

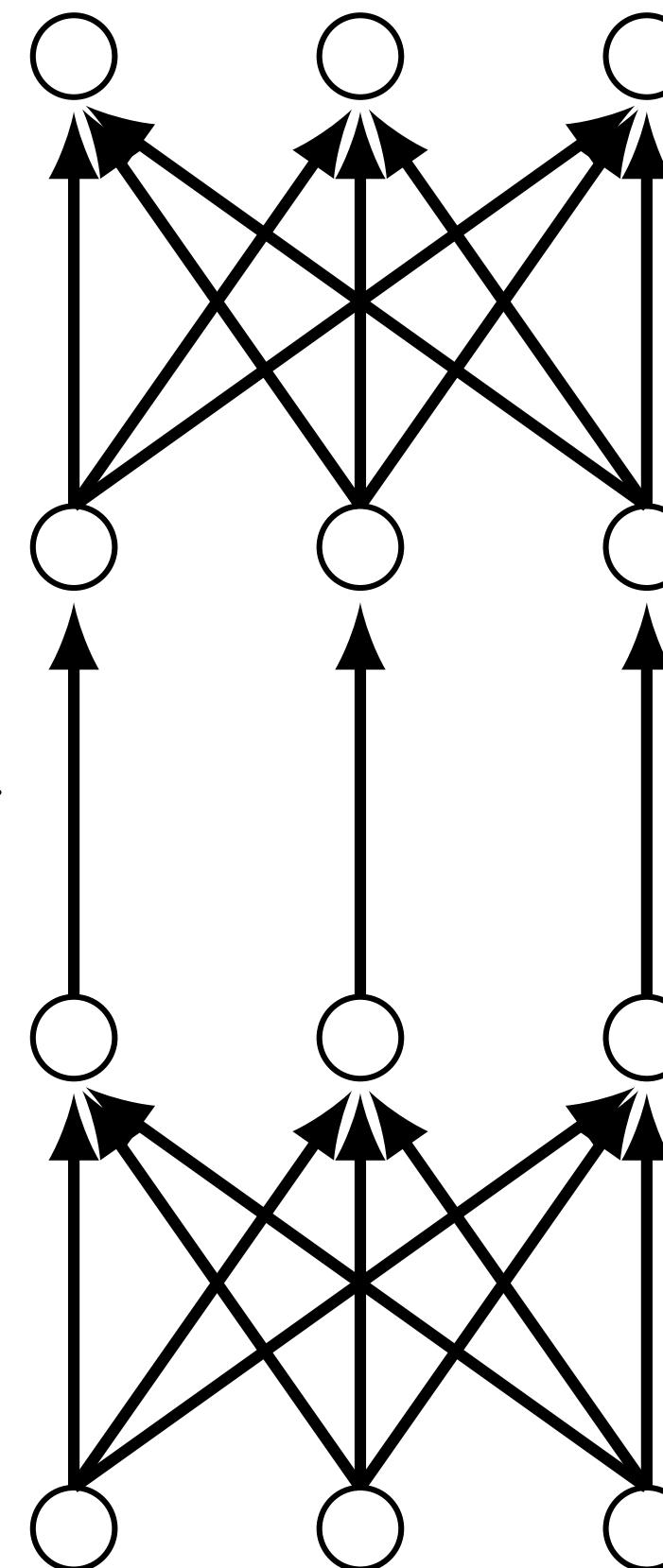
# Lecture 4: Architectures for Grids

Speaker: Sara Beery



# Multilayer Perceptron

linear comb. of neurons ▷  
neuron-wise nonlinearity ▷  
linear comb. of neurons ▷

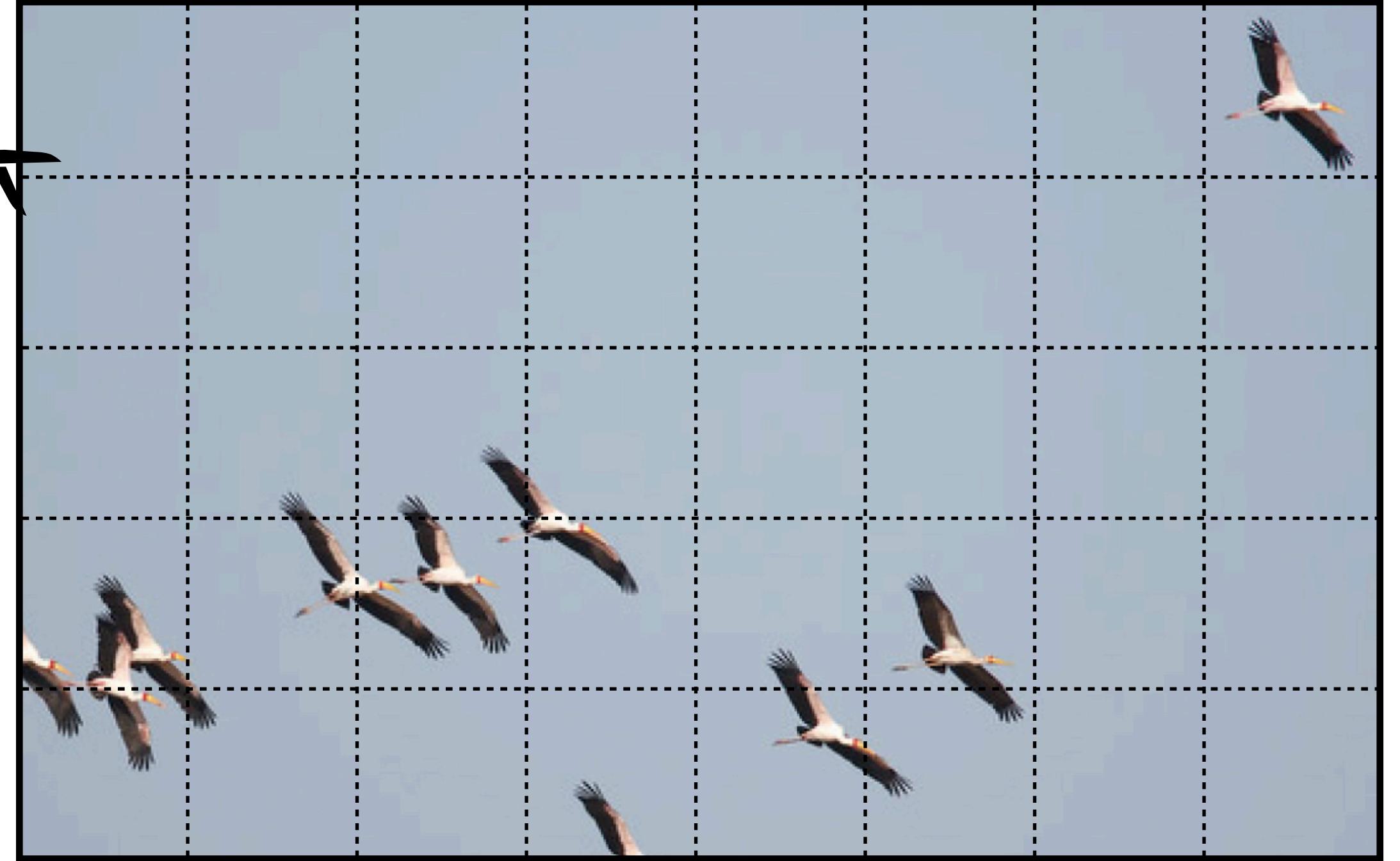


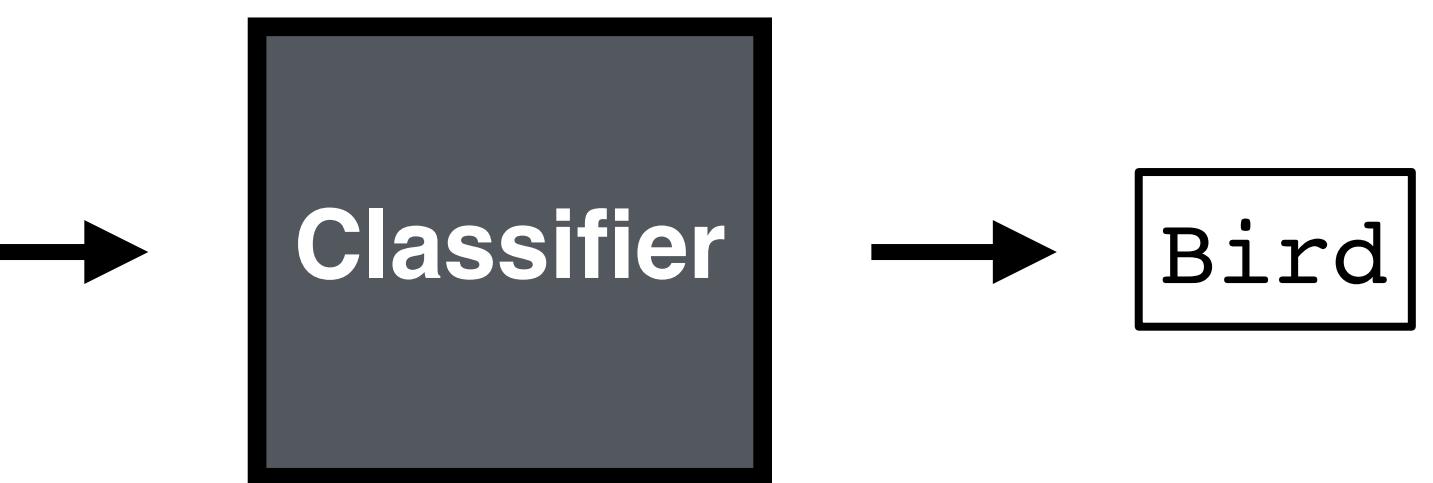
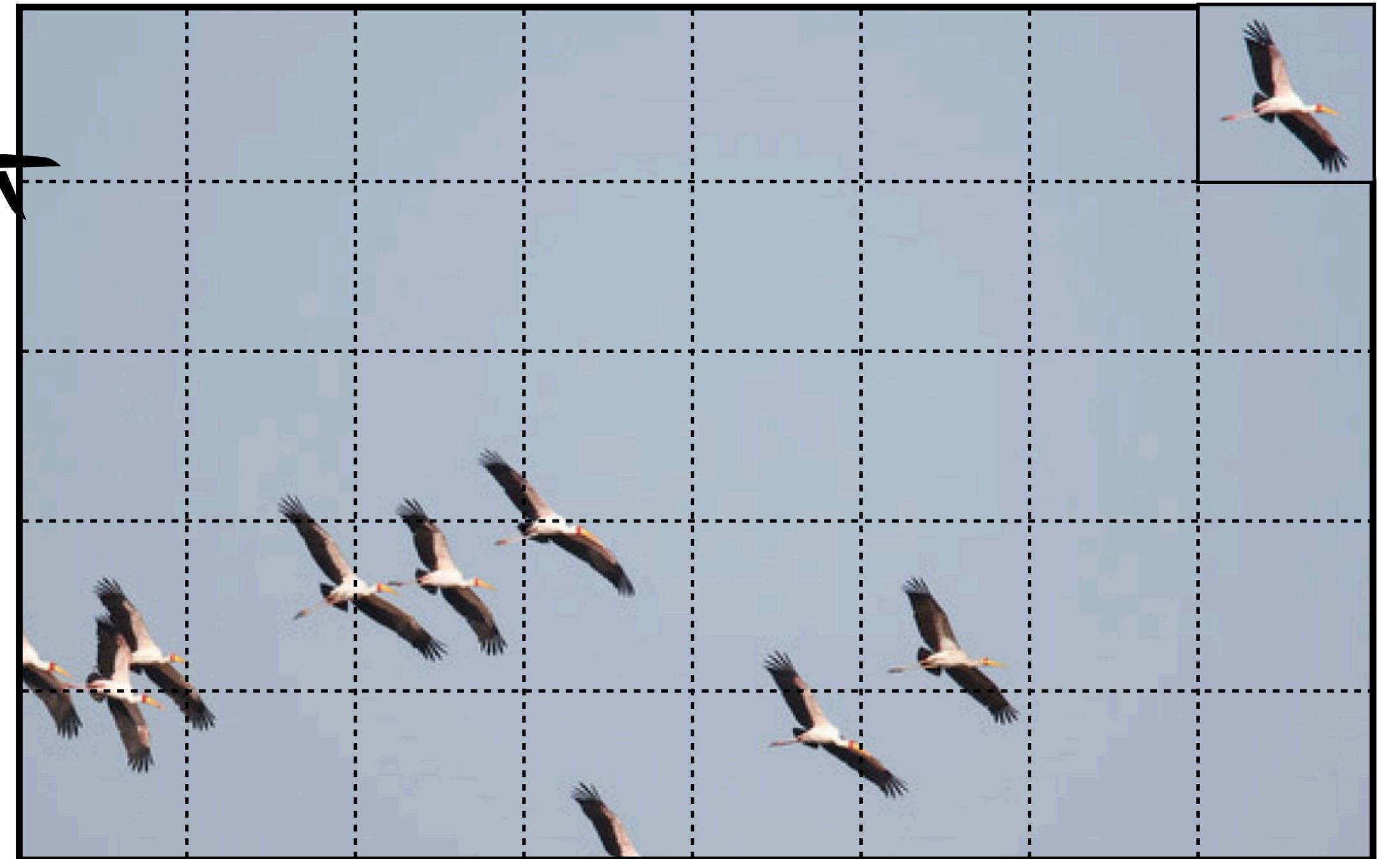
- + Universal
- + Simple (elegant theory)
- + Embarassingly parallel
  
- Weak inductive biases
- Sample inefficient / data hungry
- Dense (fully-connected) linear layers take a lot of compute

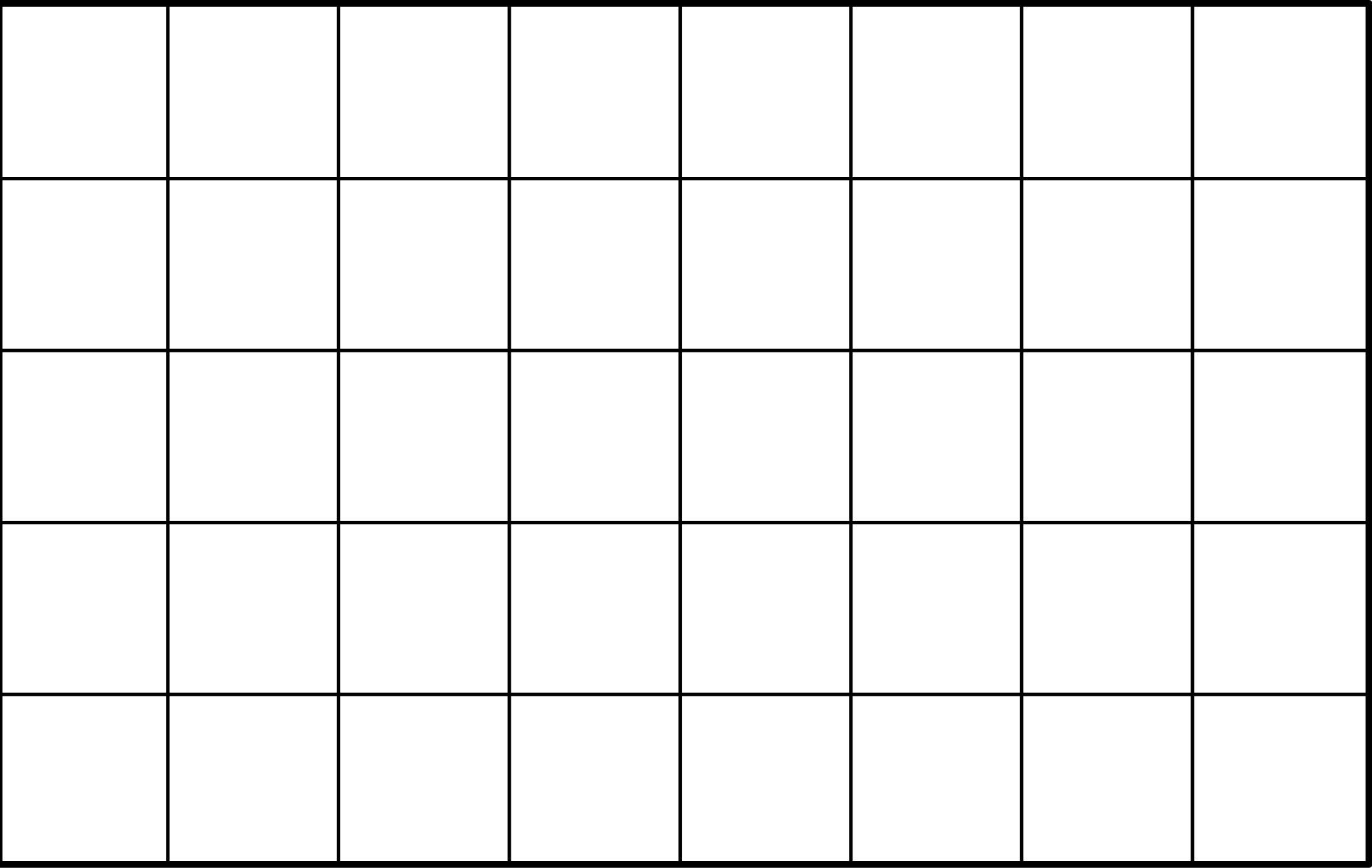
# Convolutional Neural Networks

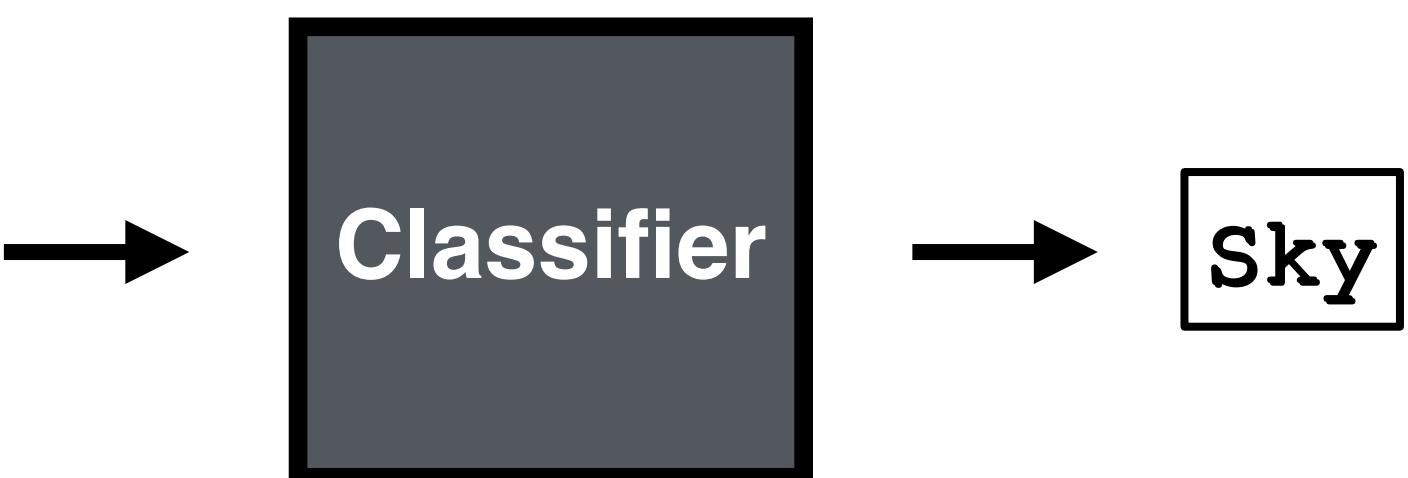
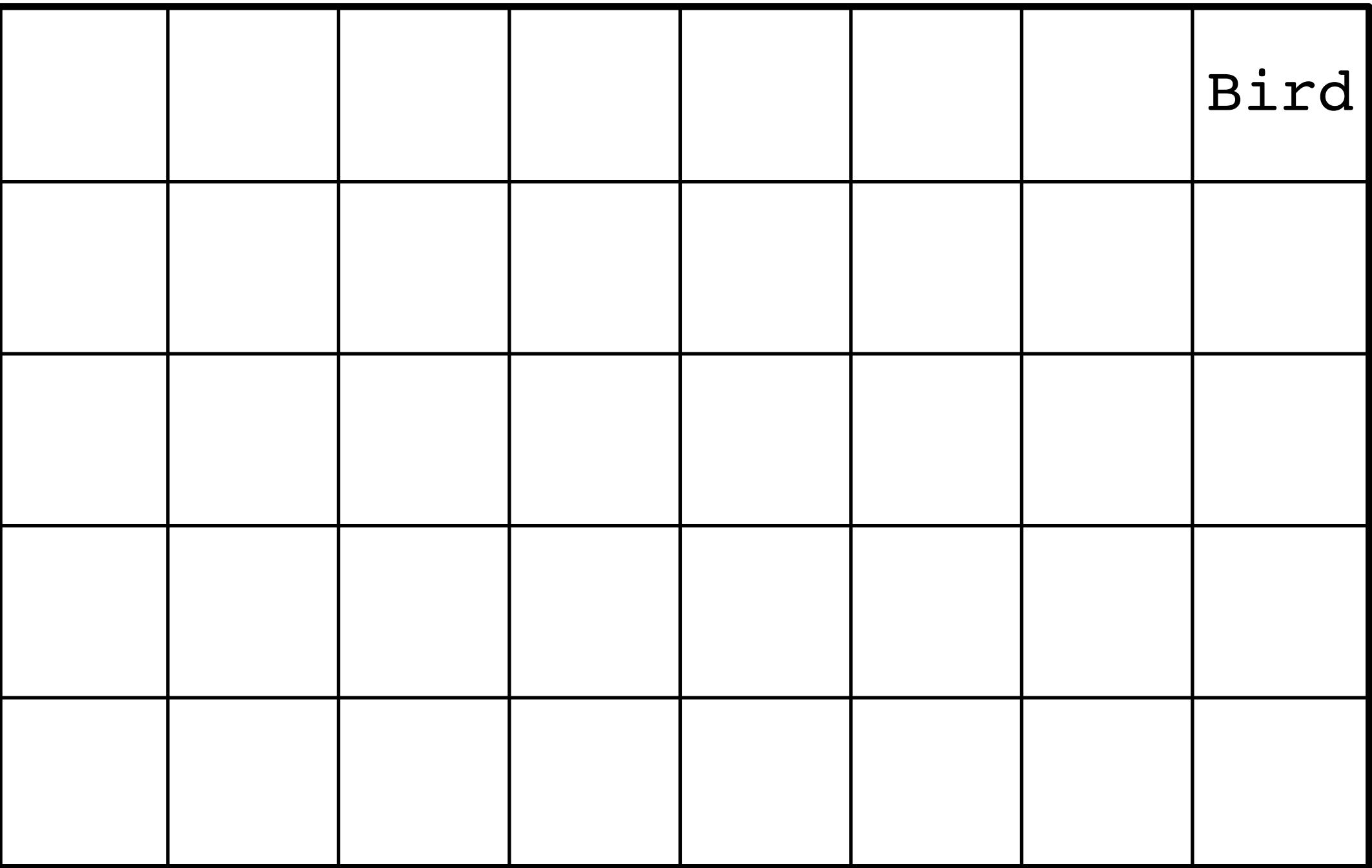


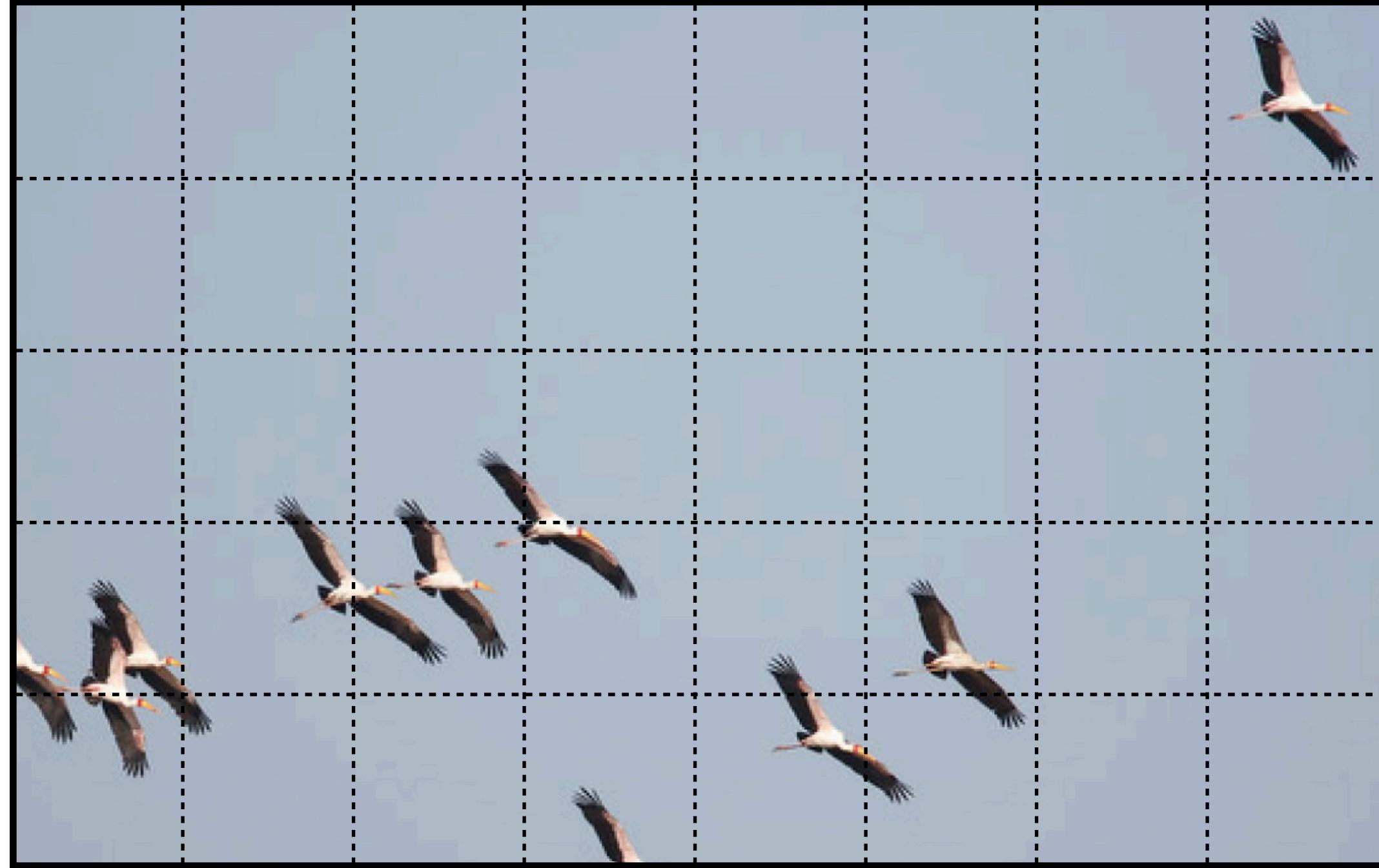
Photo credit: Fredo Durand











Sky	Sky	Sky	Sky	Sky	Sky	Sky	Bird
Sky	Sky	Sky	Sky	Sky	Sky	Sky	Sky
Sky	Sky	Sky	Sky	Sky	Sky	Sky	Sky
Bird	Bird	Bird	Sky	Bird	Sky	Sky	Sky
Sky	Sky	Sky	Bird	Sky	Sky	Sky	Sky



sky	sky	sky	sky	sky	sky	sky	bird
sky	sky	sky	sky	sky	sky	sky	sky
sky	sky	sky	sky	sky	sky	sky	sky
bird	bird	bird	bird	sky	bird	sky	sky
sky	sky	sky	bird	bird	sky	sky	sky

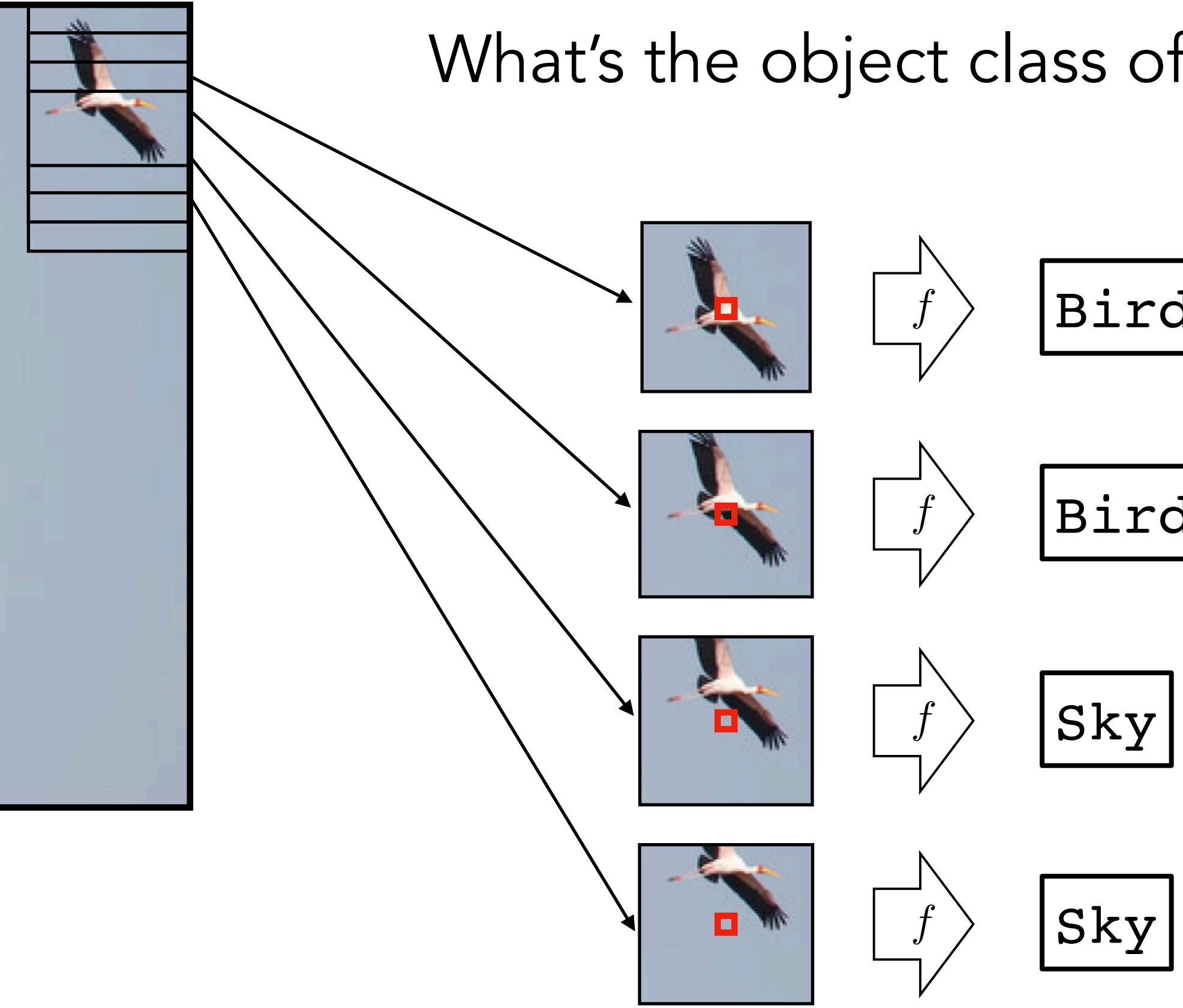
## Problem:

What if objects don't fit neatly into these patches?

How to increase the resolution of the output map?

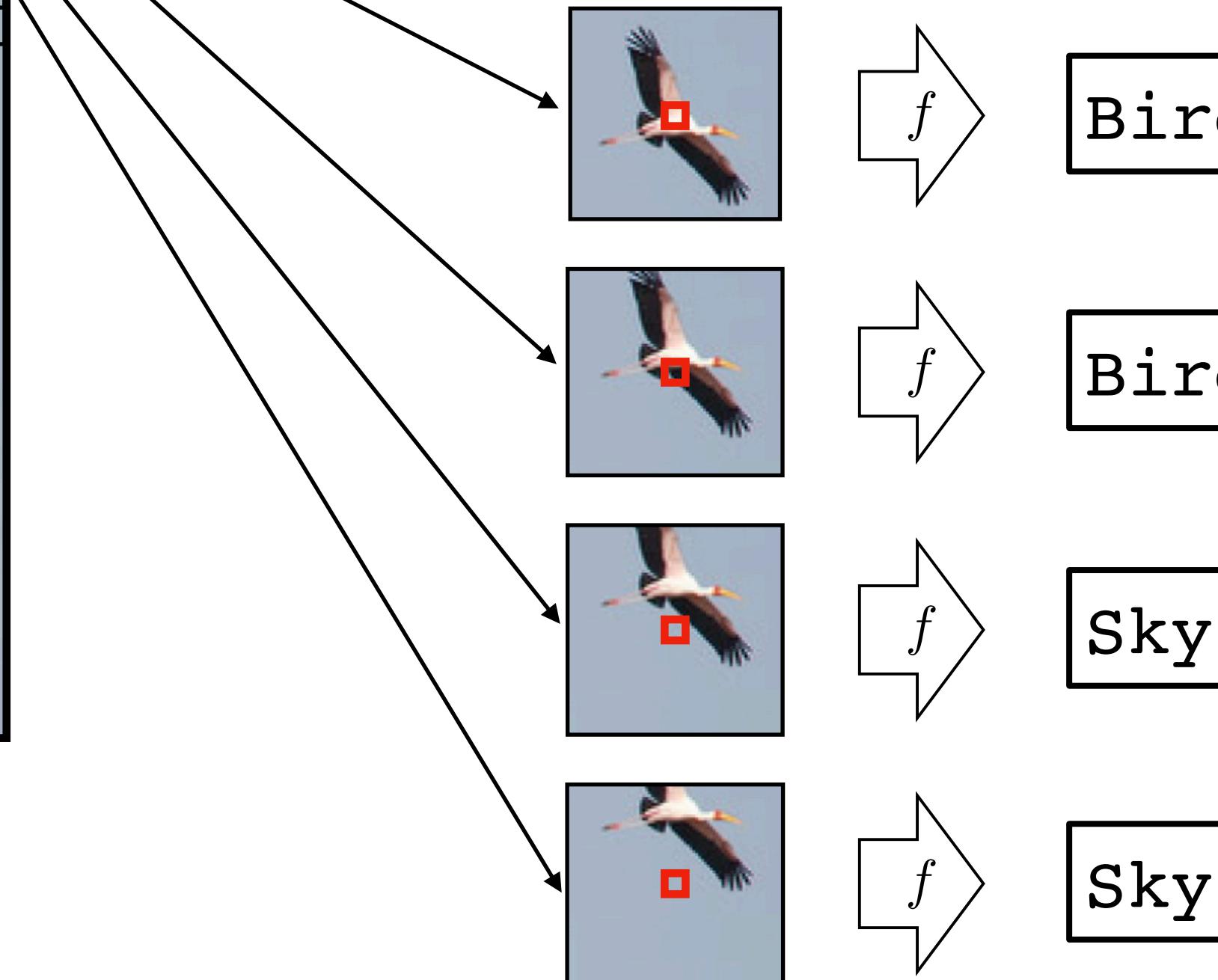
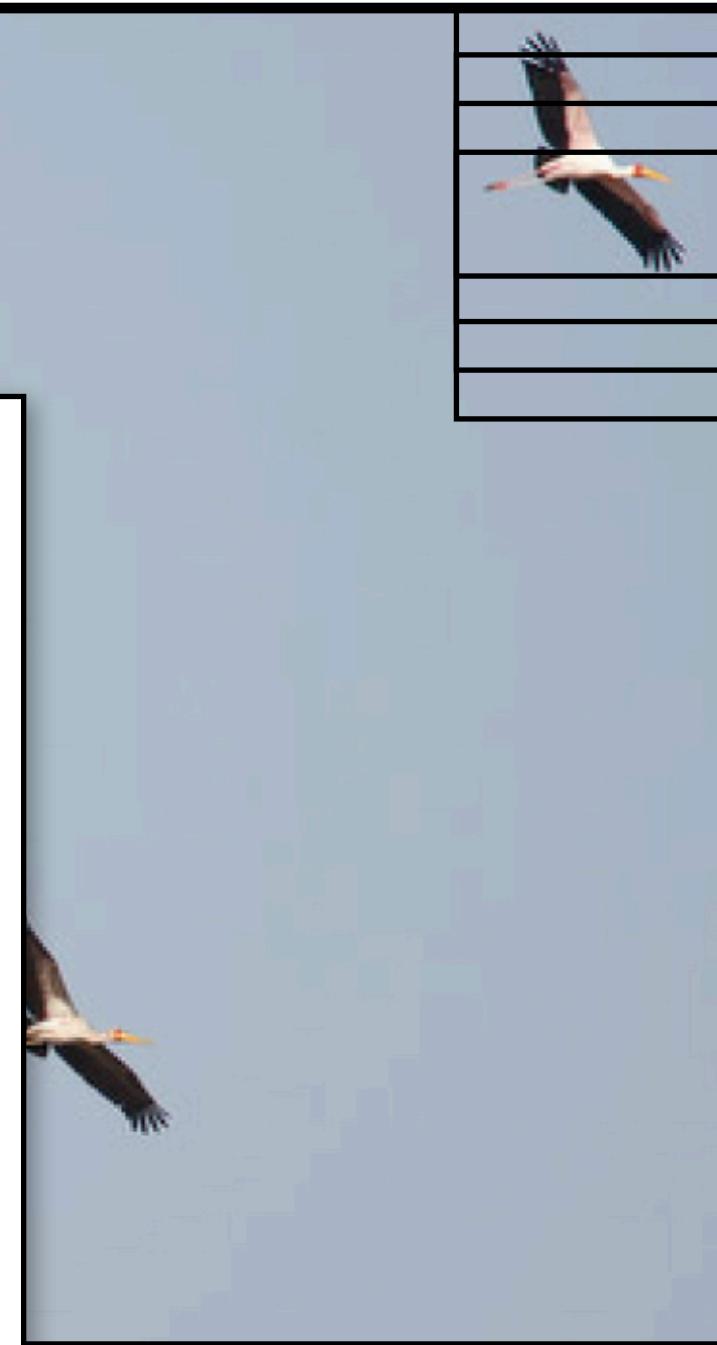


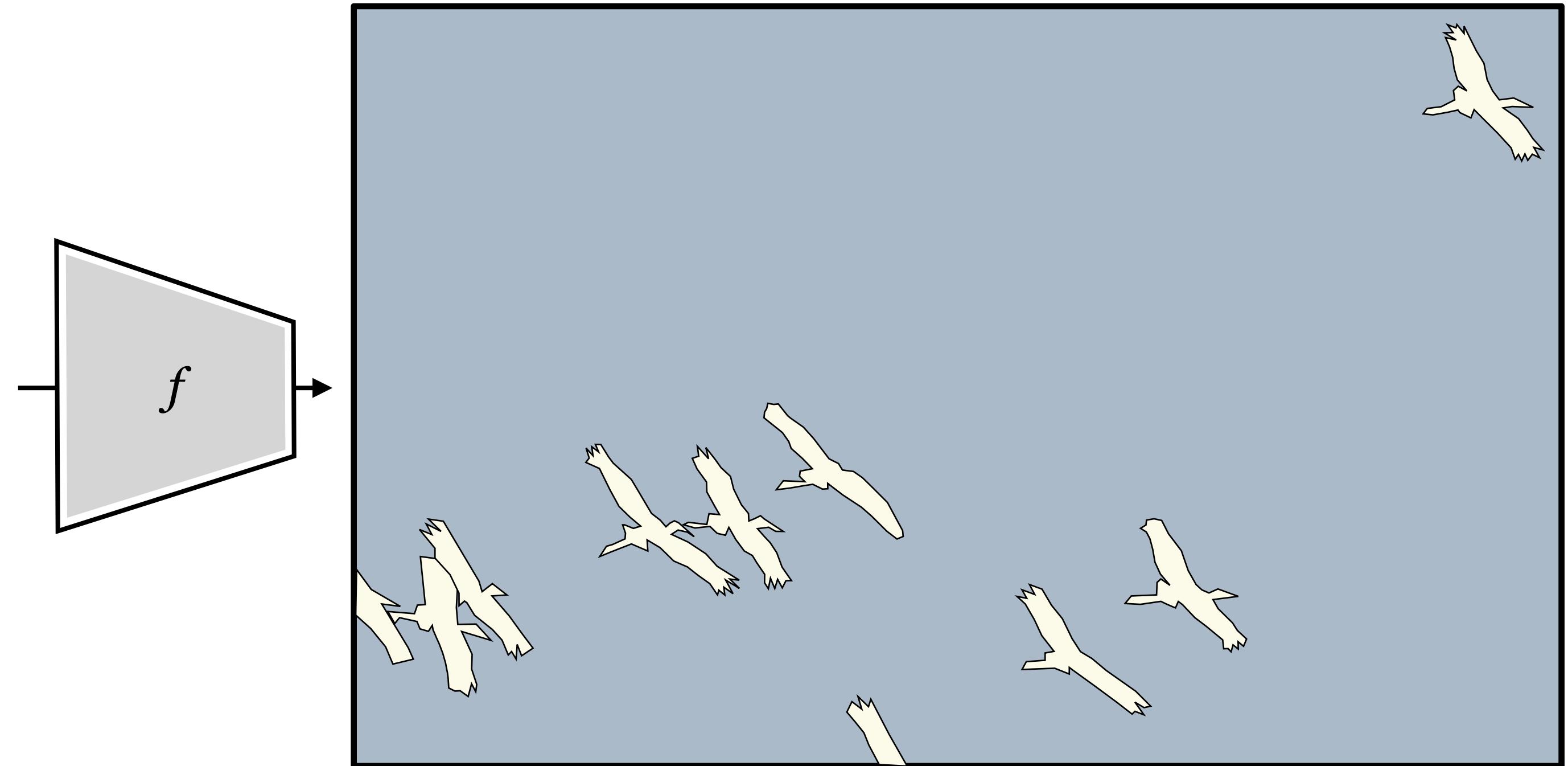
What's the object class of the center pixel?



What's the object class of the center pixel?

<i>Training data</i>	
$x$	$y$
{	Bird
{	Bird
{	Sky
:	



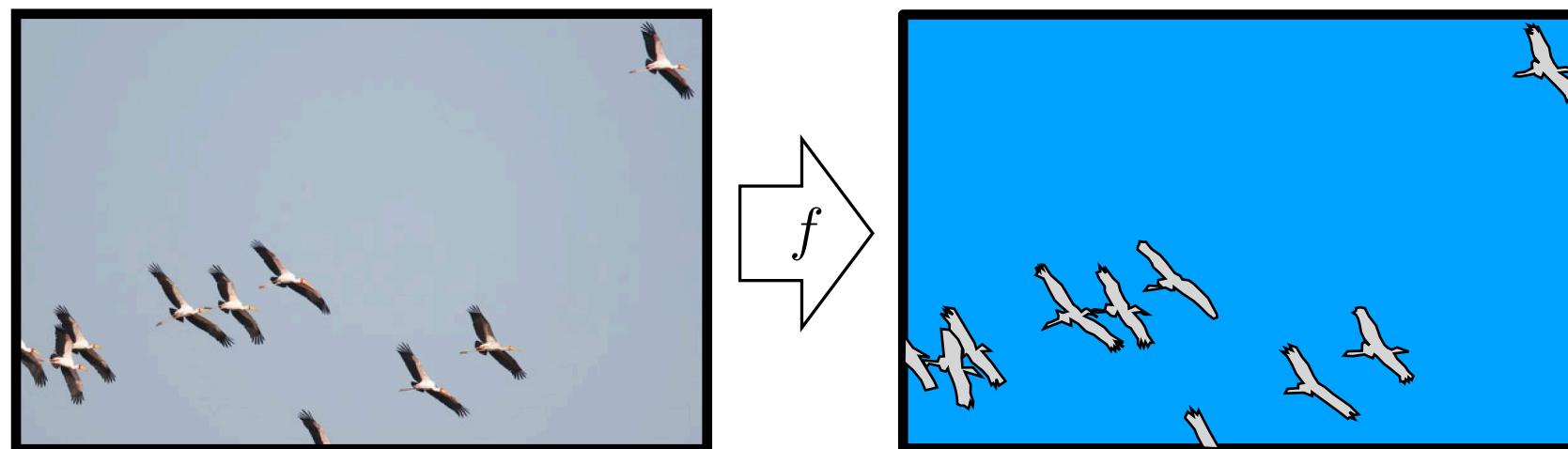
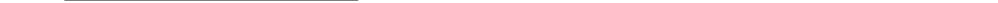
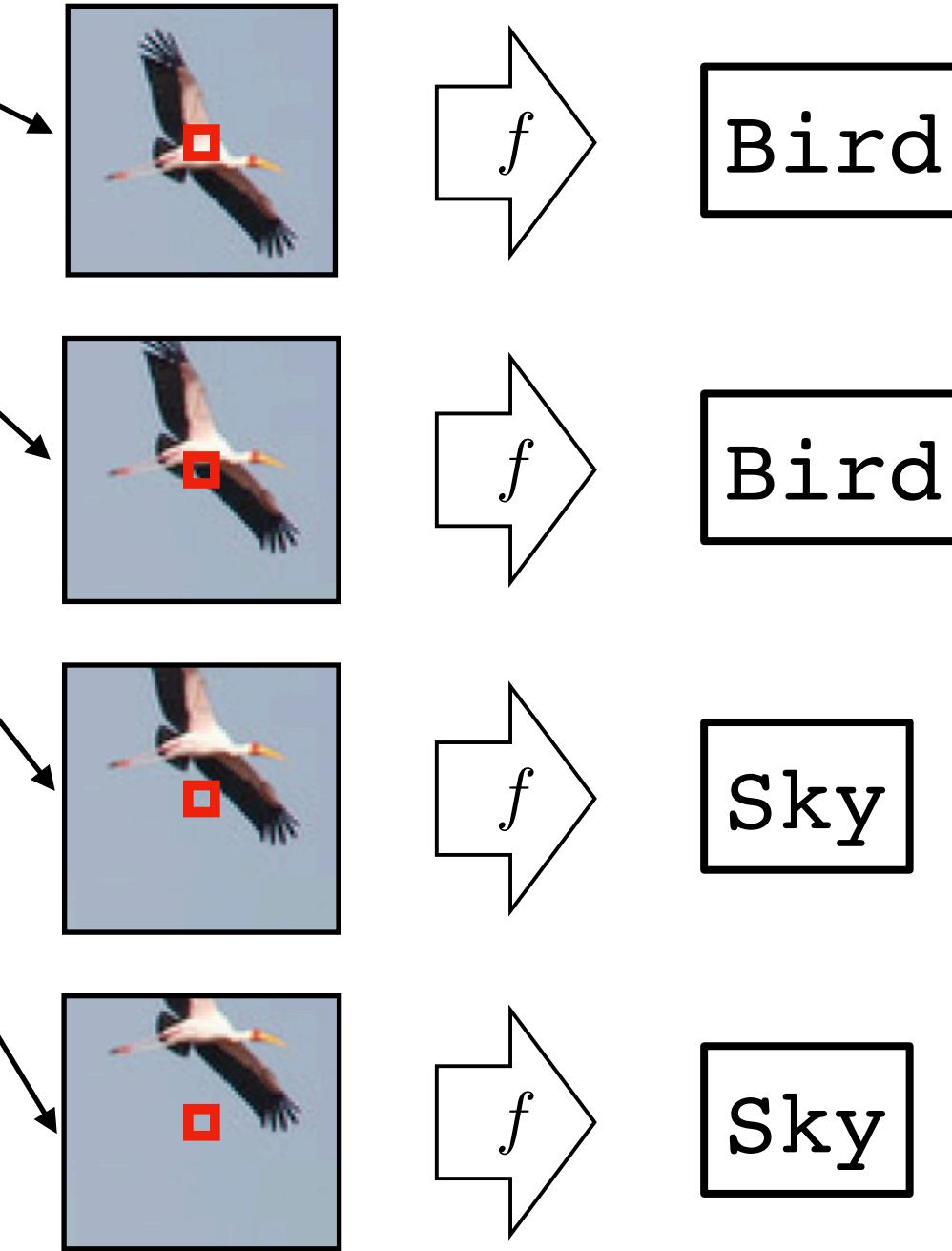


(Colors represent one-hot codes)

This problem is called **semantic segmentation**



What's the object class of the center pixel?

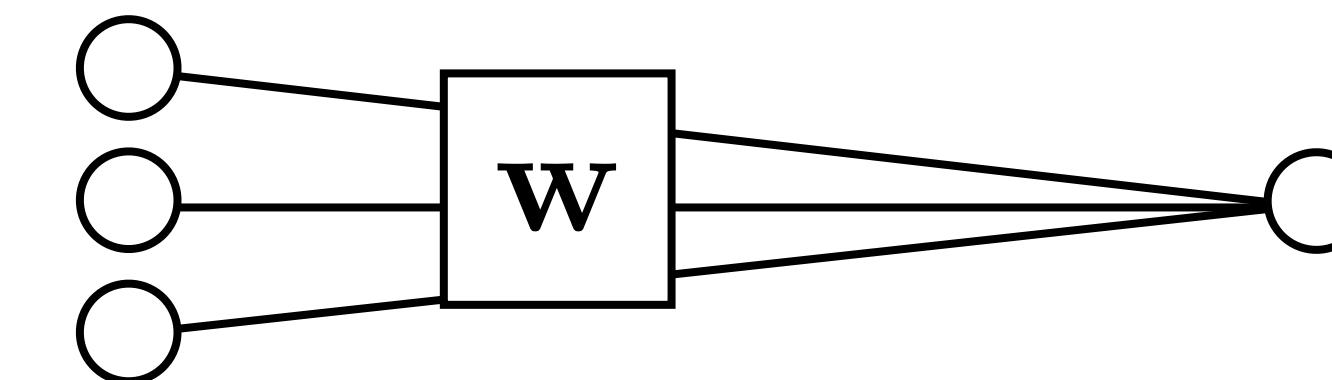


Translation invariance: process each patch in the same way.

An equivariant mapping:

$$f(\text{translate}(x)) = \text{translate}(f(x))$$

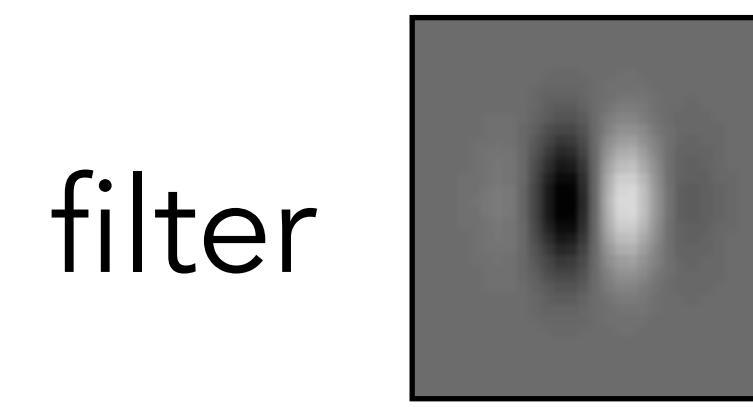
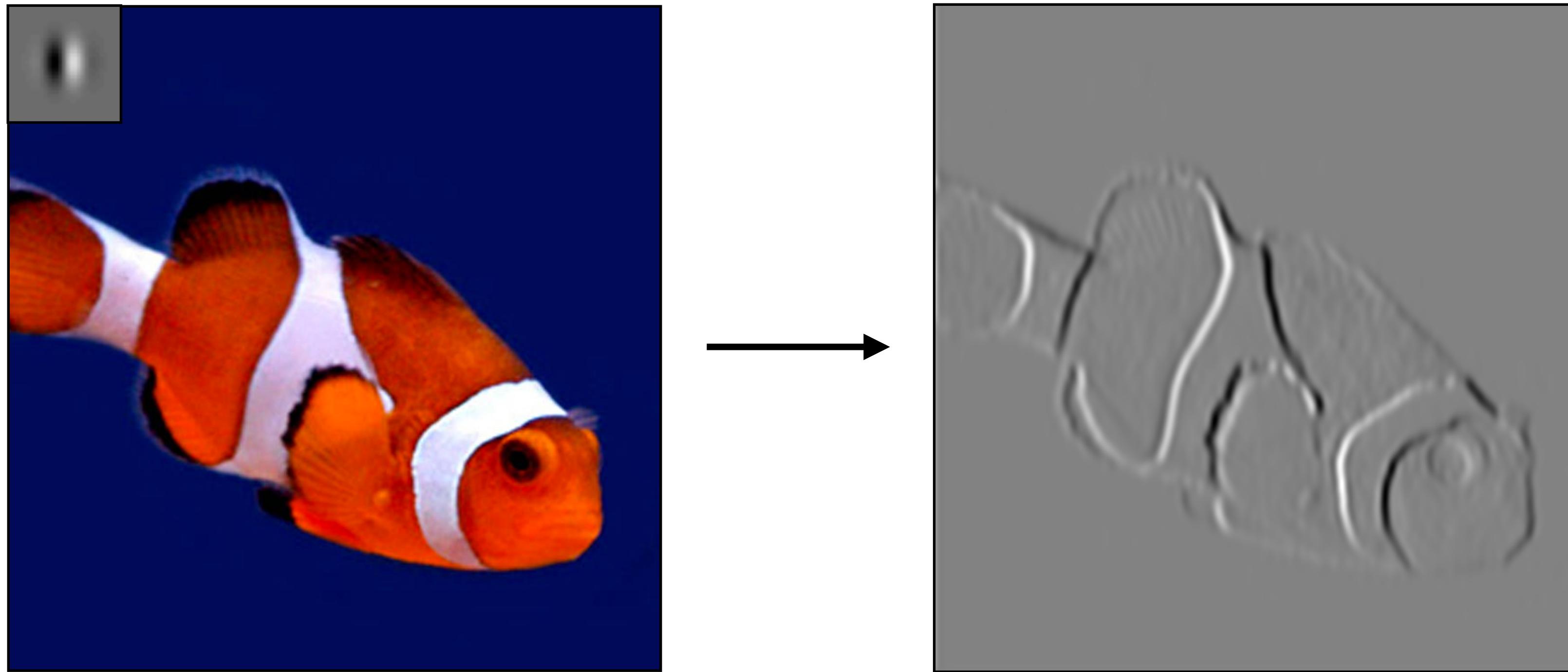
$W$  computes a weighted sum of all pixels in the patch



$W$  is a **convolutional kernel** applied to the full image!

# Convolution

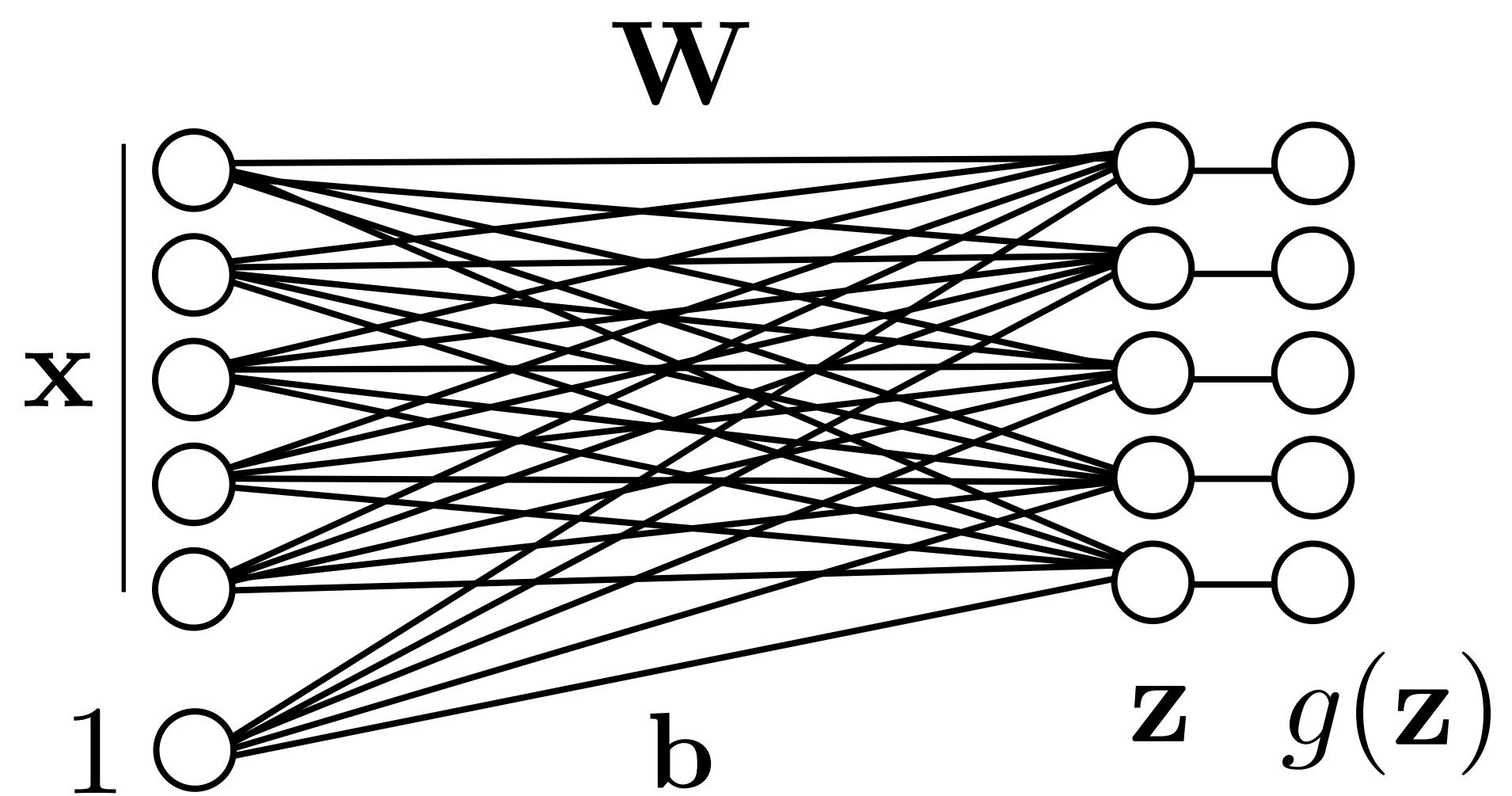
Linear, shift-invariant transformation



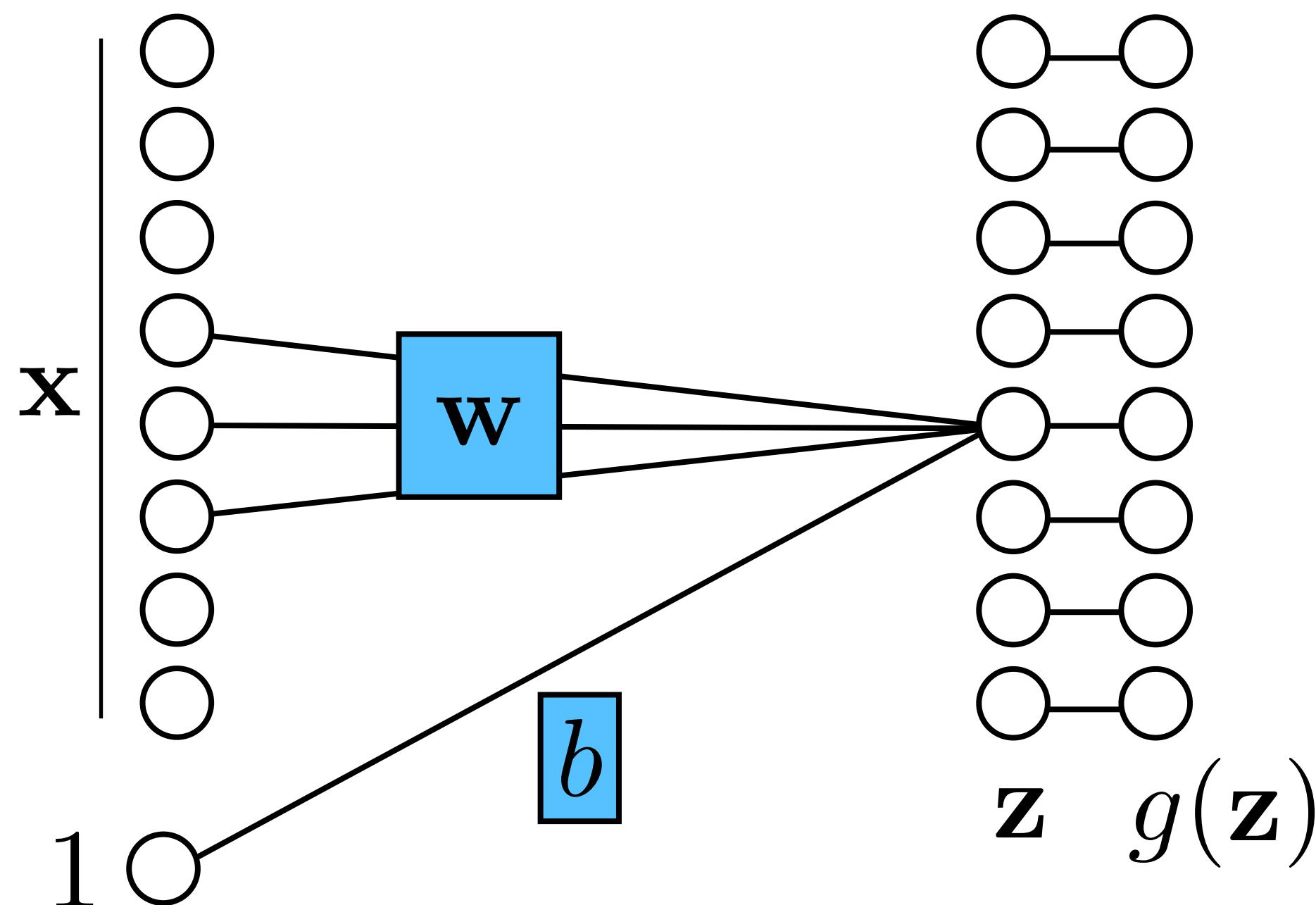
$$x_{\text{out}}[n, m] = b + \sum_{k_1, k_2=-K}^K w[k_1, k_2] x_{\text{in}}[n + k_1, m + k_2]$$

# Fully-connected network

## Fully-connected (fc) layer



# Locally connected network

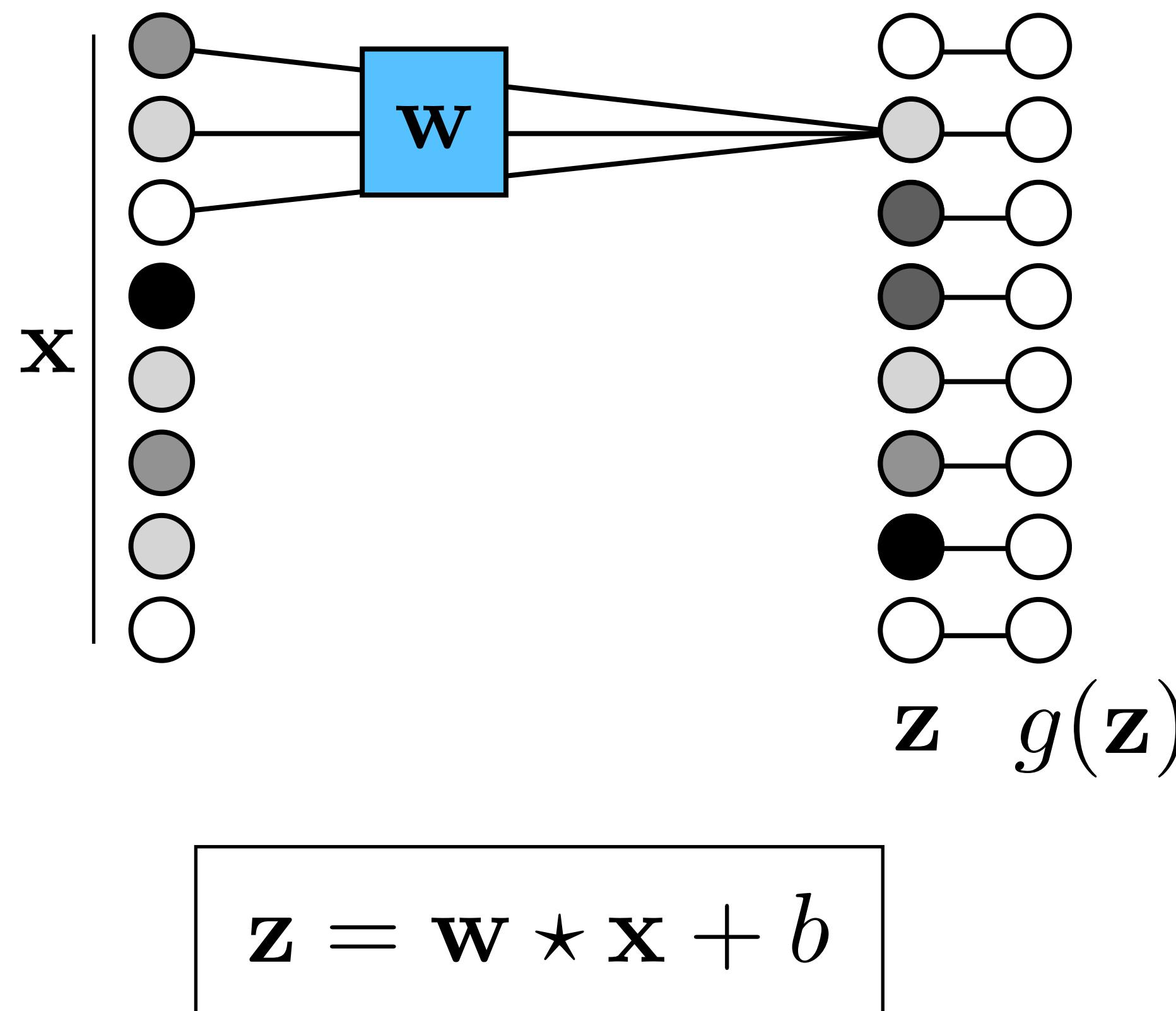


Often, we can assume output is a **local** function of input.

If we use the same weights (**weight sharing**) to compute each local function, we get a convolutional neural network.

# Convolutional neural network

## Conv layer

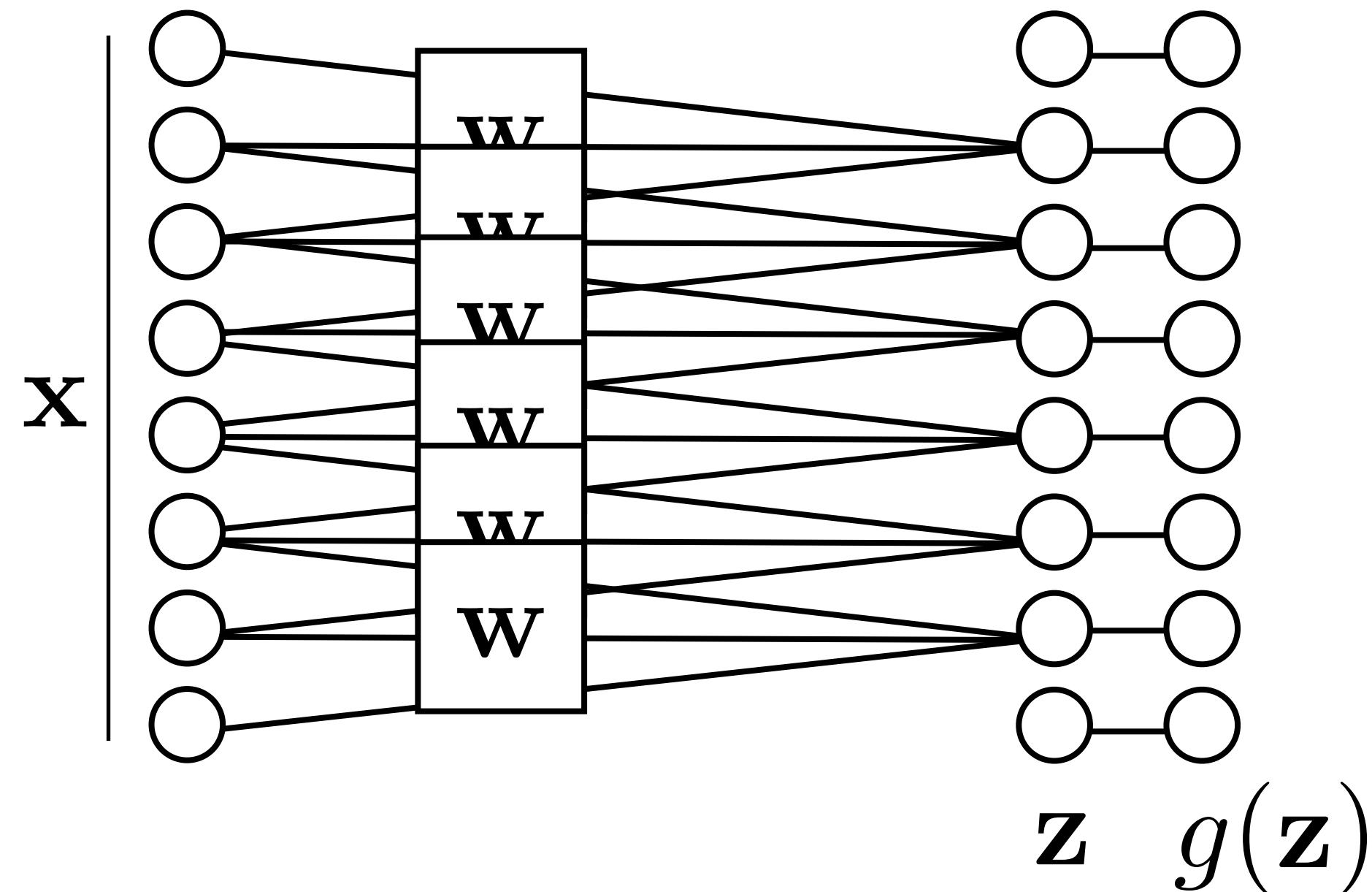


Often, we assume output is a **local** function of input.

If we use the same weights (**weight sharing**) to compute each local function, we get a convolutional neural network.

# Weight sharing

## Conv layer



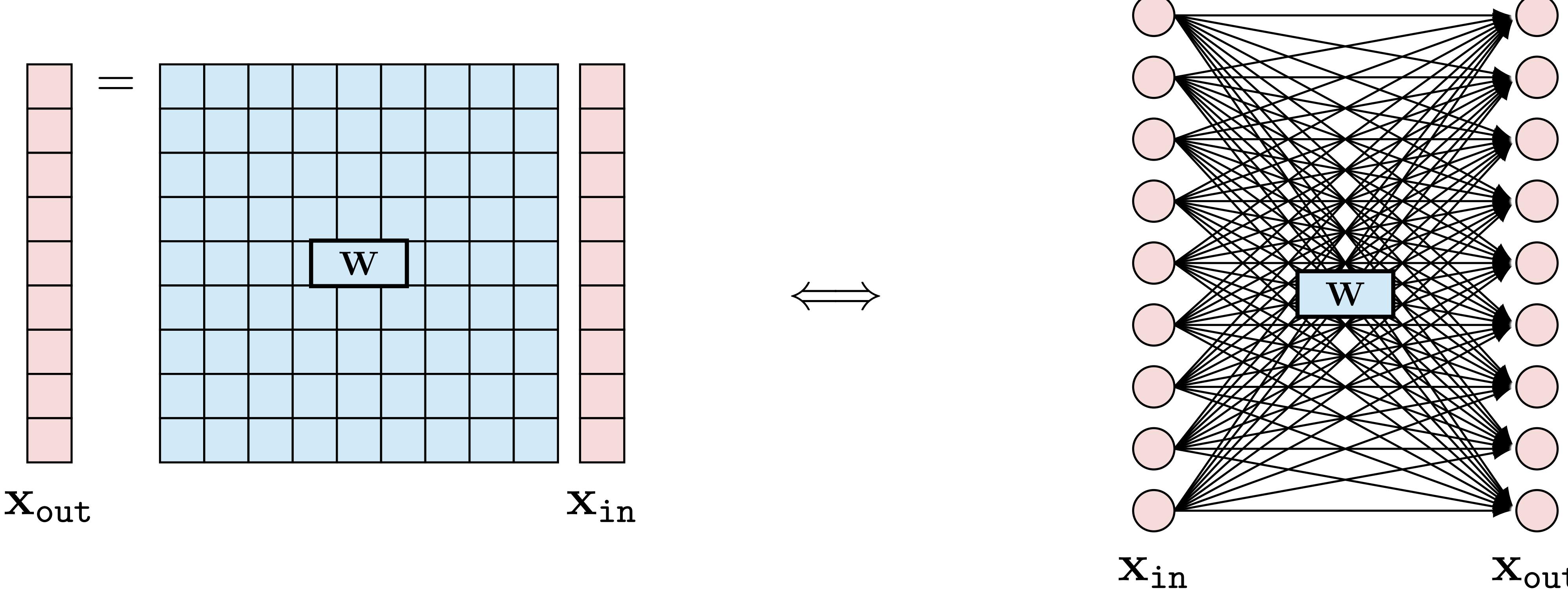
$$\mathbf{z} = \mathbf{w} \star \mathbf{x} + b$$

Often, we assume output is a **local** function of input.

If we use the same weights (**weight sharing**) to compute each local function, we get a convolutional neural network.

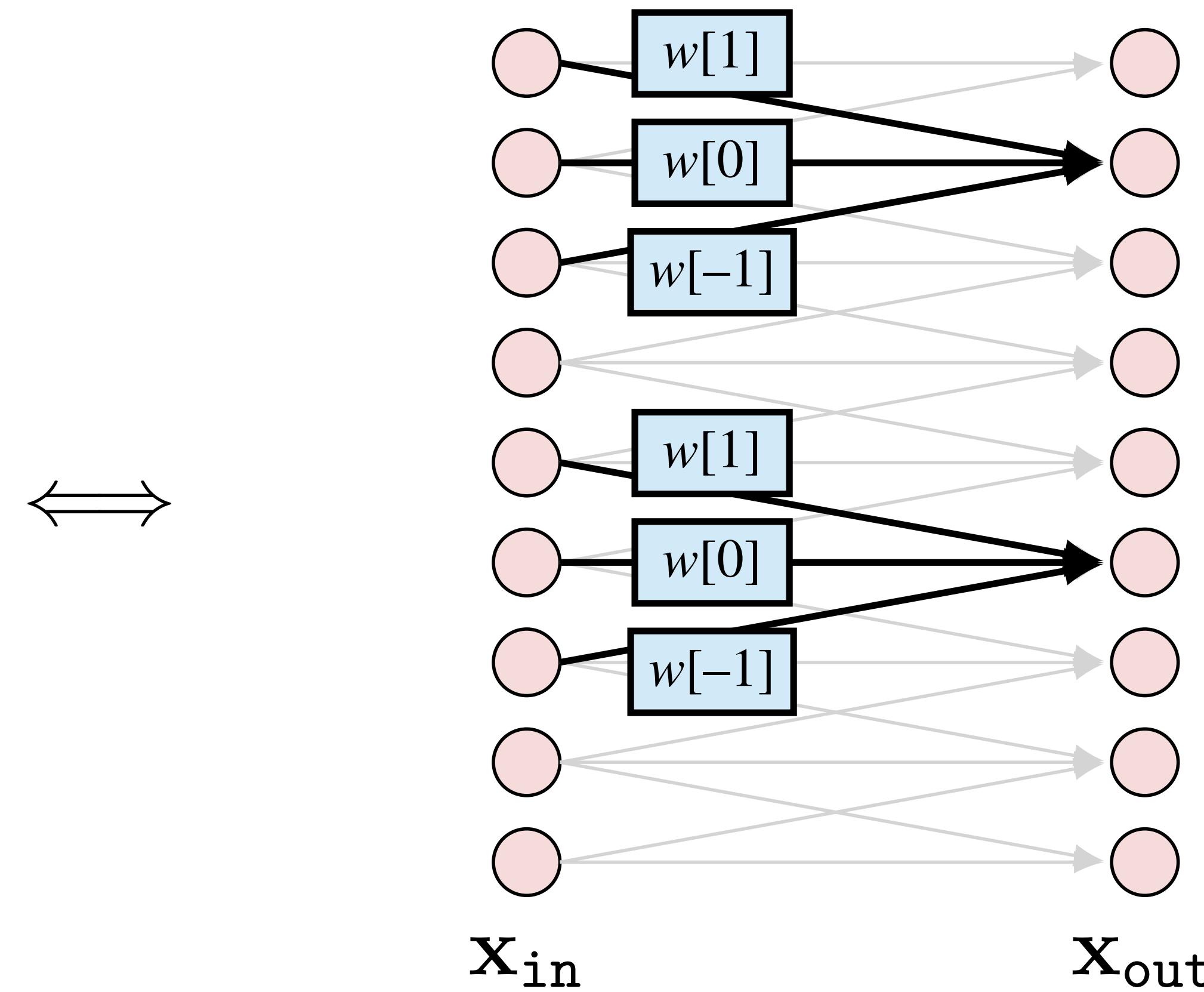
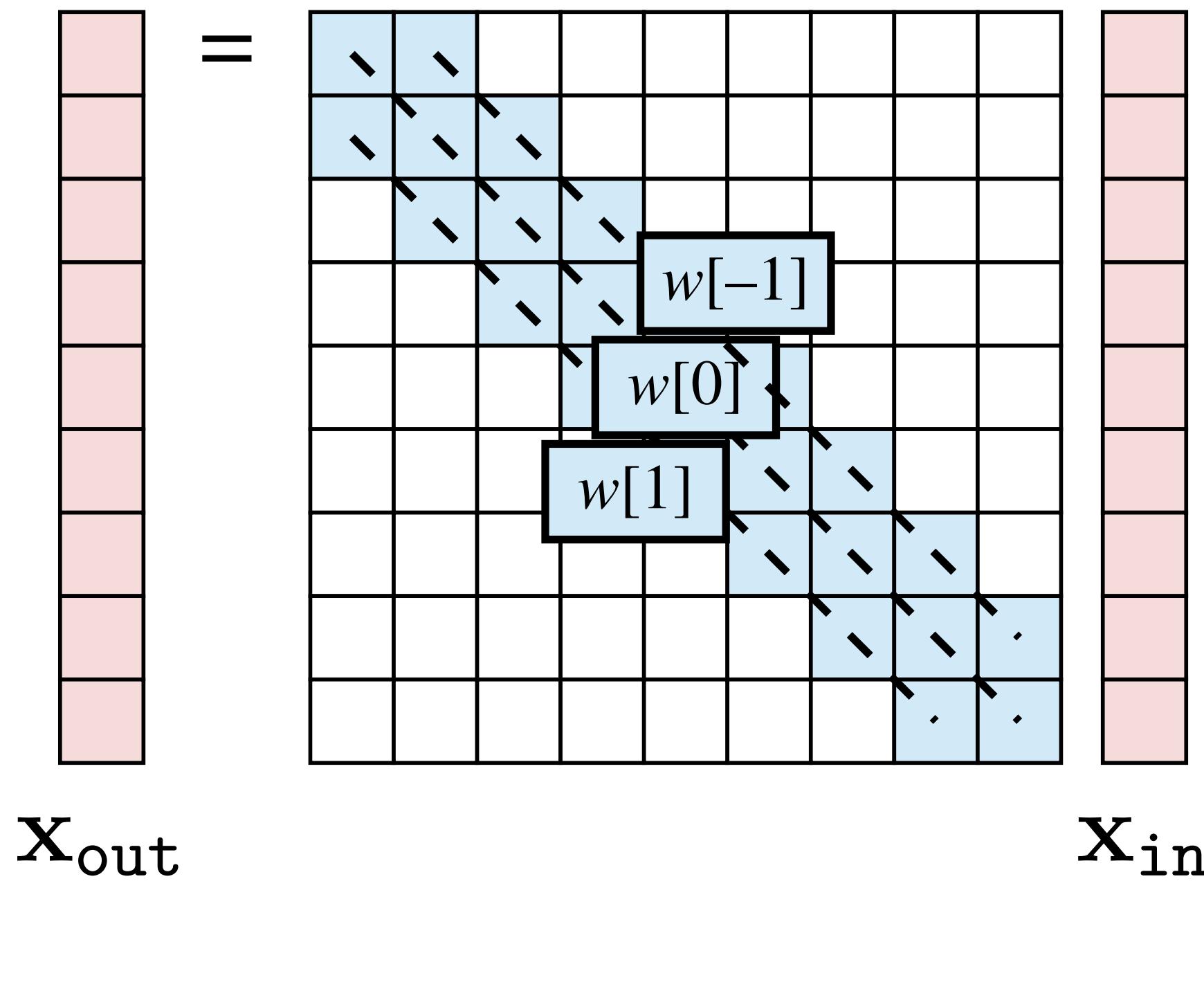
# (Fully-connected) linear layer

$$\mathbf{x}_{\text{out}} = \mathbf{W}\mathbf{x}_{\text{in}} + \mathbf{b}$$



# Convolutional layer

$$\mathbf{x}_{\text{out}} = \mathbf{w} \star \mathbf{x}_{\text{in}} + b$$



## Toeplitz matrix

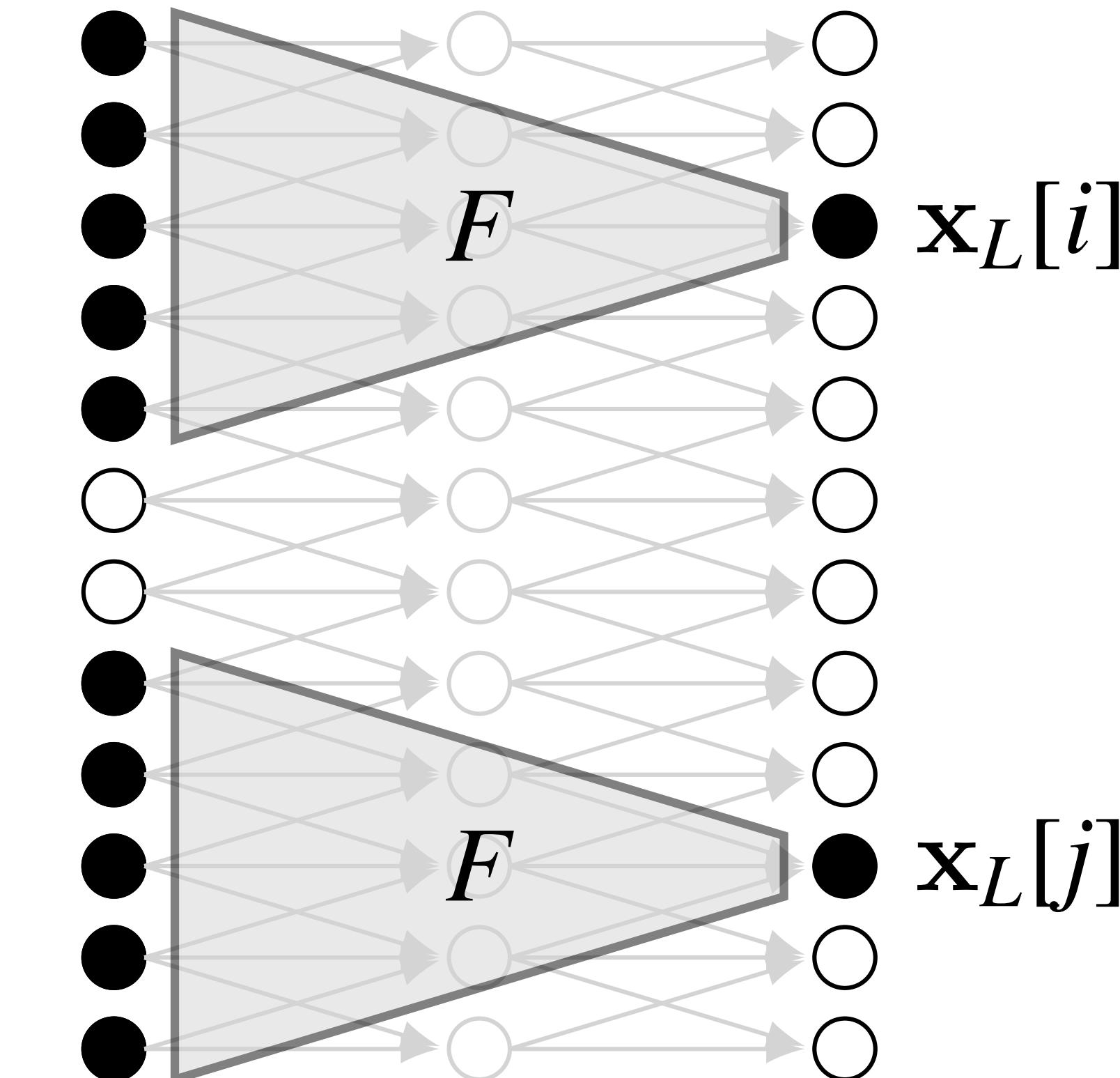
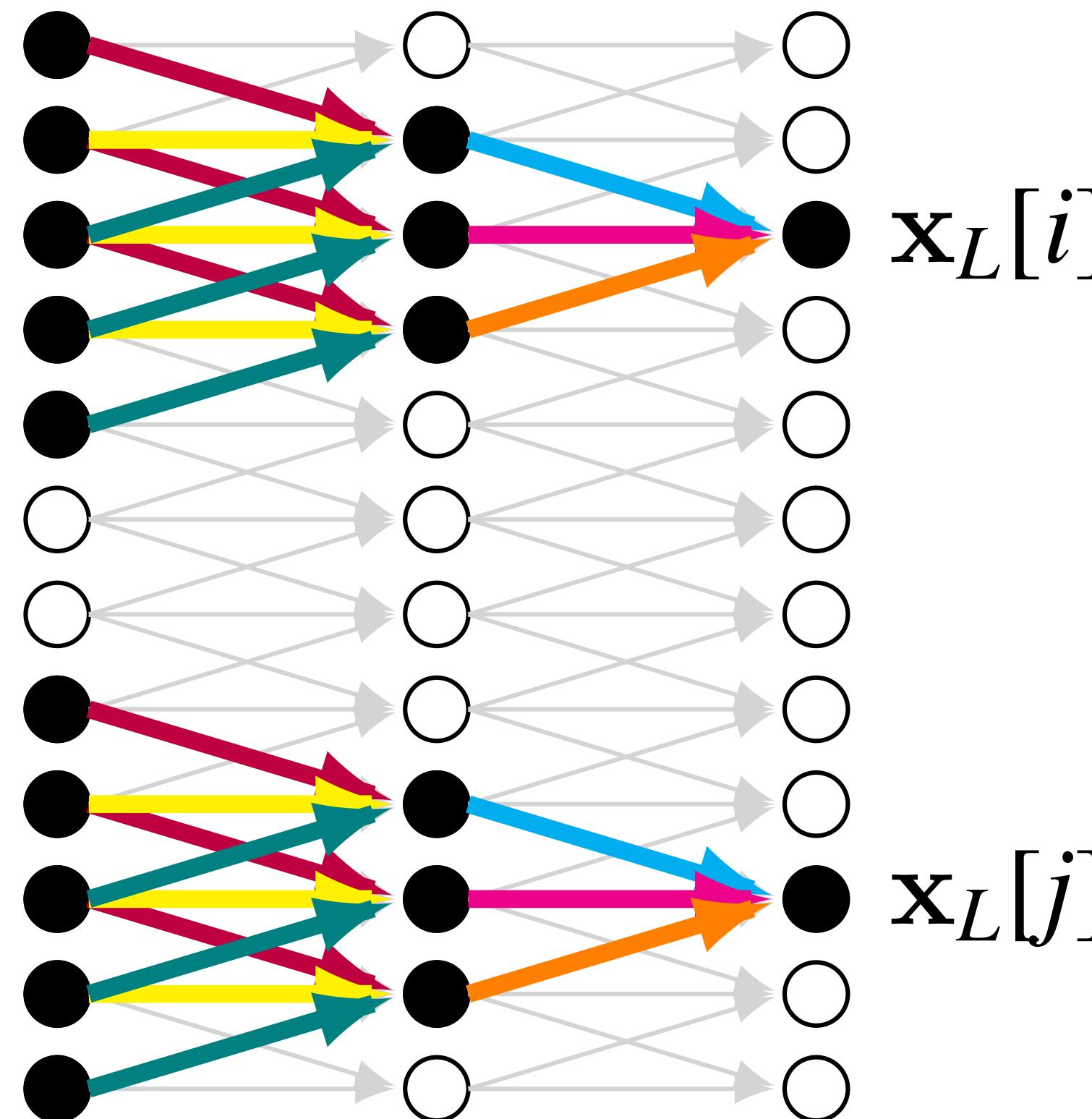
$$\begin{pmatrix} a & b & c & d & e \\ f & a & b & c & d \\ g & f & a & b & c \\ h & g & f & a & b \\ i & h & g & f & a \end{pmatrix}$$

$$\begin{array}{c|c|c} & & \\ \textbf{y} & = & \textbf{A} & \textbf{x} \\ & & \end{array}$$


e.g., pixel image

- Constrained linear layer
- Fewer parameters —> easier to learn, less overfitting

# What happens when you stack convolutional layers?



The whole CNN acts like a (nonlinear) convolutional filter!

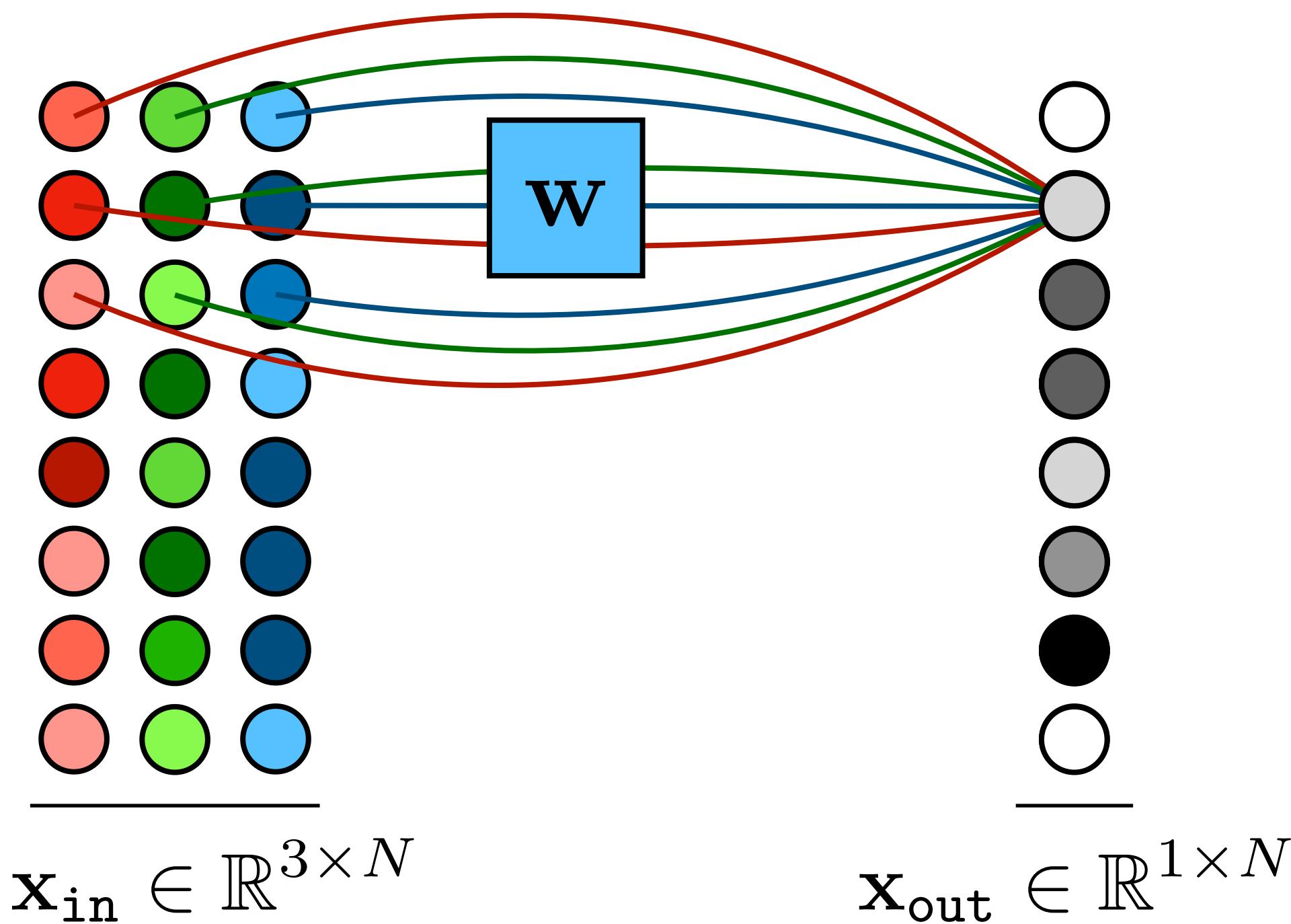
# What if we have color?

(aka multiple input channels?)

# Multichannel inputs

## Conv layer

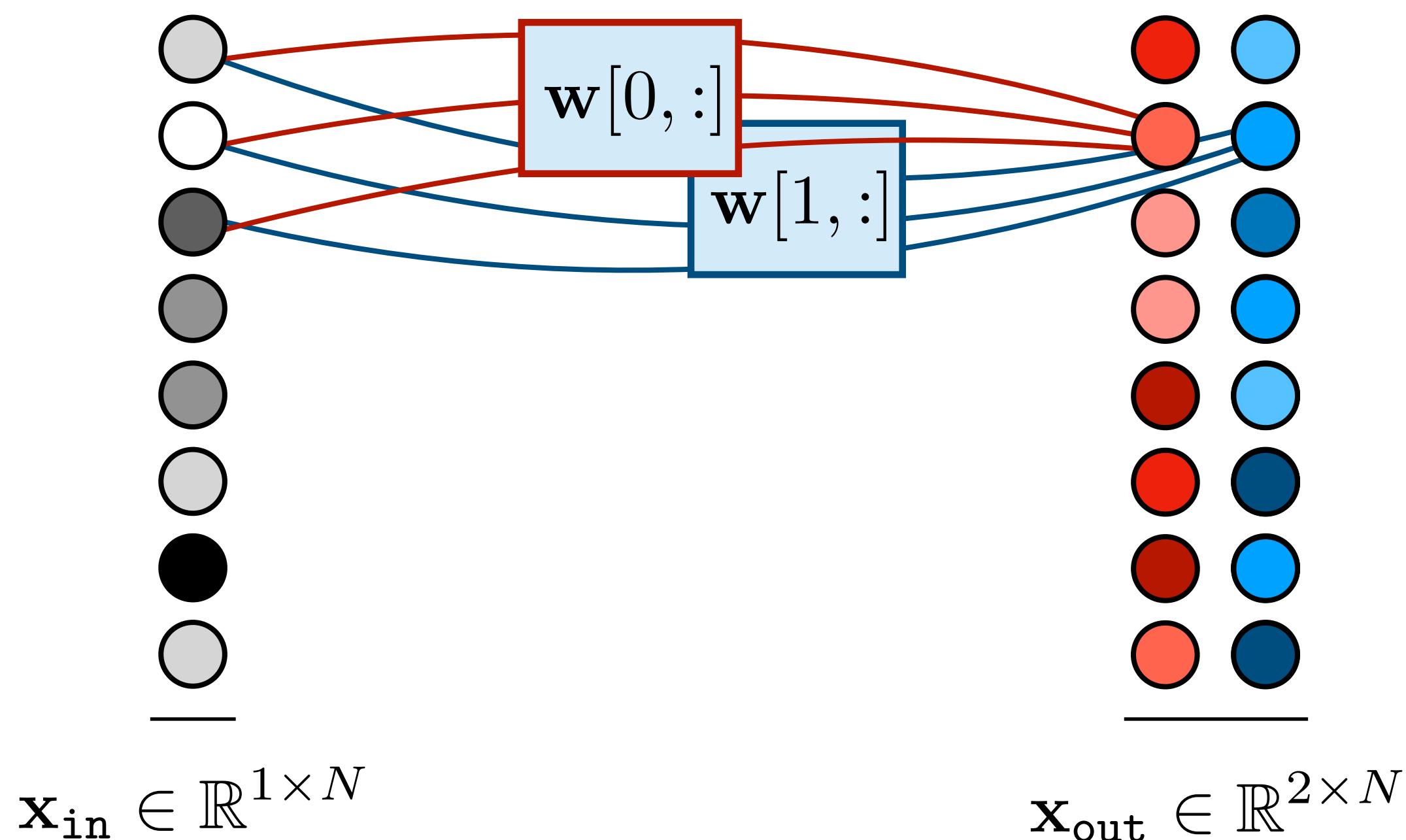
$$\mathbf{w} \in \mathbb{R}^{3 \times 3}$$



$$\mathbf{x}_{\text{out}} = \sum_c \mathbf{w}[c, :] \star \mathbf{x}_{\text{in}}[c, :] + b[c]$$

# Multichannel outputs

## Conv layer



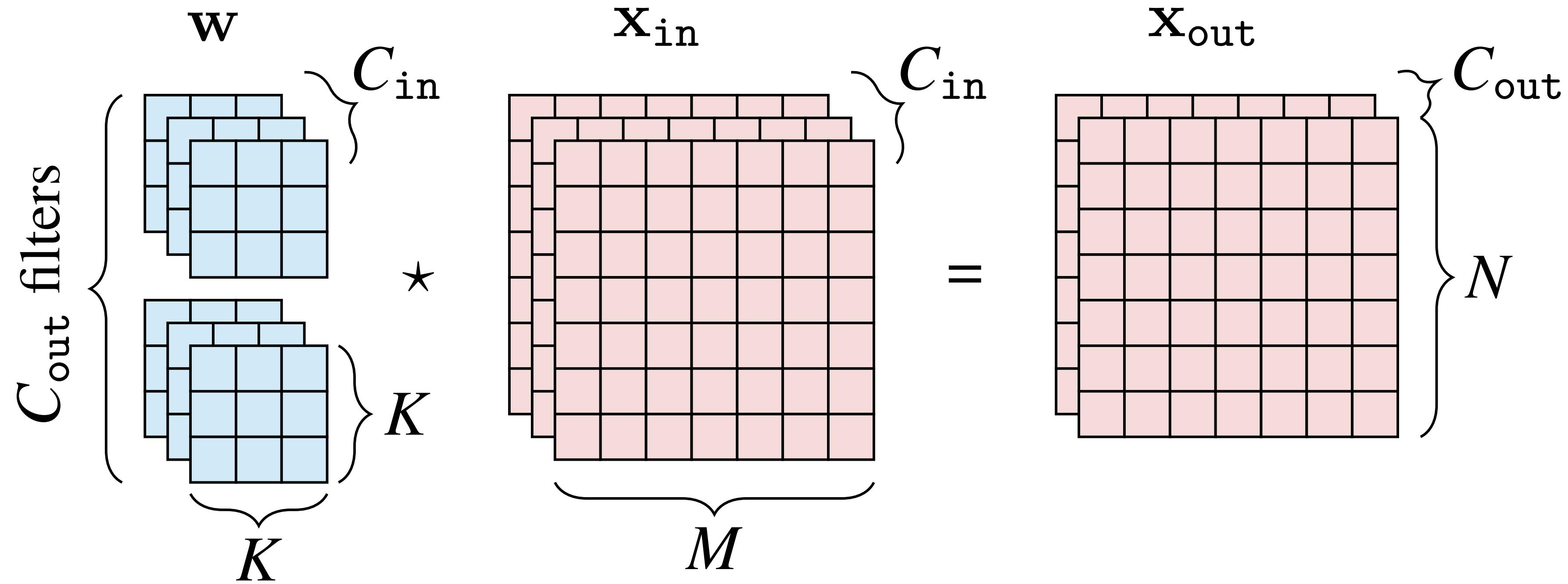
Filter bank of  $C$  filters

$$\mathbf{x}_{\text{out}}[0, :] = \mathbf{w}[0, :] \star \mathbf{x}_{\text{in}} + b[0]$$

⋮

$$\mathbf{x}_{\text{out}}[C, :] = \mathbf{w}[C - 1, :] \star \mathbf{x}_{\text{in}} + b[C - 1]$$

# General Convolutional Layer Form: Multi-Input, Multi-Output



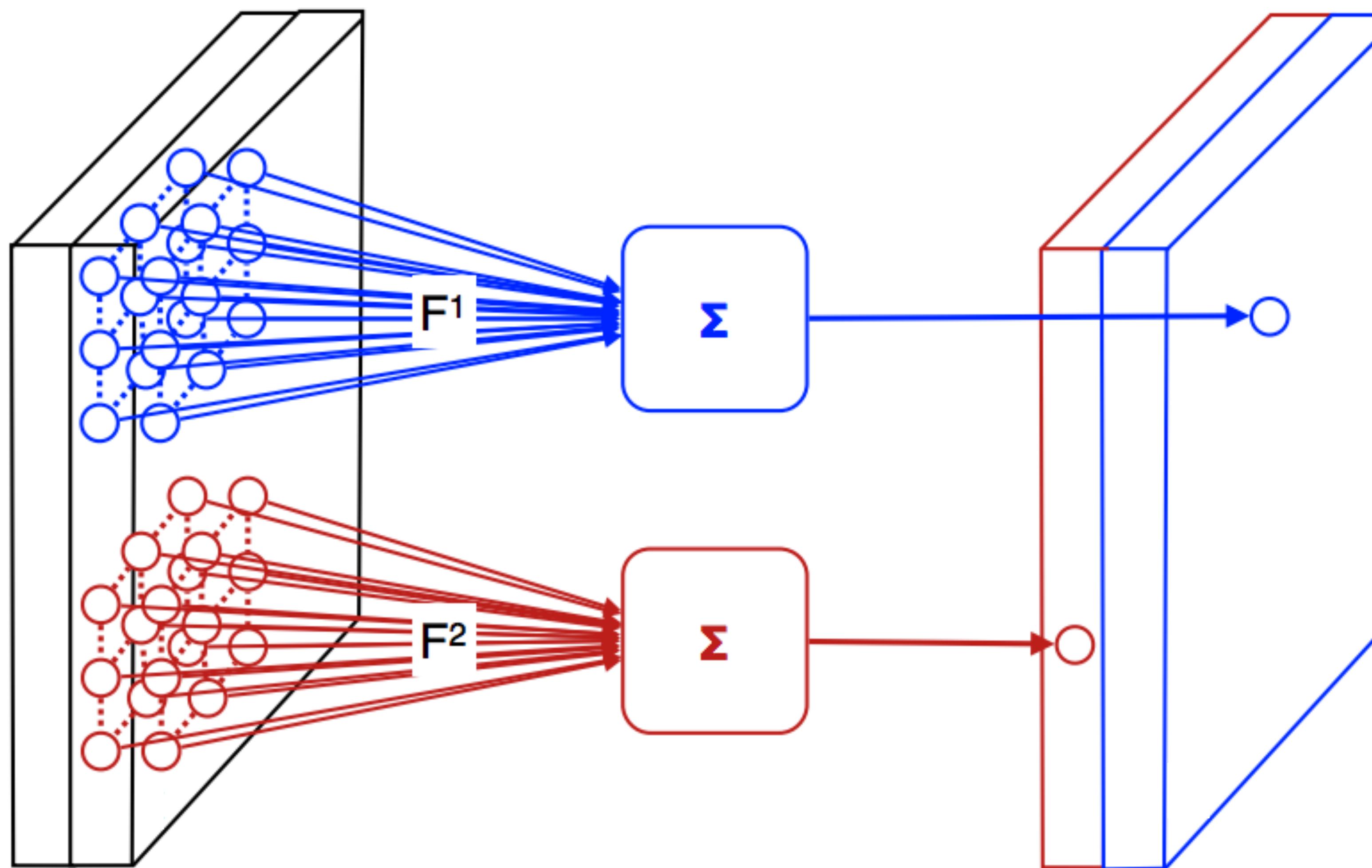
$$\mathbf{x}_{\text{out}}[c_2, :, :] = \sum_{c_1=1}^{C_{\text{in}}} \mathbf{w}[c_1, c_2, :, :] \star \mathbf{x}_{\text{in}}[c_1, :, :] + b[c_2]$$

Input channel                      Output channel

Input features

A bank of 2 filters

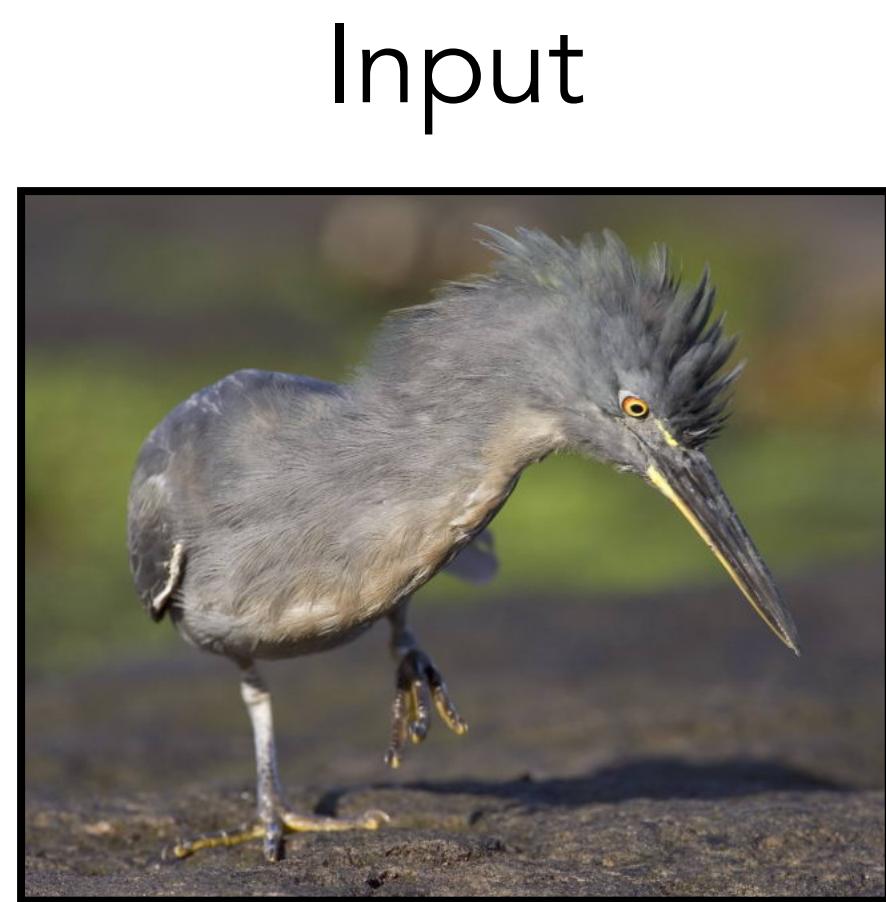
2-dimensional output  
**feature maps**



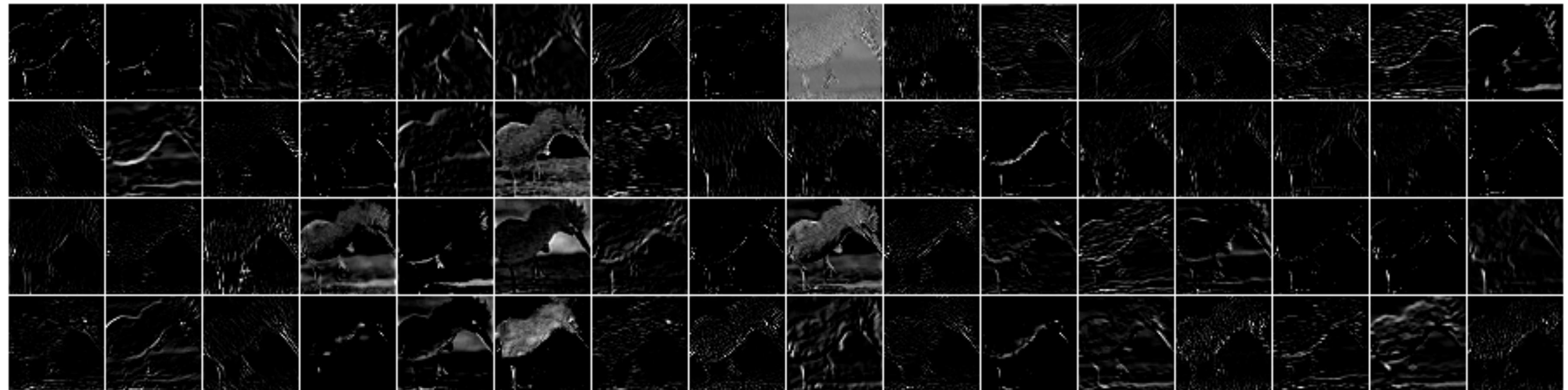
$$\mathbf{x}_{\text{in}} \in \mathbb{R}^{C_{\text{in}} \times H \times W} \rightarrow \mathbf{x}_{\text{out}} \in \mathbb{R}^{C_{\text{out}} \times H \times W}$$

[Figure modified from Andrea Vedaldi]

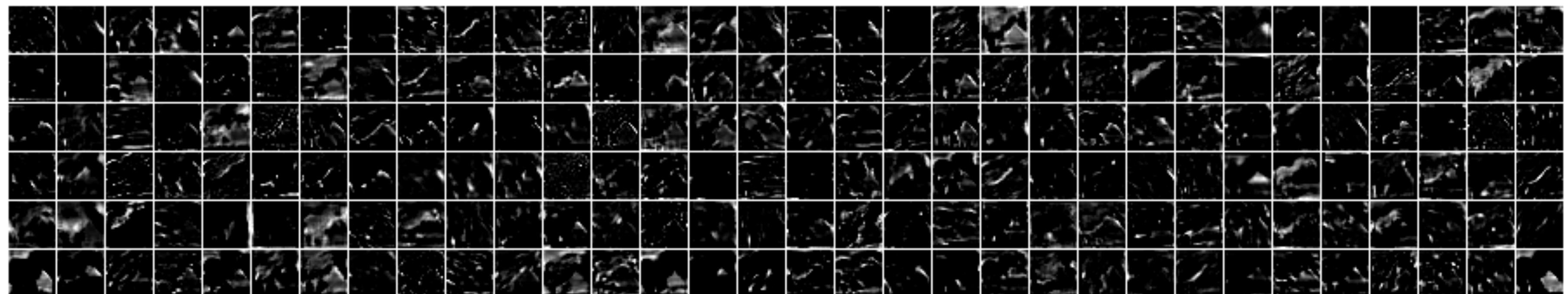
# Feature maps



conv1 (after first conv layer)



conv2 (after second conv layer)



- Each layer can be thought of as a set of C **feature maps** aka **channels**
- Each feature map is an NxM image

# Filter sizes

When mapping from

$$\mathbf{x}_l \in \mathbb{R}^{C_l \times N \times M} \rightarrow \mathbf{x}_{(l+1)} \in \mathbb{R}^{C_{(l+1)} \times N \times M}$$

using an filter of spatial extent  $K_1 \times K_2$

Number of parameters per filter:  $K_1 \times K_2 \times C_l$

Number of filters:  $C_{(l+1)}$

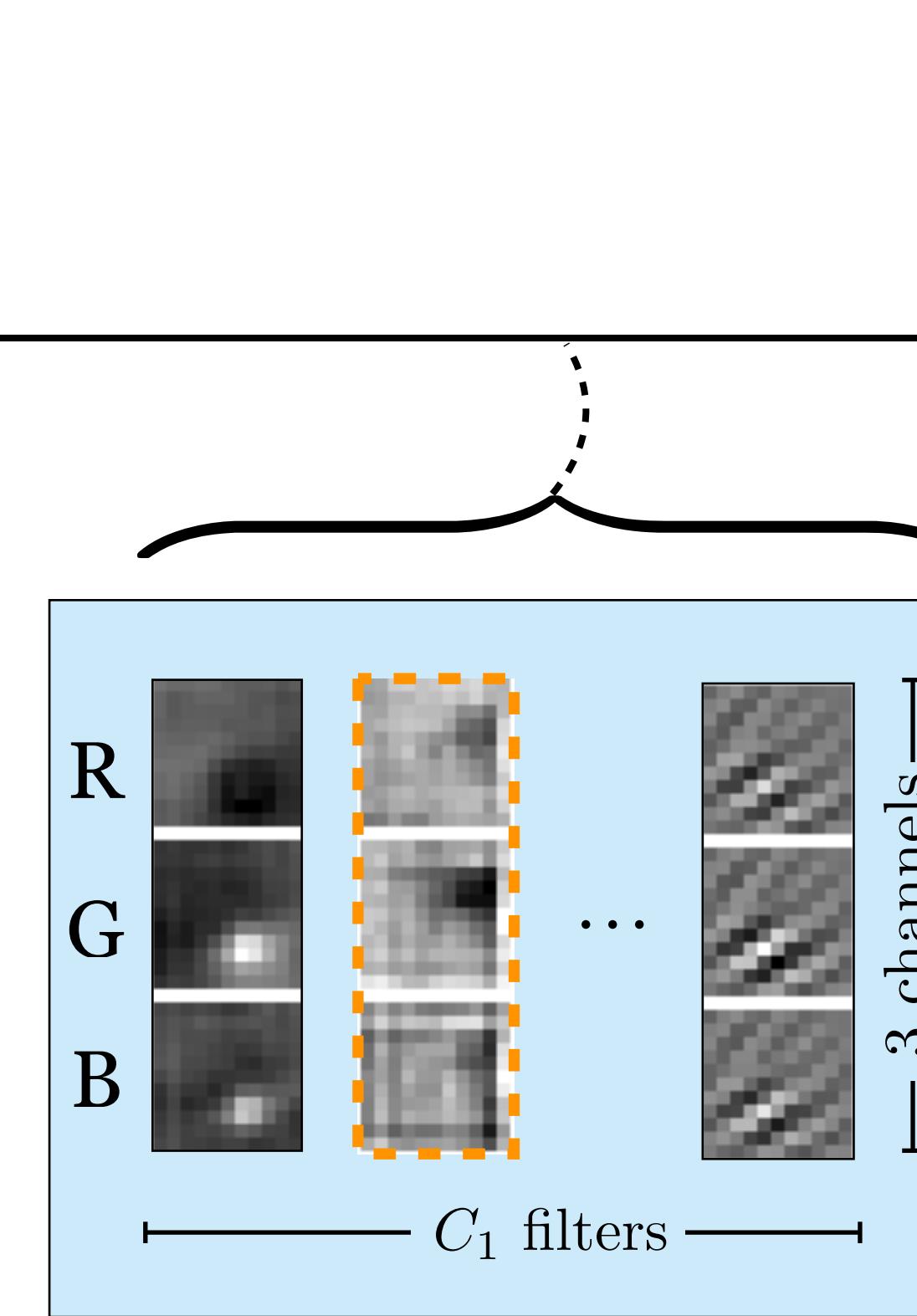
Input image (RGB)

$[H \times W \times 3]$

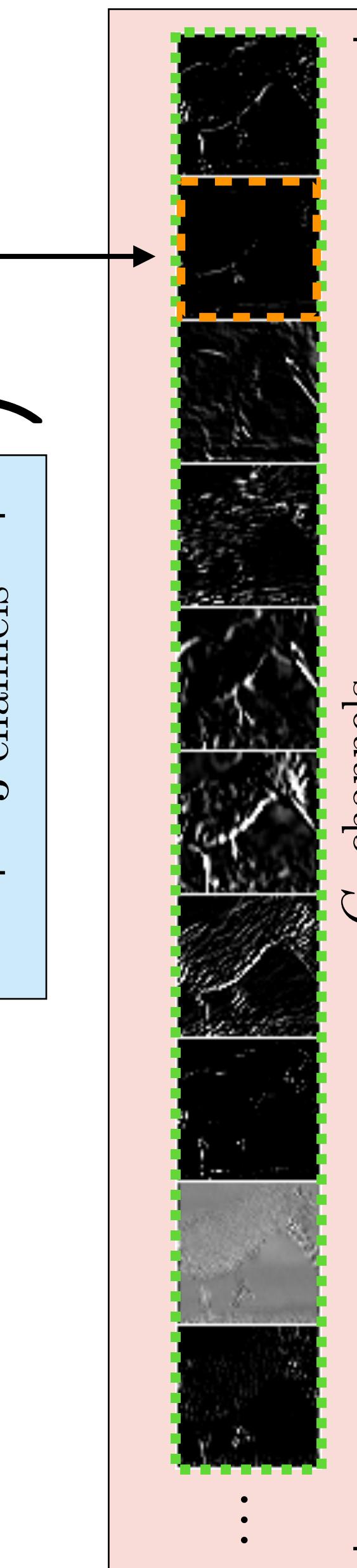


Layer 1 feature maps

$[H/4 \times W/4 \times C_1]$

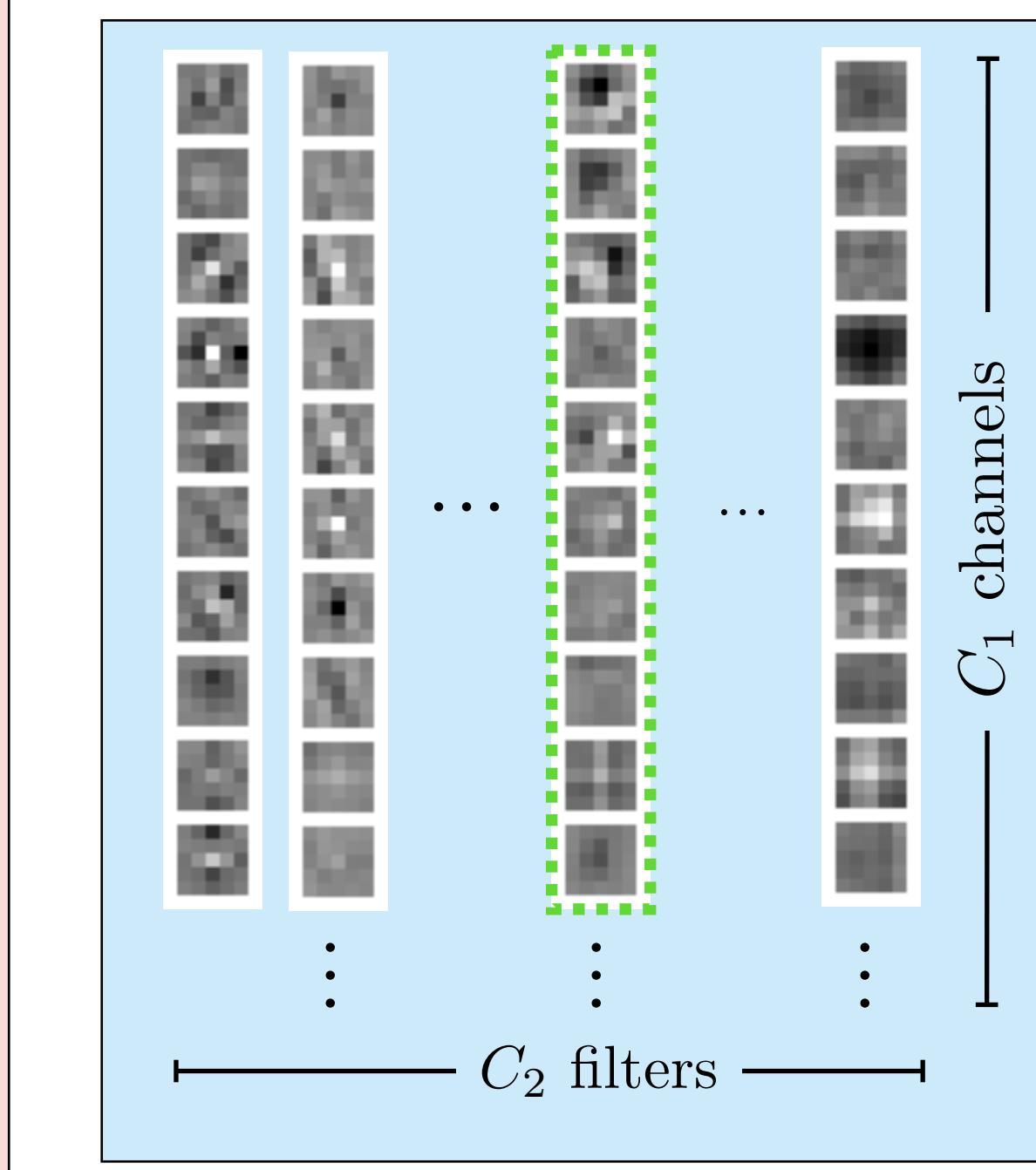


Layer 1  
filters  
(4x zoom)

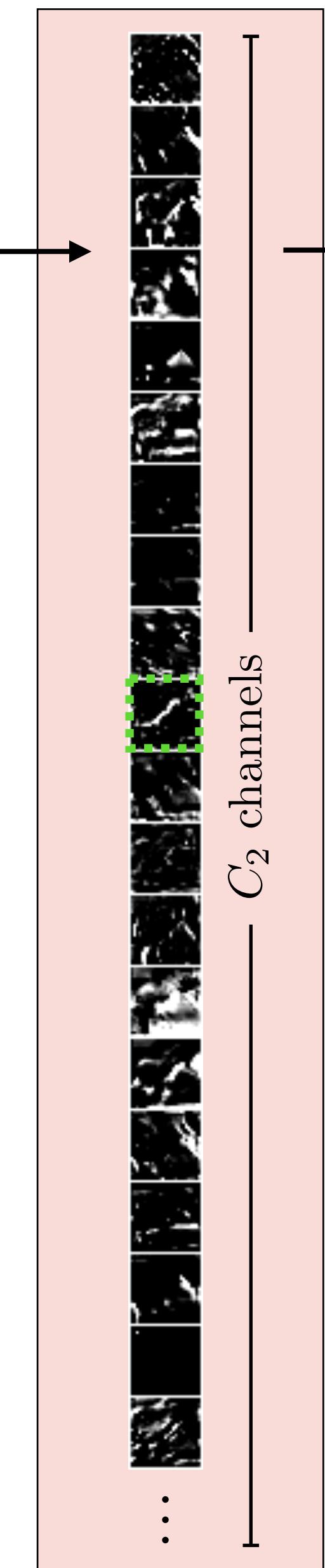


Layer 2 feature maps

$[H/8 \times W/8 \times C_2]$

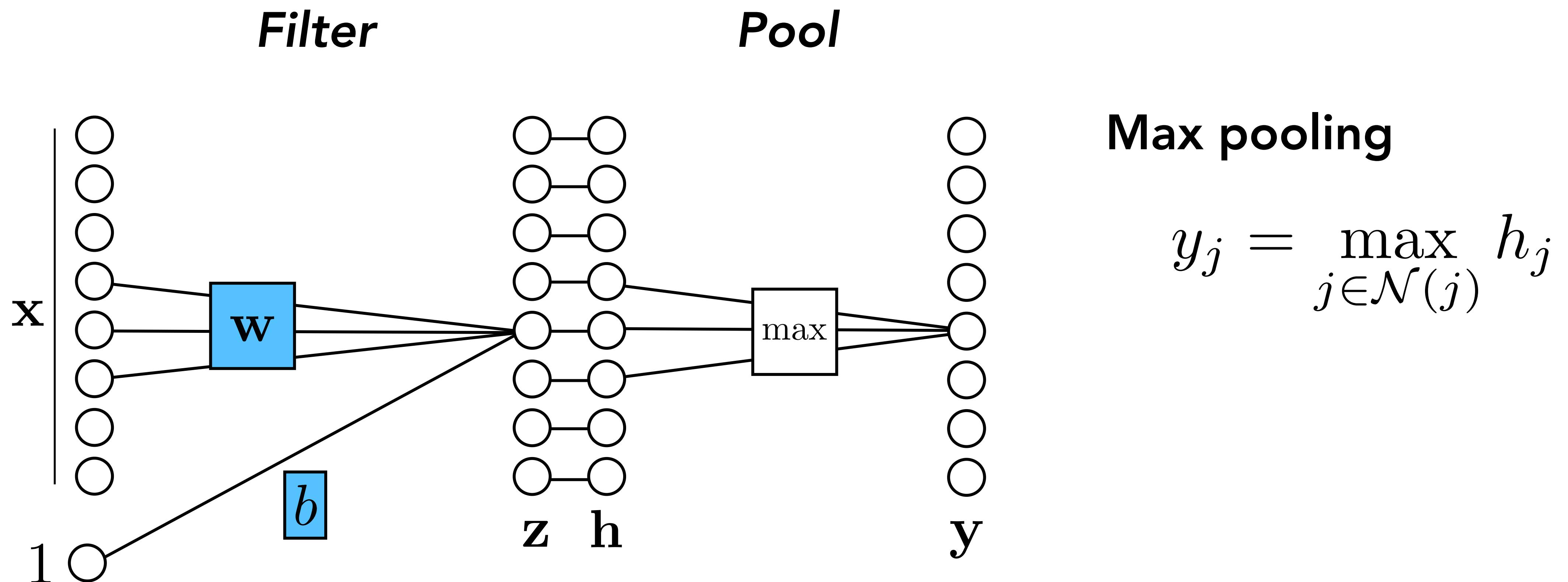


Layer 2  
filters  
(4x zoom)

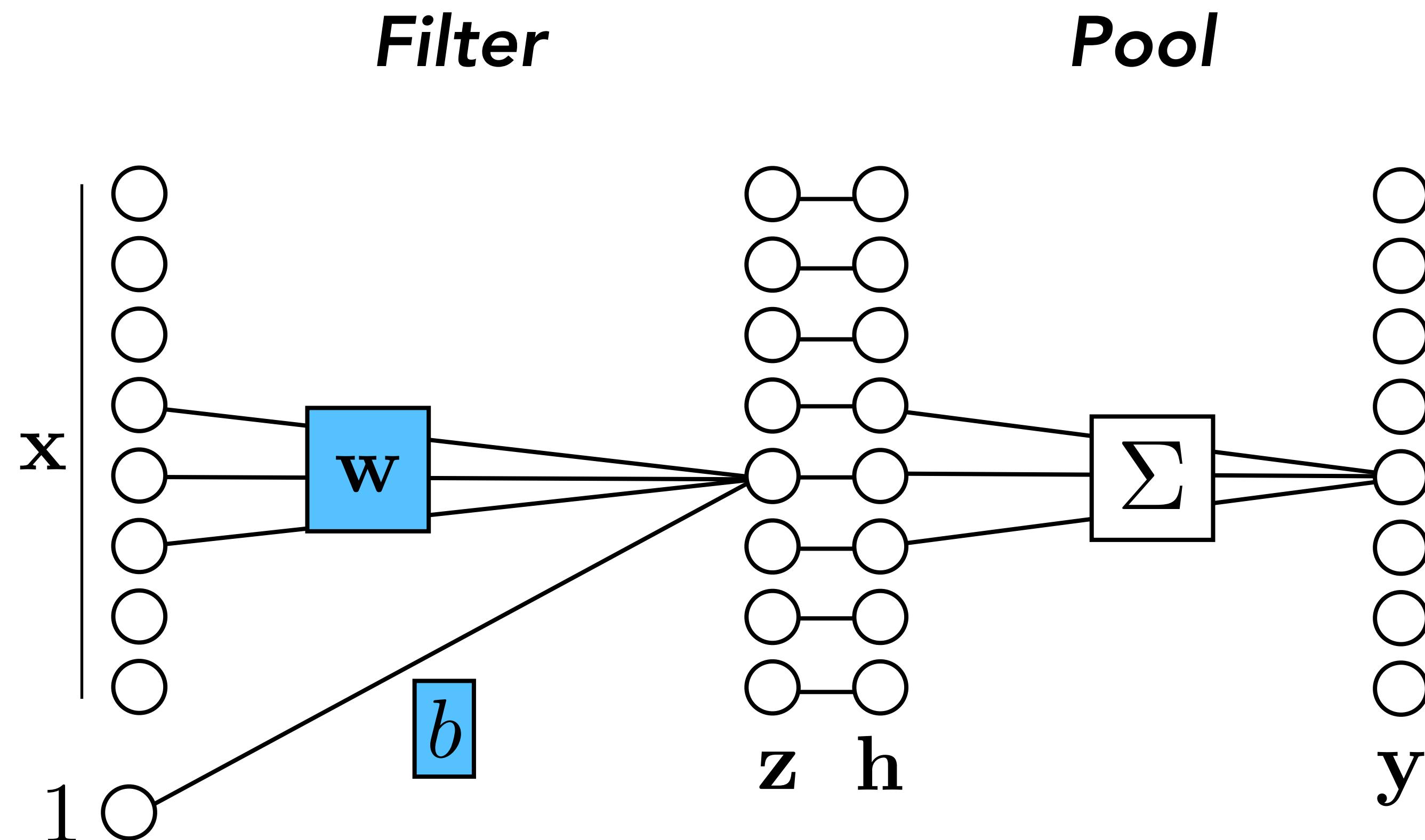


$\dots$

# Pooling



# Pooling



**Max pooling**

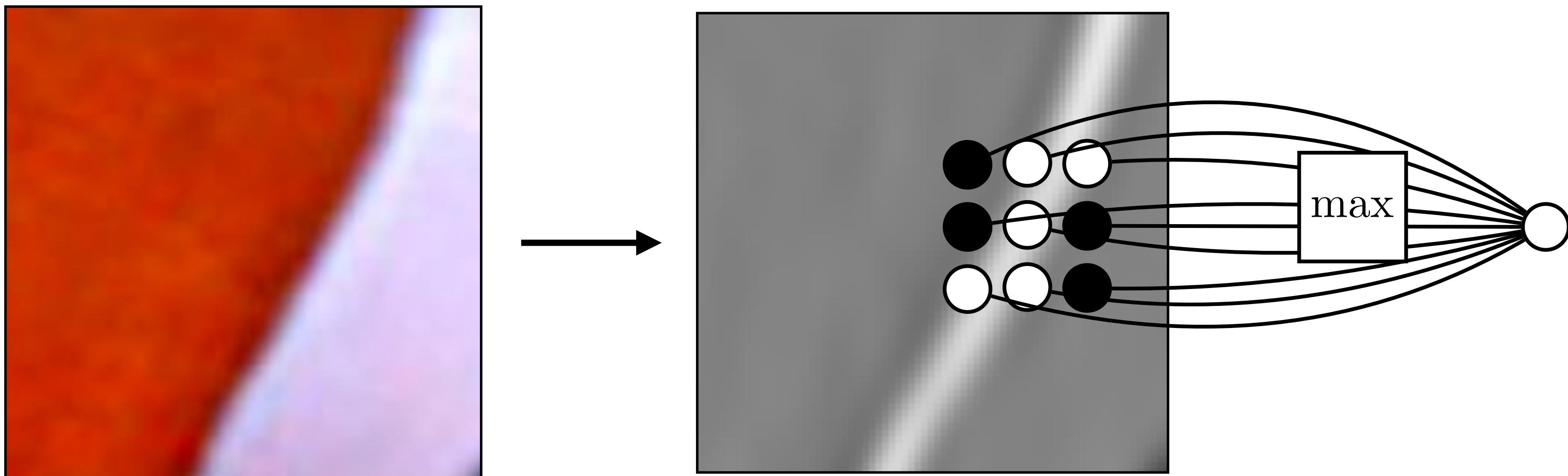
$$y_j = \max_{j \in \mathcal{N}(j)} h_j$$

**Mean pooling**

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

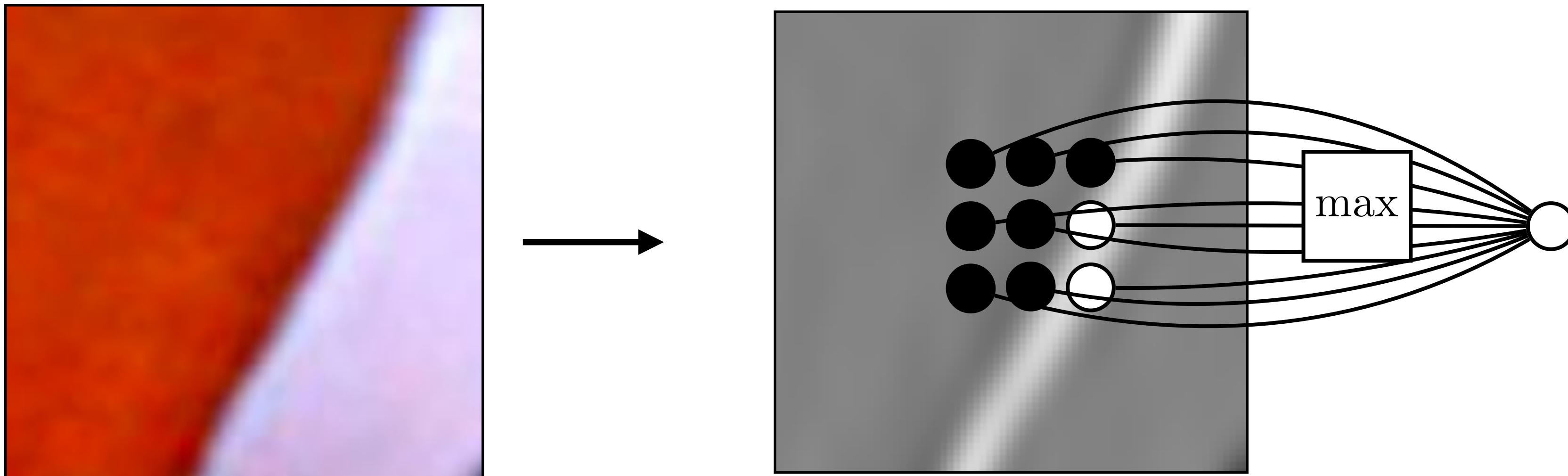
# Pooling — Why?

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



# Pooling — Why?

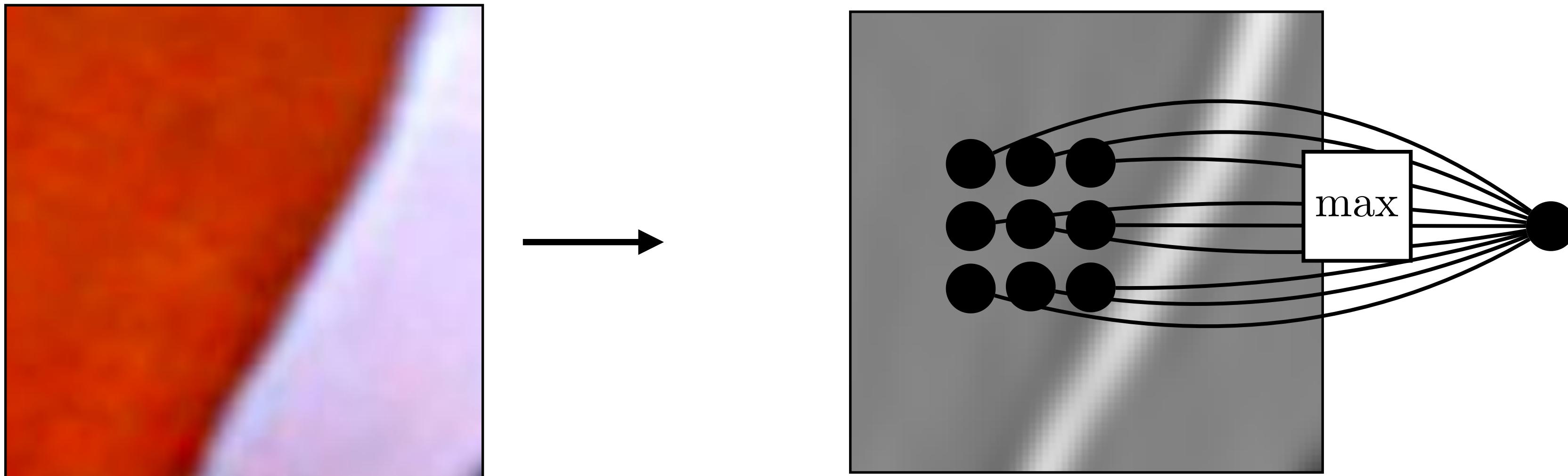
Pooling across spatial locations achieves stability w.r.t. small translations:



large response  
regardless of exact  
position of edge

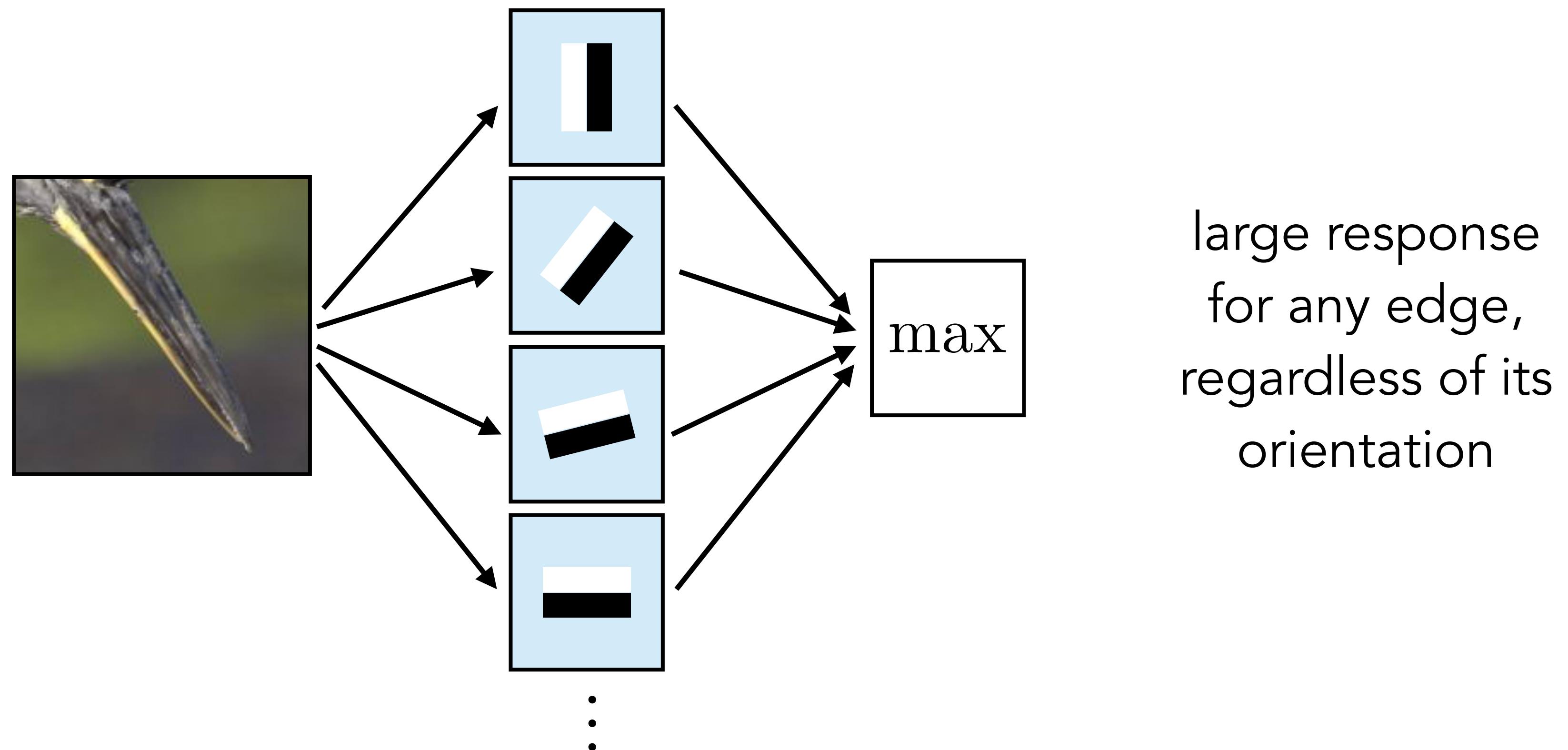
# Pooling — Why?

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

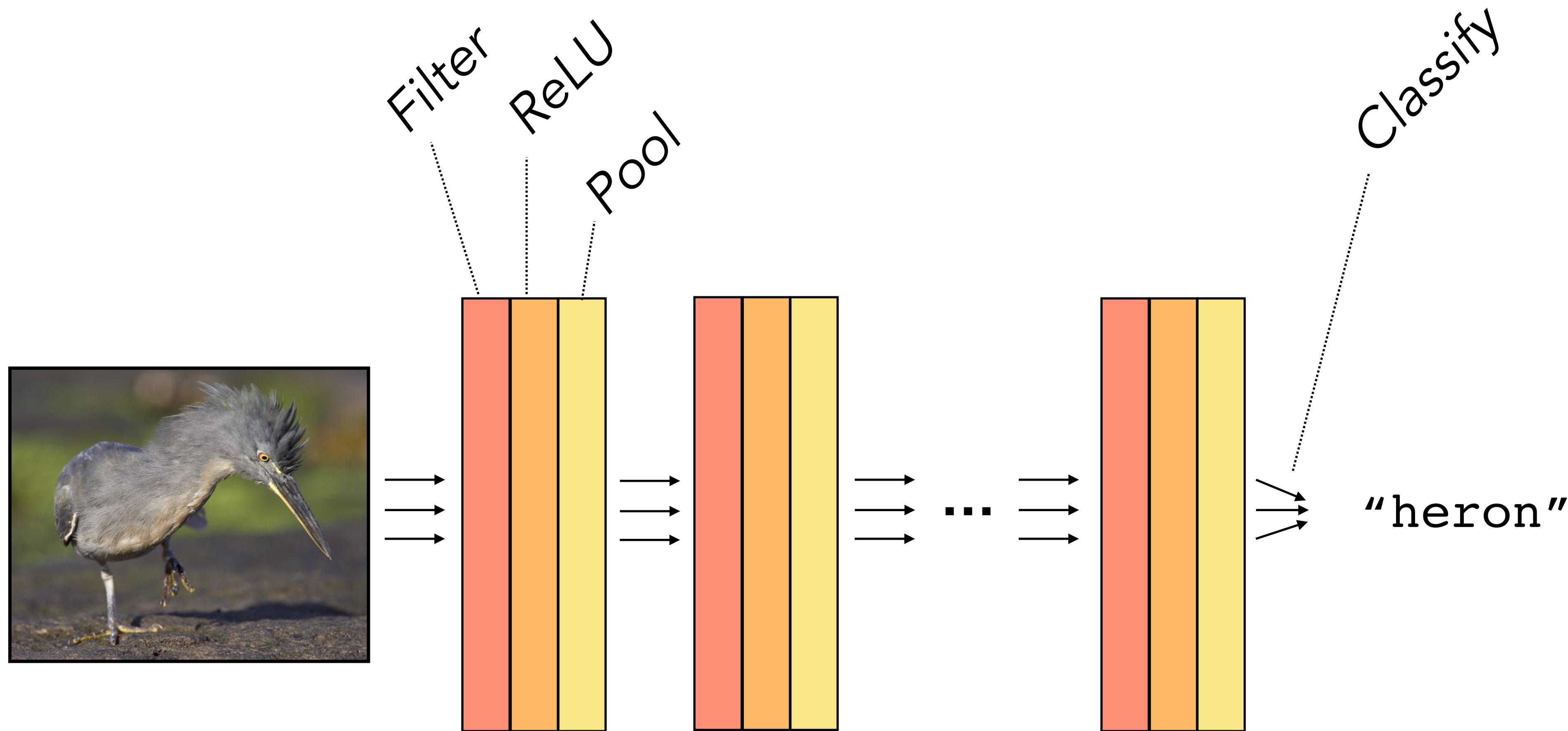


# Pooling across channels — Why?

Pooling across feature channels (filter outputs)  
can achieve other kinds of invariances:

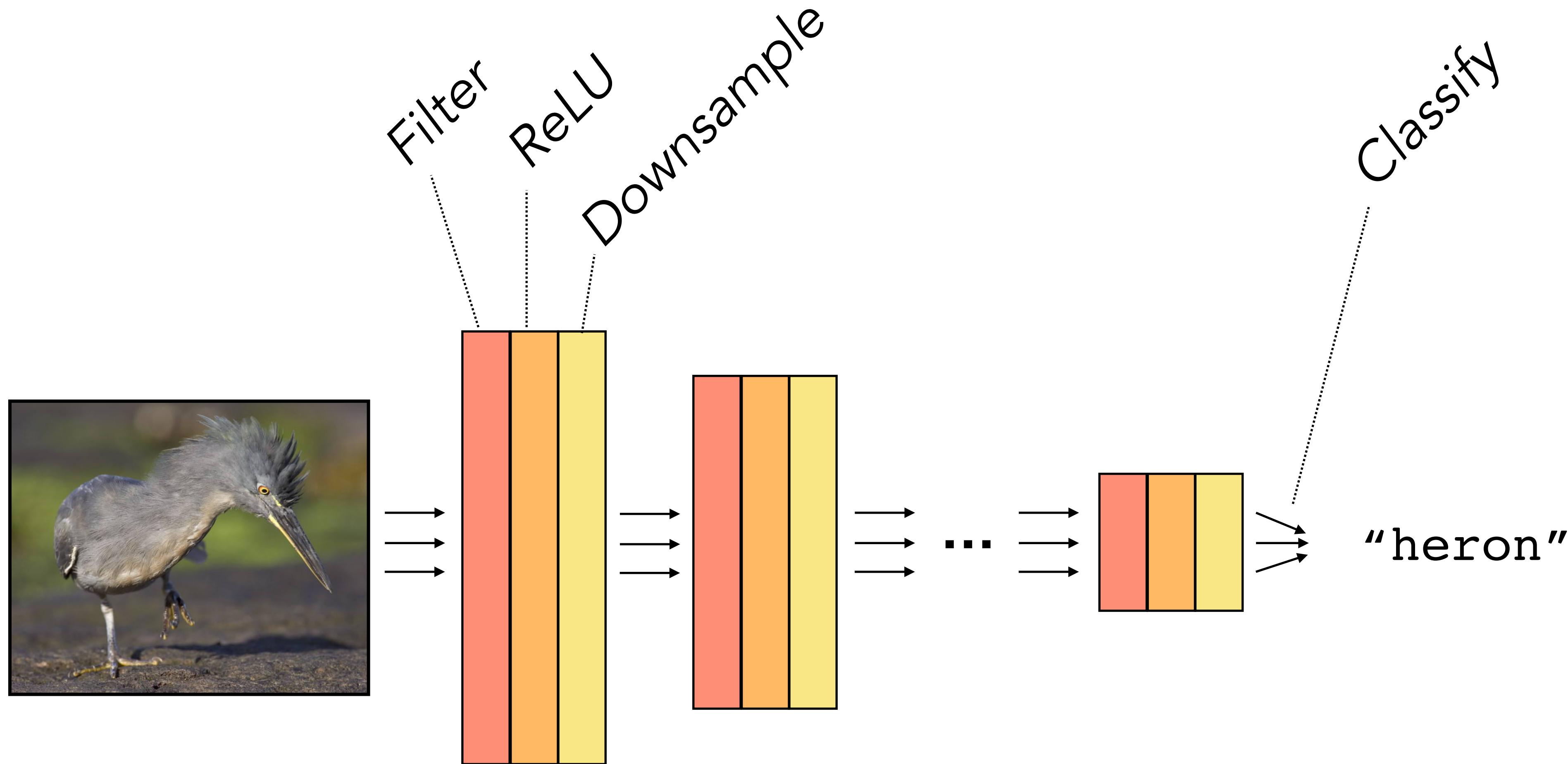


# Computation in a neural net



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

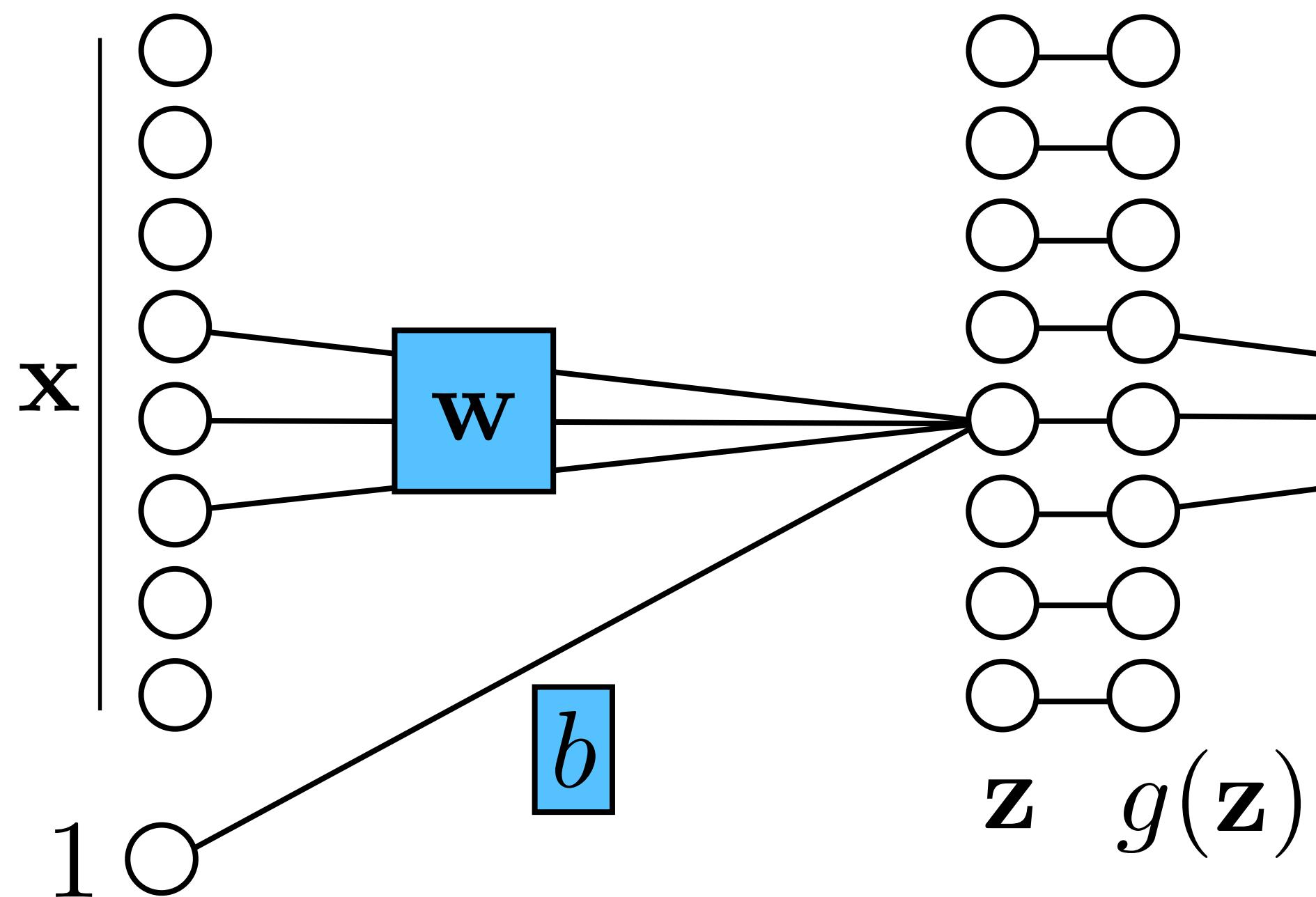
# Computation in a neural net



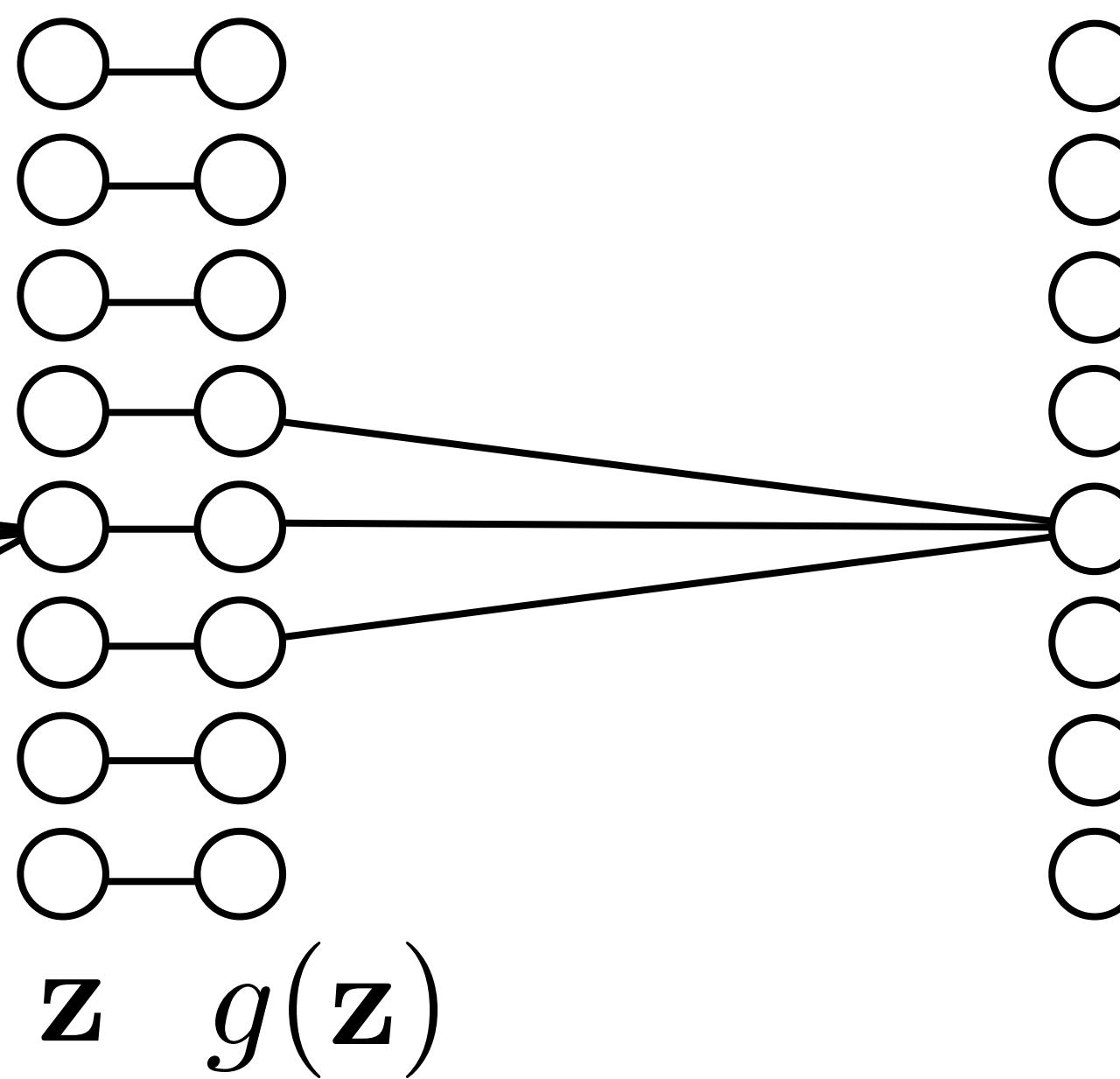
$$f(\mathbf{x}) = f_L(\dots f_2(f_1(\mathbf{x})))$$

# Downsampling

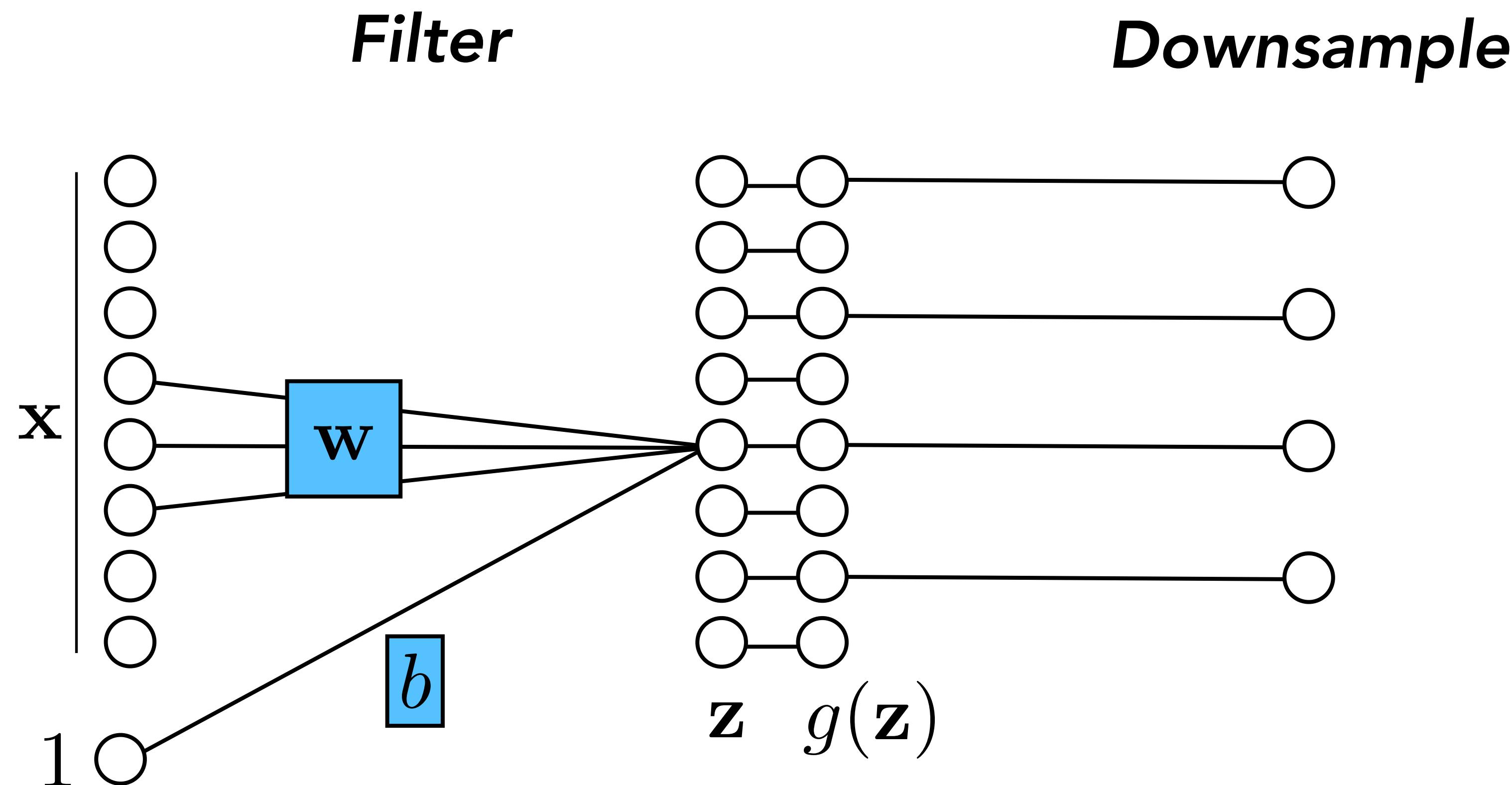
*Filter*



*Pool and downsample*

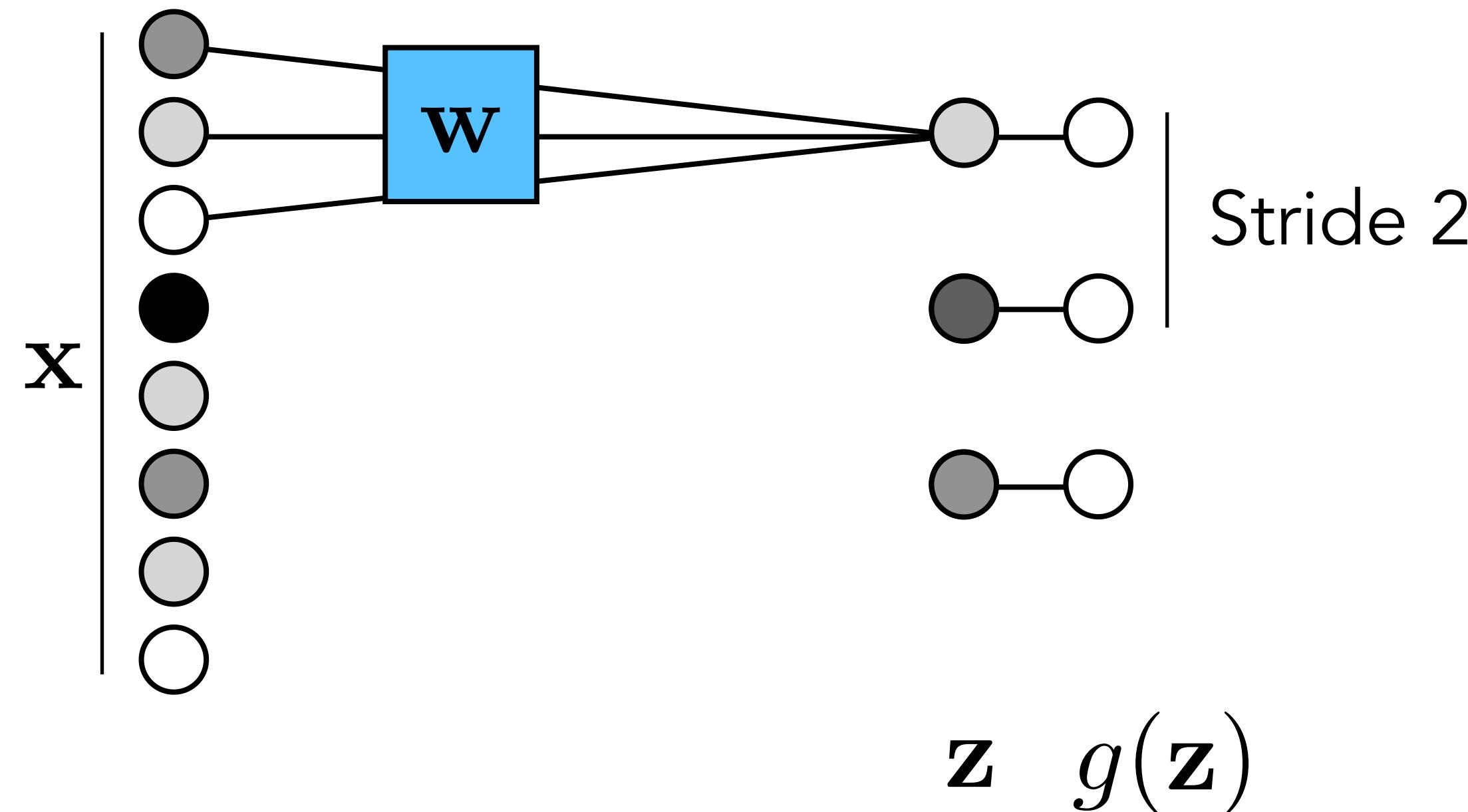


# Downsampling



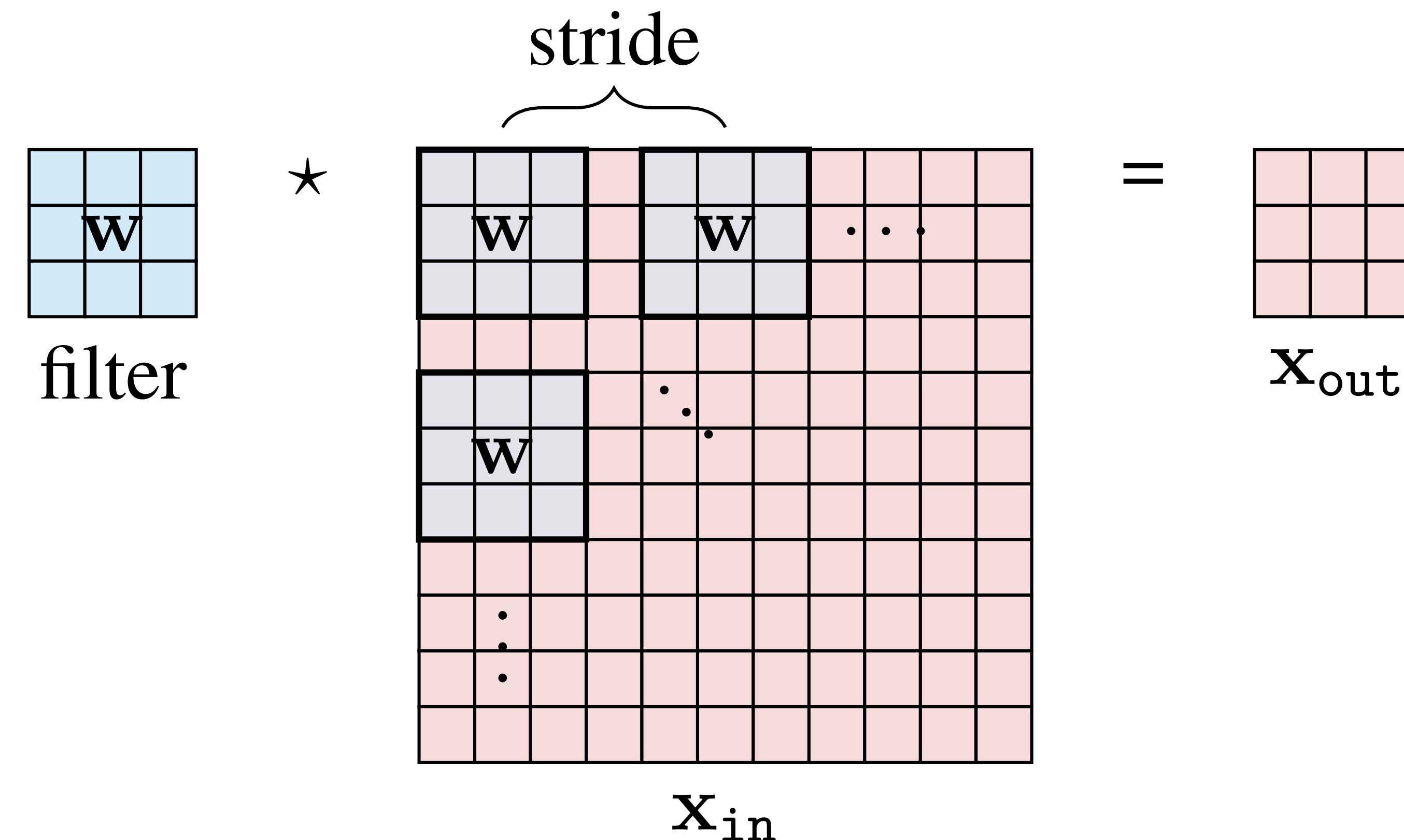
# Strided operations

## Conv layer

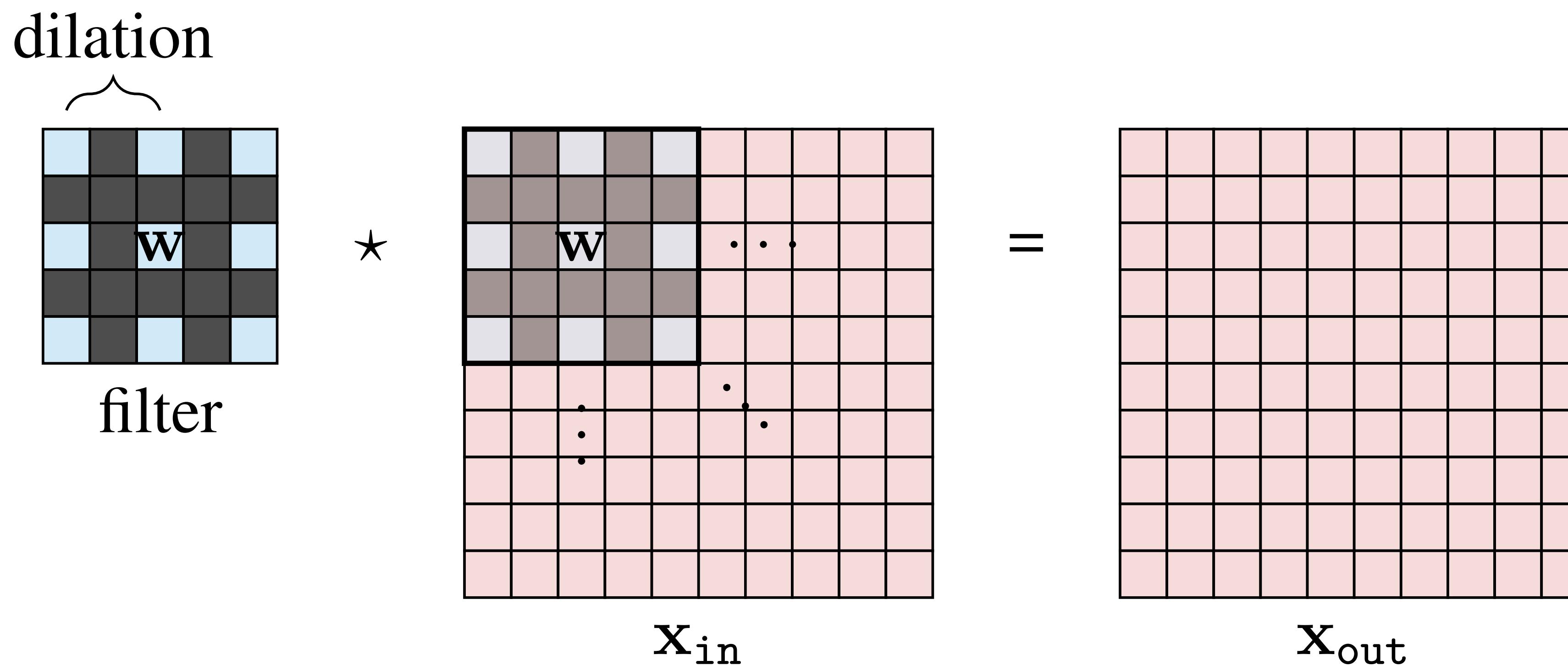


**Strided operations** combine a given operation (convolution or pooling) and downsampling into a single operation.

# Strided operations (2D)

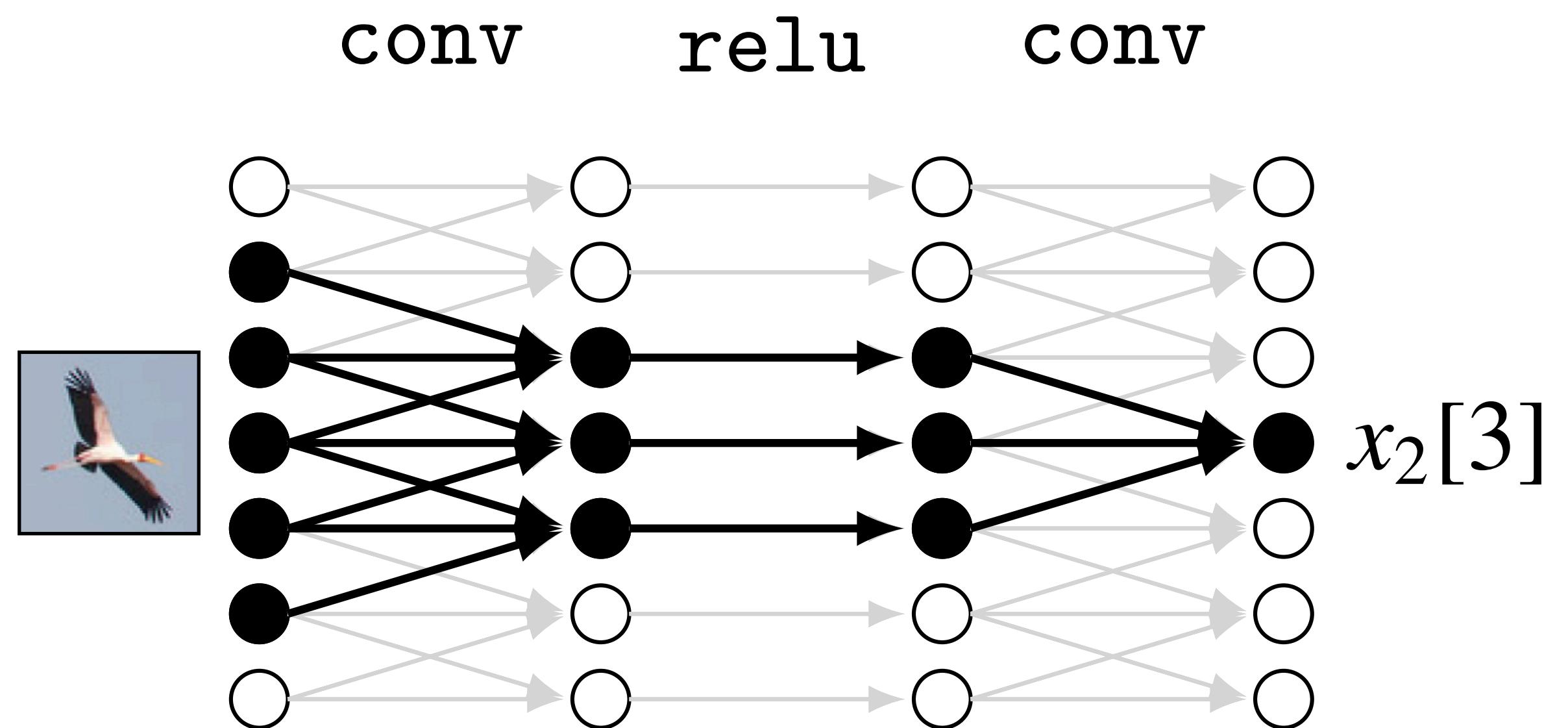


# Dilated filter

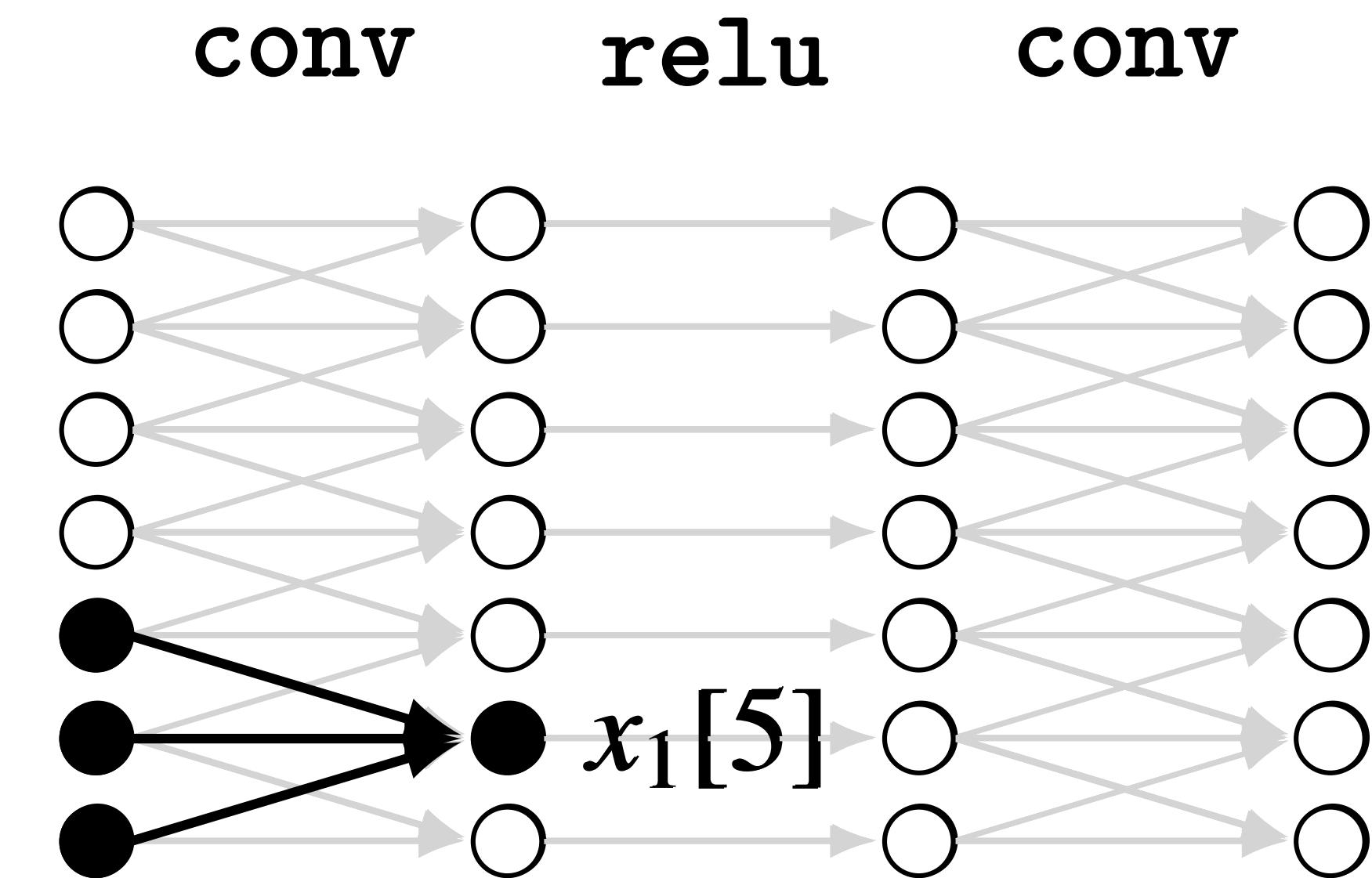


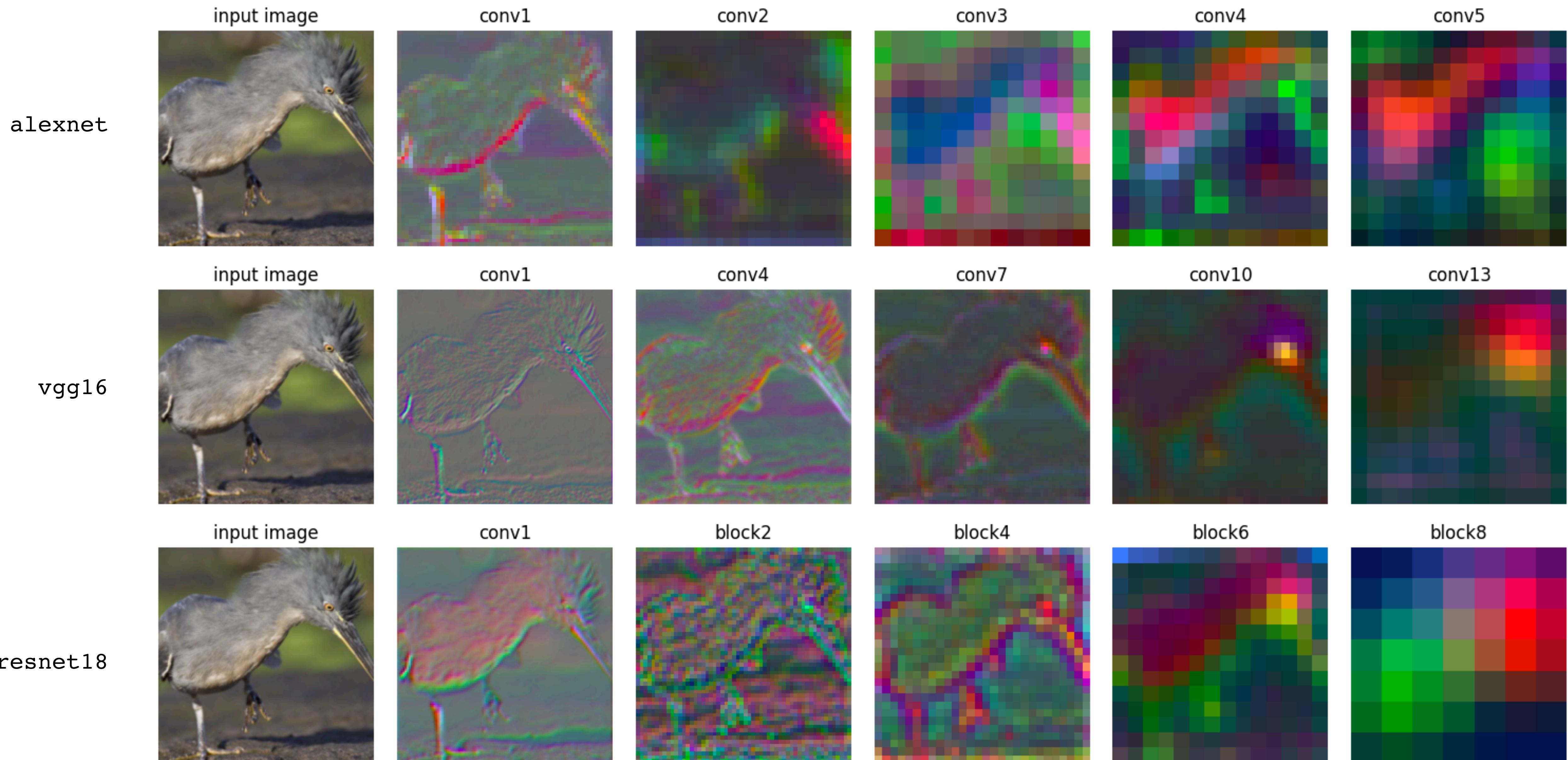
Covers a large receptive field with fewer parameters.

# Receptive fields



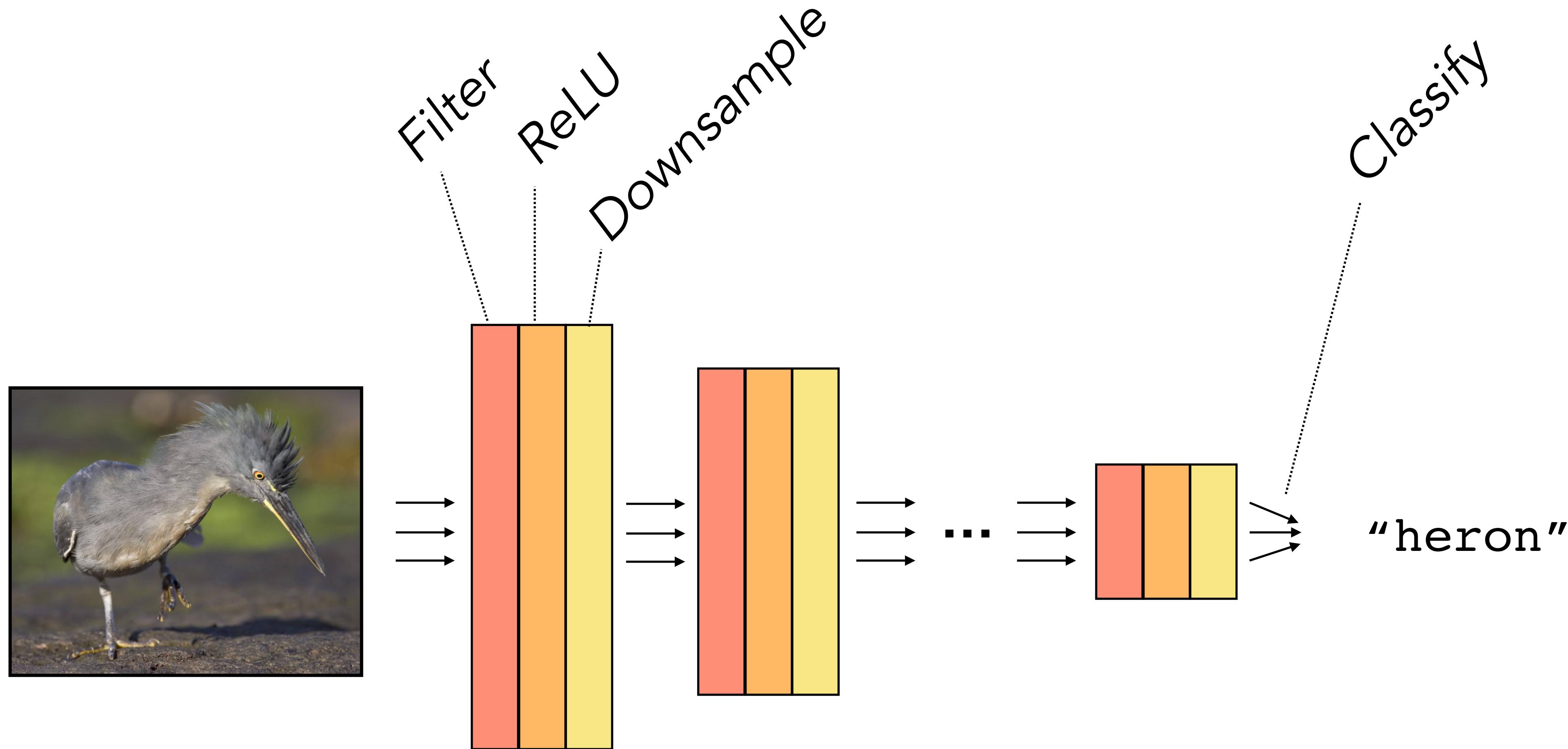
Sky
Sky
Bird
Sky



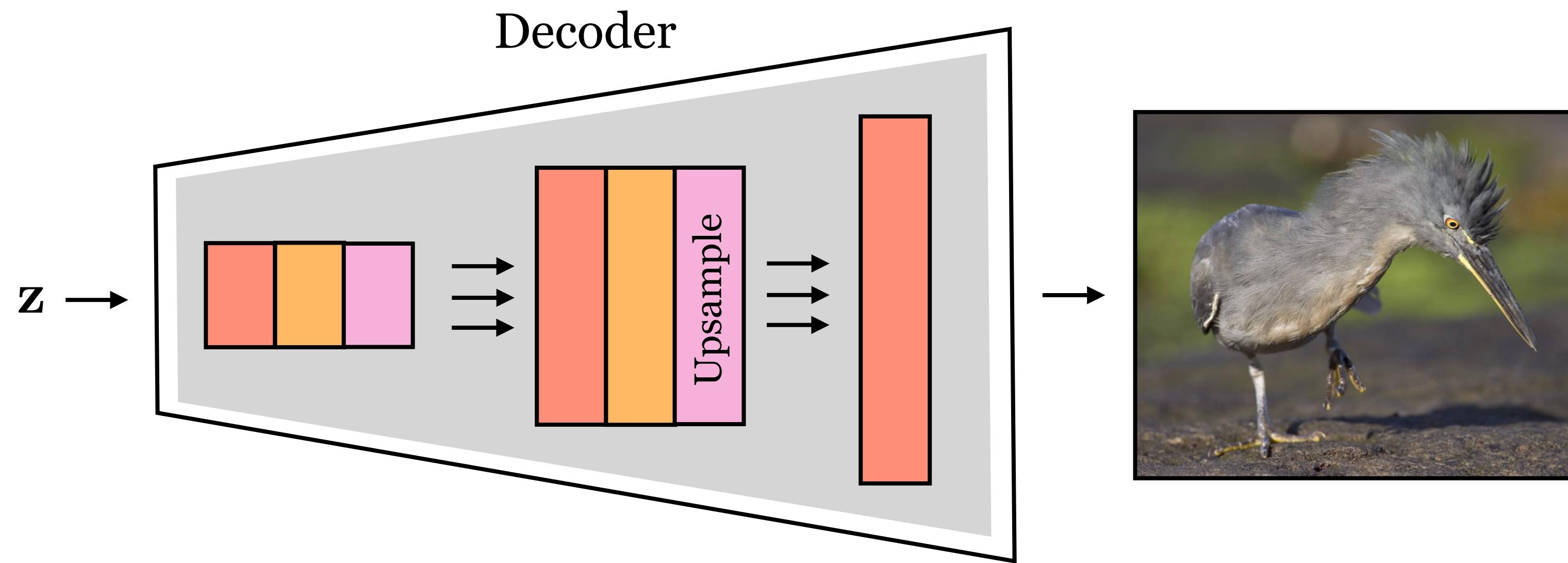
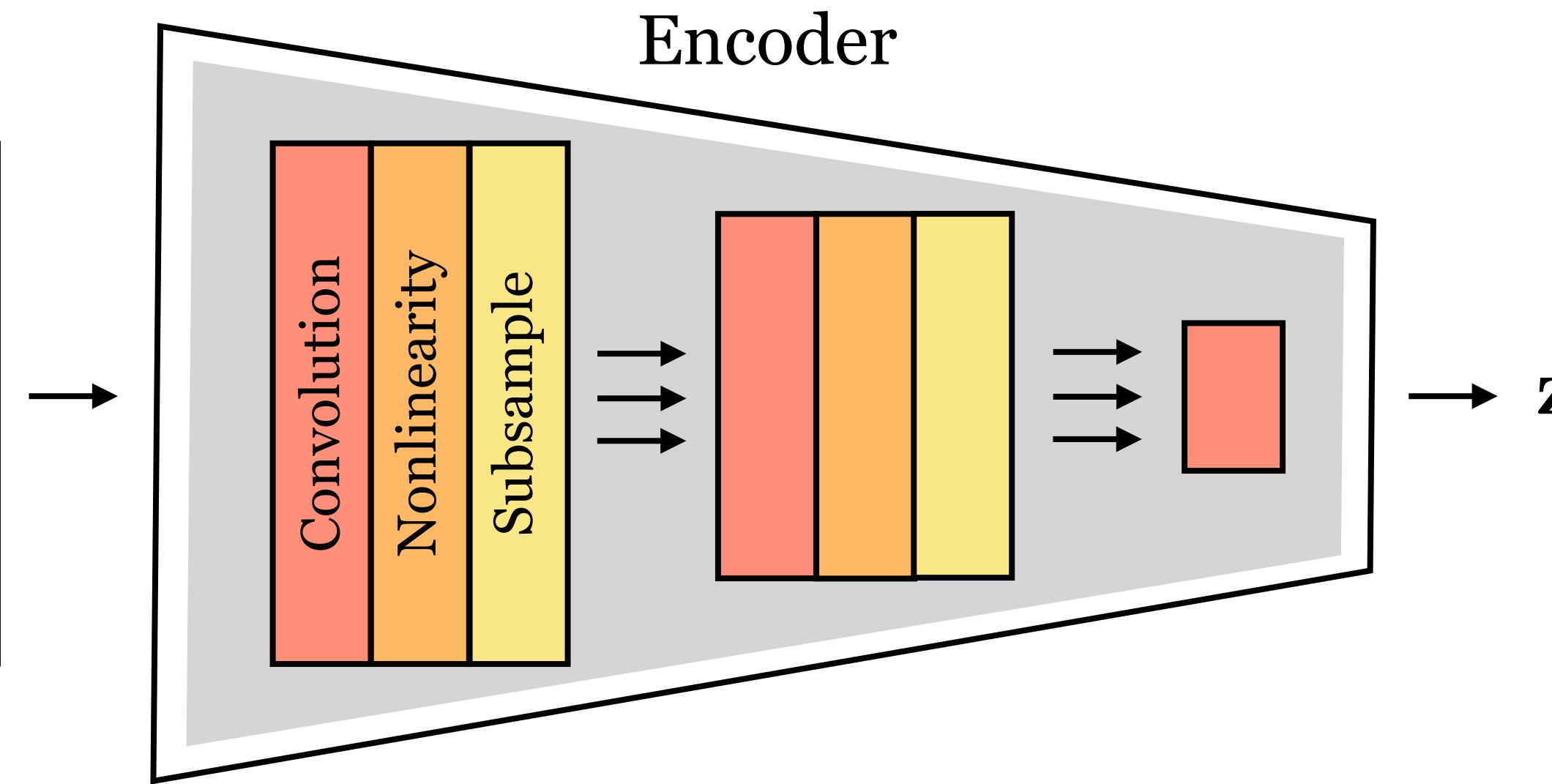


# Popular CNN Architectures

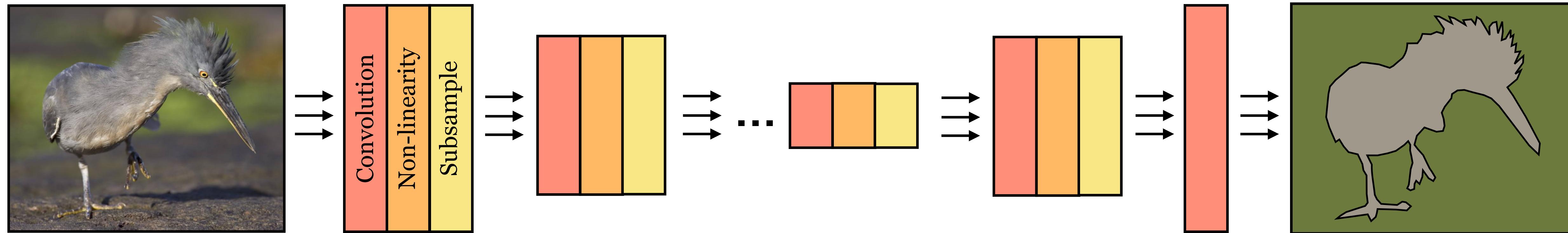
# Computation in a neural net



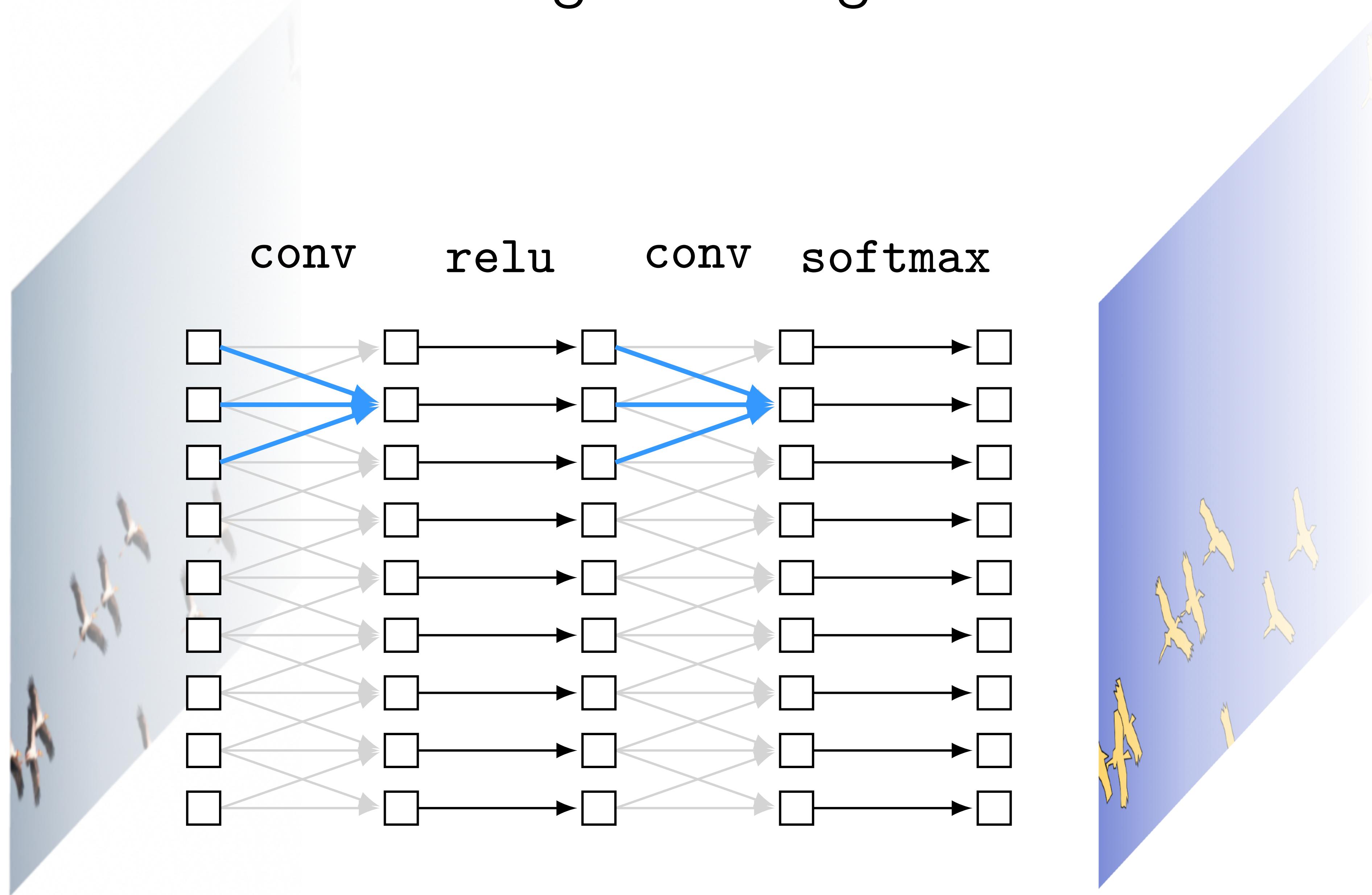
$$f(\mathbf{x}) = f_L(\dots f_2(f_1(\mathbf{x})))$$



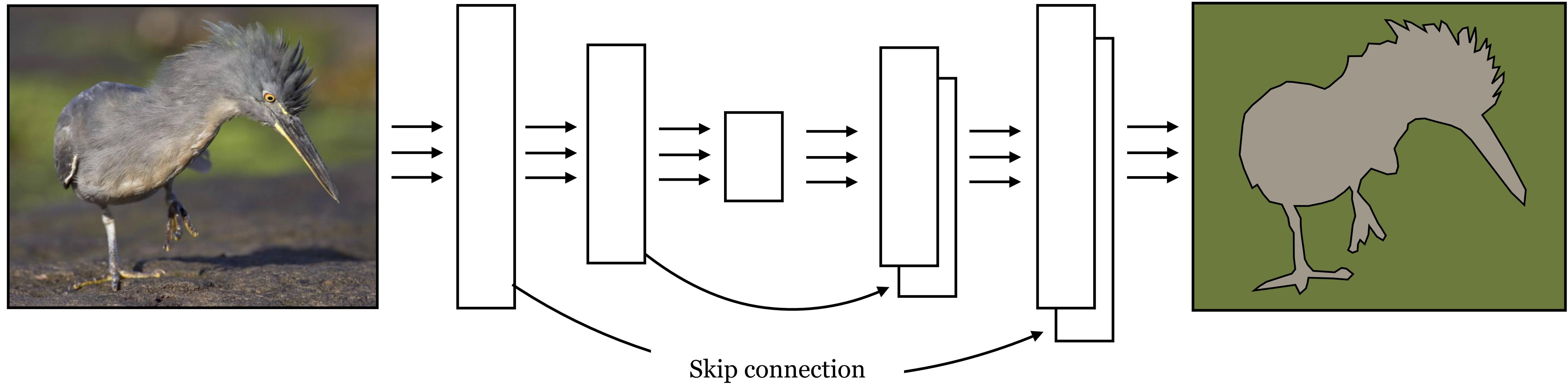
# Image-to-image



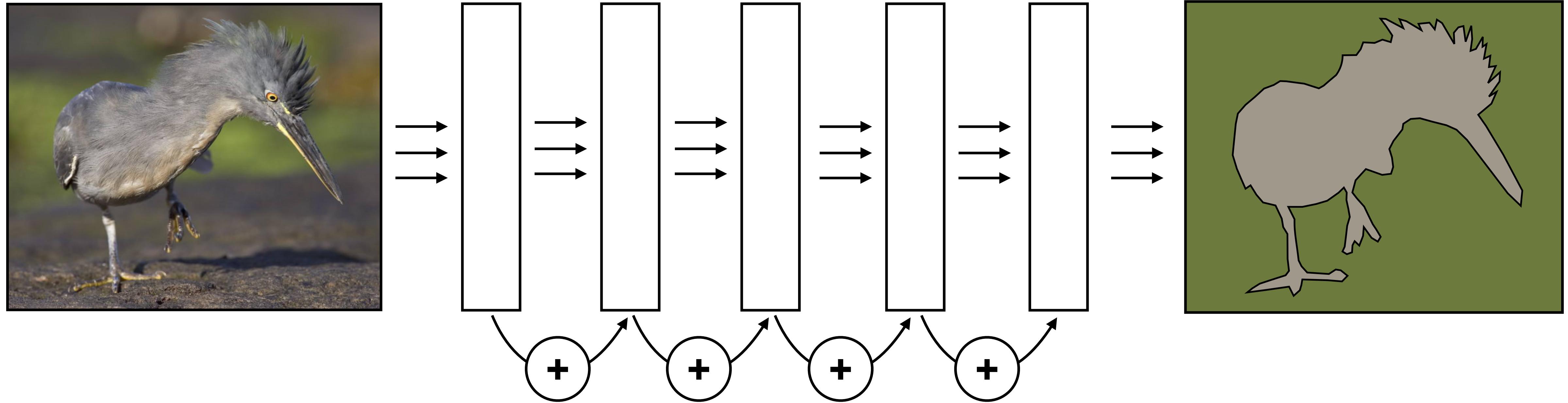
# Image-to-image



# U-net



# ResNet



**Residual connection:**  $\mathbf{x}_{\text{out}} = F(\mathbf{x}_{\text{in}}) + \mathbf{x}_{\text{in}}$

Or, if you want to change dimensionality:  $\mathbf{x}_{\text{out}} = F(\mathbf{x}_{\text{in}}) + \mathbf{W}\mathbf{x}_{\text{in}}$