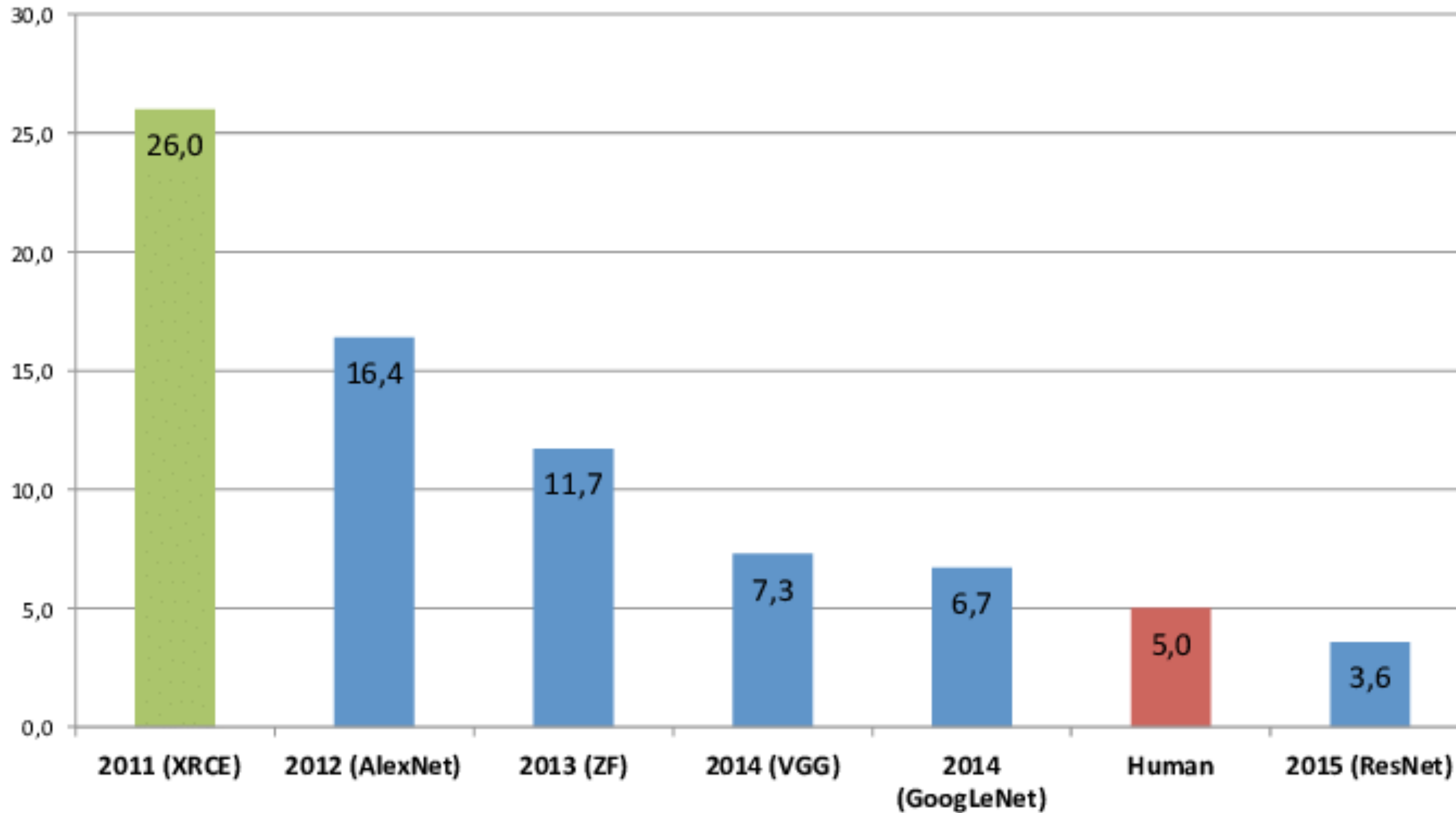


Split channels into N groups, and process separately with N convolution layers.

## ImageNet Classification Error (Top 5)



# ImageNet Classification Error (Top 5)



2016: ResNet  
>100 conv. layers

Error: 3.6%

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



2016: ResNet  
>100 conv. layers

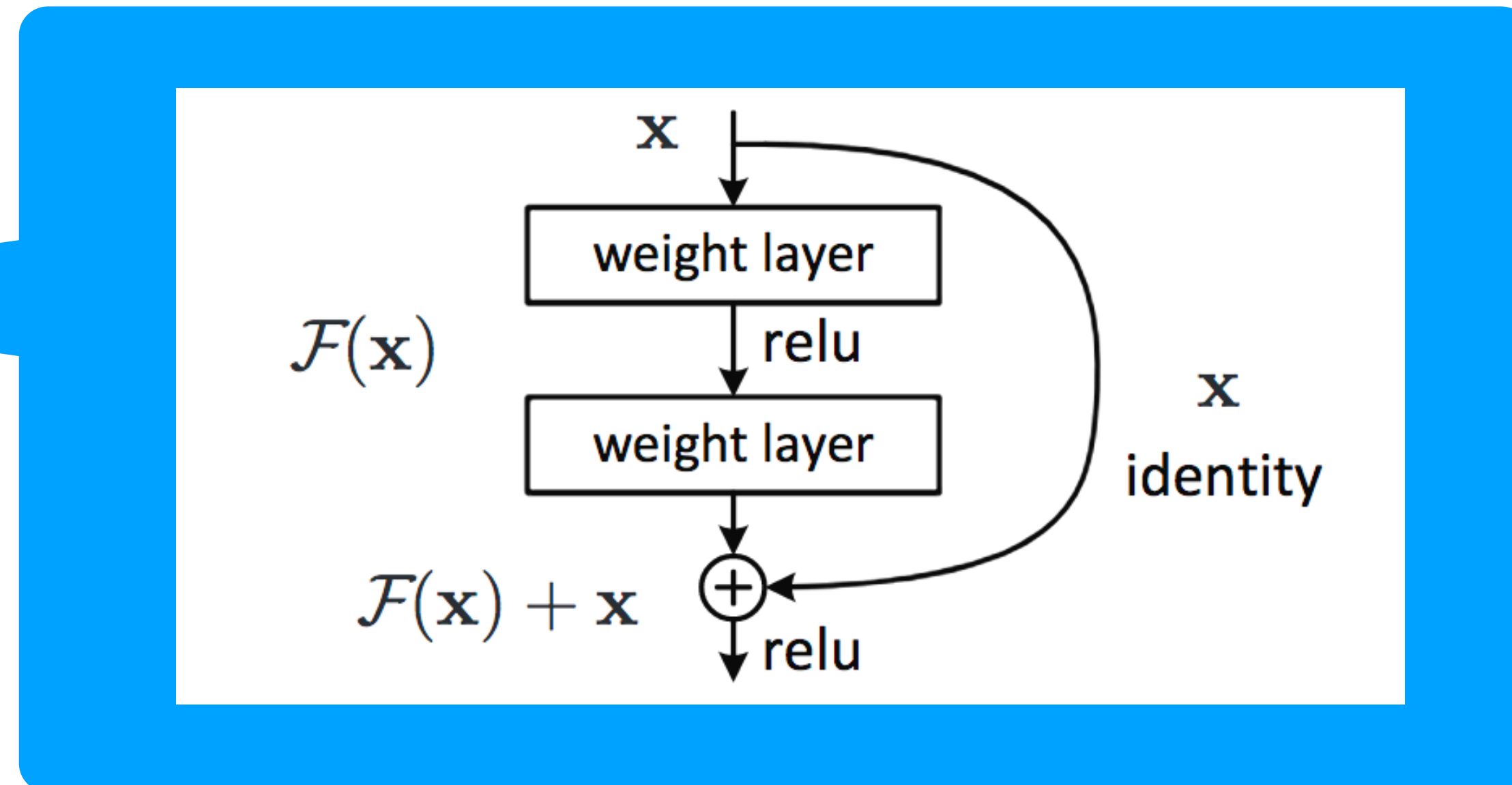
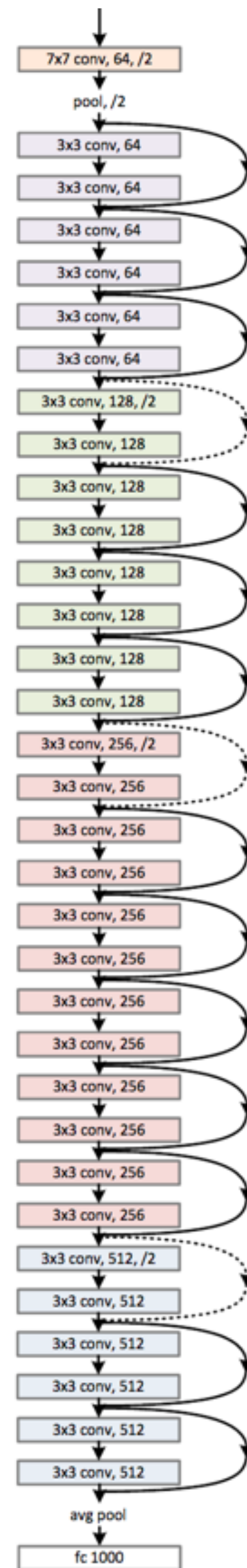
# ResNet [He et al, 2016]

## Main developments

- Increased depth possible through residual blocks



Error: 3.6%



# 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:**

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

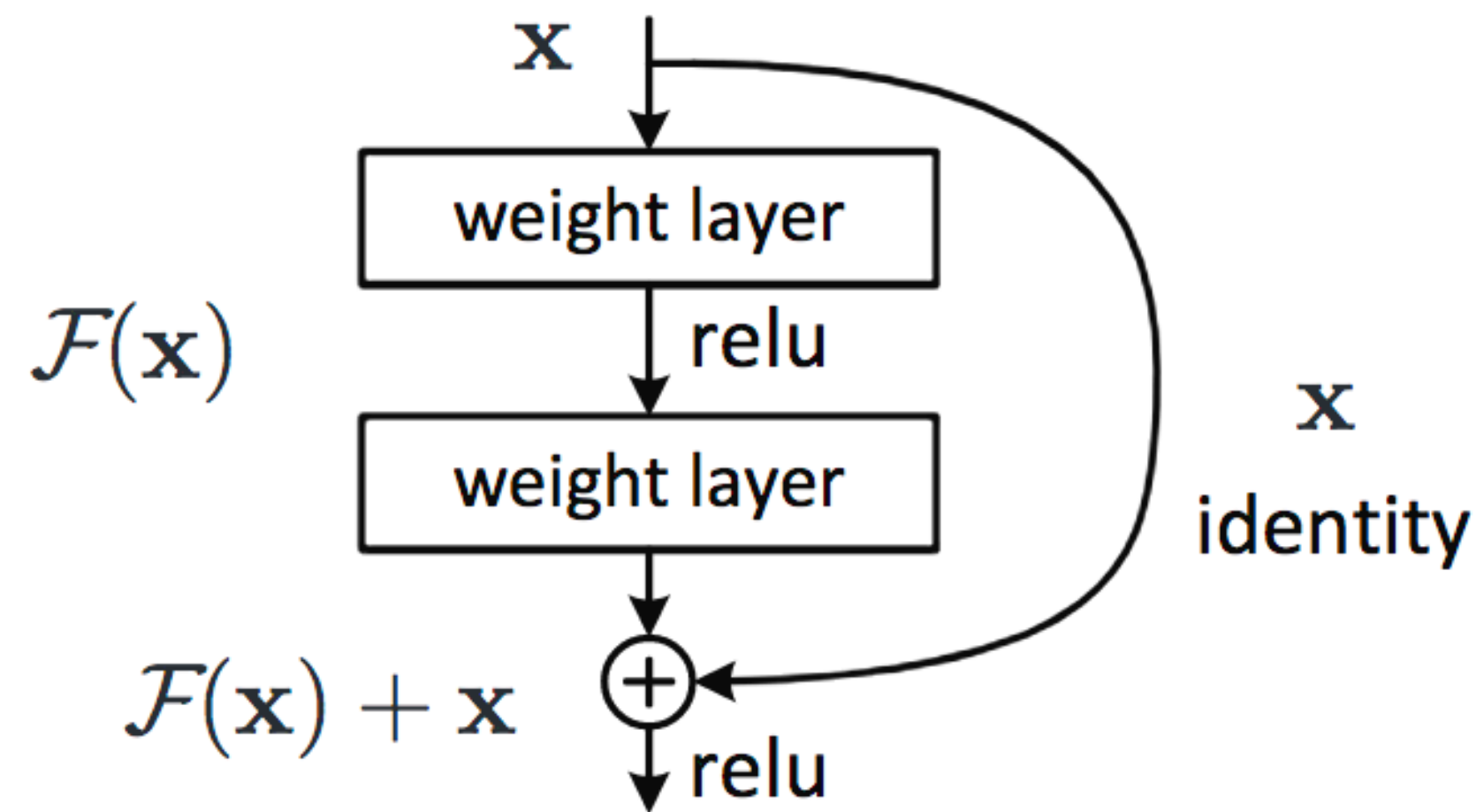
**Residual connection**

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

**In general, do multiple convolutions (with nonlinearities) before summing:**

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

# Residual Blocks



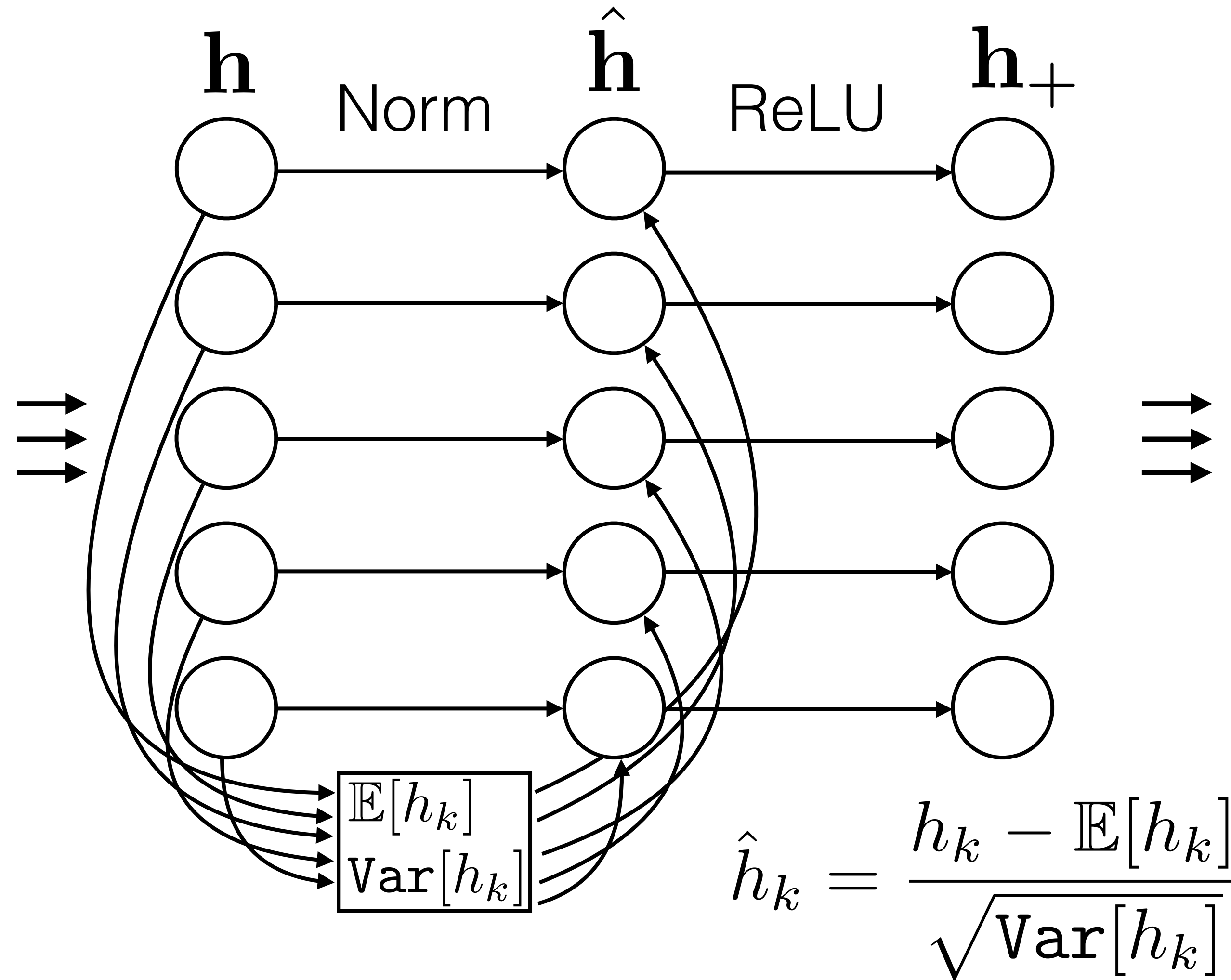
## Why do they work?

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

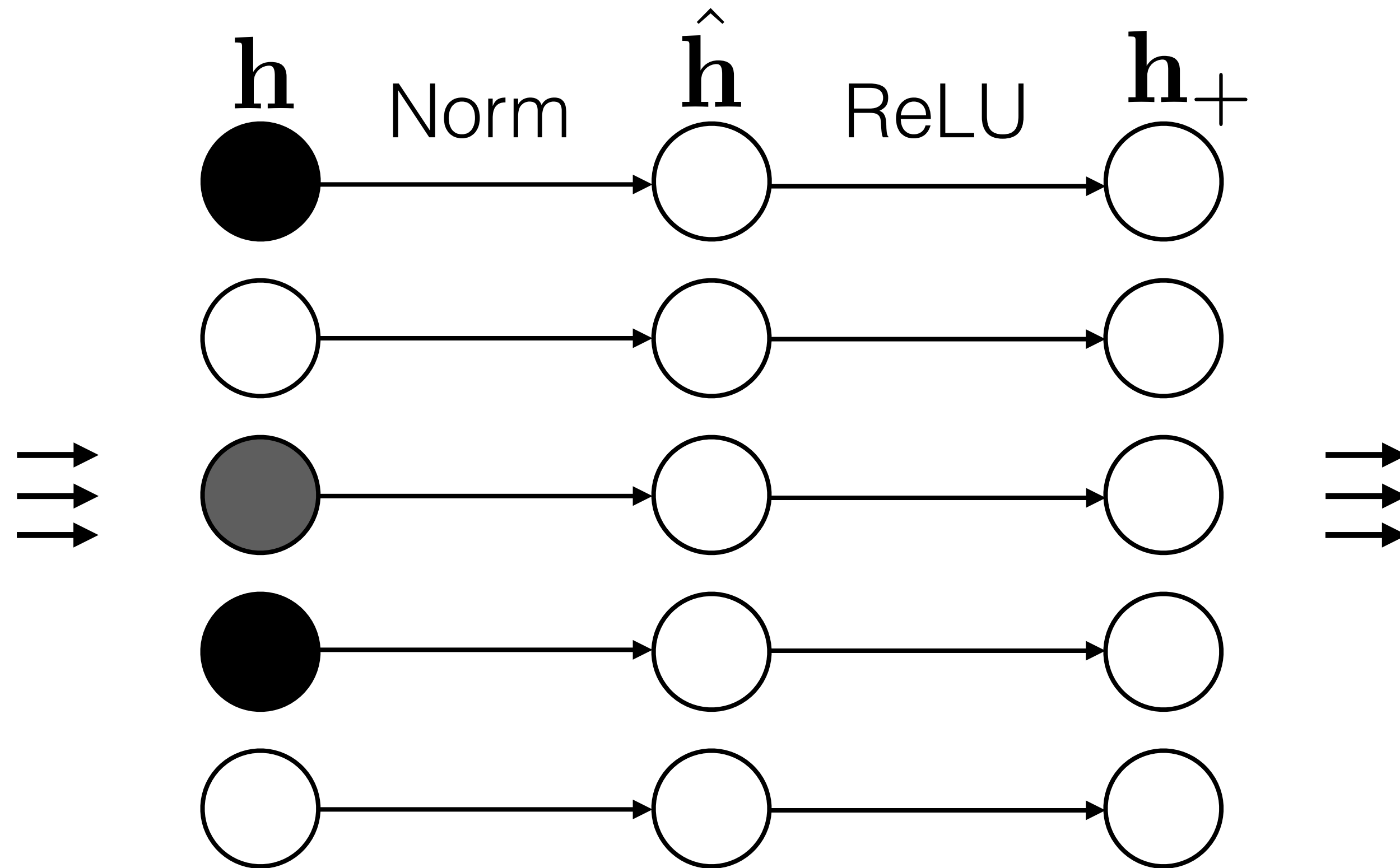
# Normalization layers

- 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

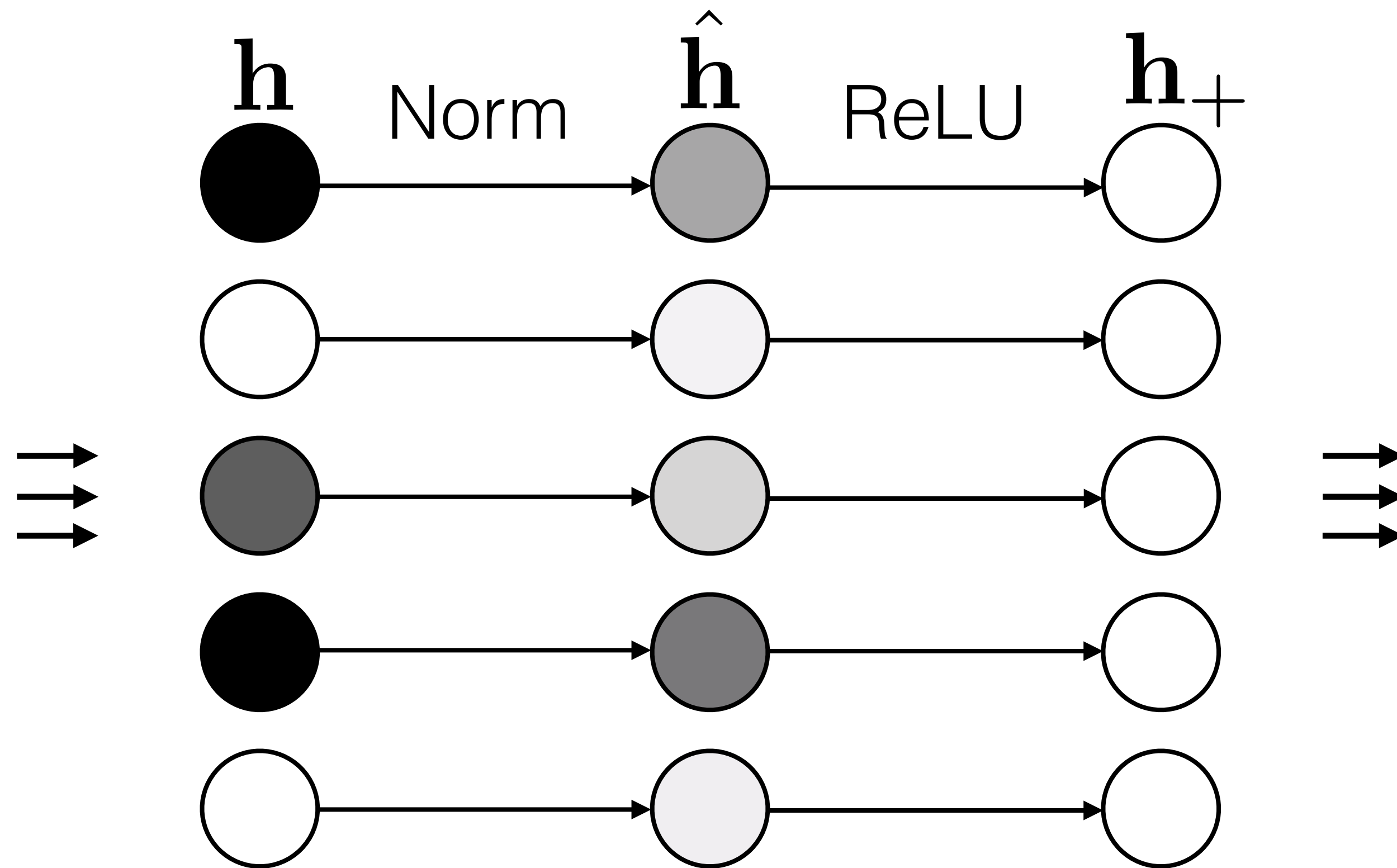


# Normalization layers



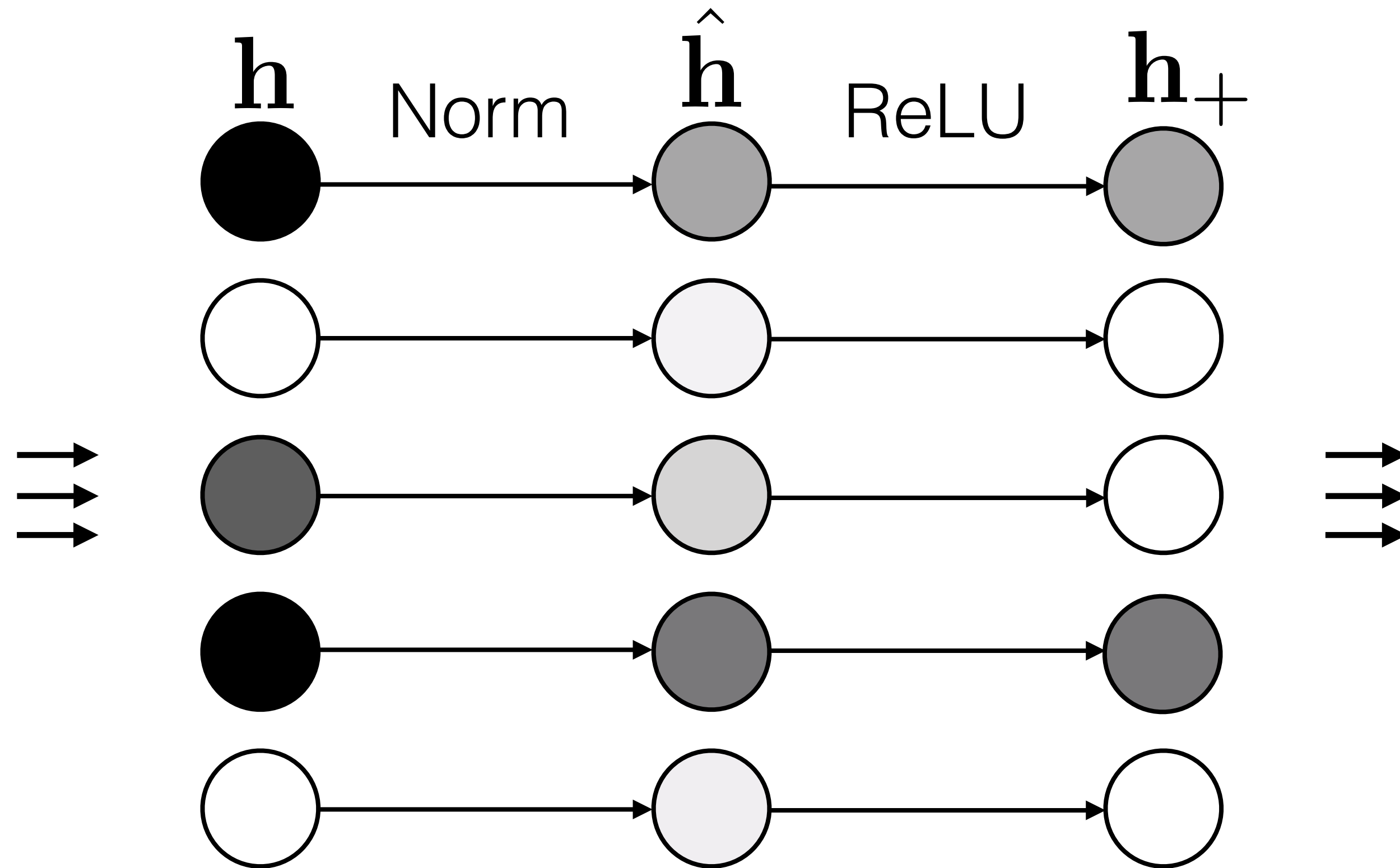
$$\hat{h}_k = \frac{h_k - \mathbb{E}[h_k]}{\sqrt{\text{Var}[h_k]}}$$

# Normalization layers



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

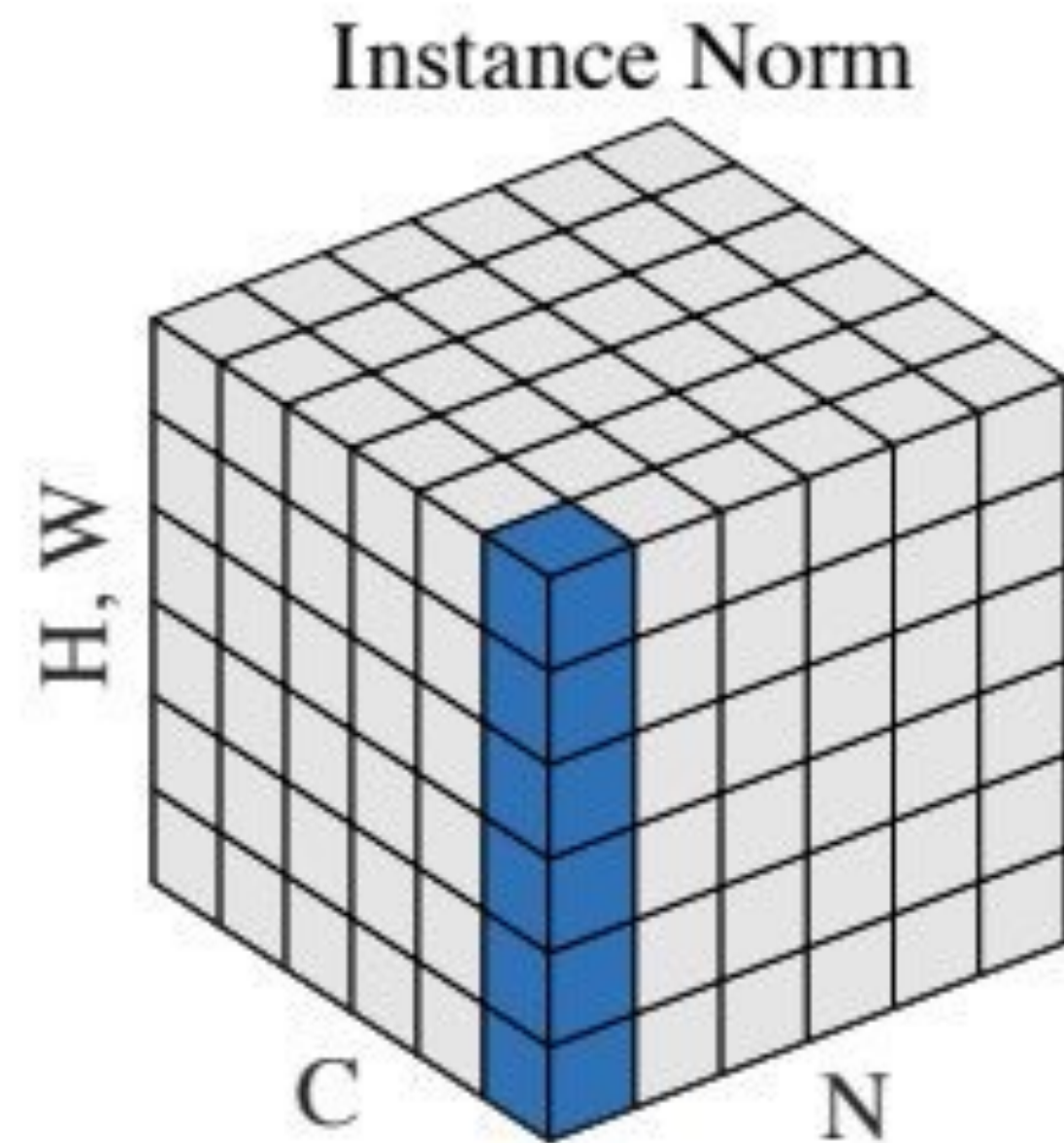
# Normalization layers



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



# Instance normalization



Filter value in a CNN layer

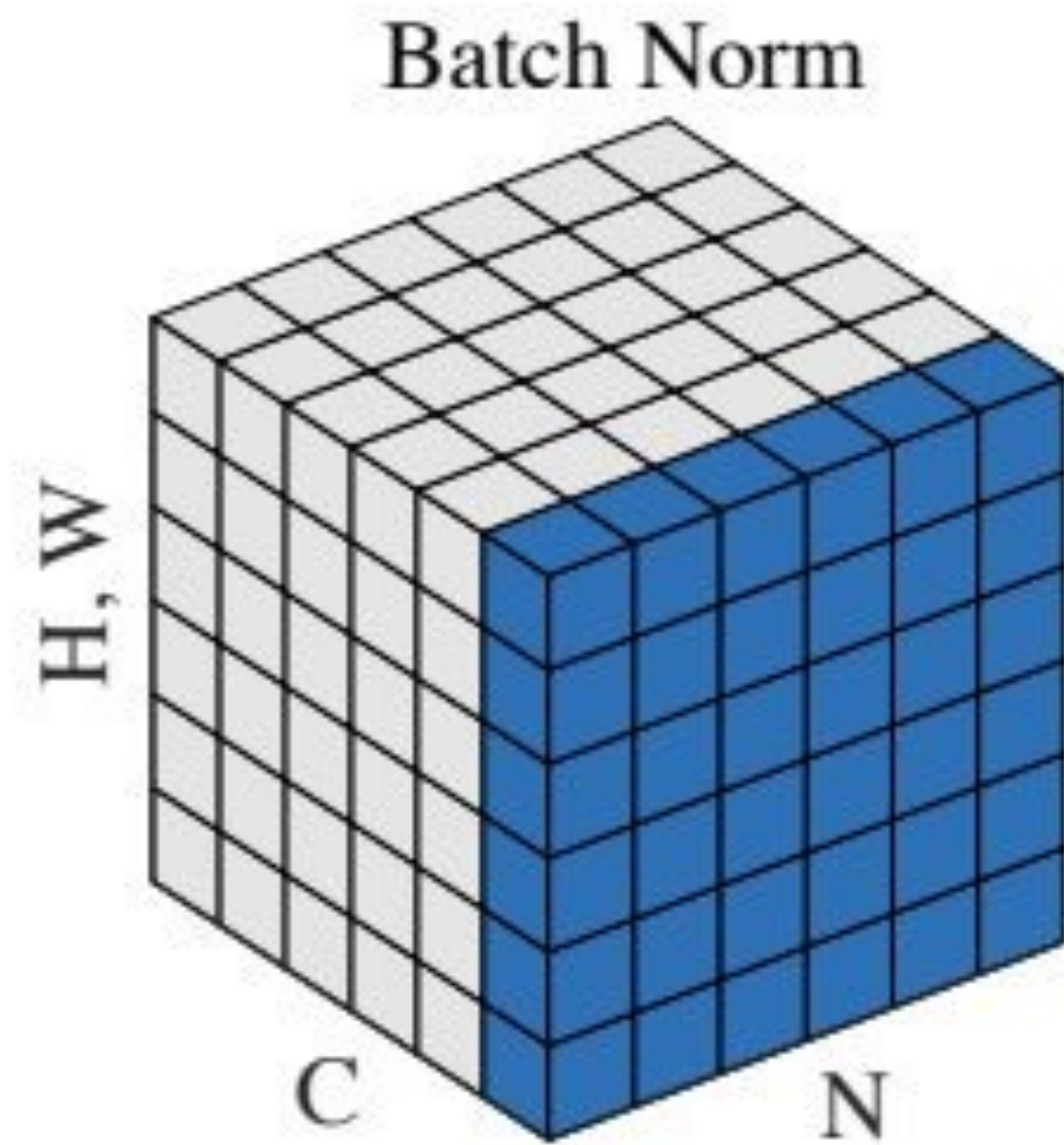
Average and standard deviation each filter response, estimated over **instance**.

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

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

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

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

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]

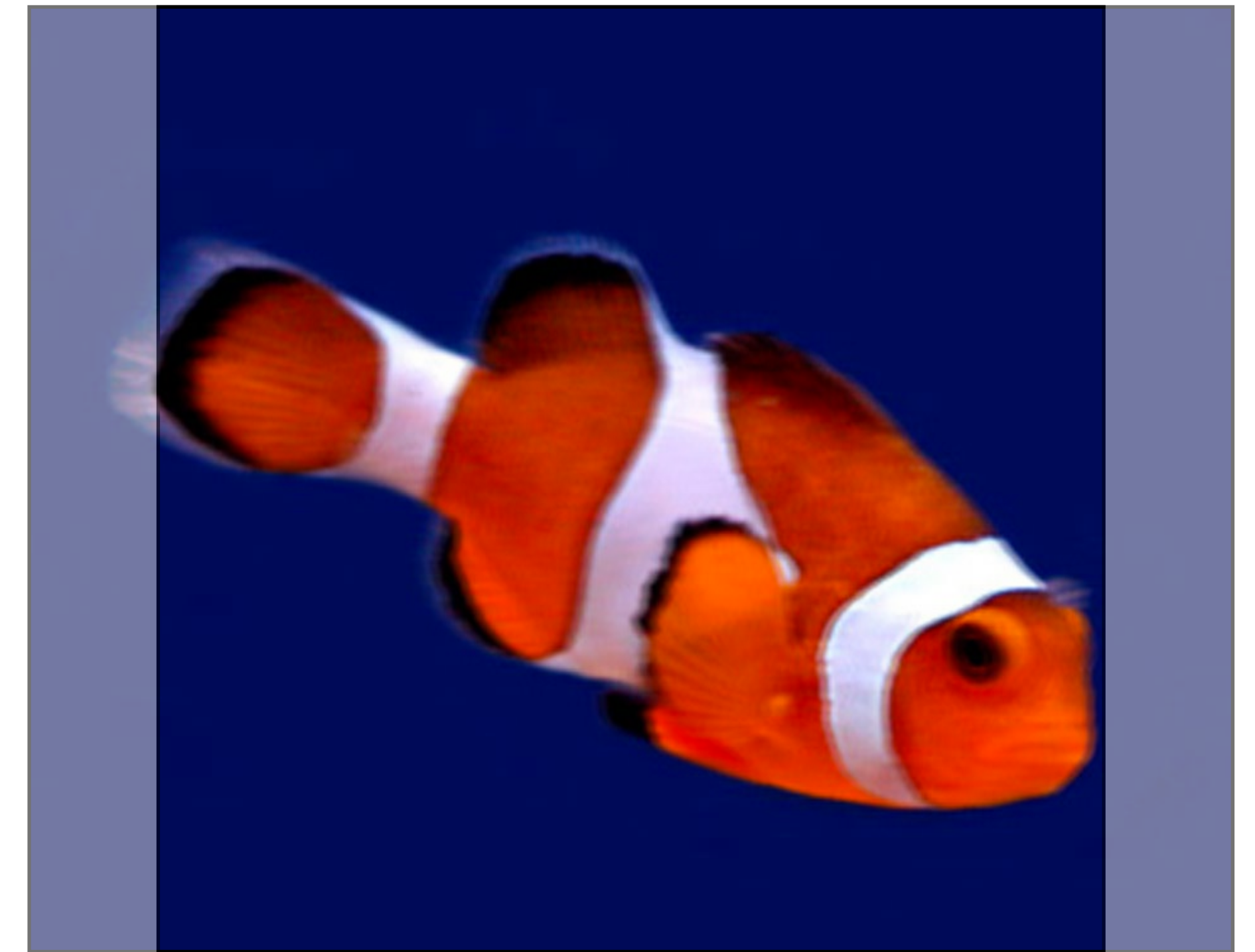
# A few practical issues



# 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

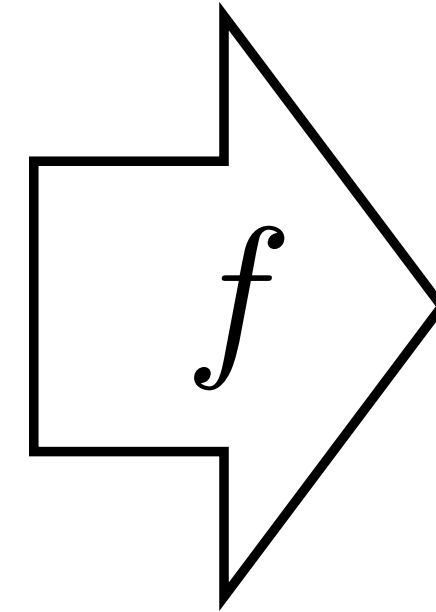


Color jittering

# Beyond image labeling

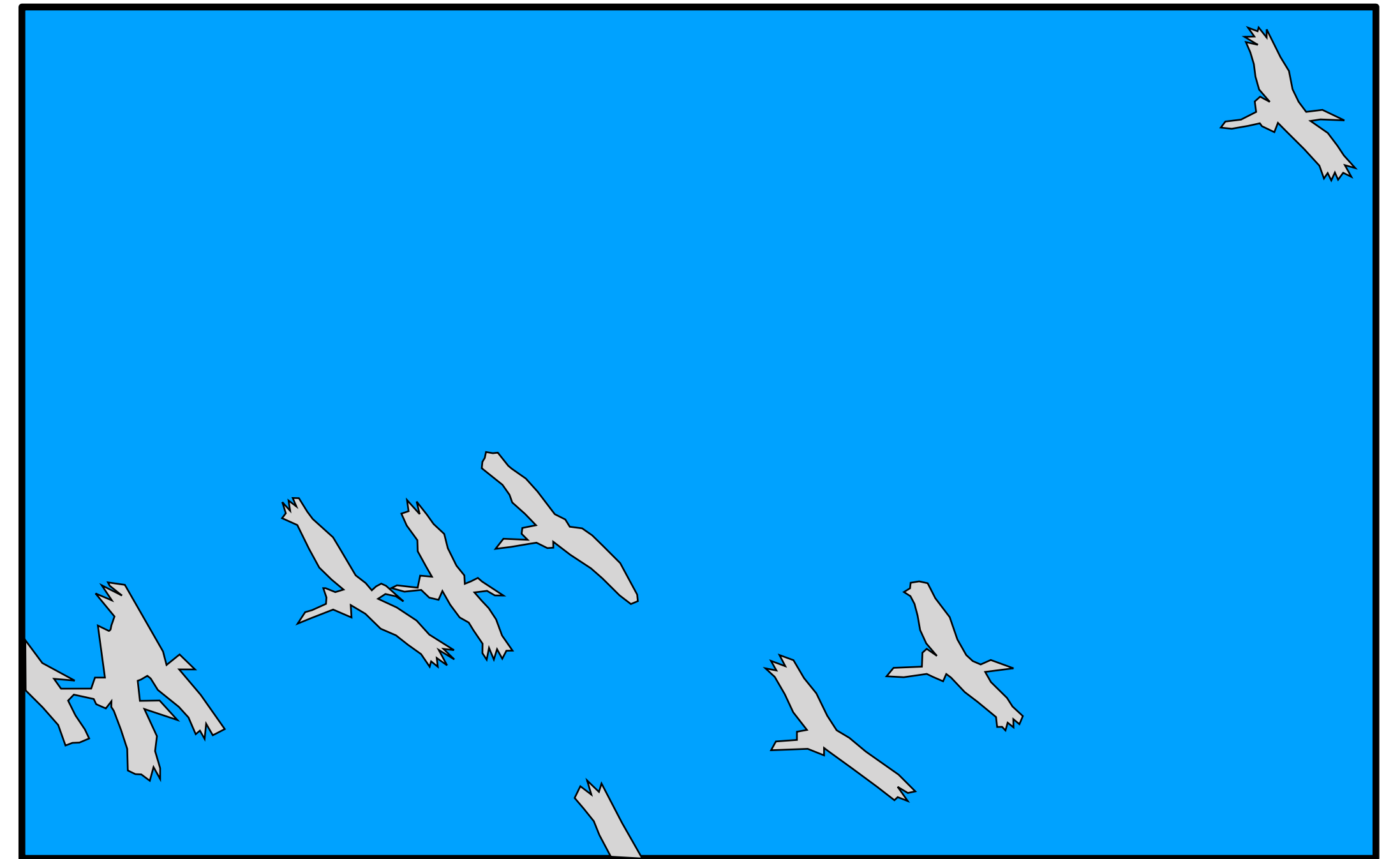
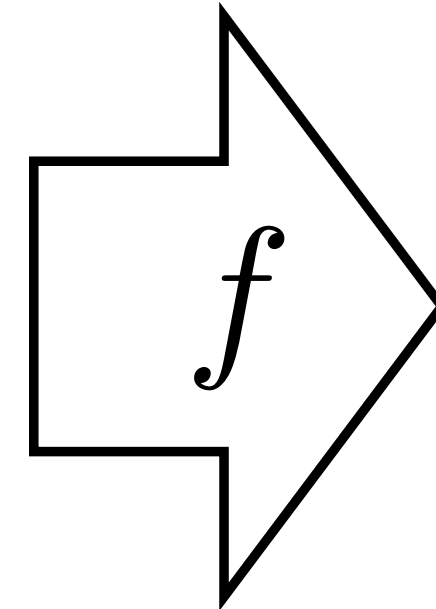


# Object recognition: what objects are in the image?



“Birds”

# Semantic segmentation

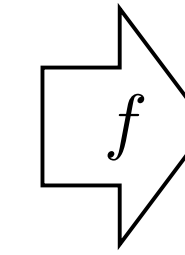
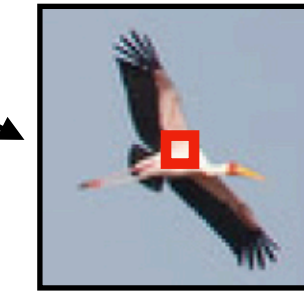
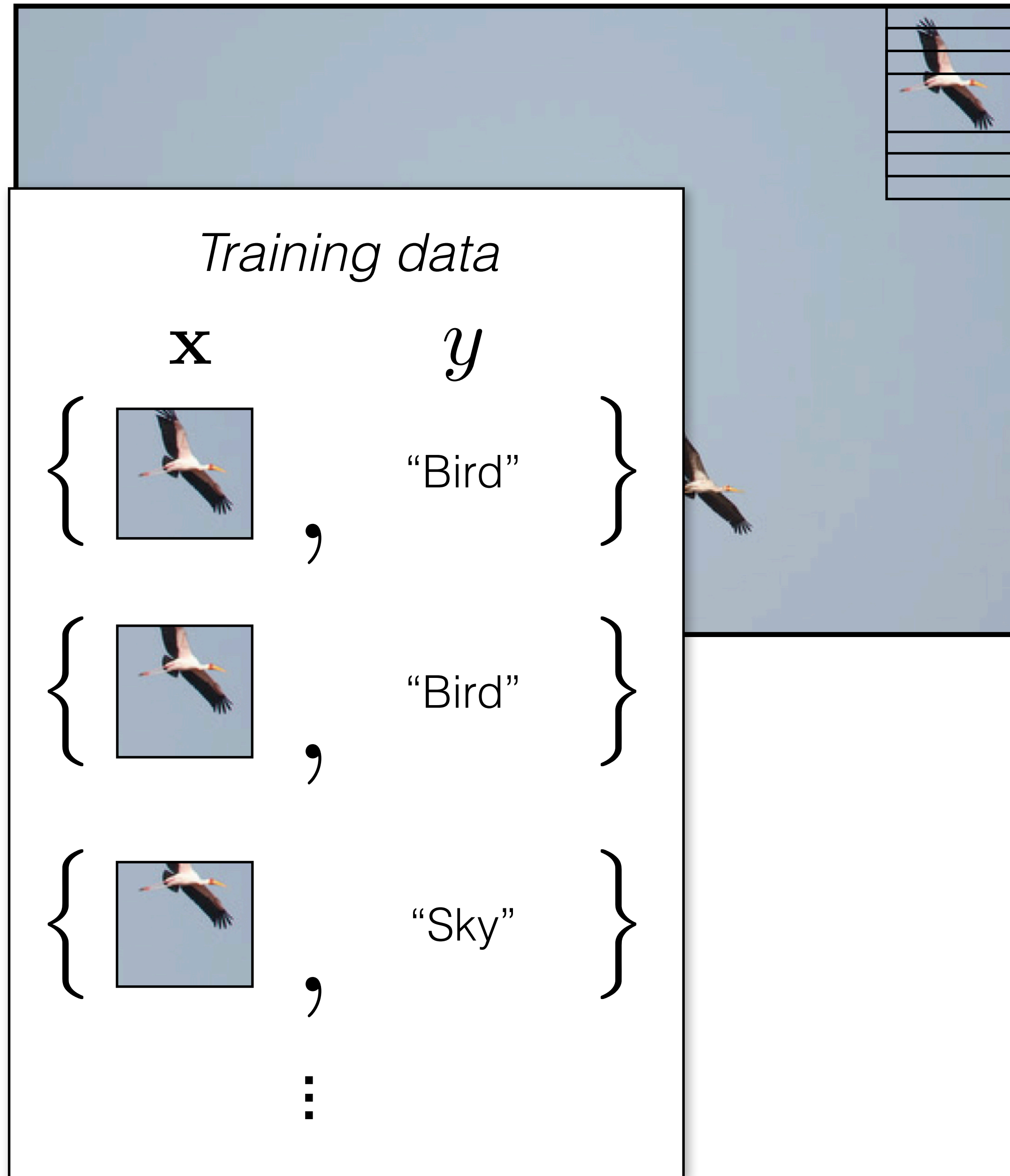


(Colors represent categories)

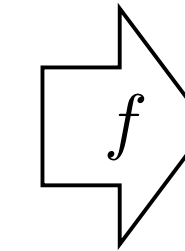
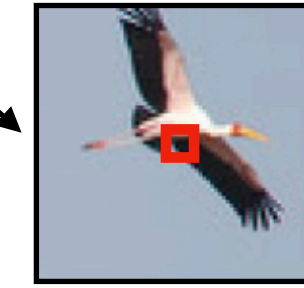
General technique: predict<sup>77</sup> something at every pixel!

**Idea #1: Independently classify windows**

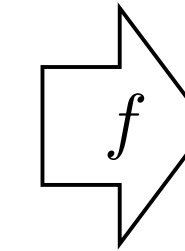
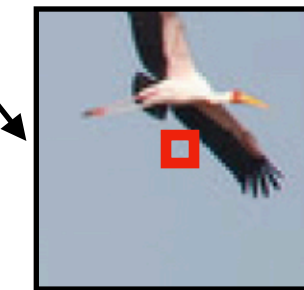
What's the object class of the center pixel?



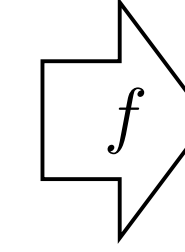
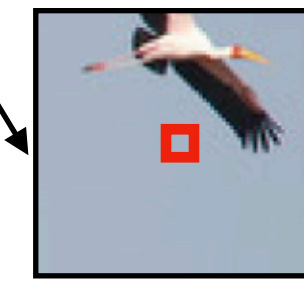
“Bird”



“Bird”



“Sky”

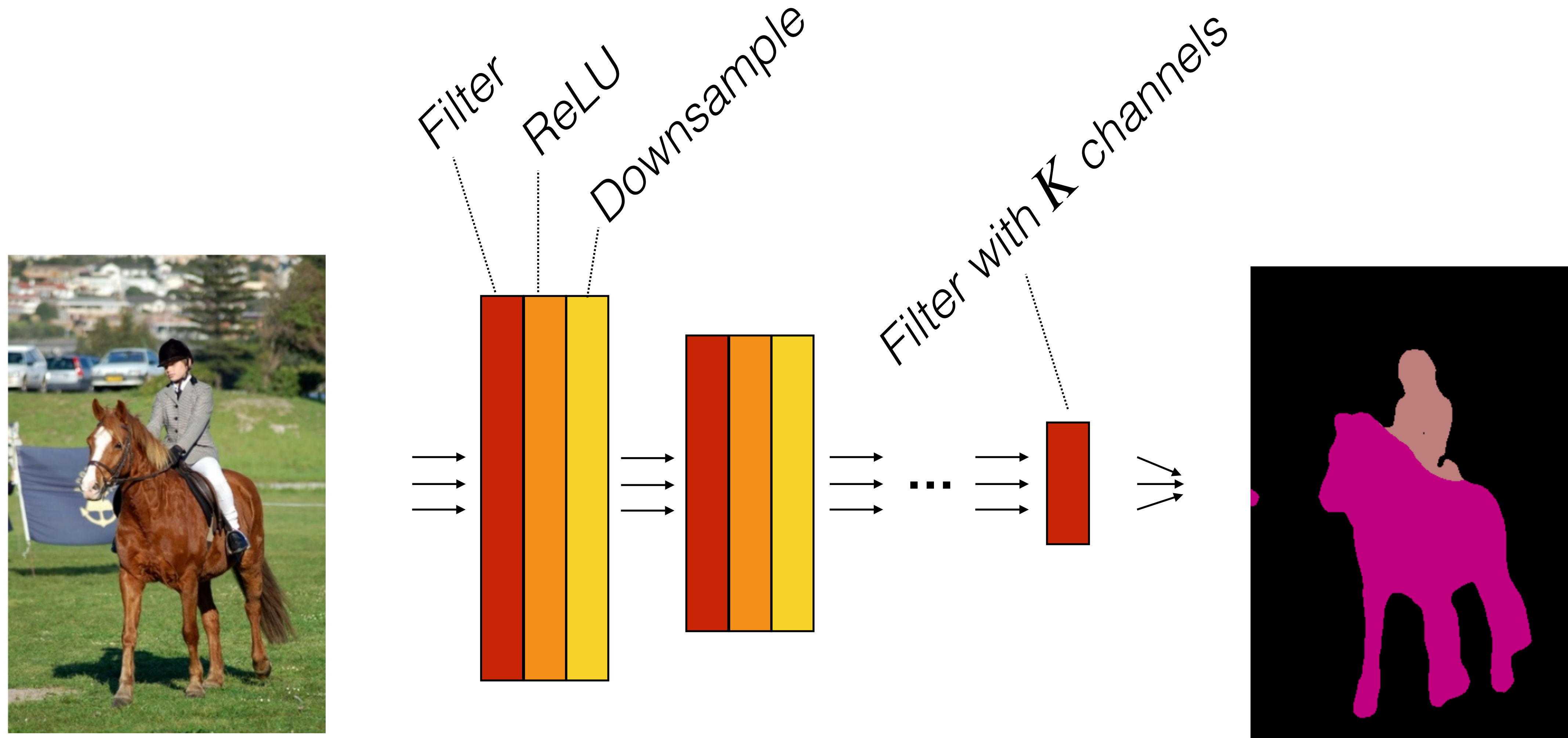


“Sky”

**Idea #2: Fully convolutional networks**



# Fully convolutional network



Classification problem with  $K$  classes

# Idea #3: Dilated convolutions



# Dilated convolutions

3x3

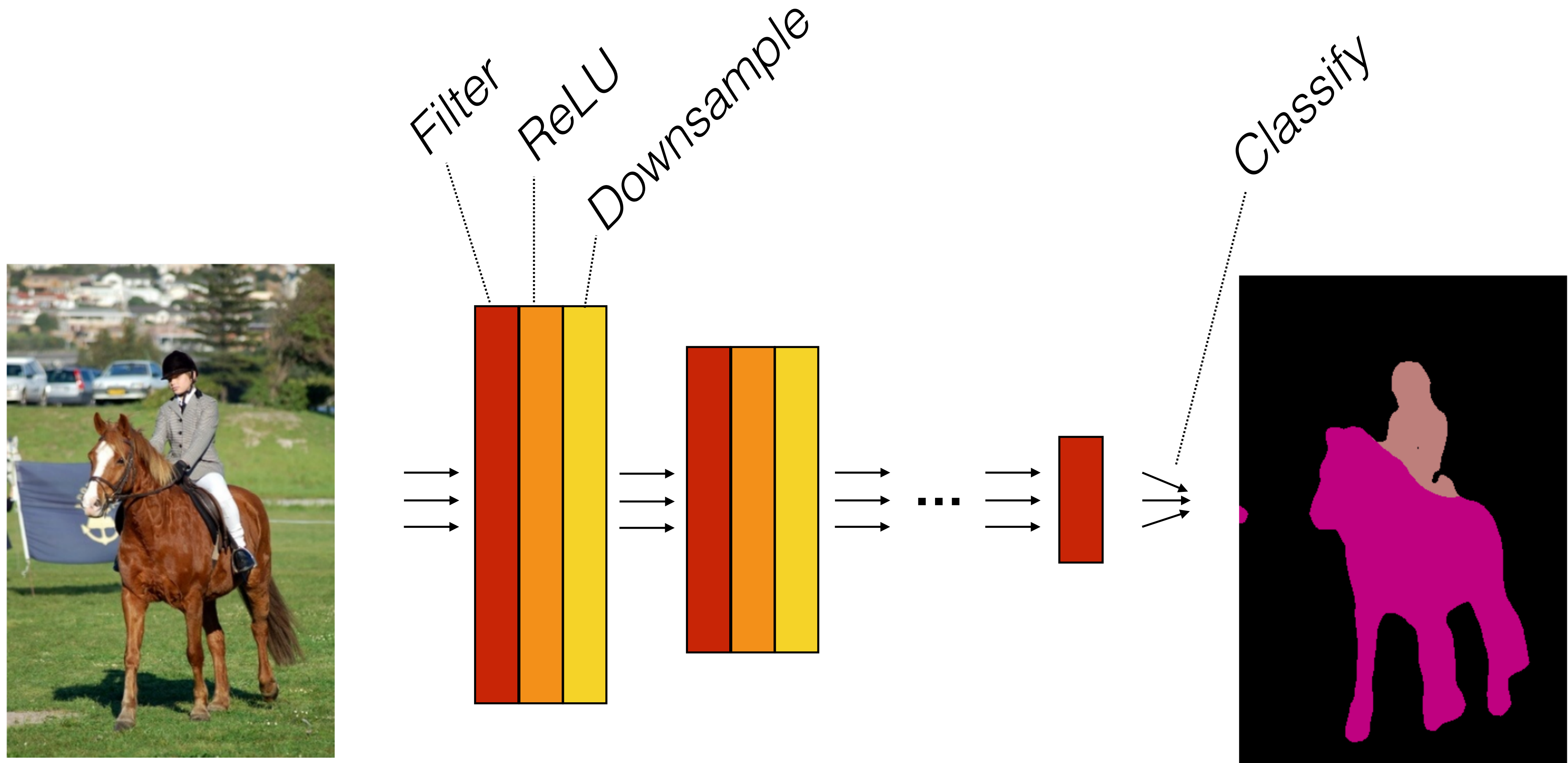
<b>a</b>	<b>b</b>	<b>c</b>
<b>d</b>	<b>e</b>	<b>f</b>
<b>g</b>	<b>h</b>	<b>i</b>

5x5

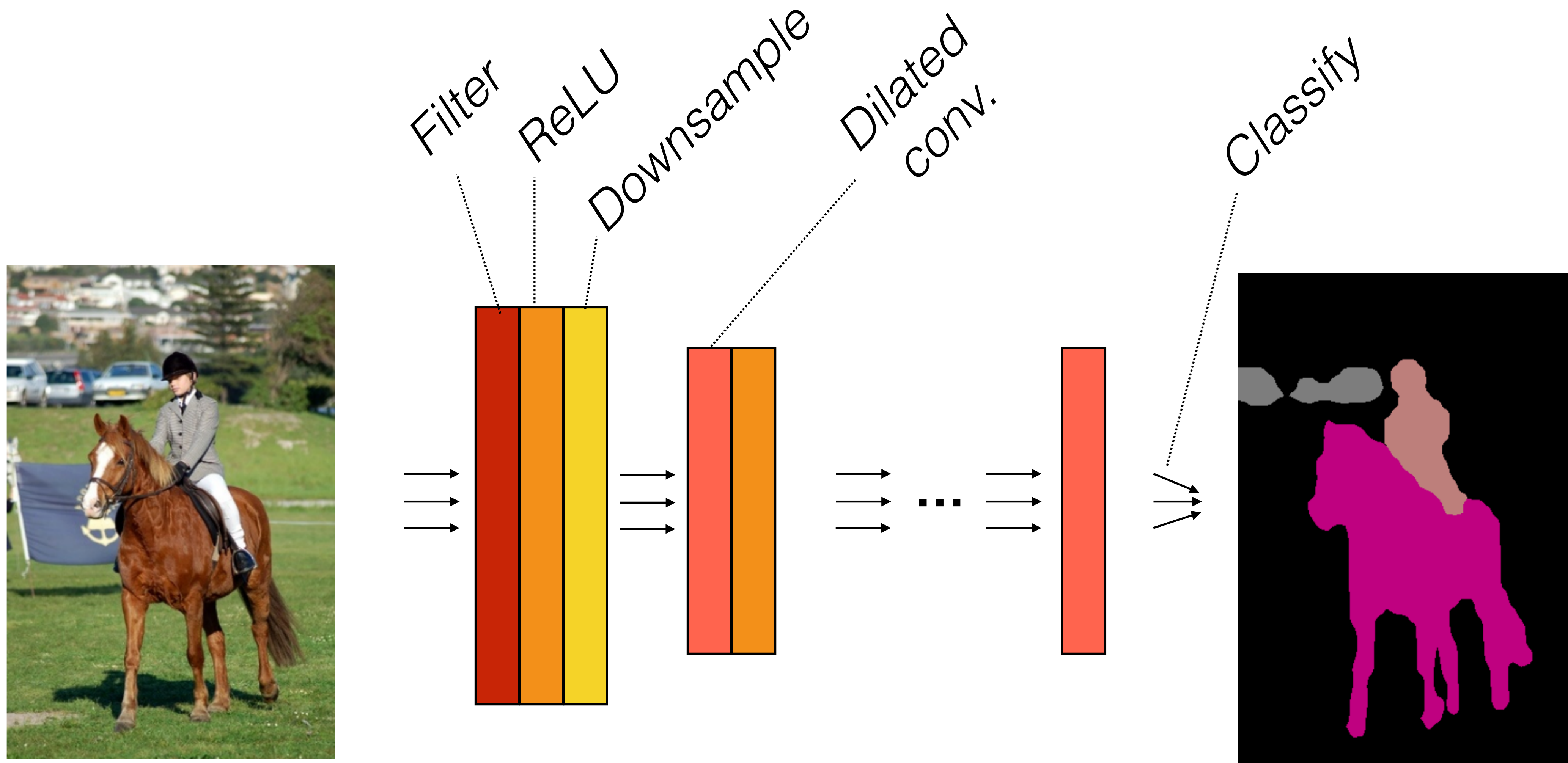
<b>a</b>	0	<b>b</b>	0	<b>c</b>
0	0	0	0	0
<b>d</b>	0	<b>e</b>	0	<b>f</b>
0	0	0	0	0
<b>g</b>	0	<b>h</b>	0	<b>i</b>

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

# CNN without dilated convolutions



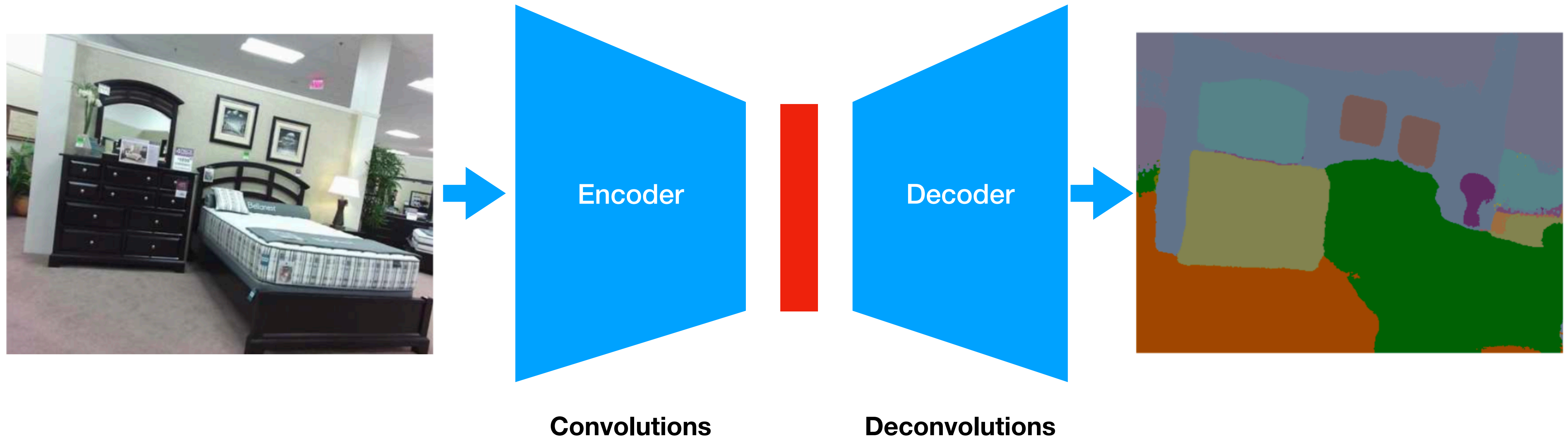
# CNN with dilated convolutions



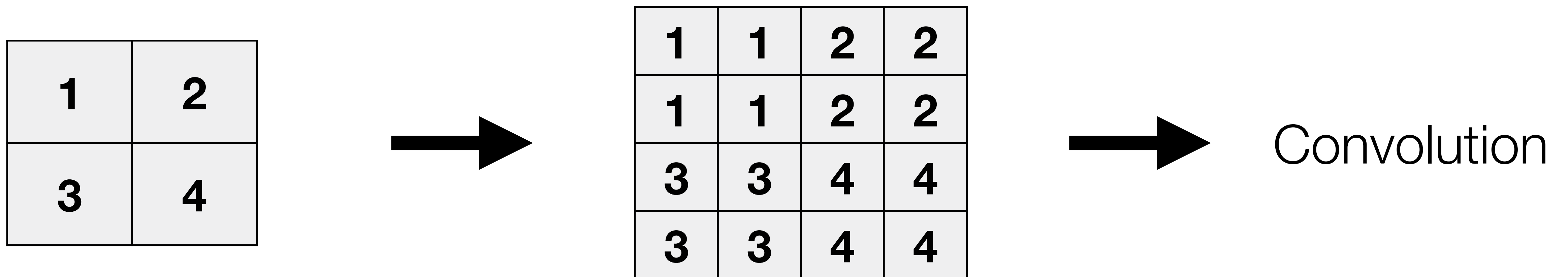
# Idea #4: Encoder-decoder models



# Encoder-decoder architectures

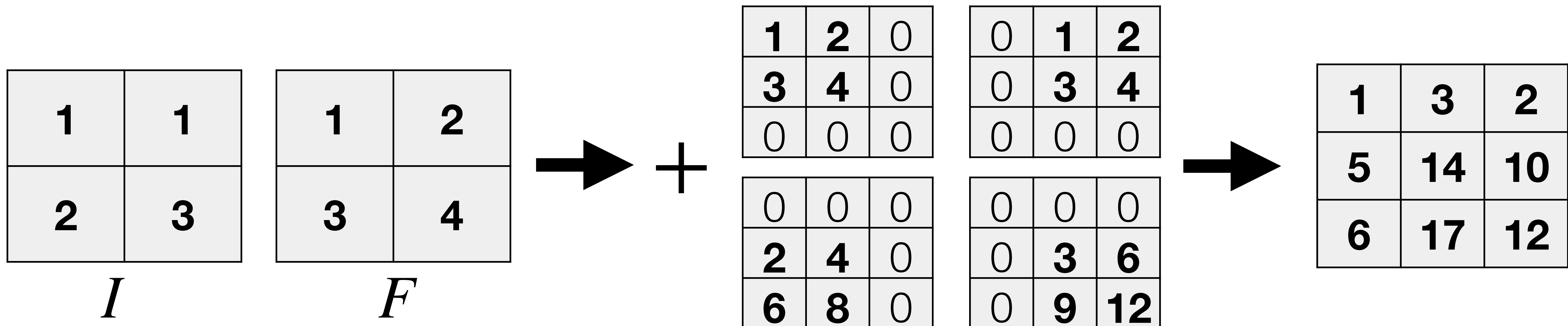


# Upsampling



- 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

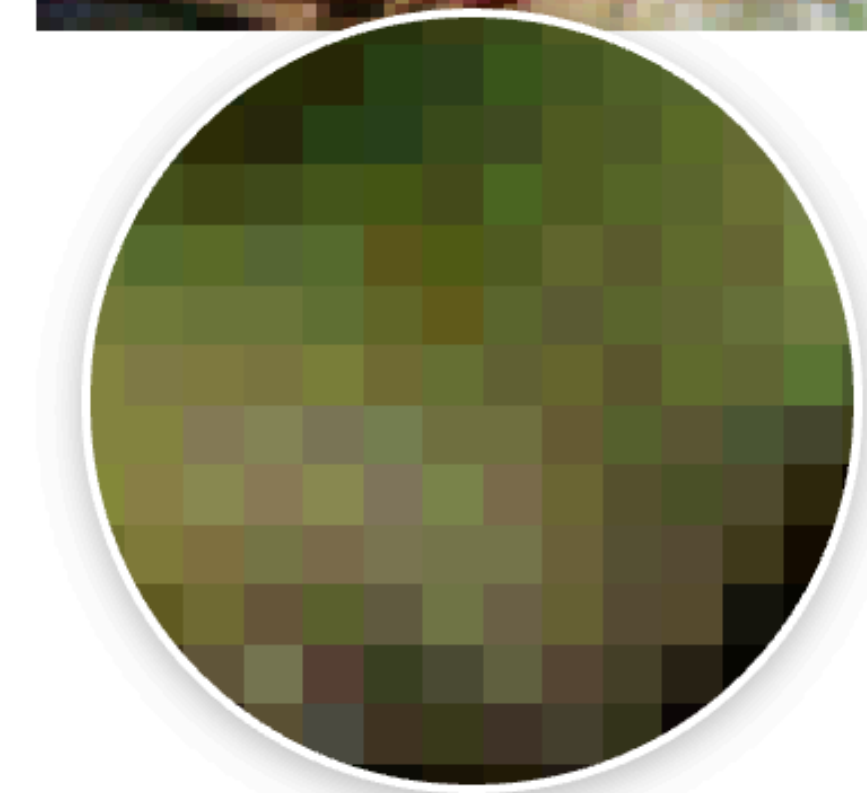
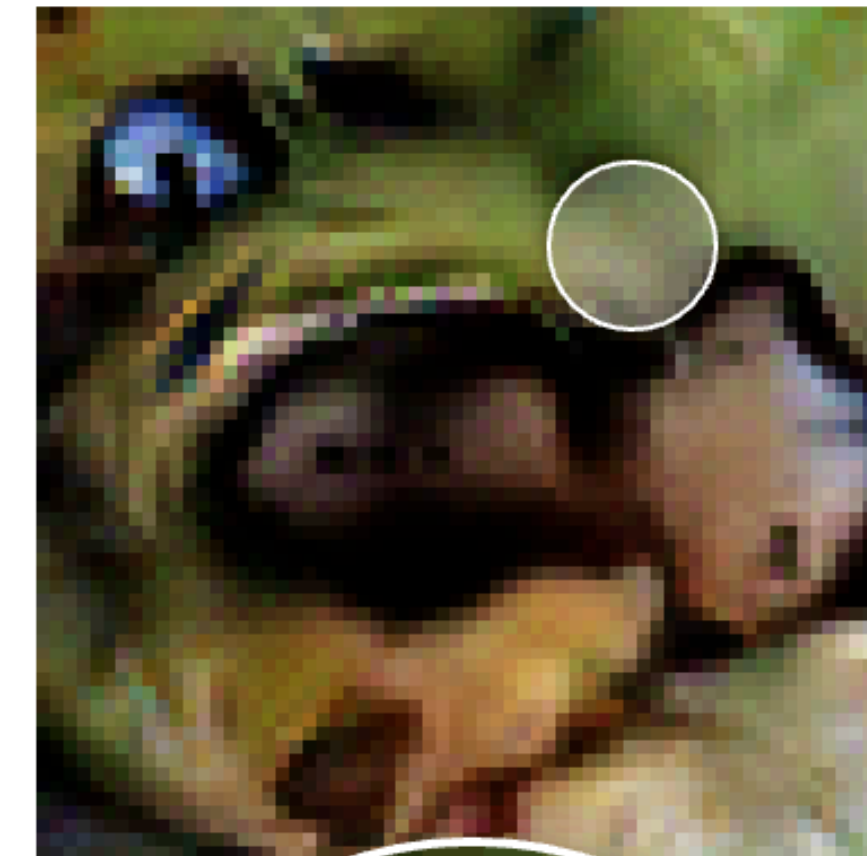
<b>1</b>	<b>1</b>
<b>2</b>	<b>3</b>

*I*

<b>1</b>	<b>2</b>
<b>3</b>	<b>4</b>

*F*

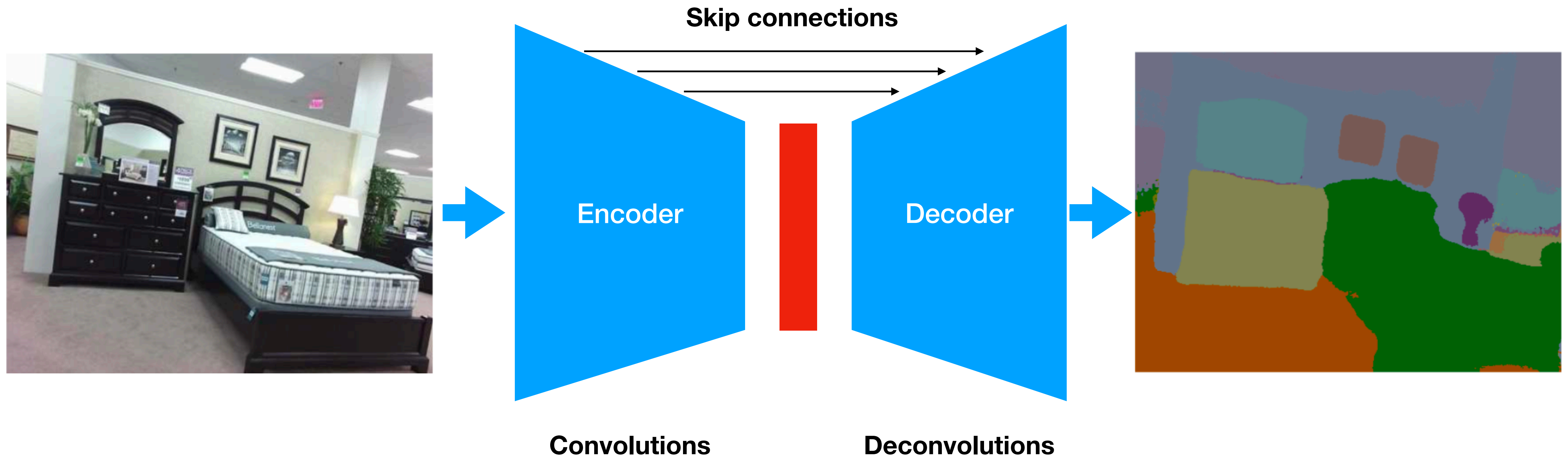
Can lead to “checkerboard” artifacts.



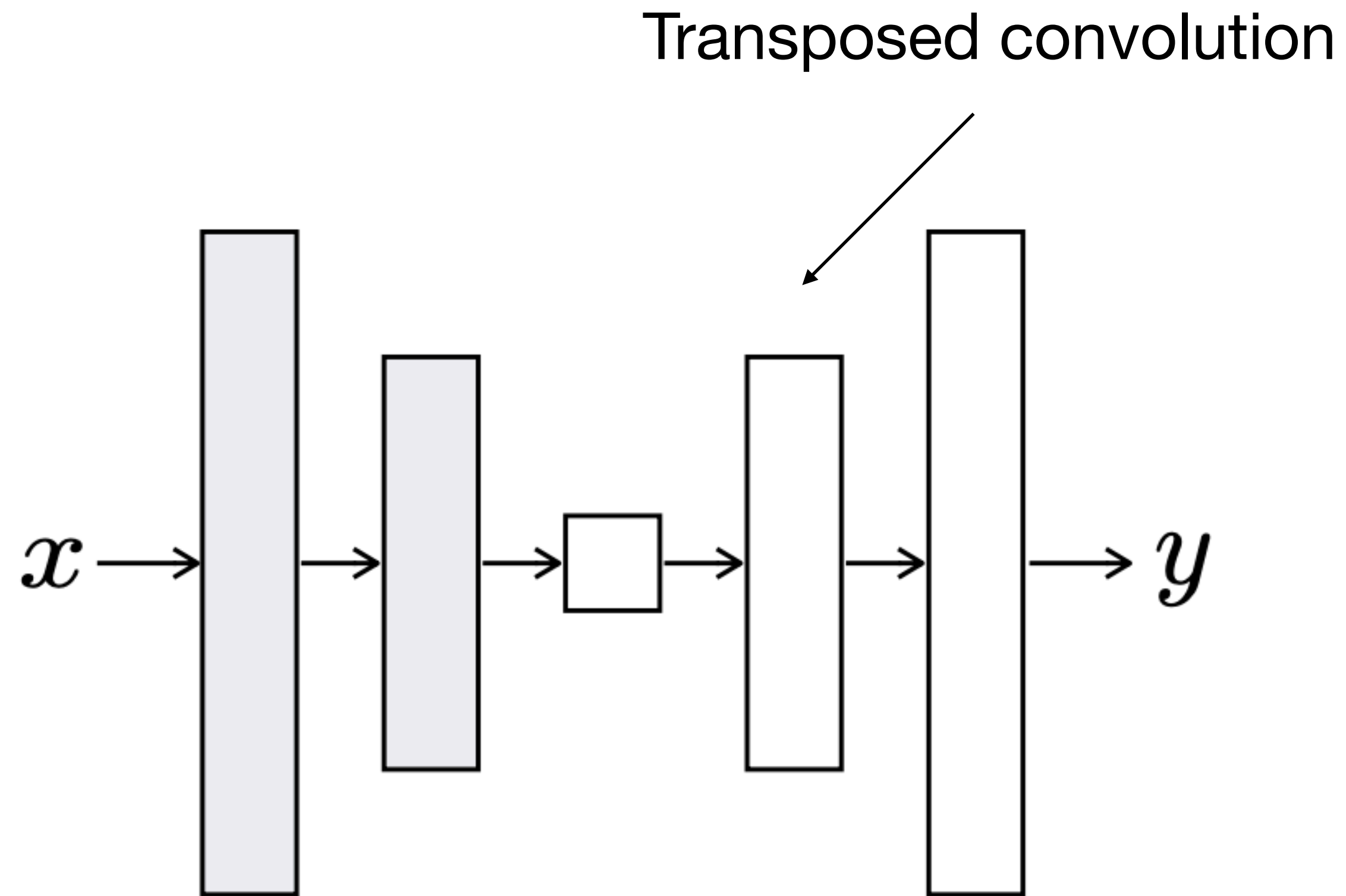
Donahue, et al., 2016 [3]

[Odena et al. Distill article]

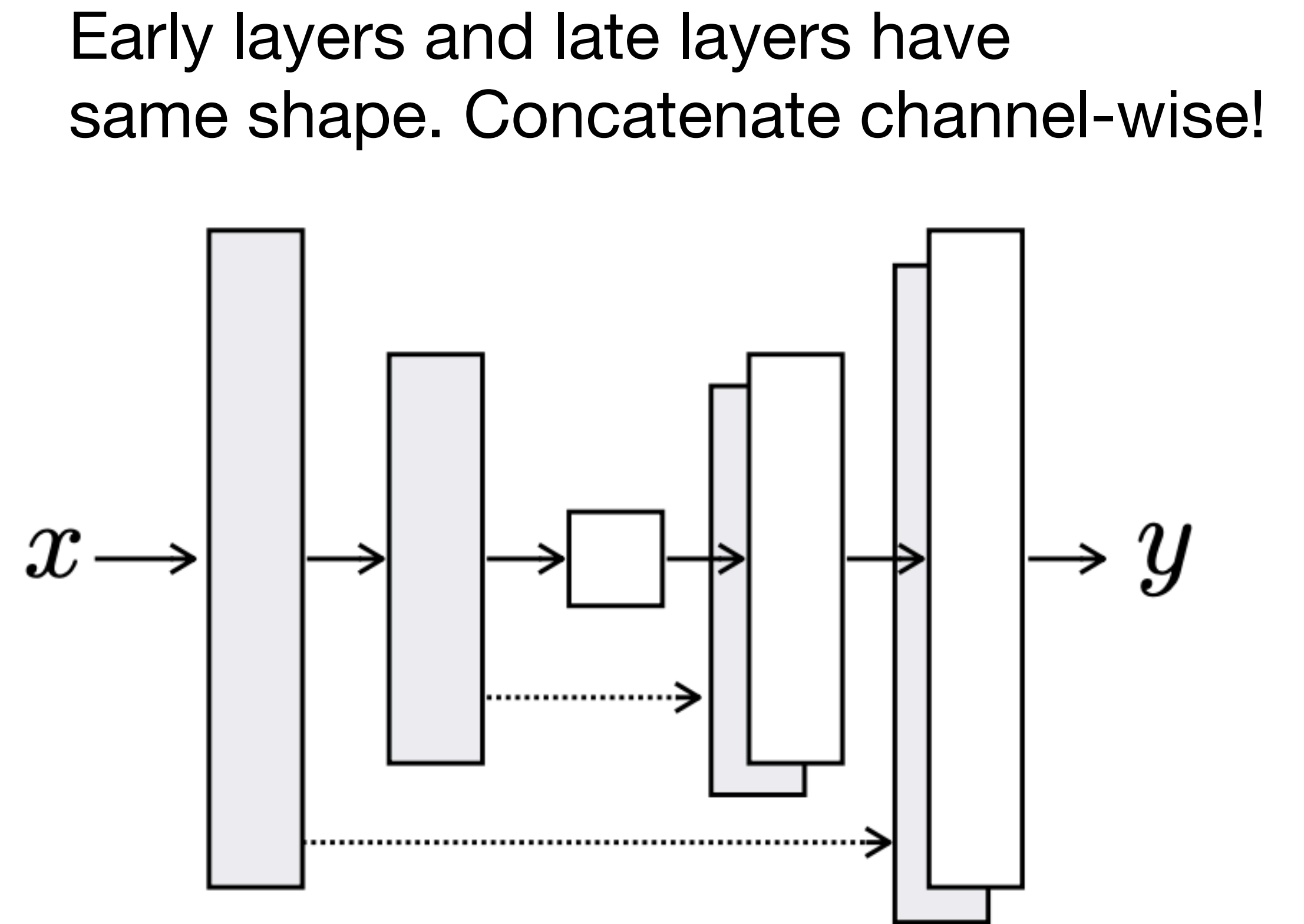
# Encoder-decoder architectures



# Encoder-decoder architectures

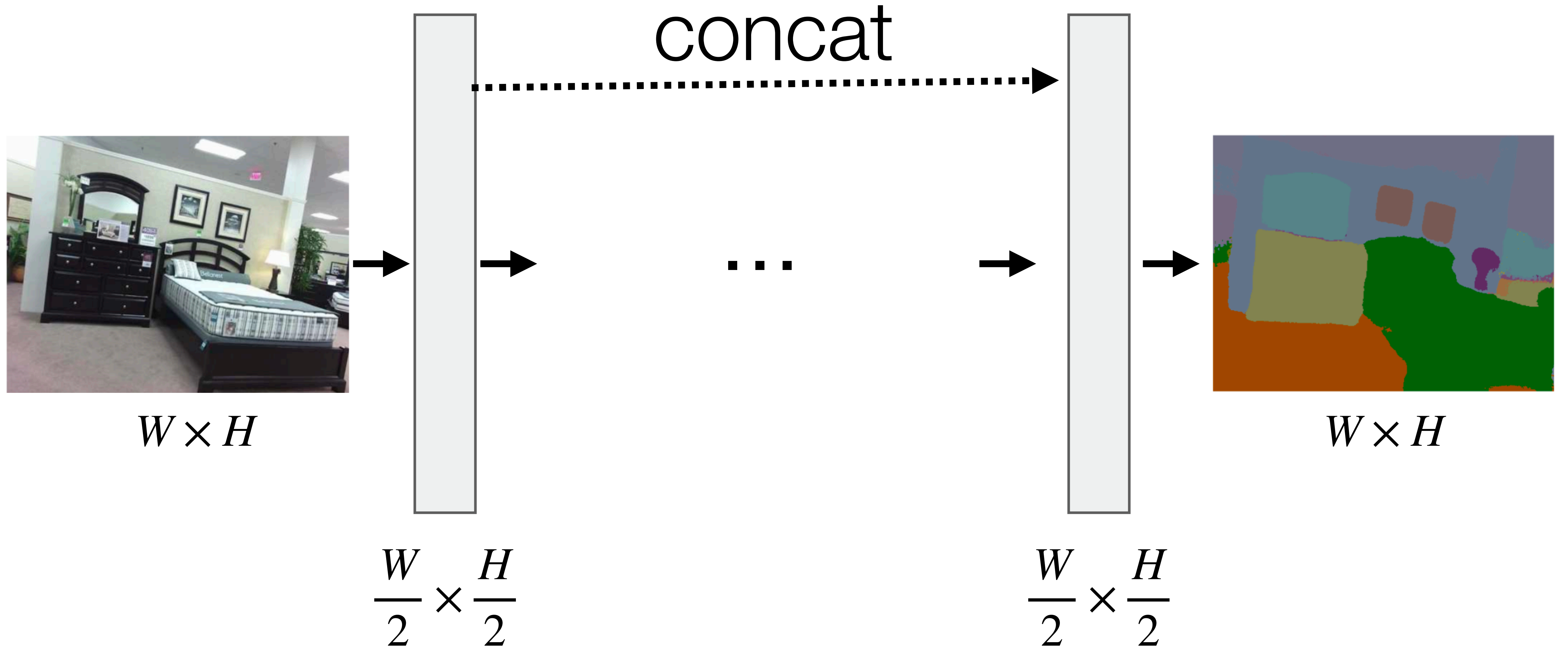


“Vanilla” encoder-decoder architecture

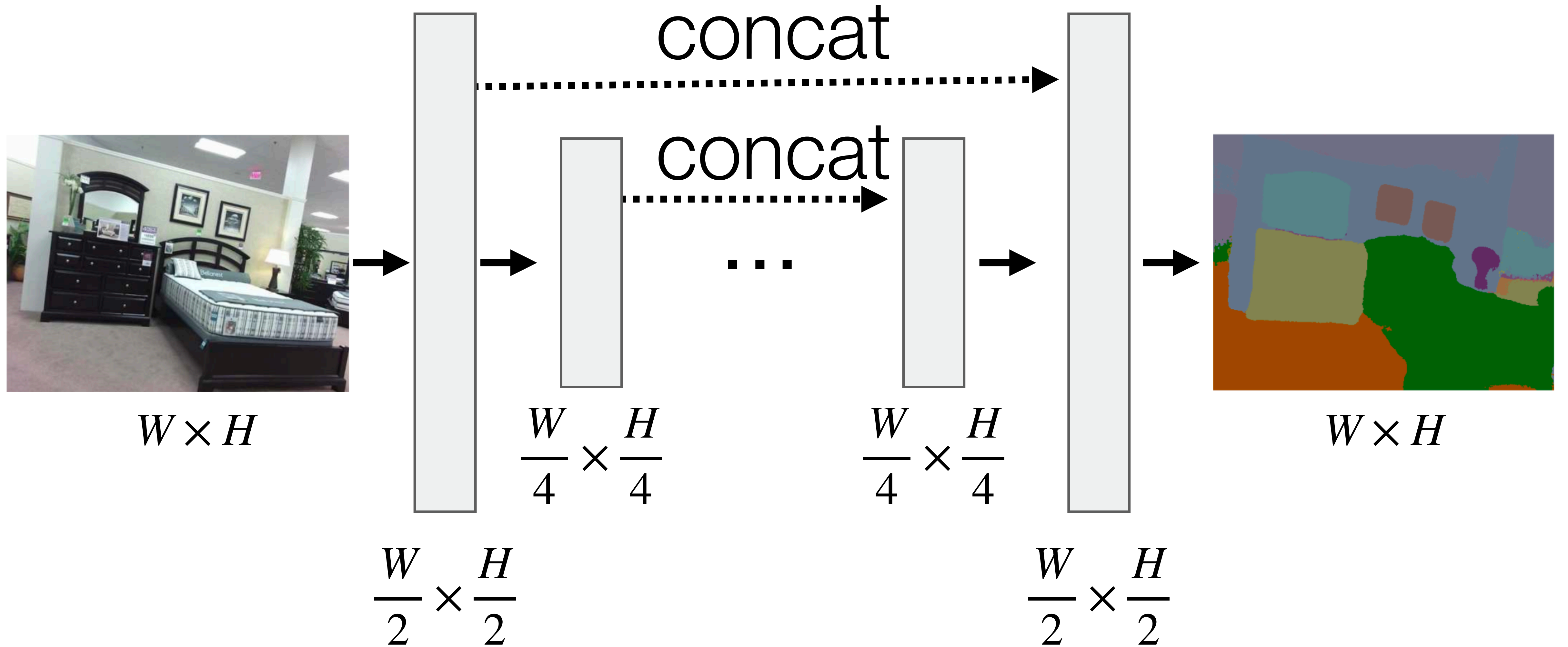


U-Net

# U-net



# U-net





# Encoder-decoder architectures

## SegNet: A Deep Convolutional Encoder-Decoder Architecture for Image Segmentation

Riccardo Cipolla, Senior Member, IEEE,

This paper presents a network architecture for semantic pixel-wise segmentation. In the encoder network, a corresponding decoder network followed by a skip connection. The encoder feature maps to full input resolution feature maps which the decoder upsamples its lower resolution input feature maps. The decoder upsamples its lower resolution input feature maps by a max-pooling step of the corresponding encoder to produce feature maps. The upsampled maps are sparse and are then combined with the proposed architecture with the widely adopted FCN [1]. This comparison reveals the memory versus accuracy trade-off.

The network is designed to be efficient both in terms of memory and number of trainable parameters than other competing architectures. We also performed a controlled benchmark of SegNet on semantic segmentation tasks. These quantitative assessments show that SegNet is the most efficient inference memory-wise as compared to other architectures. A web demo at <http://mi.eng.cam.ac.uk/projects/segnet/>

Outdoor Scenes, Indoor Scenes, Road Scenes, Encoder,

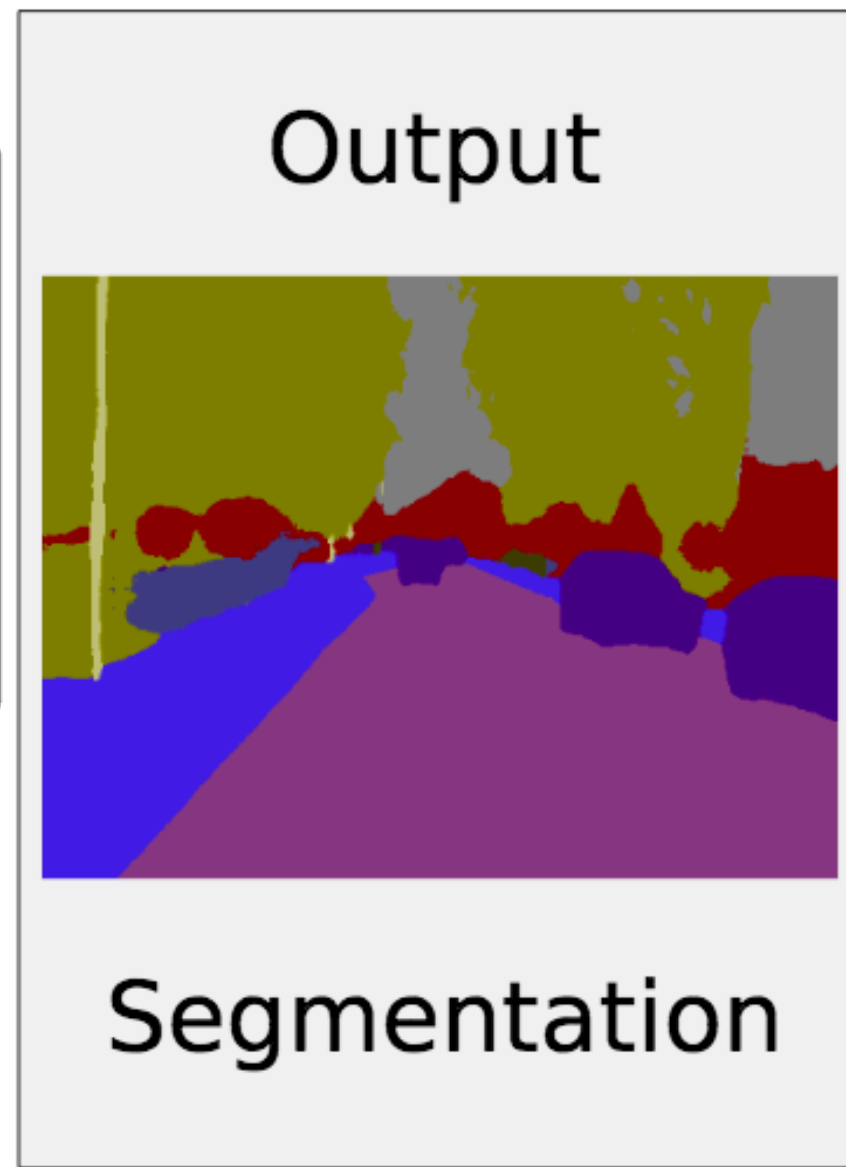
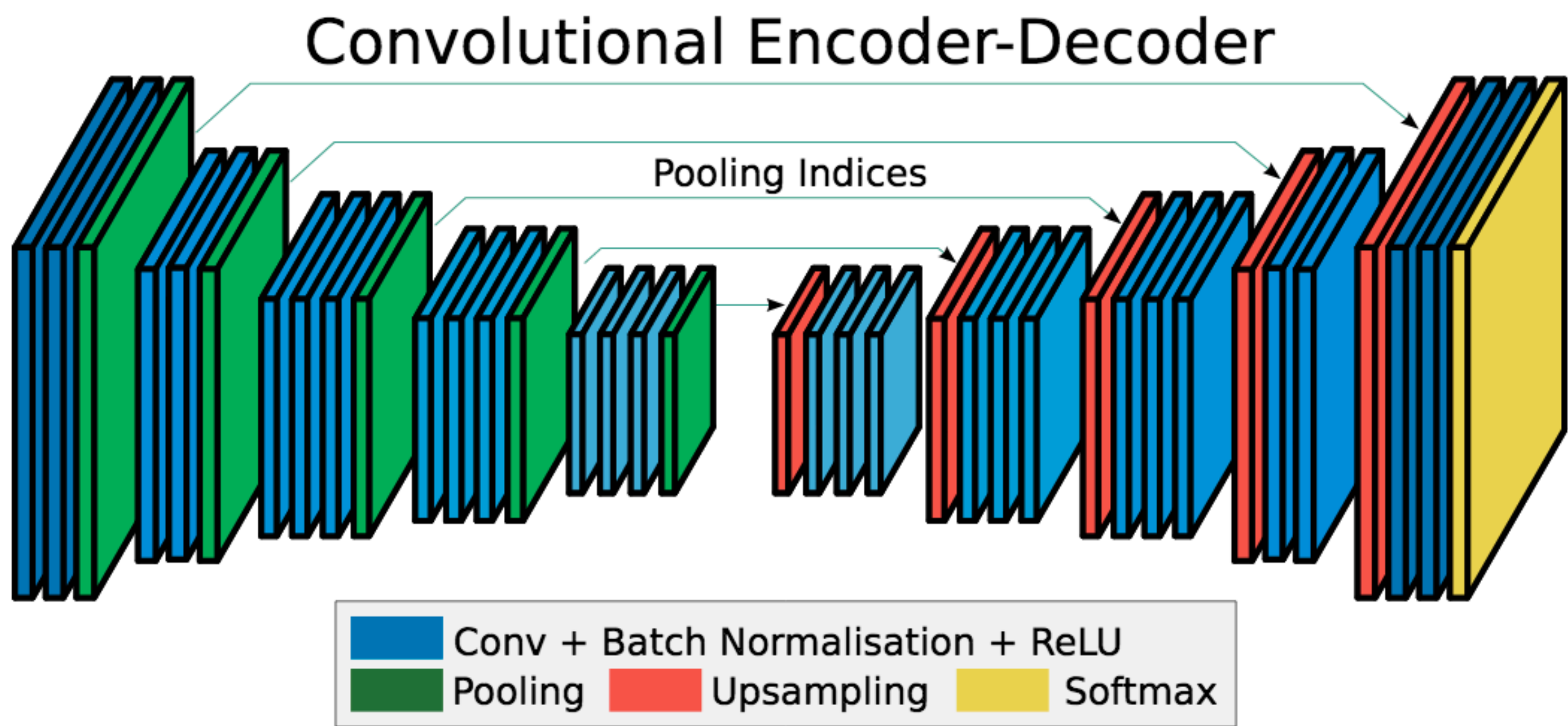
to understand the spatial-relationship (context) between different classes such as road and side-walk. In typical road scenes, the majority of the pixels belong to large classes such as building and hence the network must produce smooth boundaries. The engine must also have the ability to delineate boundaries based on their shape despite their small size. Hence it is necessary to retain boundary information in the extracted image features. From a computational perspective, it is necessary for the network to be efficient in terms of both memory and inference time during inference. The ability to train end-to-end and to jointly optimise all the weights in the network using an efficient weight update technique such as stochastic gradient descent (SGD) [17] is an additional benefit since it is more easily repeatable. The design of SegNet arose from a need to match these criteria.

This is primarily because max pooling and sub-sampling reduce feature map resolution. Our motivation to design SegNet arises from this need to map low resolution features to input resolution for pixel-wise classification. This mapping must produce features which are useful for accurate boundary localization.

Our architecture, SegNet, is designed to be an efficient architecture for pixel-wise semantic segmentation. It is primarily motivated by road scene understanding applications which require the ability to model appearance (road, building), shape (cars,

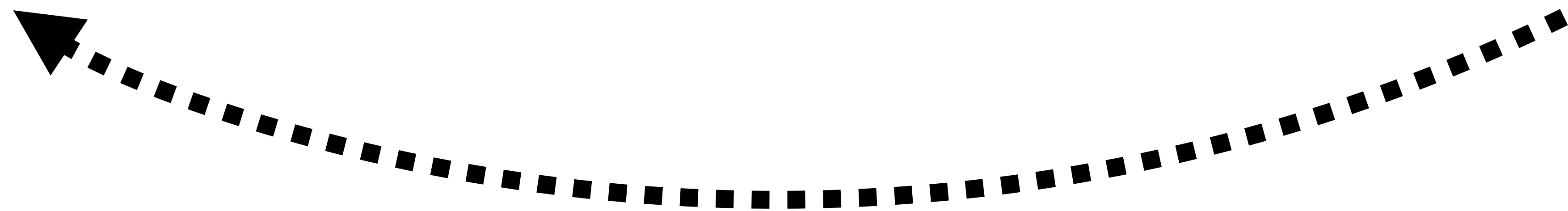
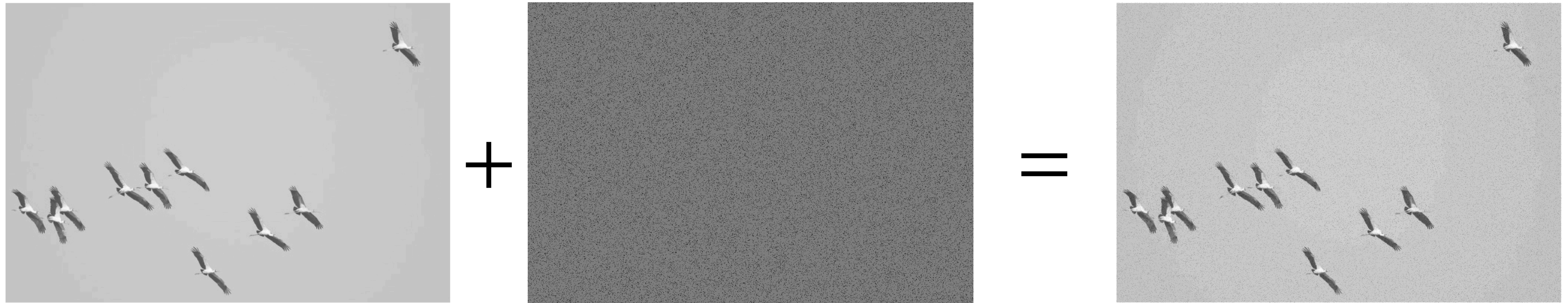
The encoder network in SegNet is topologically identical to the convolutional layers in VGG16 [1]. We remove the fully connected layers of VGG16 which makes the SegNet encoder network significantly smaller and easier to train than many other recent architectures [2], [4], [11], [18]. The key component of SegNet is the decoder network which consists of a hierarchy of decoders one corresponding to each encoder. Of these, the appropriate decoders use the max-pooling indices received from the corresponding encoder to perform non-linear upsampling of their input feature maps. This idea was inspired from an architecture designed for unsupervised feature learning [19]. Reusing max-pooling indices in the decoding process has several practical

• V. Badrinarayanan, A. Kendall, R. Cipolla are with the Machine Intelligence Lab, Department of Engineering, University of Cambridge, UK. E-mail: vb292,agk34,cipolla@eng.cam.ac.uk





# Other uses for U-nets



Goal: recover the original image

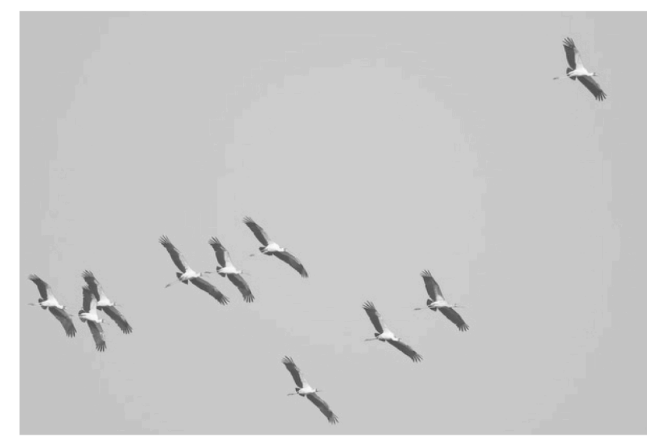
Recall: denoising problem



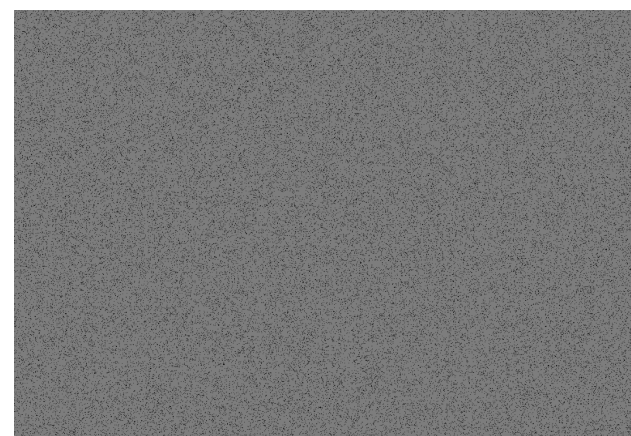
# Denoising



=

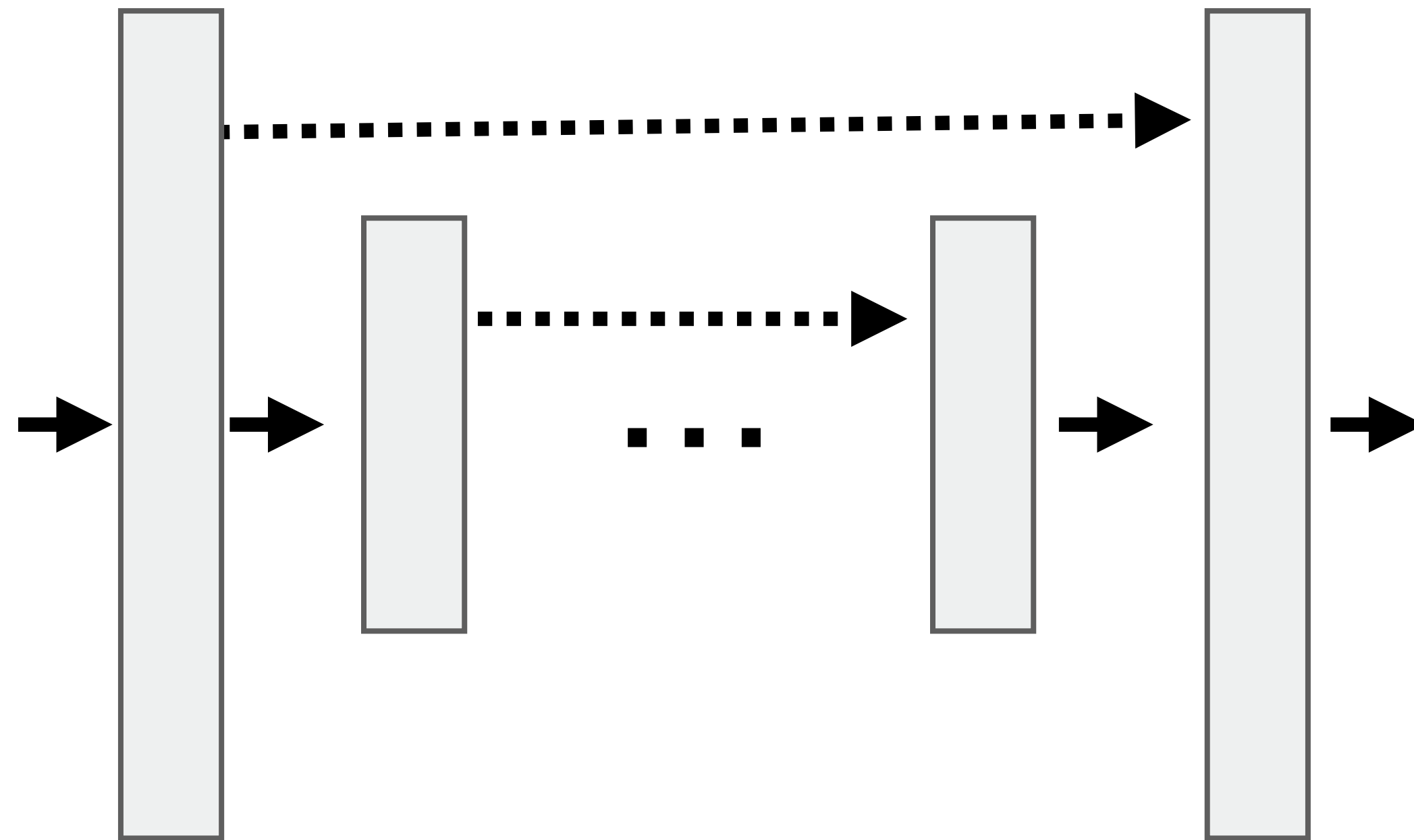


+



$\mathbf{X}_{\text{clean}}$

random noise



$\hat{\mathbf{x}}$

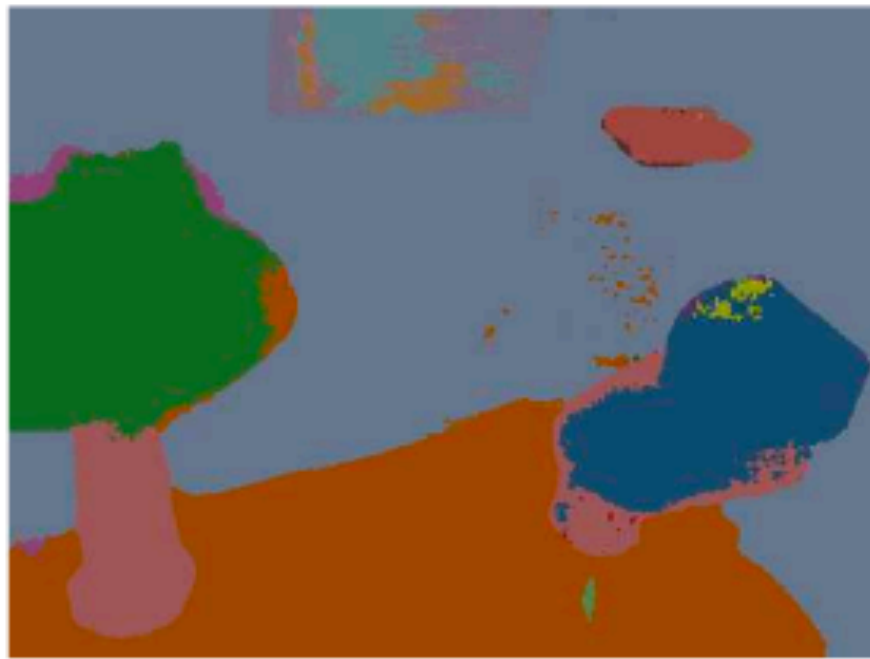
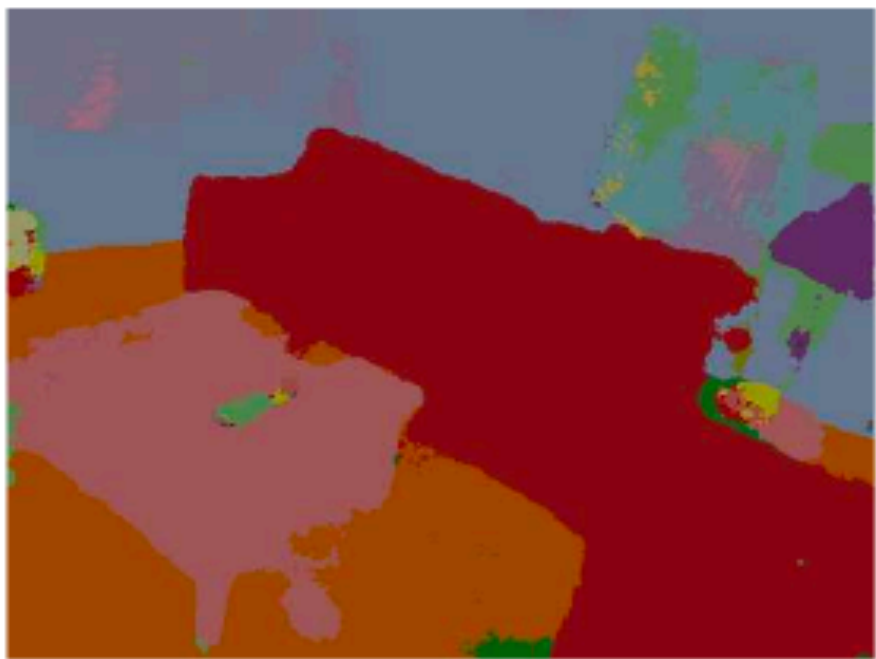
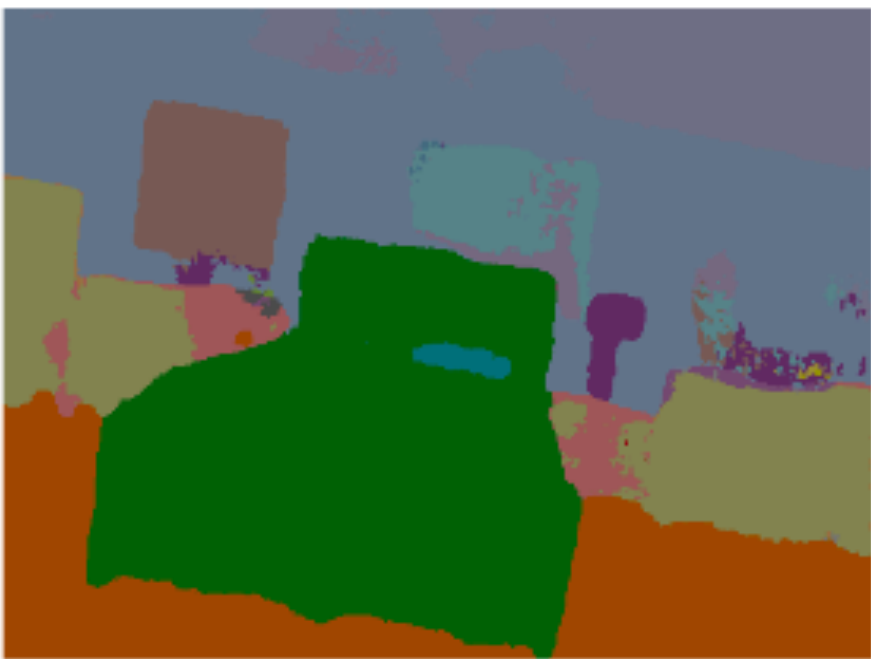
Loss:  $\|\mathbf{x}_{\text{clean}} - \hat{\mathbf{x}}\|^2$



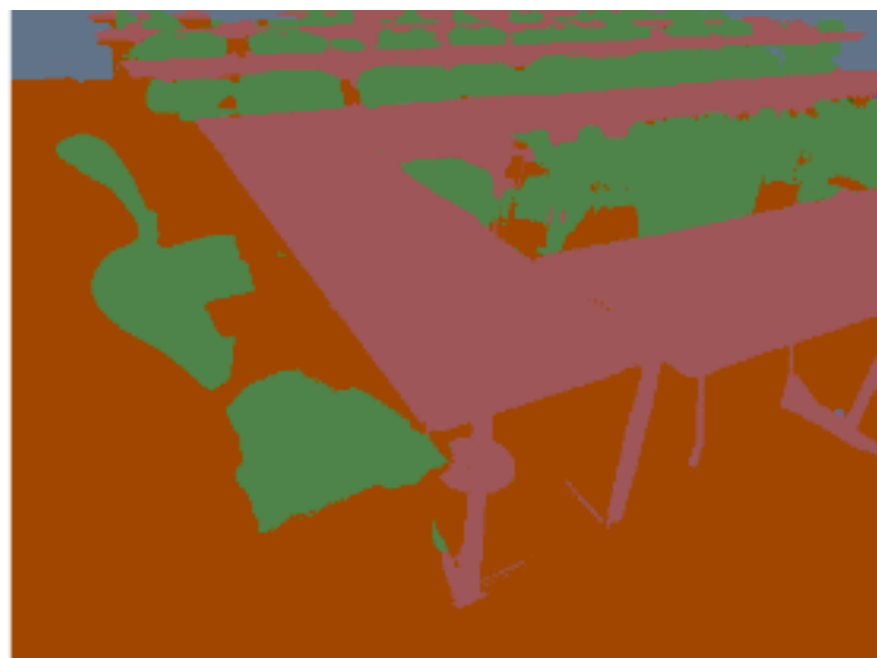
**Input**



**Segnet**



**FCN**





# PS5: CNNs



- Train a convolutional network to recognize scenes
- Use PyTorch + autodiff
- Train on GPU

**Next class:** image generation with GANs