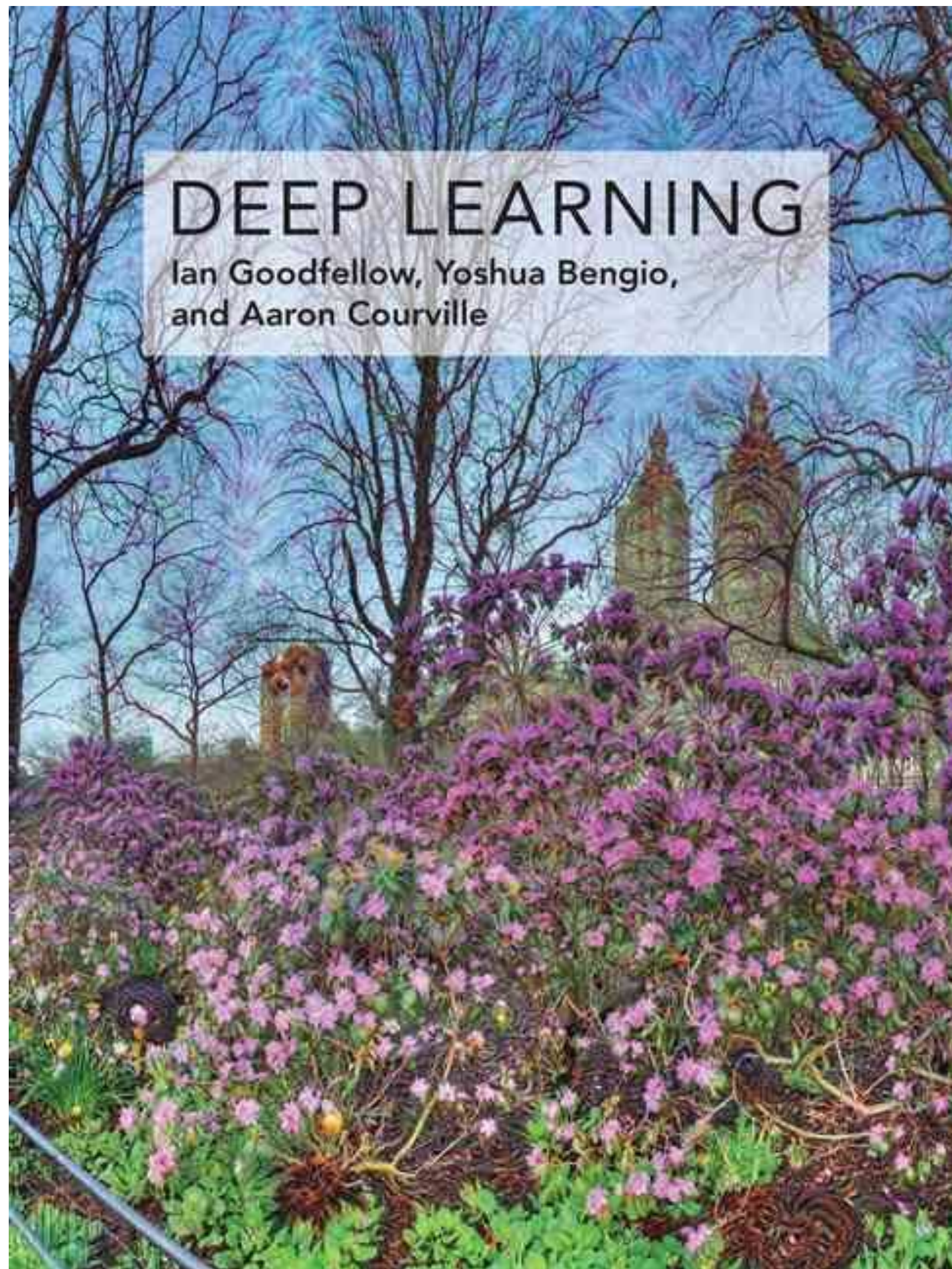


CNNs and Spatial Processing

Bill Freeman, Antonio Torralba, Phillip Isola
6.819 / 6.869

Hi!



<http://www.deeplearningbook.org/>

By Ian Goodfellow, Yoshua Bengio and Aaron Courville

November 2016

Today: parts of chapter 9

Review lectures 5 through 8 for background on signal processing, convolution, and multiscale image processing — this is the technology that underlies convnets!

Today

- How to use networks for images
- Why “C”-NNs
- Standard building blocks of CNNs
- Some important networks & their tricks
- Some debugging tools

Image classification

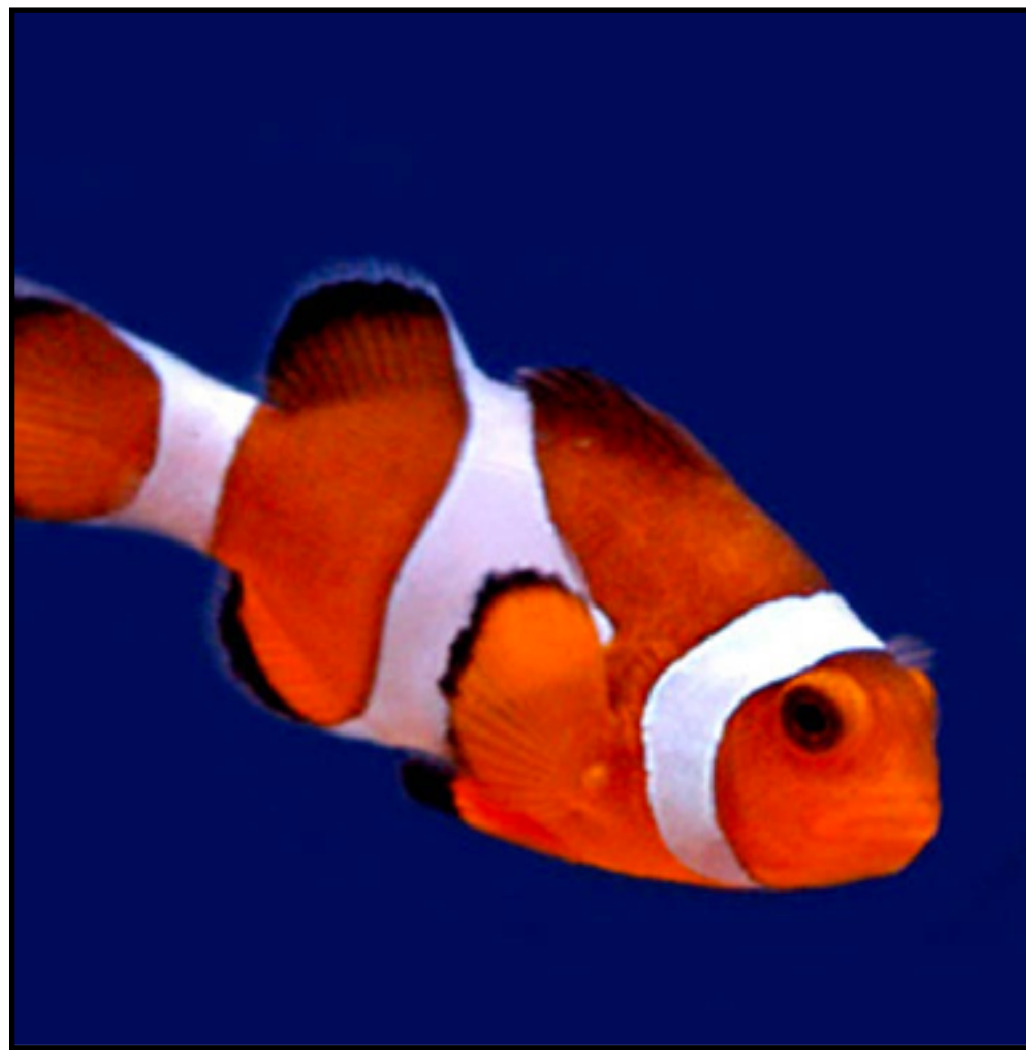


image x

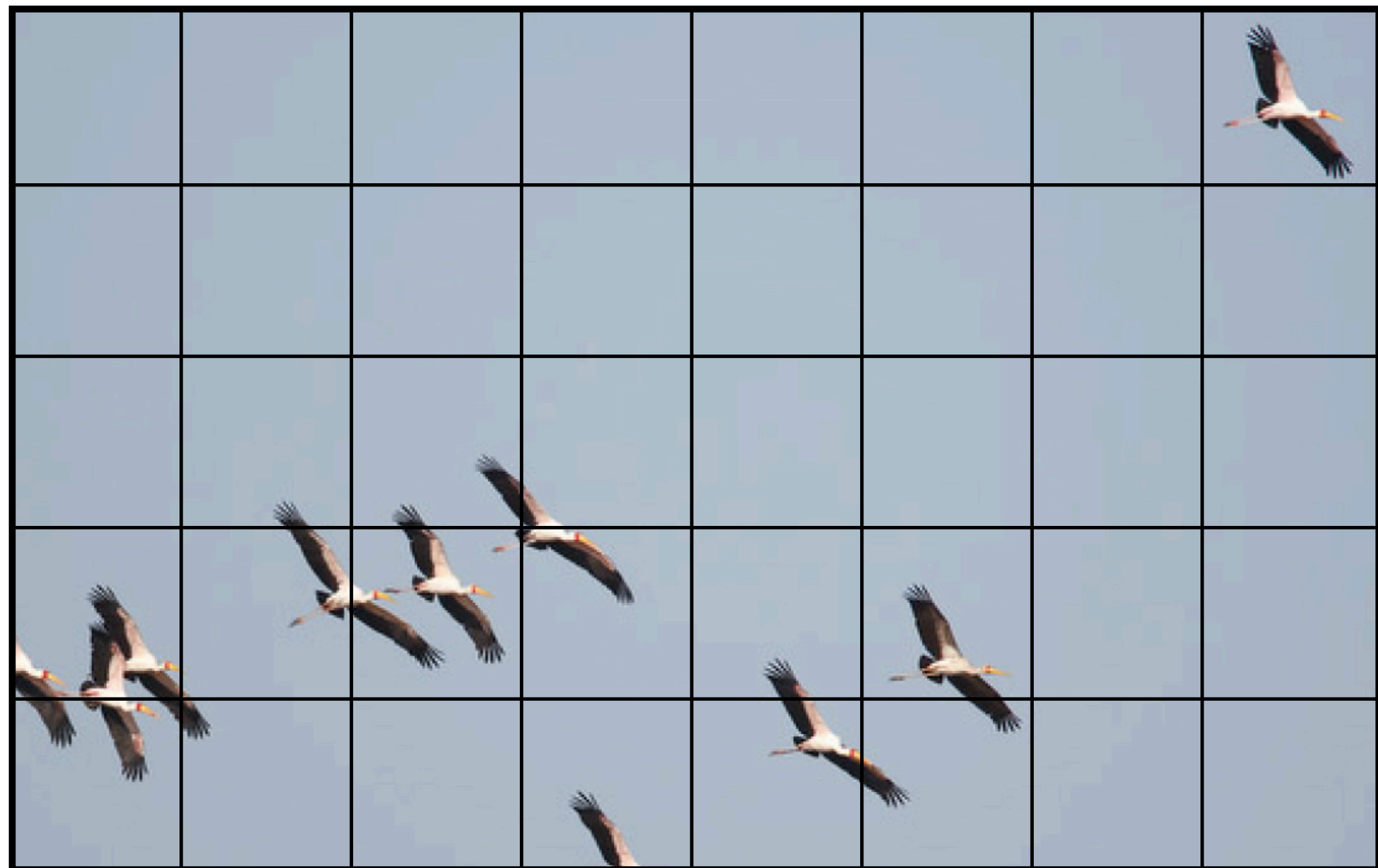


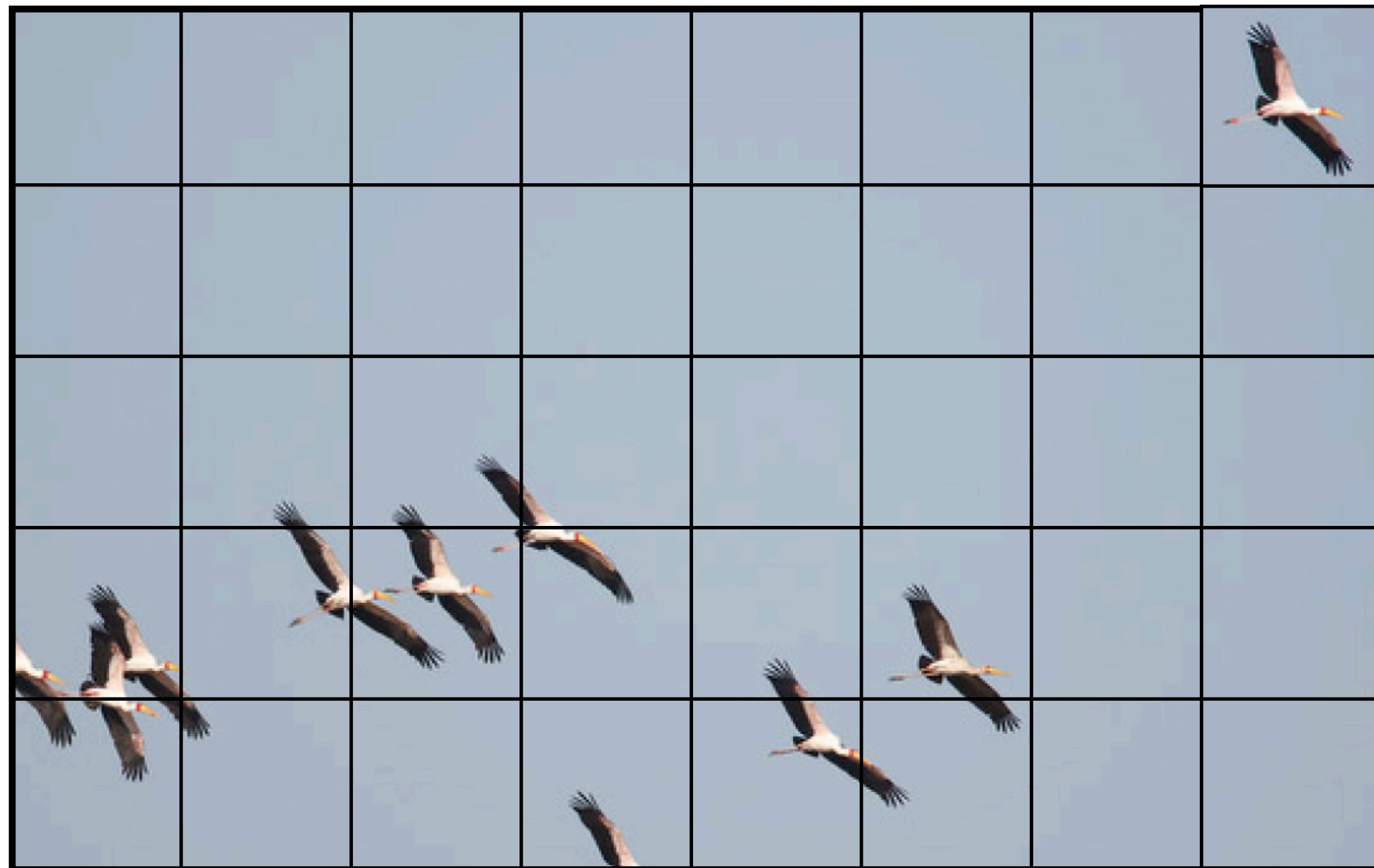
"Fish"

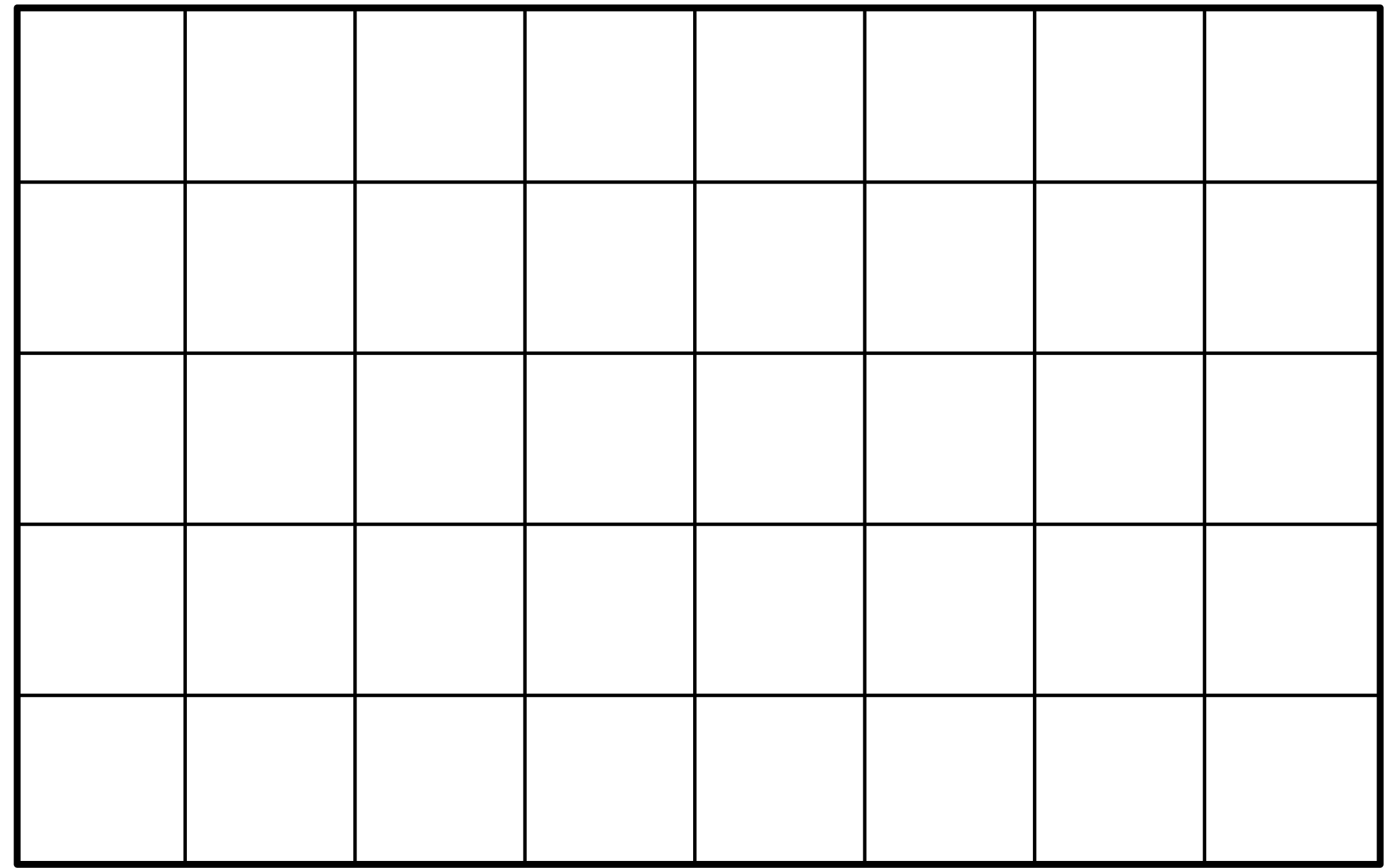
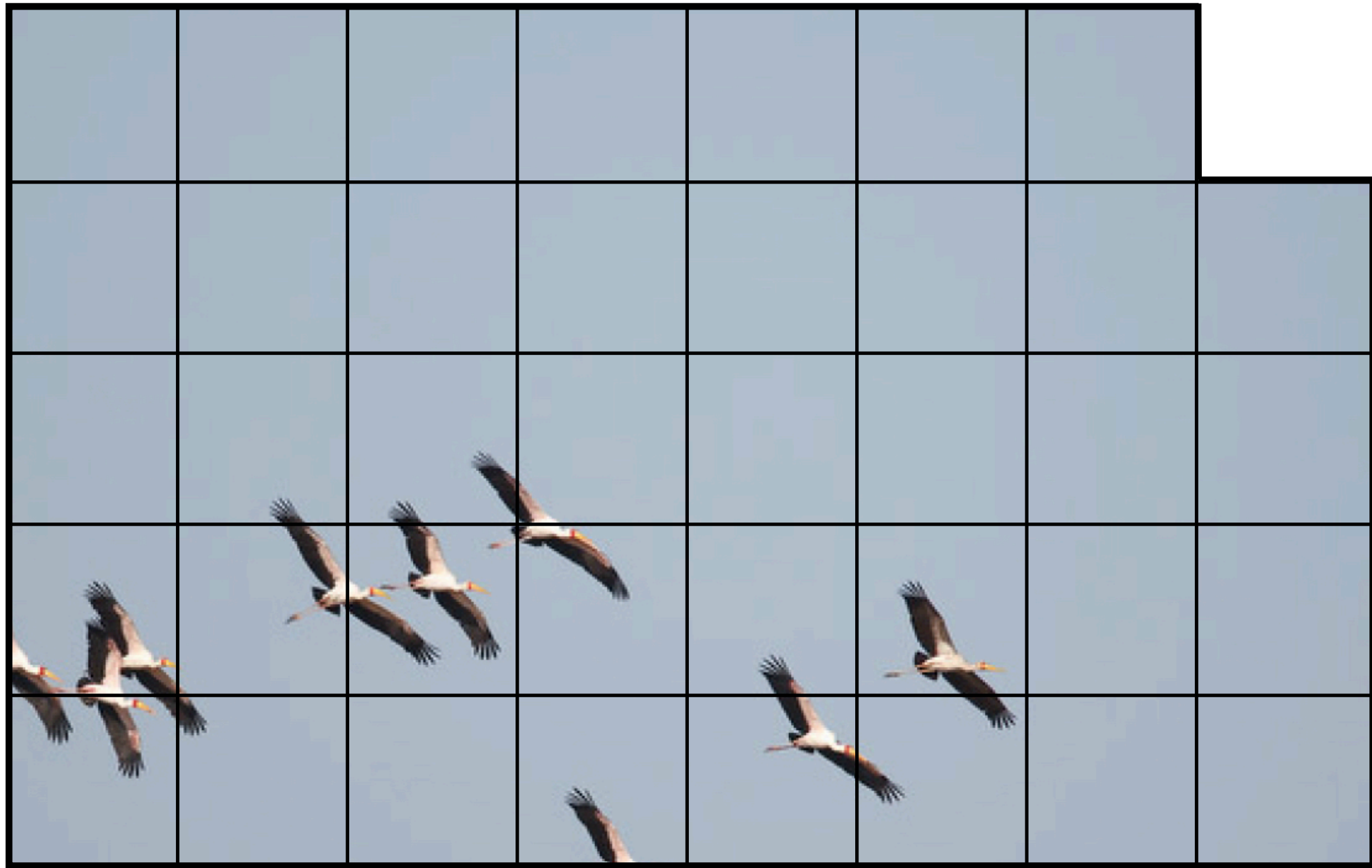
label y

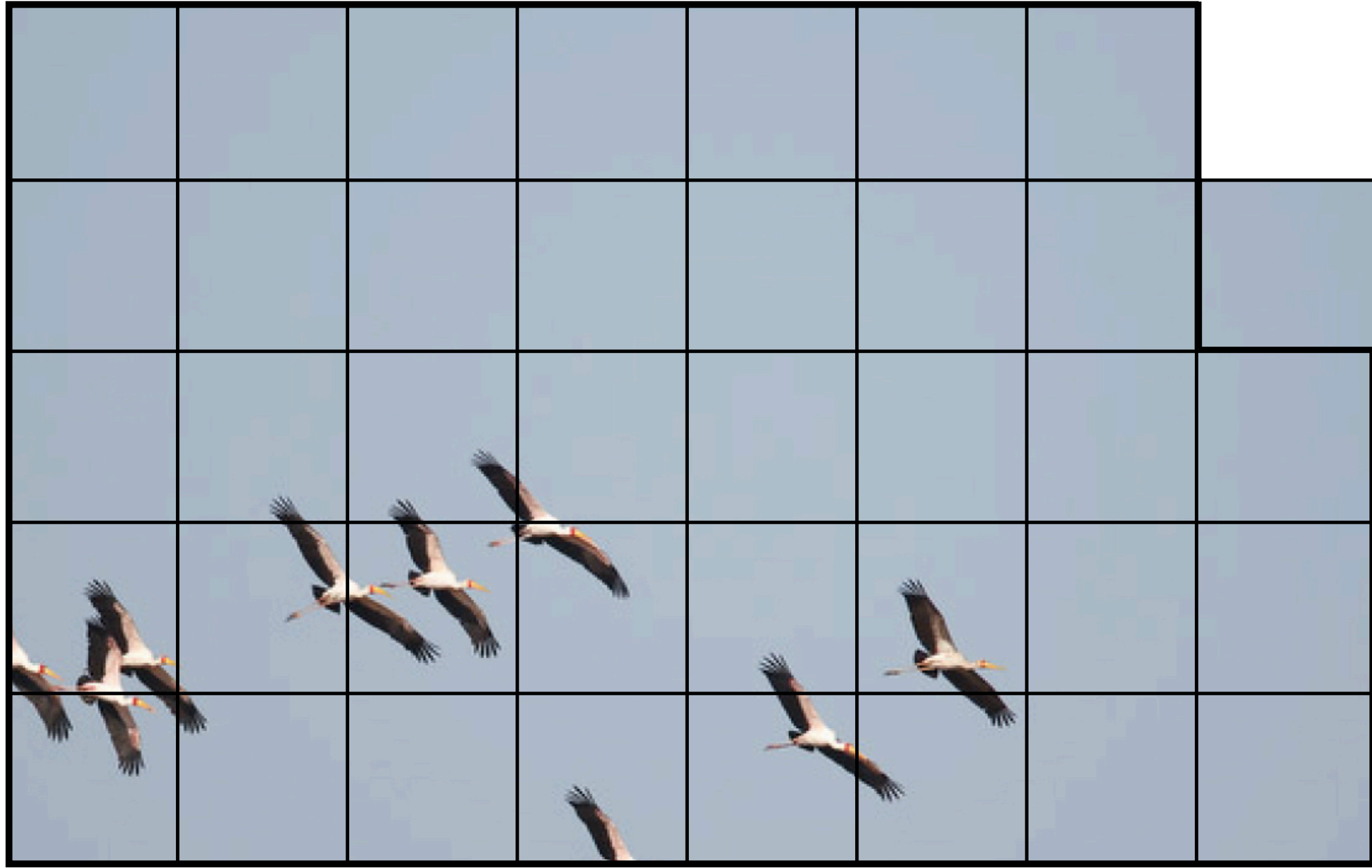


Photo credit: Fredo Durand



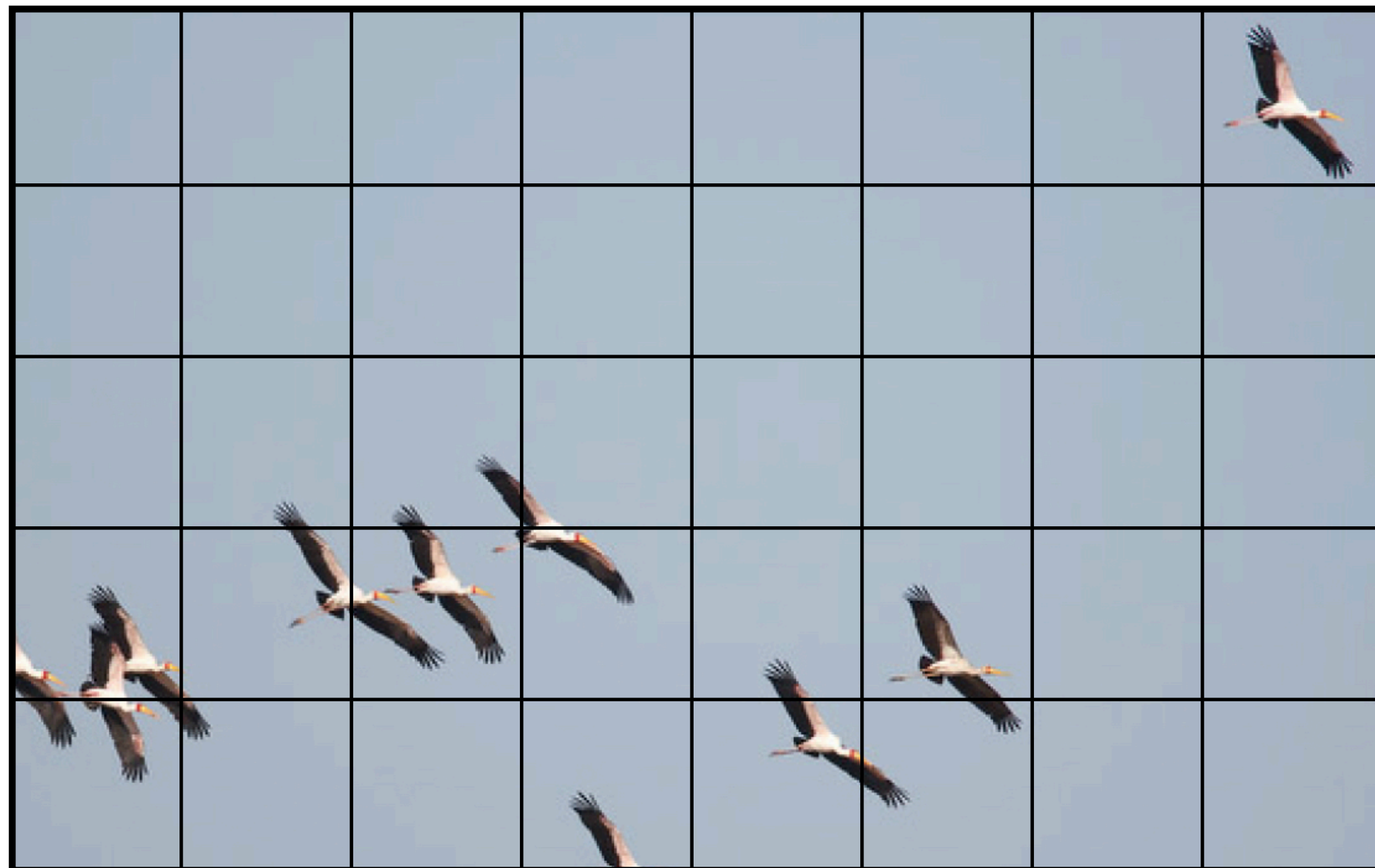




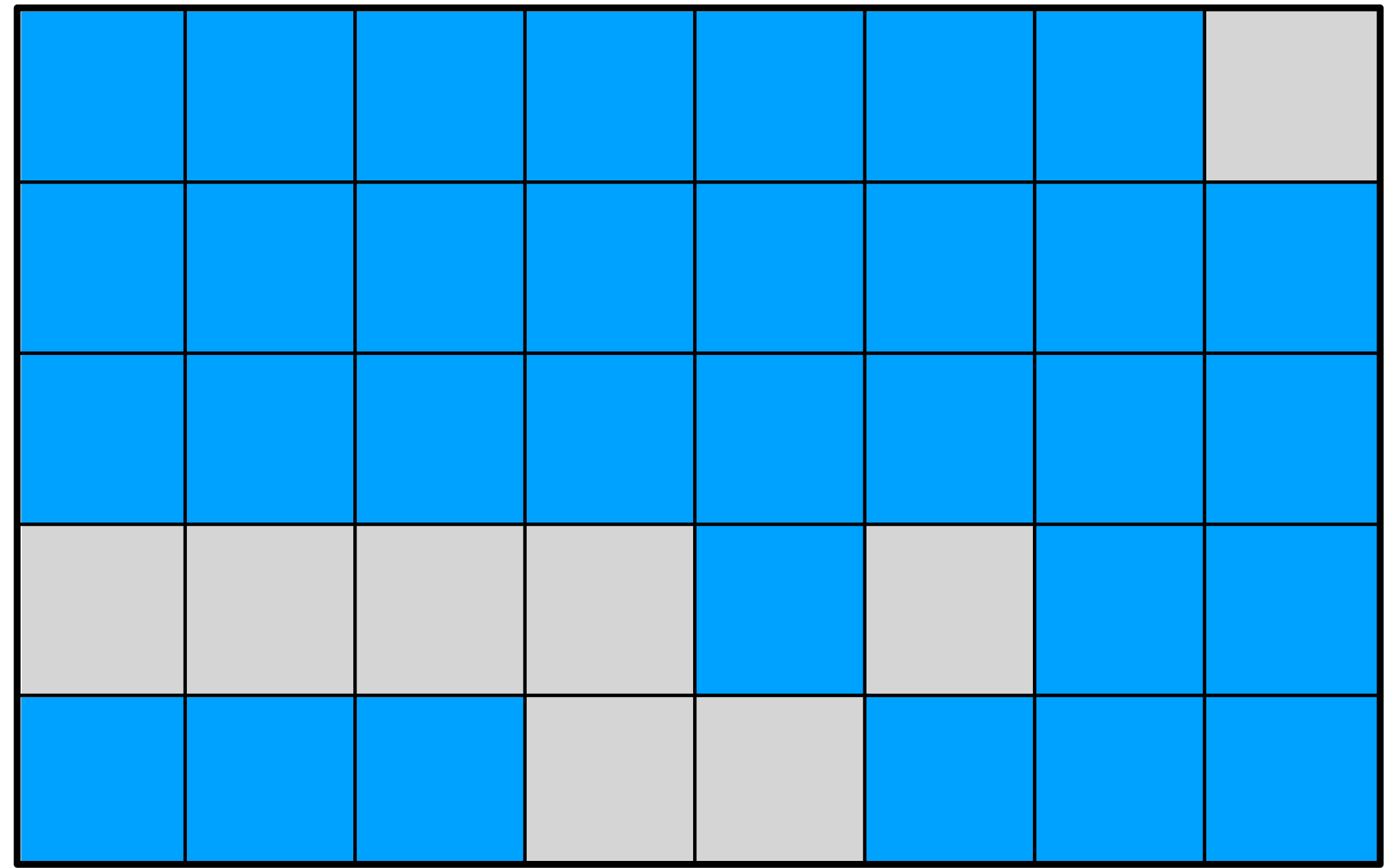
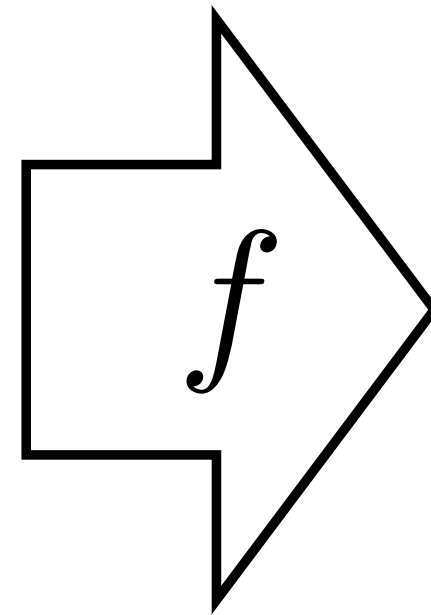
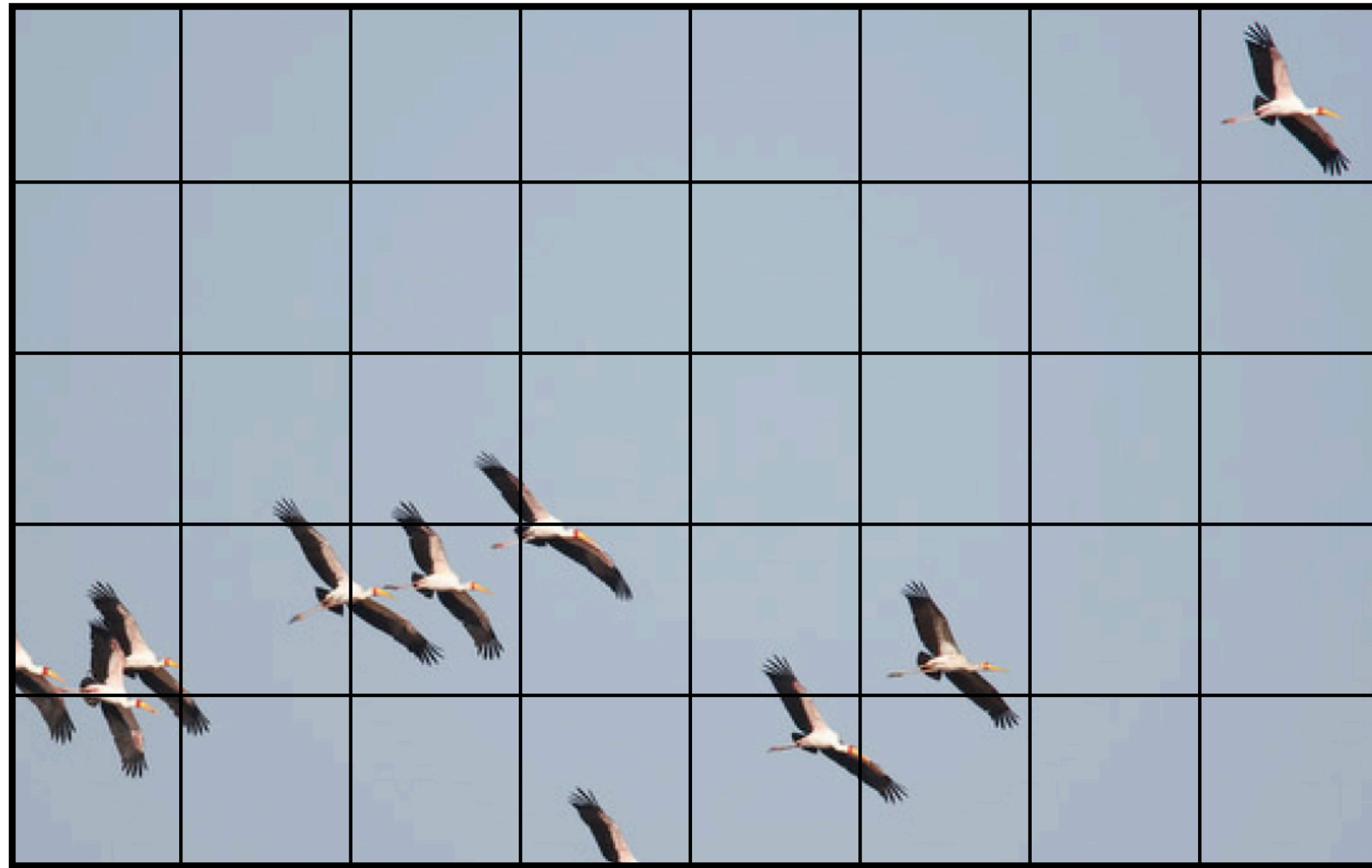


							Bird





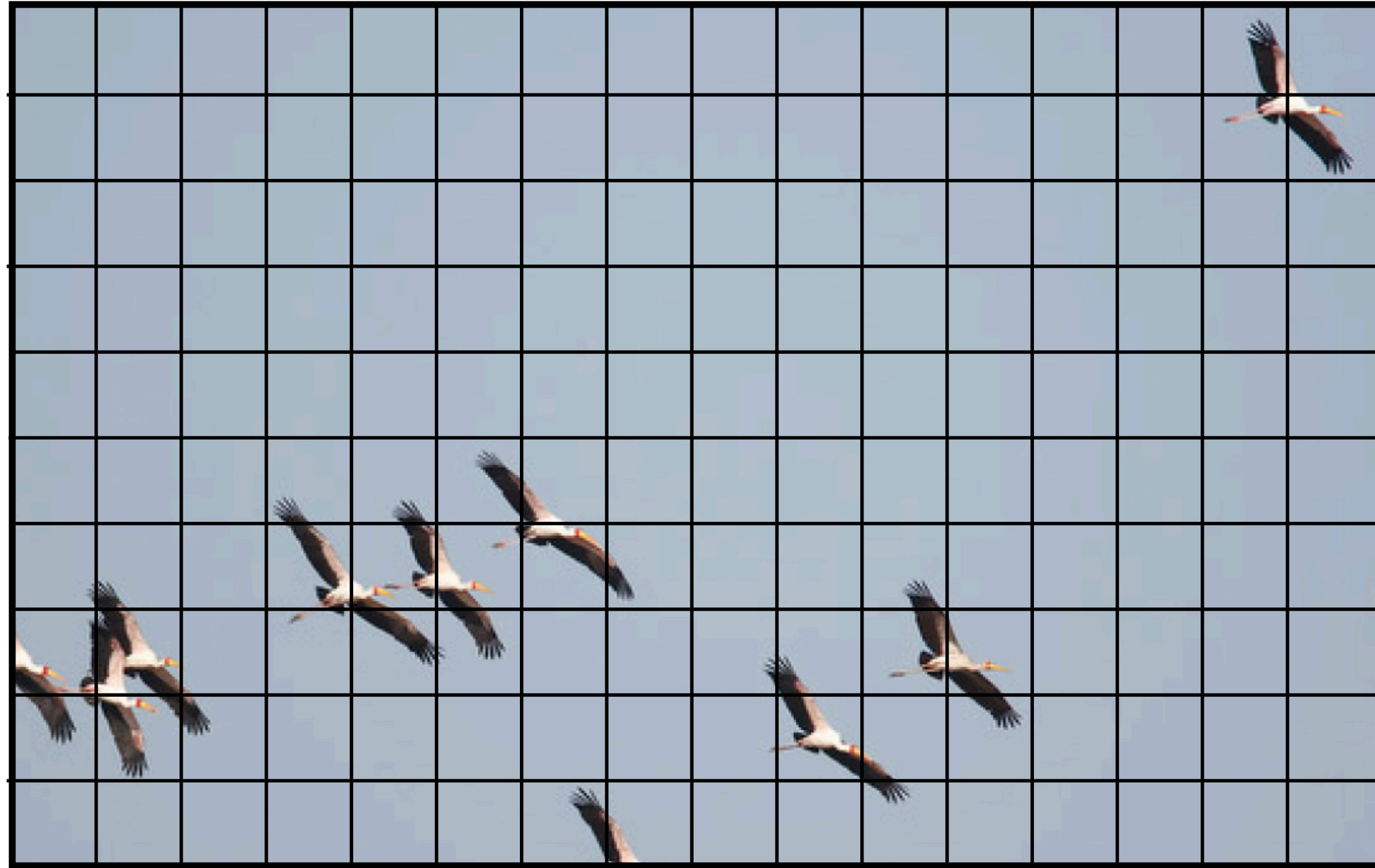
Sky	Sky	Sky	Sky	Sky	Sky	Sky	Bird
Sky	Sky	Sky	Sky	Sky	Sky	Sky	Sky
Sky	Sky	Sky	Sky	Sky	Sky	Sky	Sky
Bird	Bird	Bird	Sky	Bird	Sky	Sky	Sky
Sky	Sky	Sky	Bird	Sky	Sky	Sky	Sky



Problem:

What happens to objects that are bigger?

What if an object crosses multiple cells?



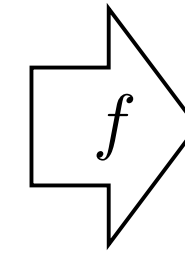
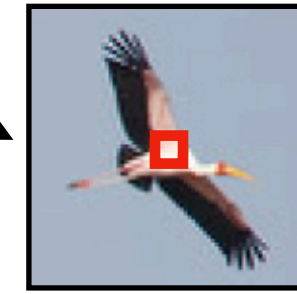
“Cell”-based approach is limited.

What can we do instead?

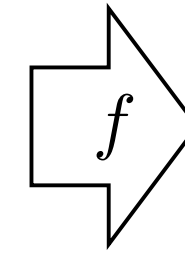
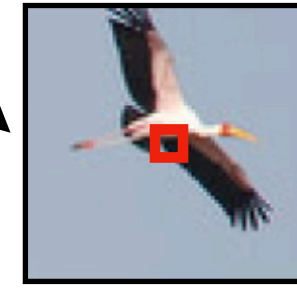




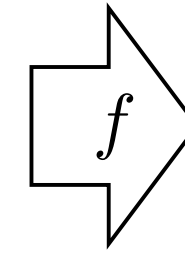
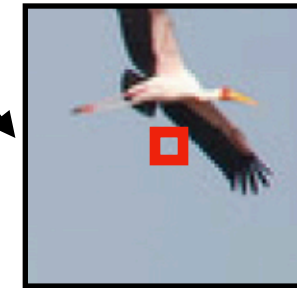
What's the object class of the center pixel?



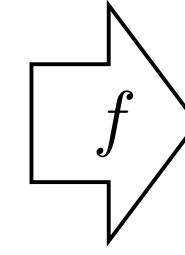
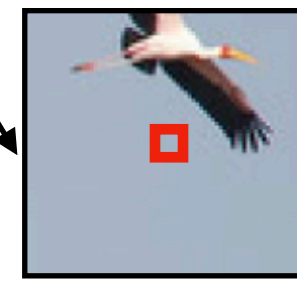
“Bird”



“Bird”

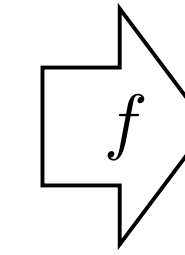
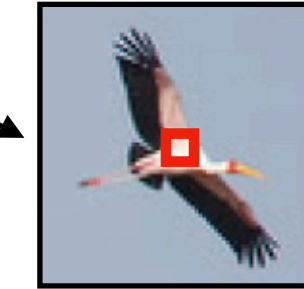
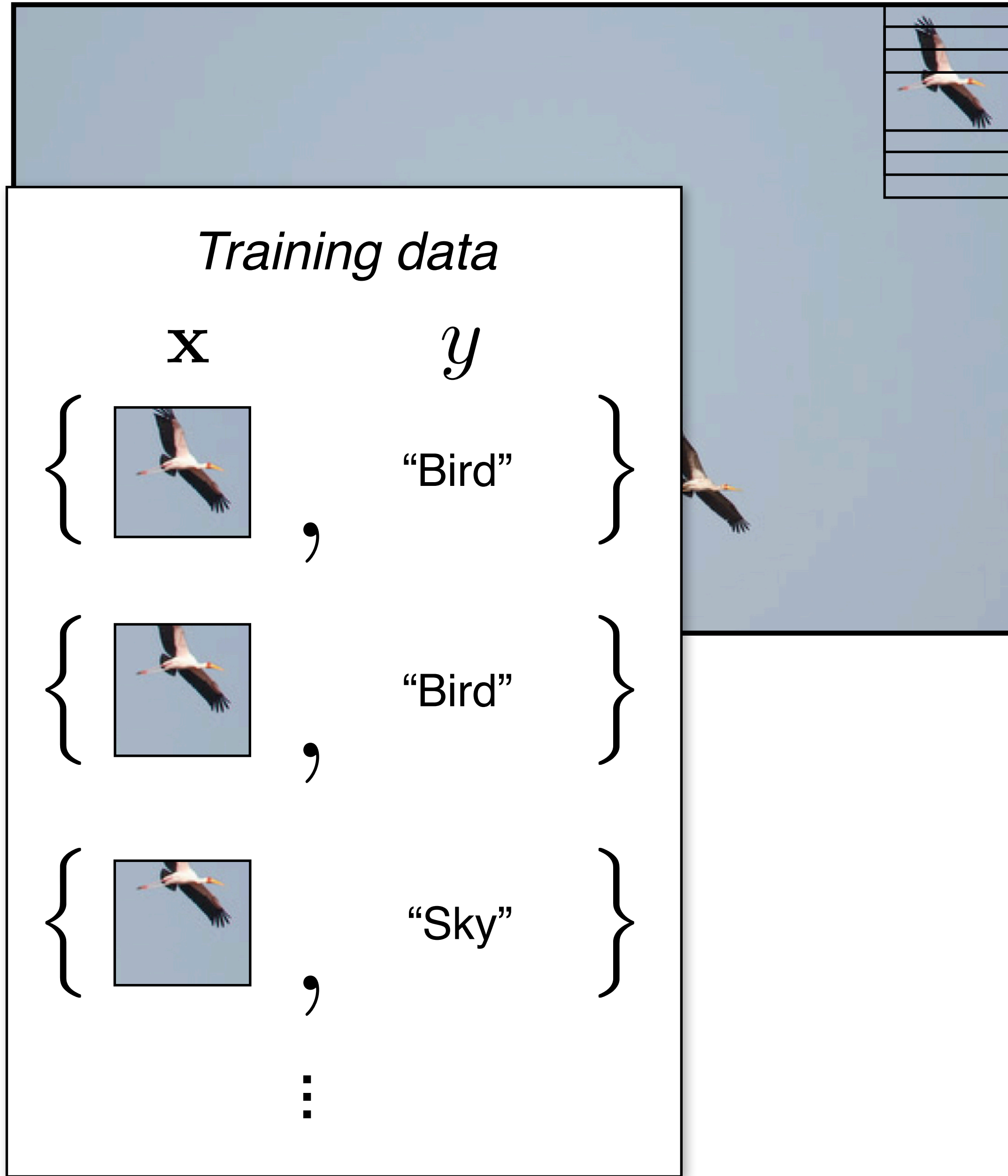


“Sky”

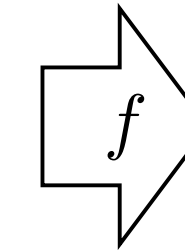
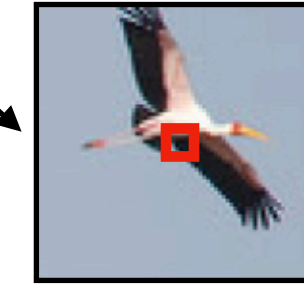


“Sky”

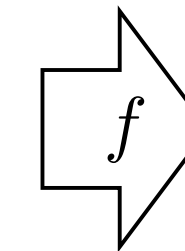
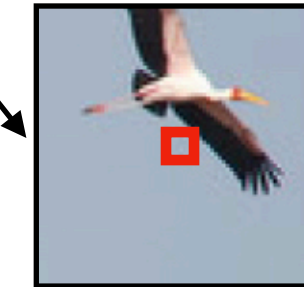
What's the object class of the center pixel?



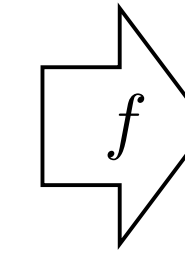
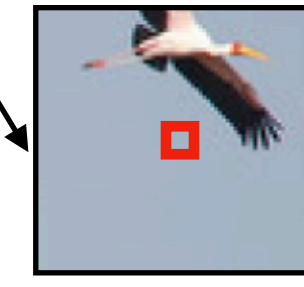
“Bird”



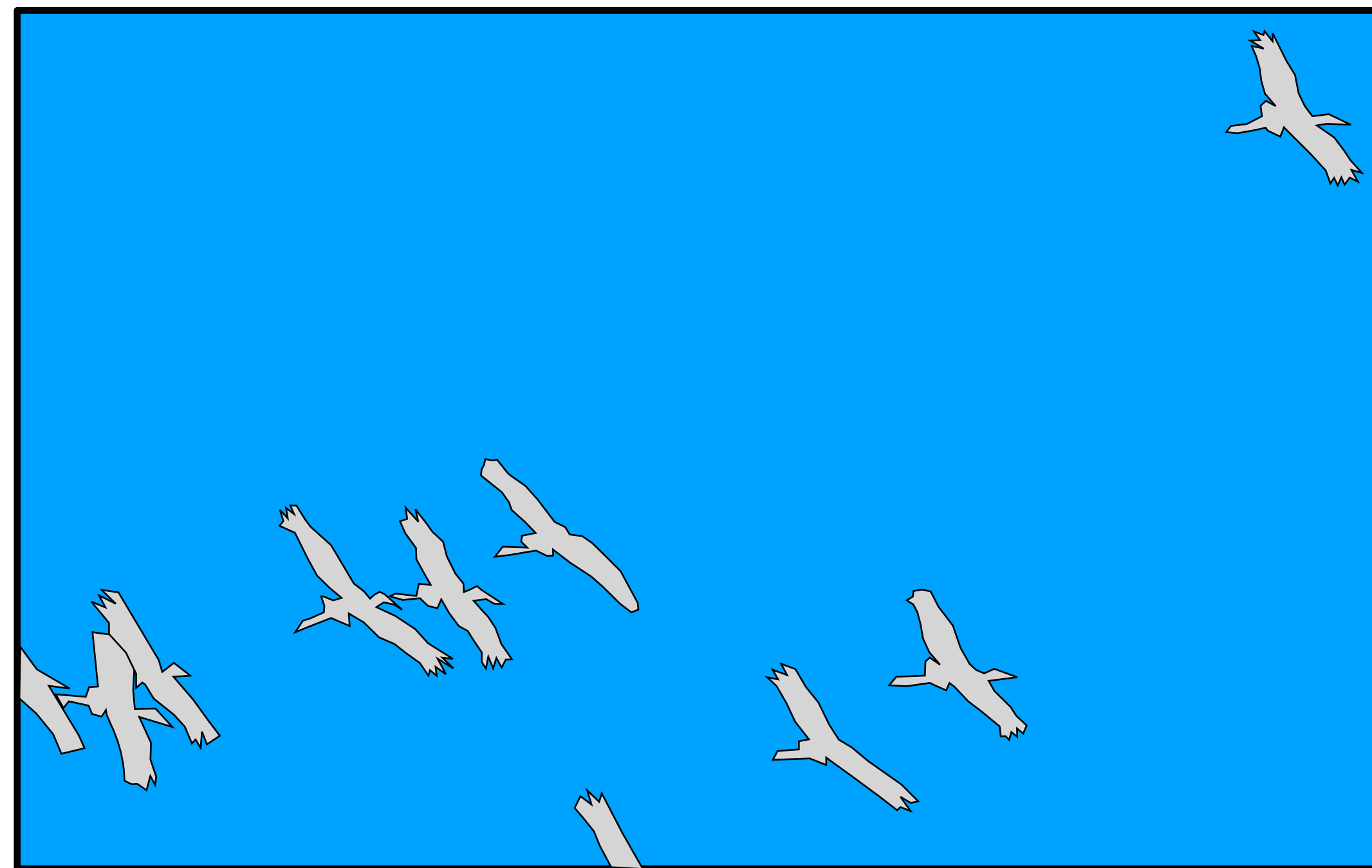
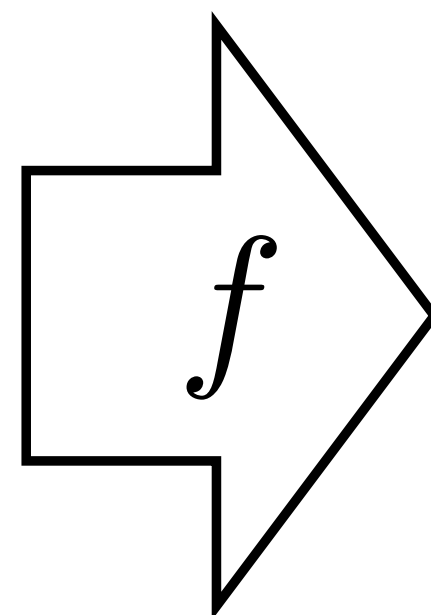
“Bird”



“Sky”

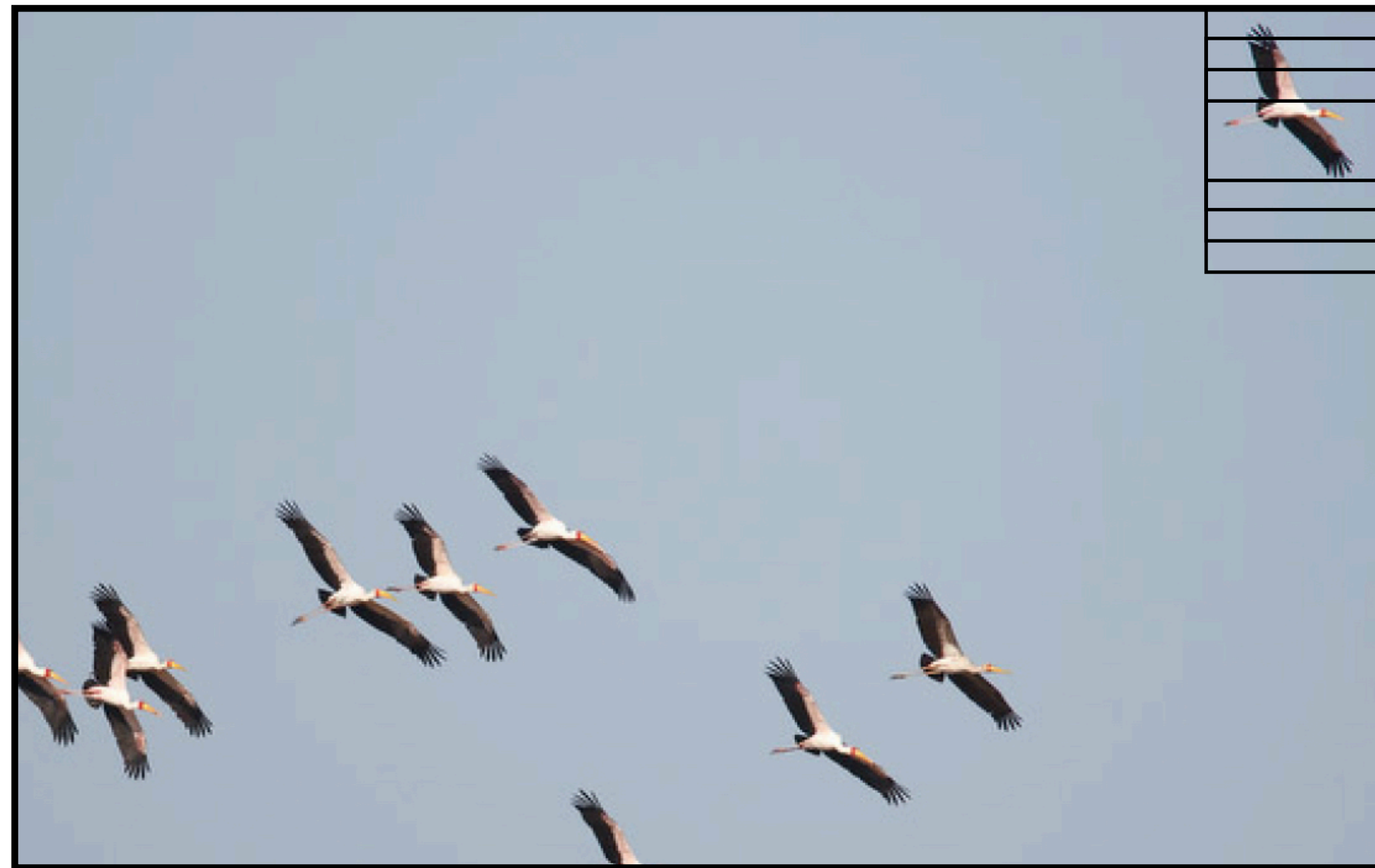


“Sky”

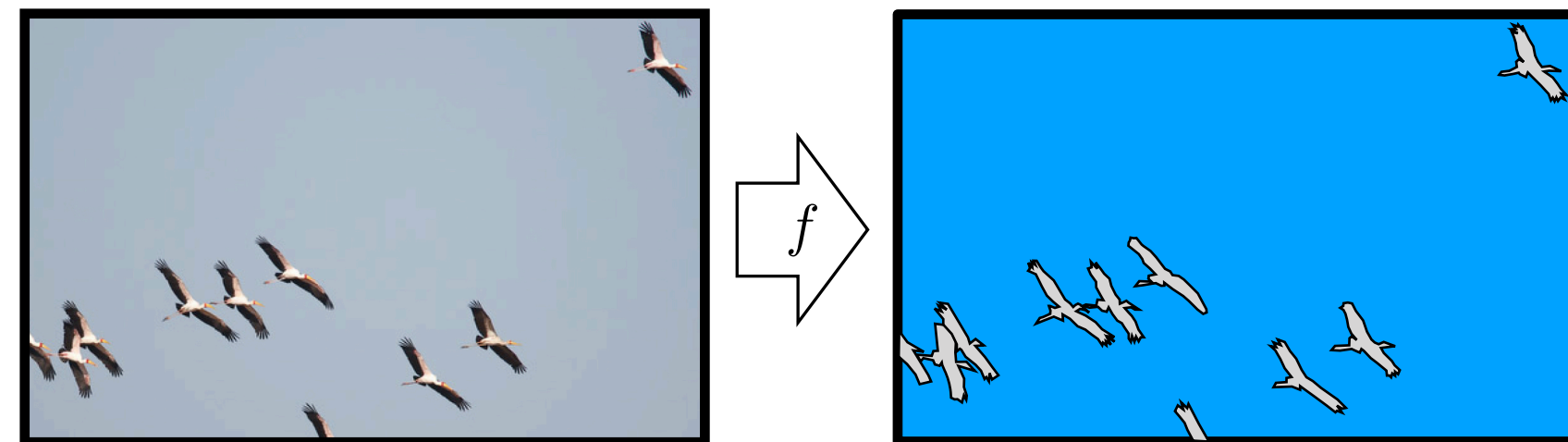
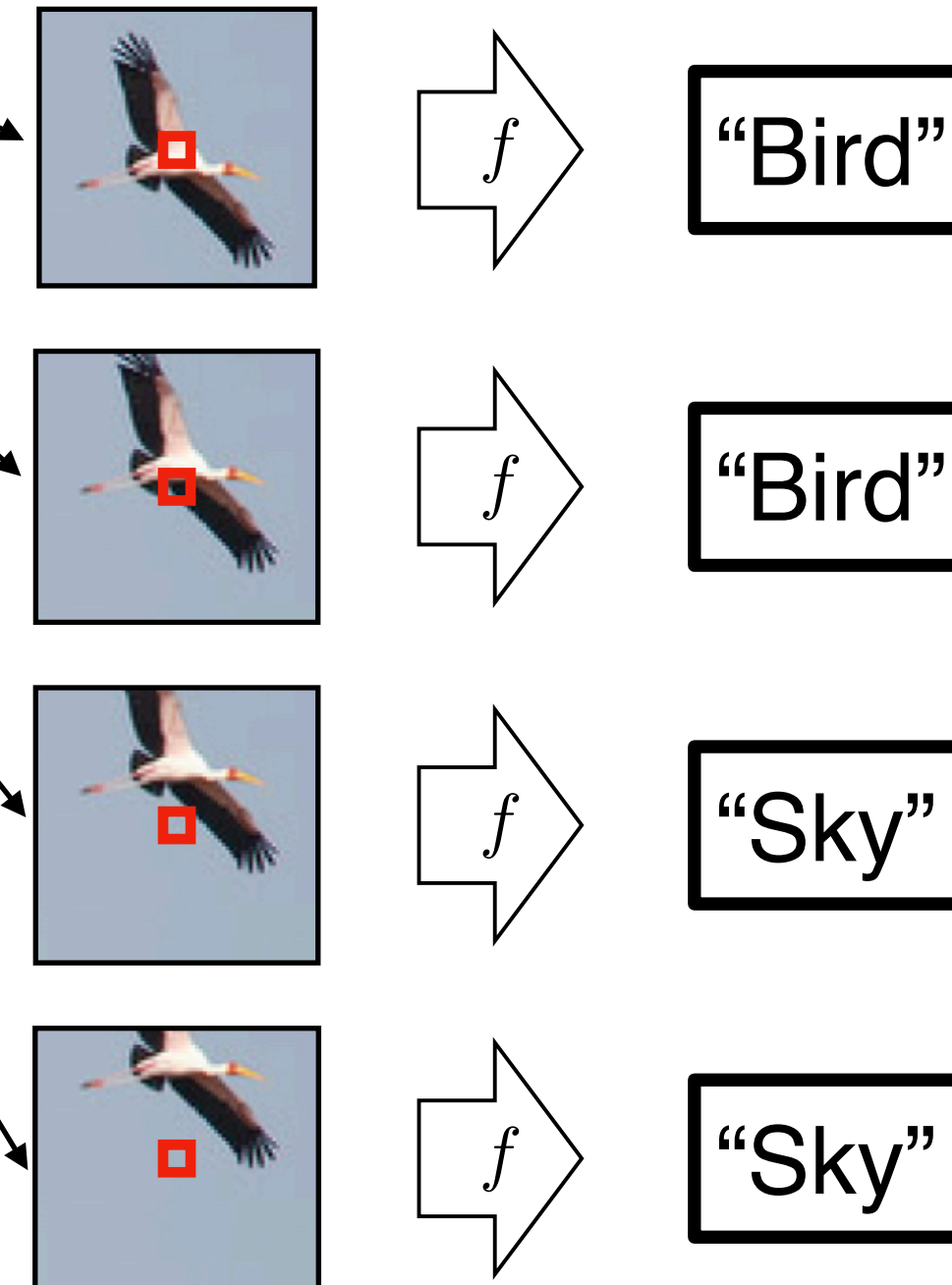


(Colors represent one-hot codes)

This problem is called **semantic segmentation**



What's the object class of the center pixel?

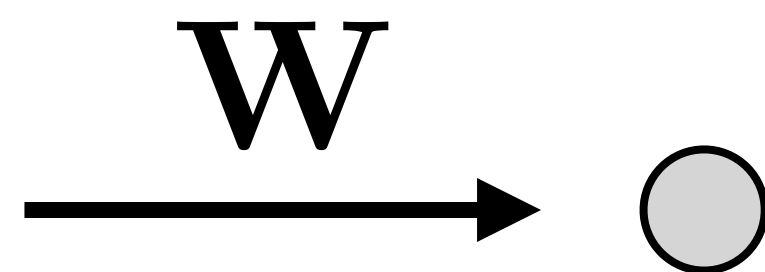
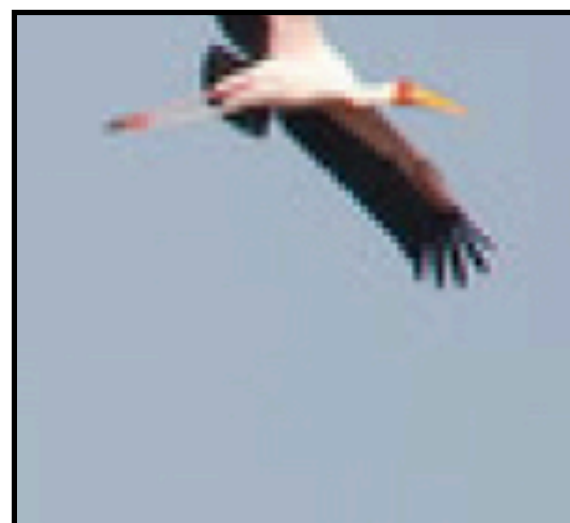
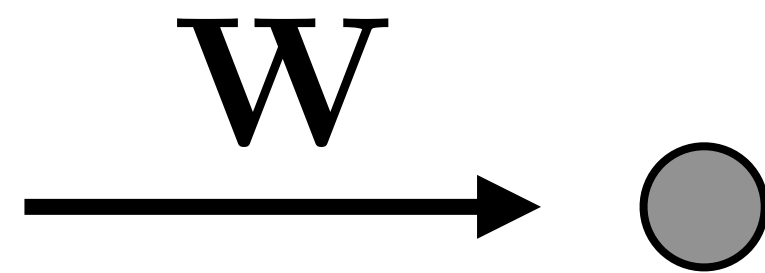
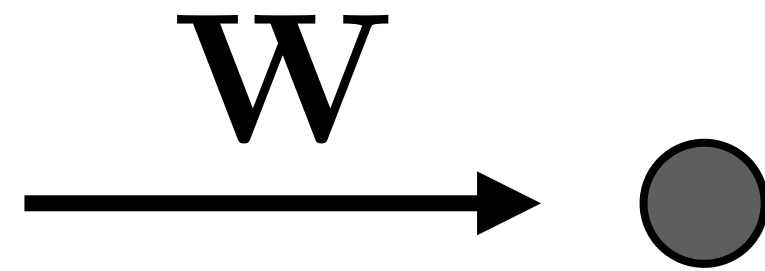
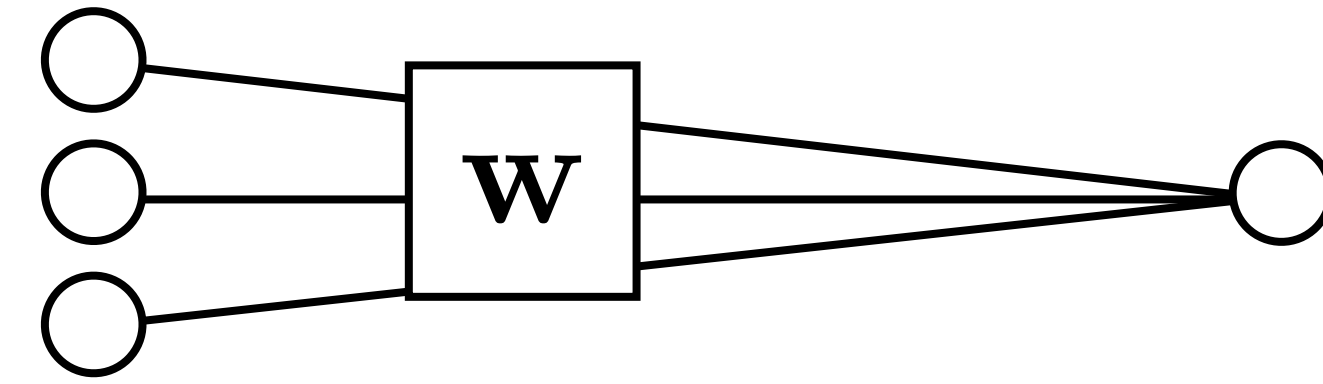


Translation invariance: process each patch in the same way.

An *equivariant* mapping:

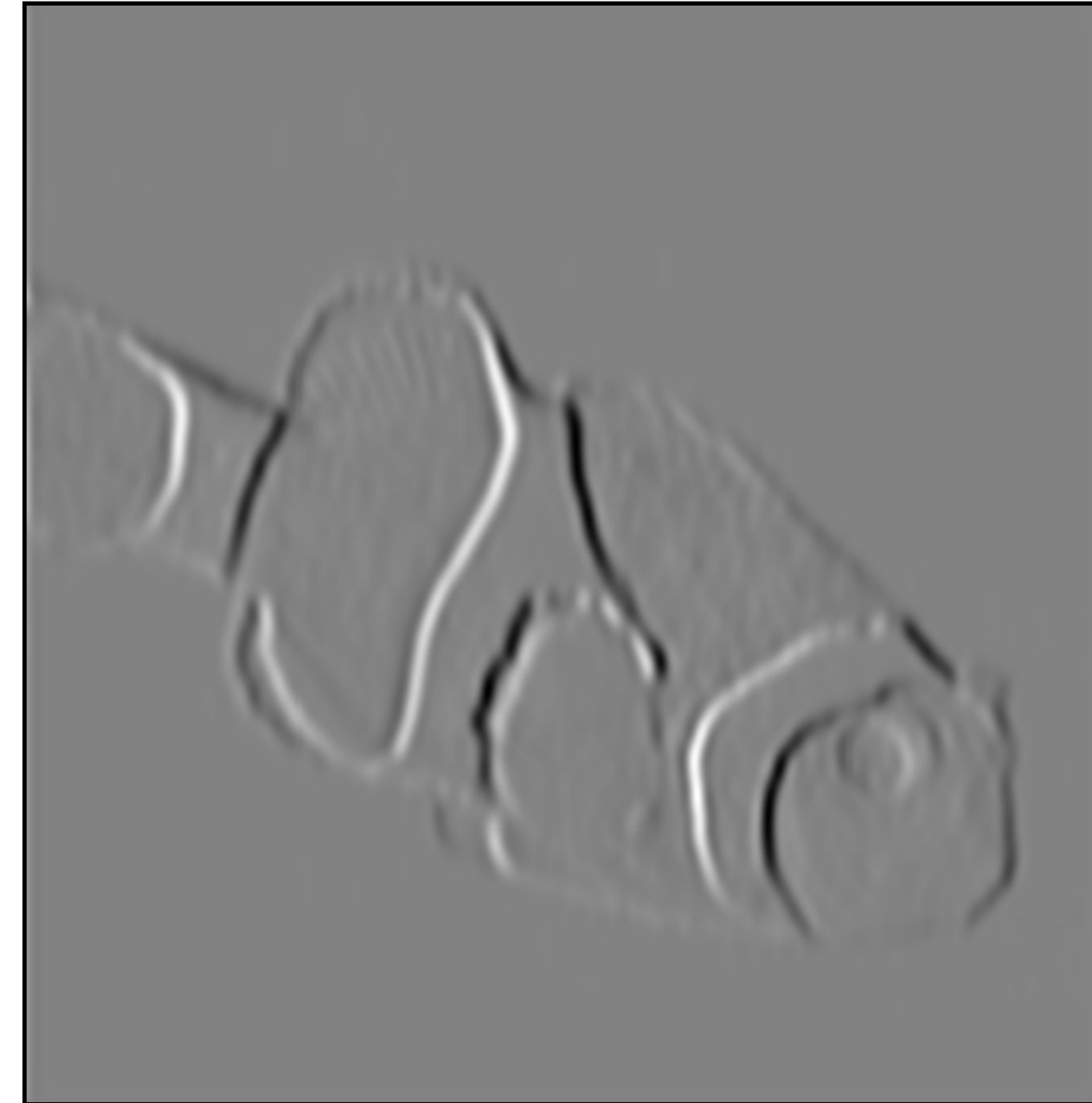
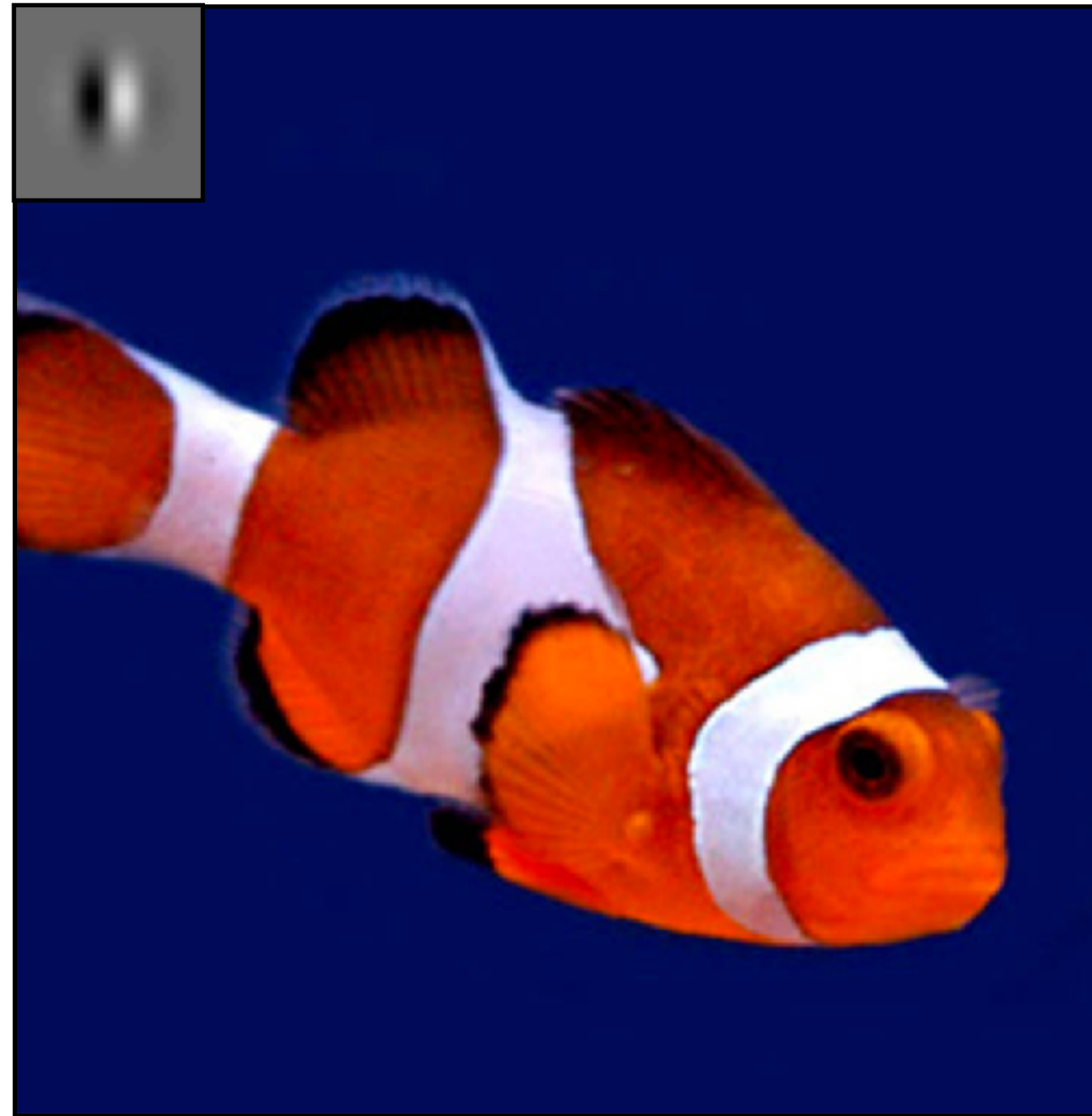
$$f(\text{translate}(x)) = \text{translate}(f(x))$$

W computes a weighted sum of all pixels in the patch

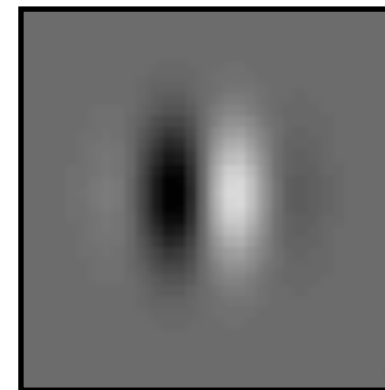


W is a convolutional kernel applied to the full image!

Convolution

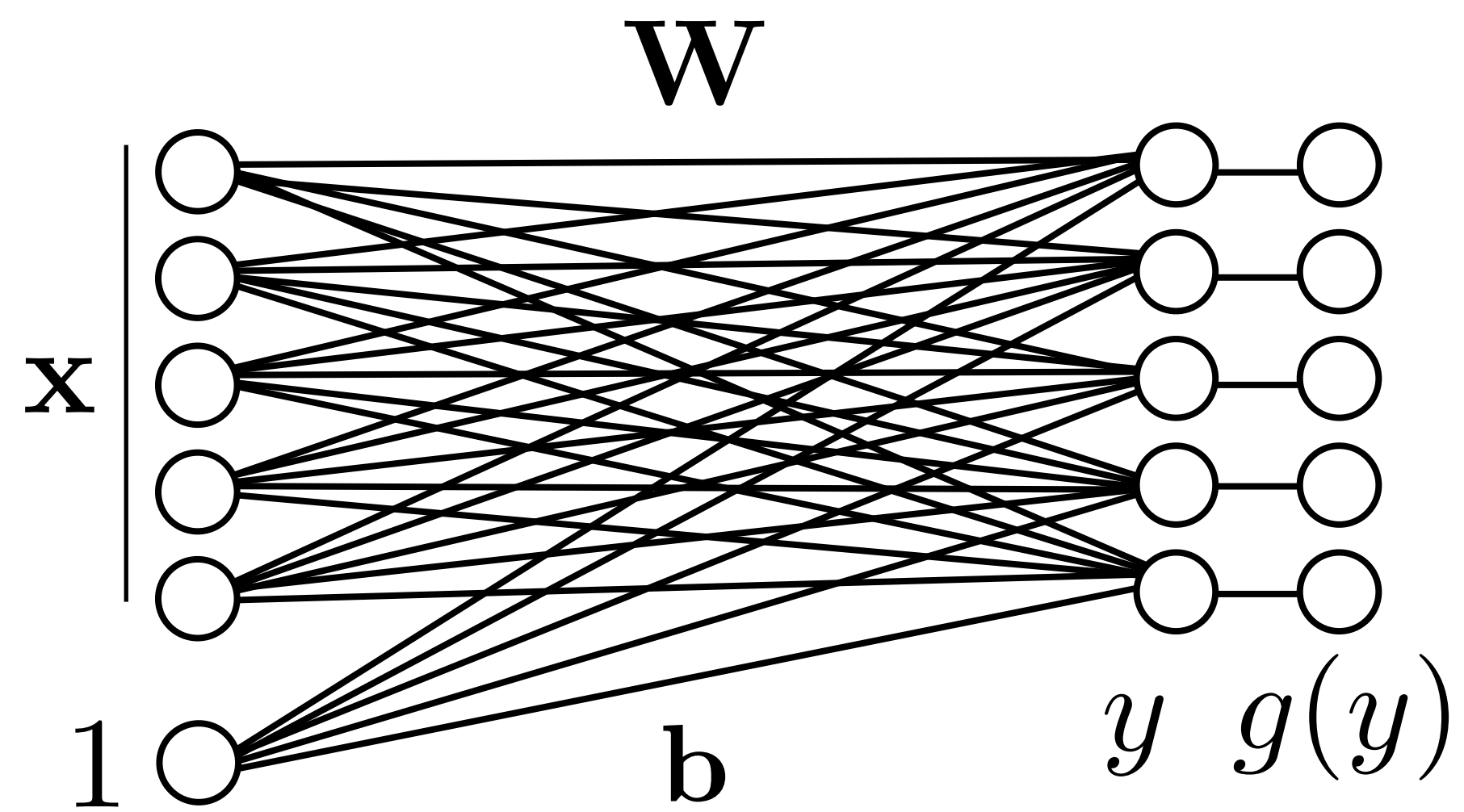


filter

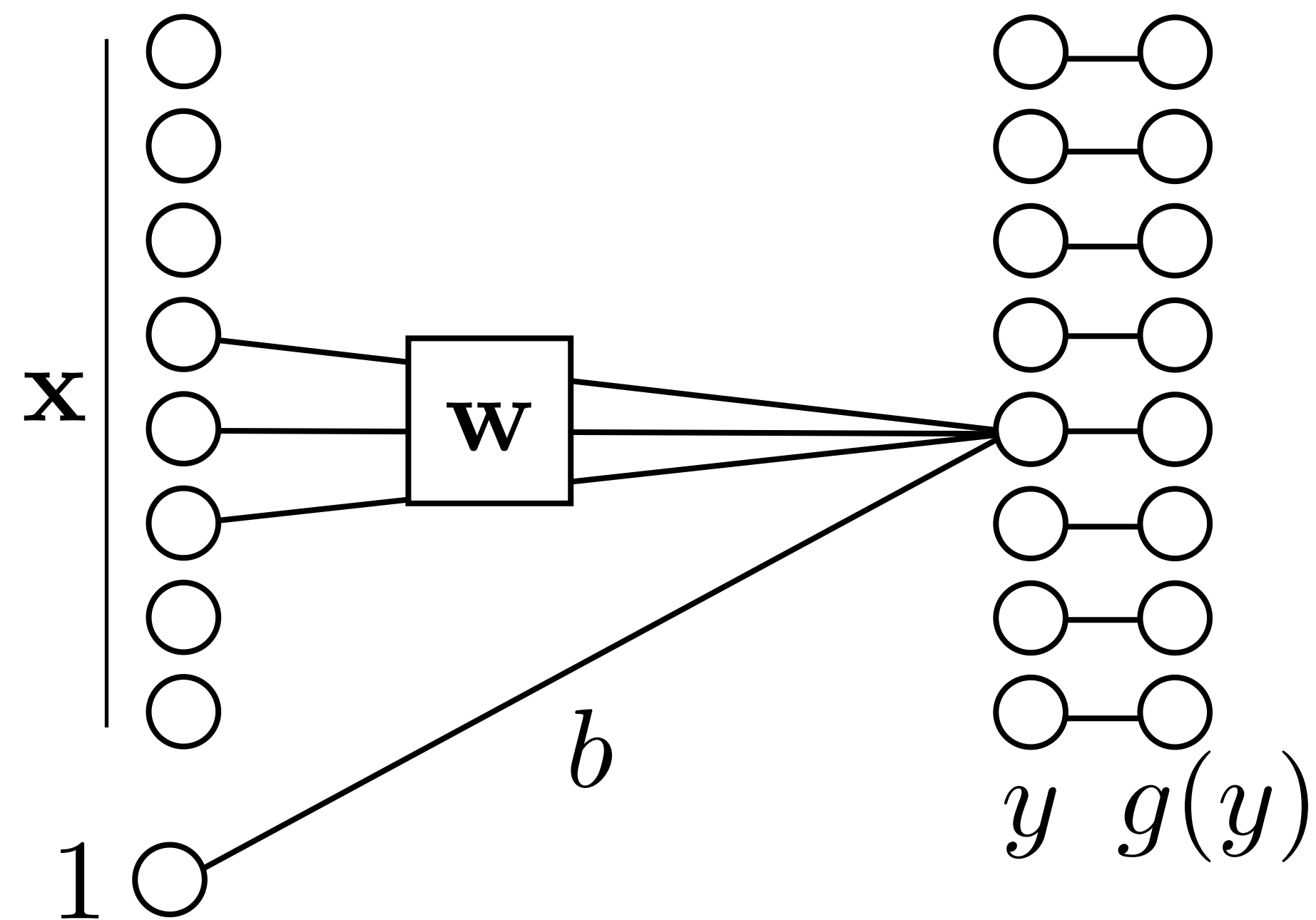


Fully-connected network

Fully-connected (fc) layer



Locally connected network

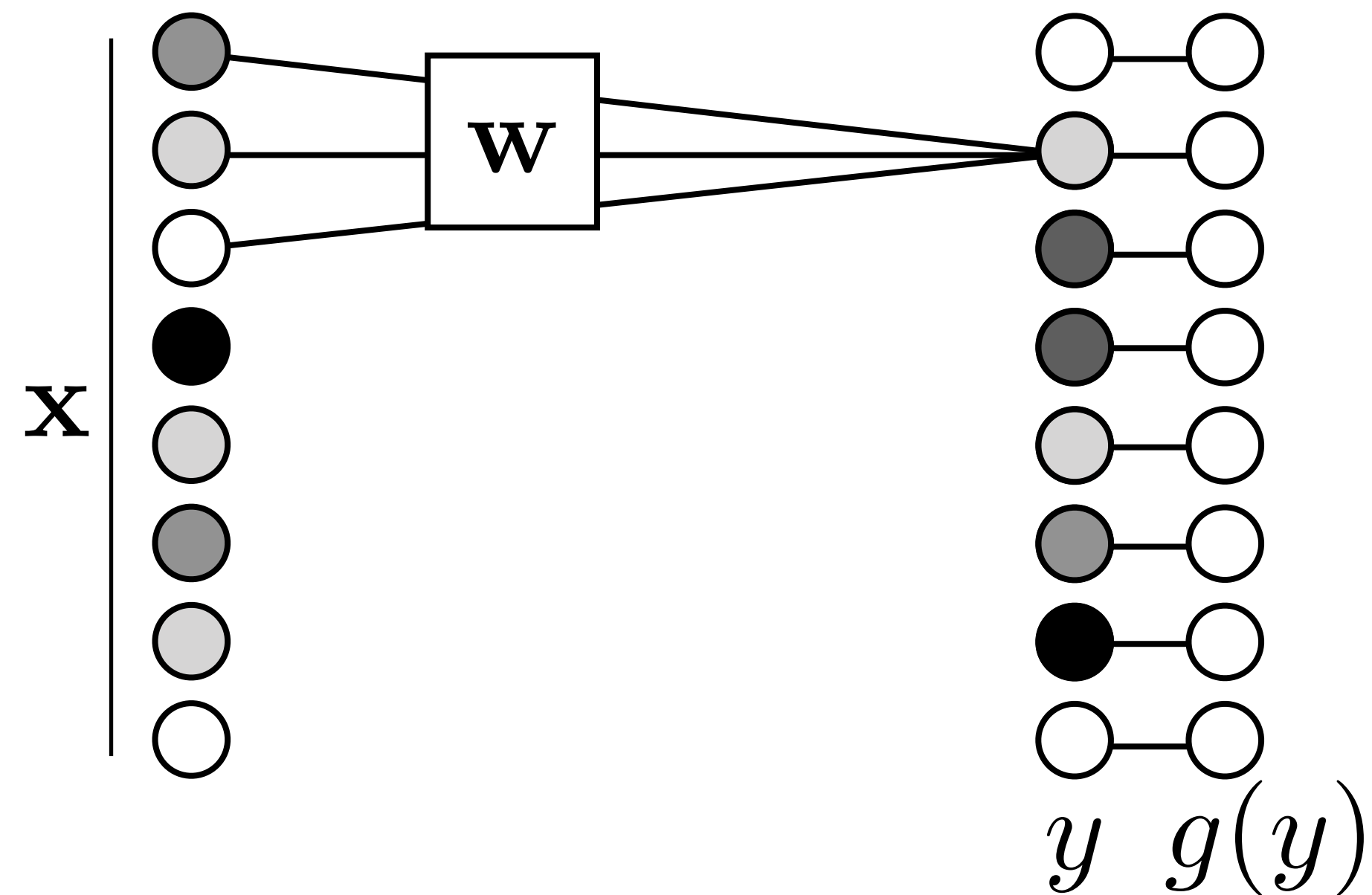


Often, we assume output is a **local** function of input.

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

Convolutional neural network

Conv layer

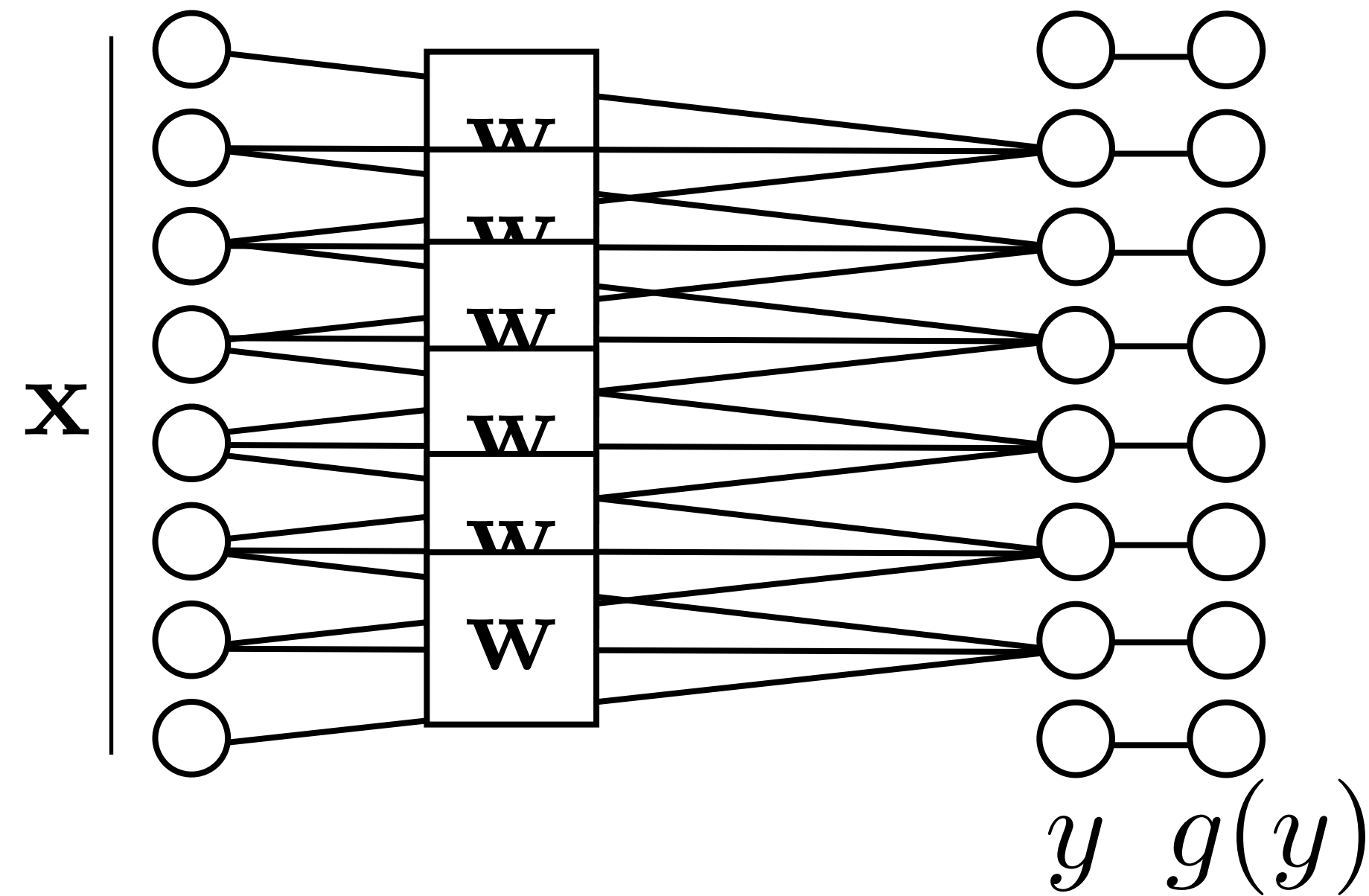


Often, we assume output is a **local** function of input.

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

Weight sharing

Conv layer

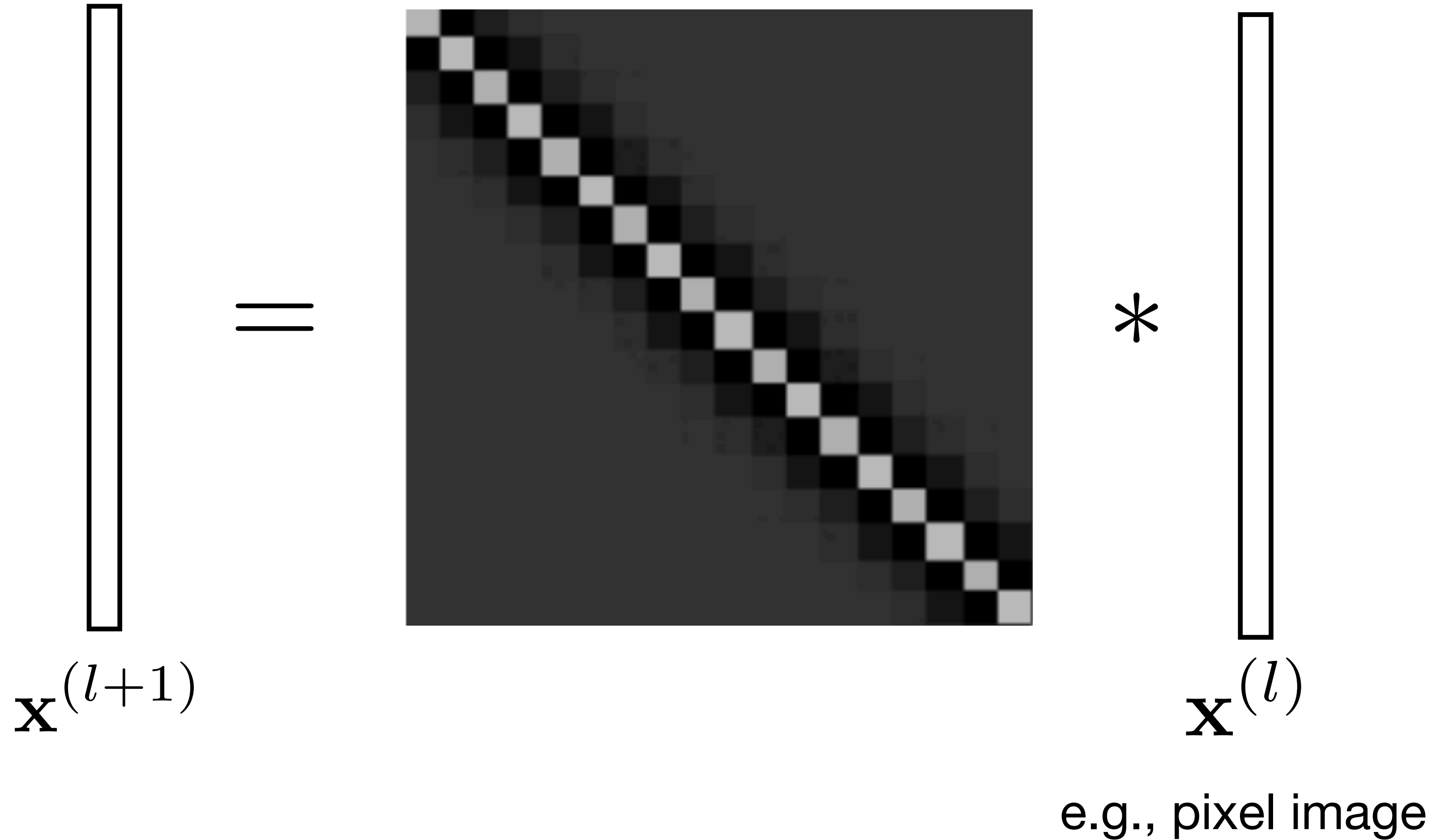


Often, we assume output is a **local** function of input.

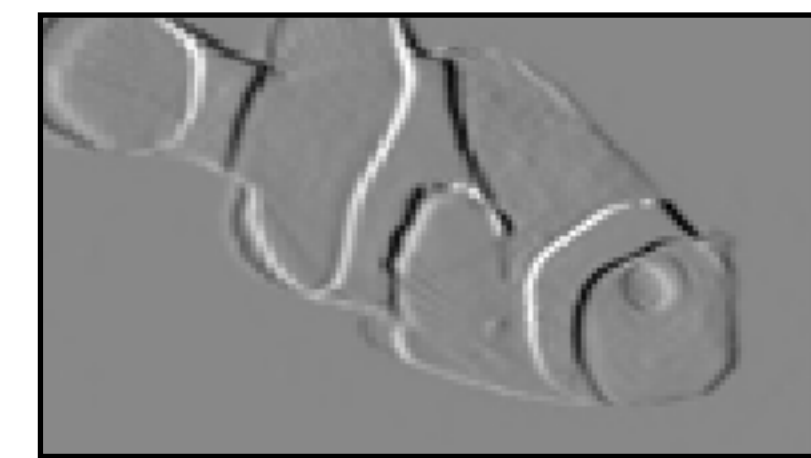
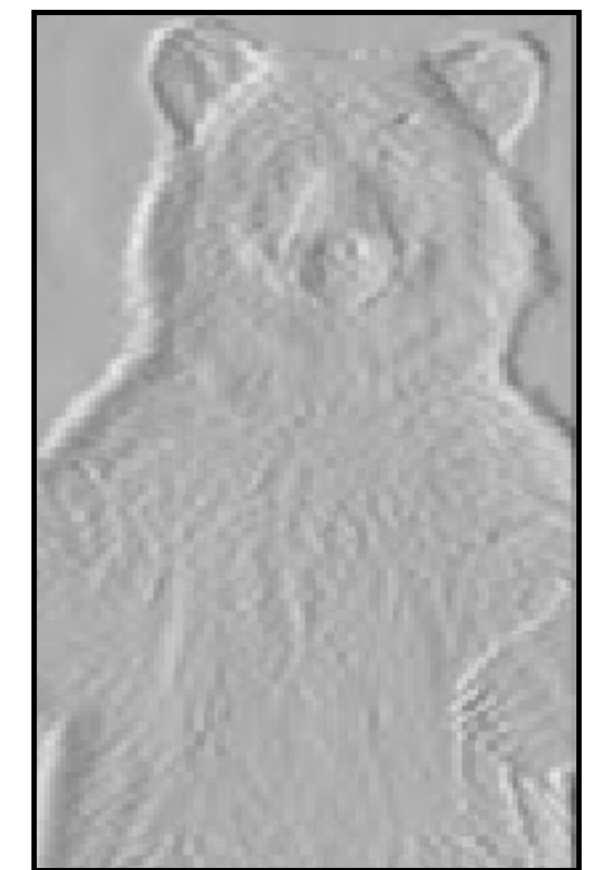
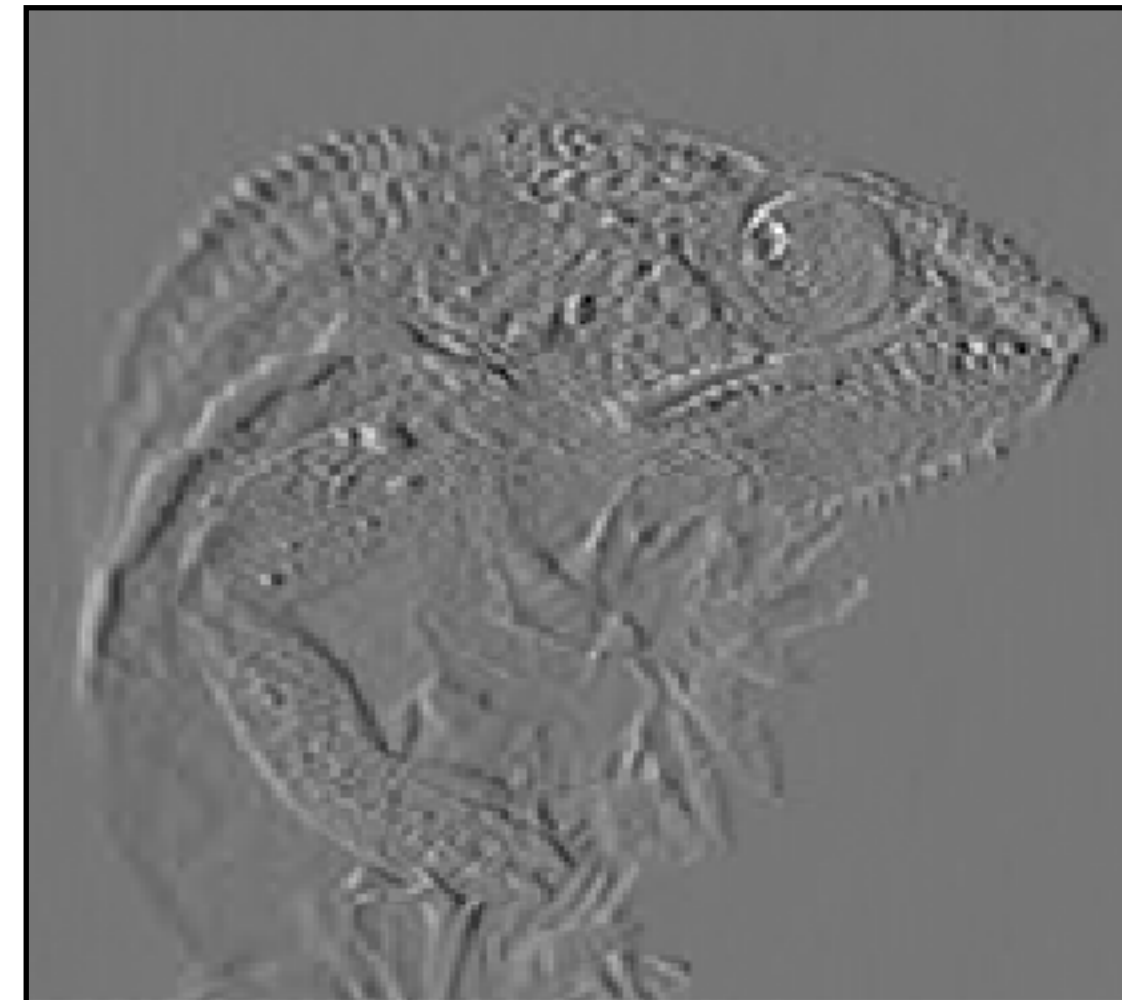
If we use the same weights (**weight sharing**) to compute each local function, we get a convolutional neural network.

Toeplitz matrix

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



- Constrained linear layer
- Fewer parameters \rightarrow easier to learn, less overfitting



Conv layers can be applied to arbitrarily-sized inputs

$$\mathbf{X}^{(l+1)} = \text{matrix} * \mathbf{X}^{(l)}$$

The diagram illustrates a matrix multiplication operation. On the left is a vertical rectangle representing the vector $\mathbf{X}^{(l+1)}$. In the center is a square matrix with a dark gray background and a light gray diagonal, representing the weight matrix. On the right is another vertical rectangle representing the vector $\mathbf{X}^{(l)}$. An equals sign (=) is positioned between the first vector and the matrix, and an asterisk (*) is positioned between the matrix and the second vector.

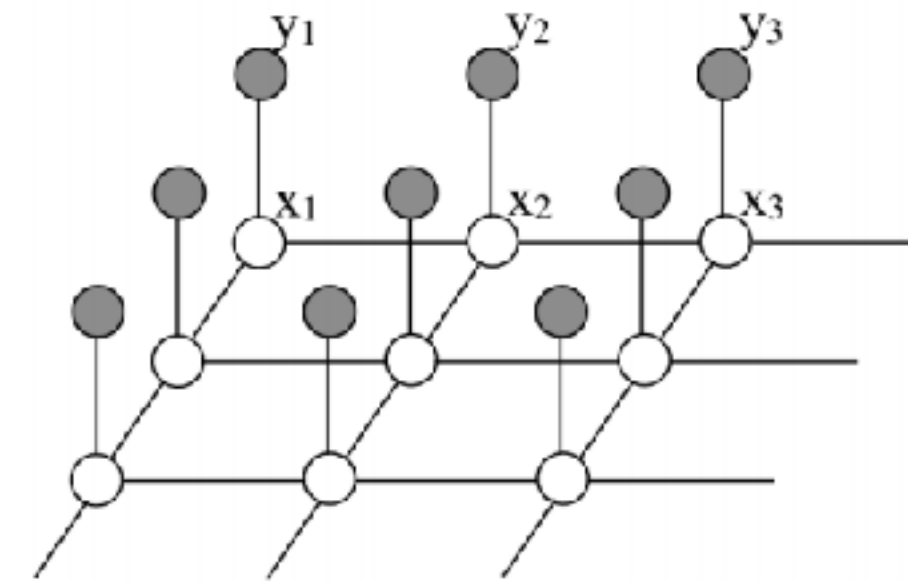
$$\mathbf{X}^{(l+1)} = \text{diag}(\mathbf{X}^{(l)}) * \mathbf{X}^{(l)}$$

The diagram illustrates the element-wise multiplication of a vector $\mathbf{X}^{(l)}$ by a diagonal matrix. On the left, a vertical rectangle represents the vector $\mathbf{X}^{(l+1)}$. In the center, a square matrix with a dark gray background and a light gray diagonal represents the diagonal matrix. On the right, another vertical rectangle represents the vector $\mathbf{X}^{(l)}$. An equals sign is placed between the vector and the matrix, and an asterisk is placed between the matrix and the vector.

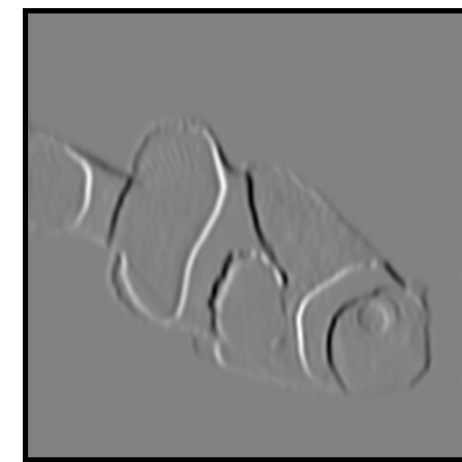
Five views on convolutional layers

1. Equivariant with translation (stationarity) $f(\text{translate}(x)) = \text{translate}(f(x))$

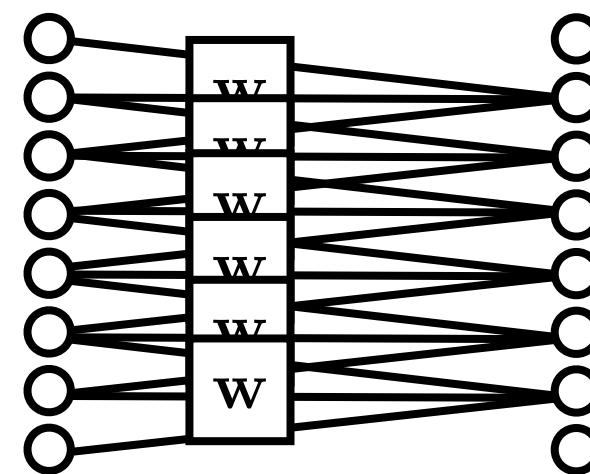
2. Patch processing (Markov assumption)



3. Image filter



4. Parameter sharing

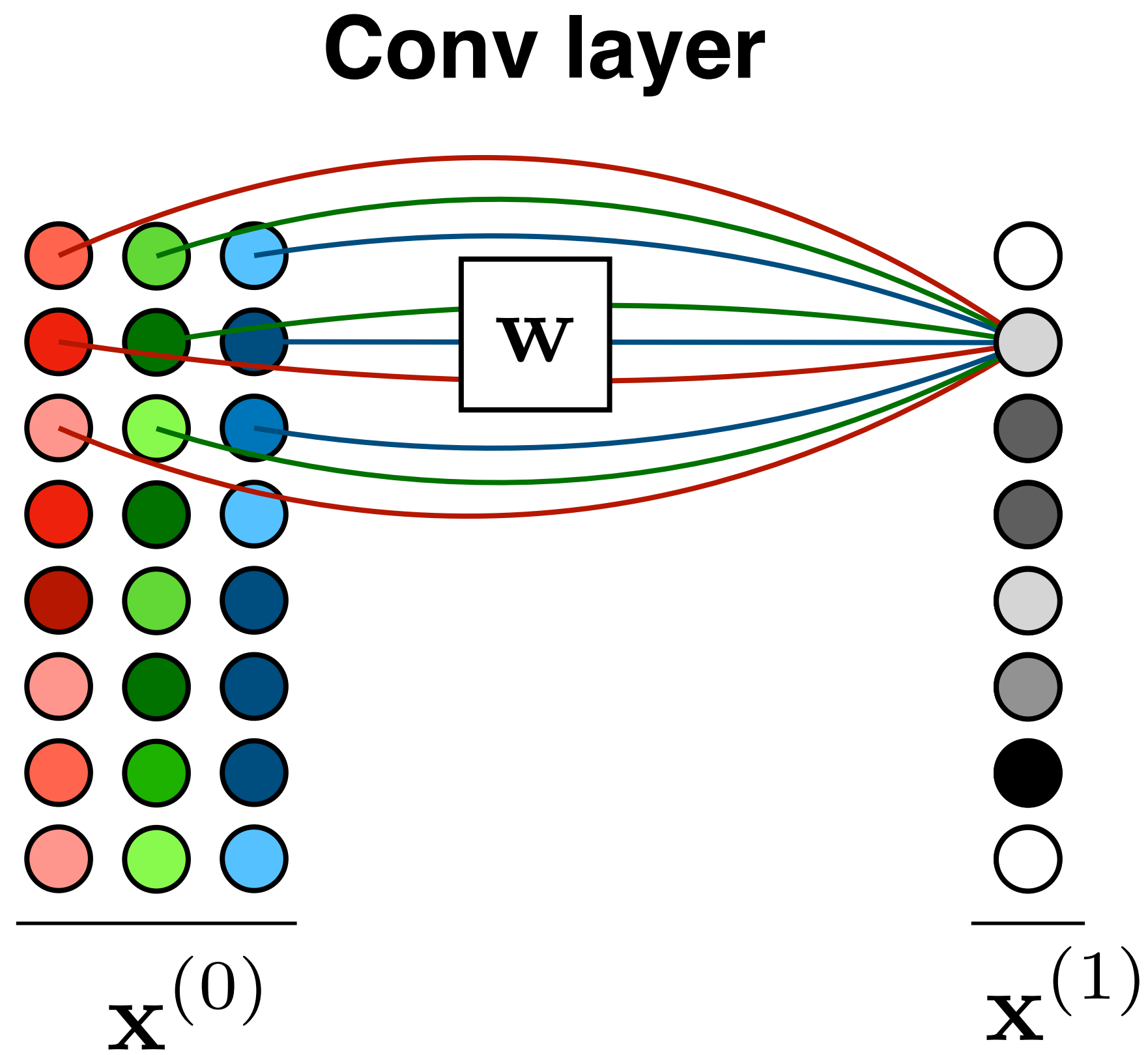


5. A way to process variable-sized tensors

What if we have color?

(aka multiple input channels?)

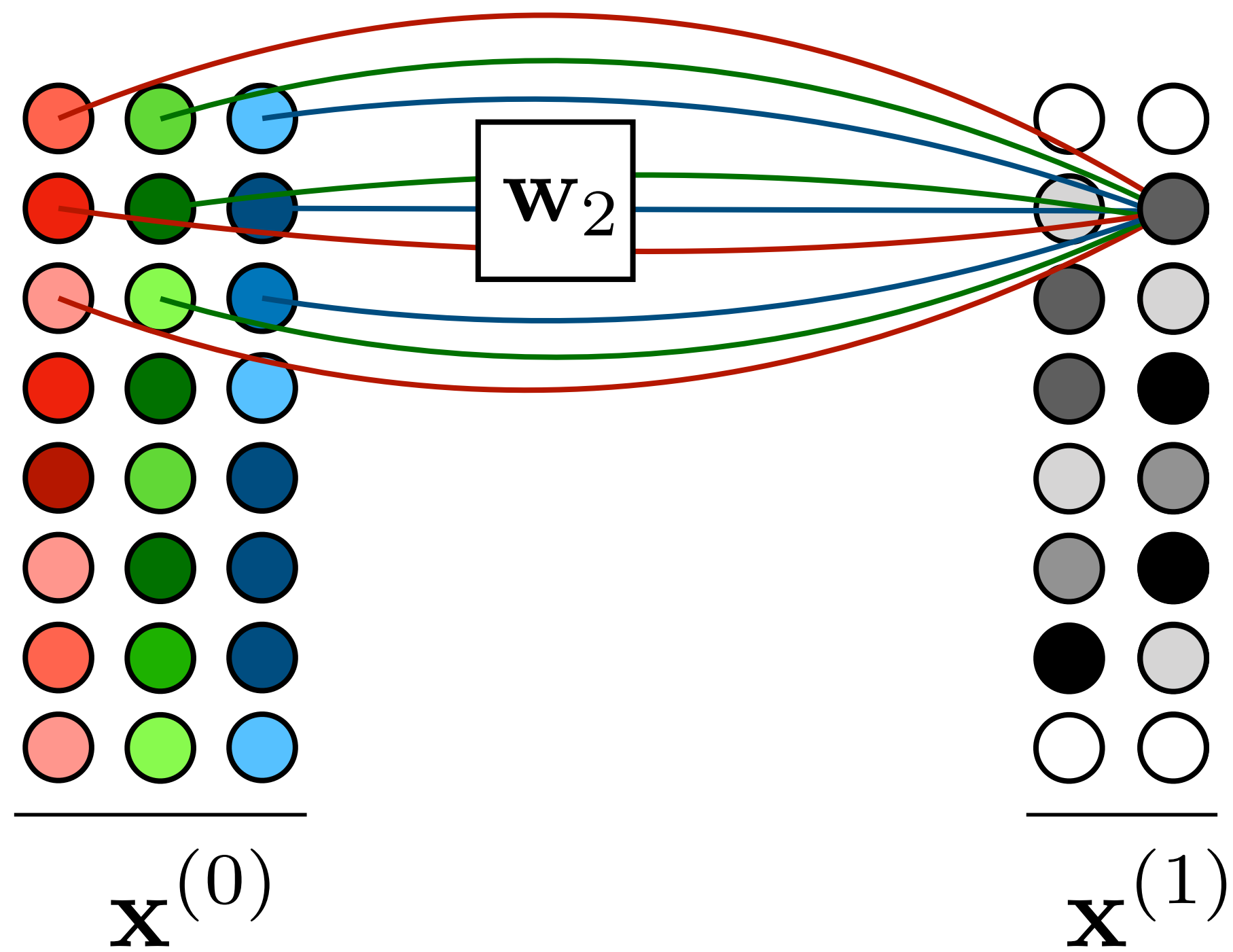
Multiple channels



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

Multiple channels

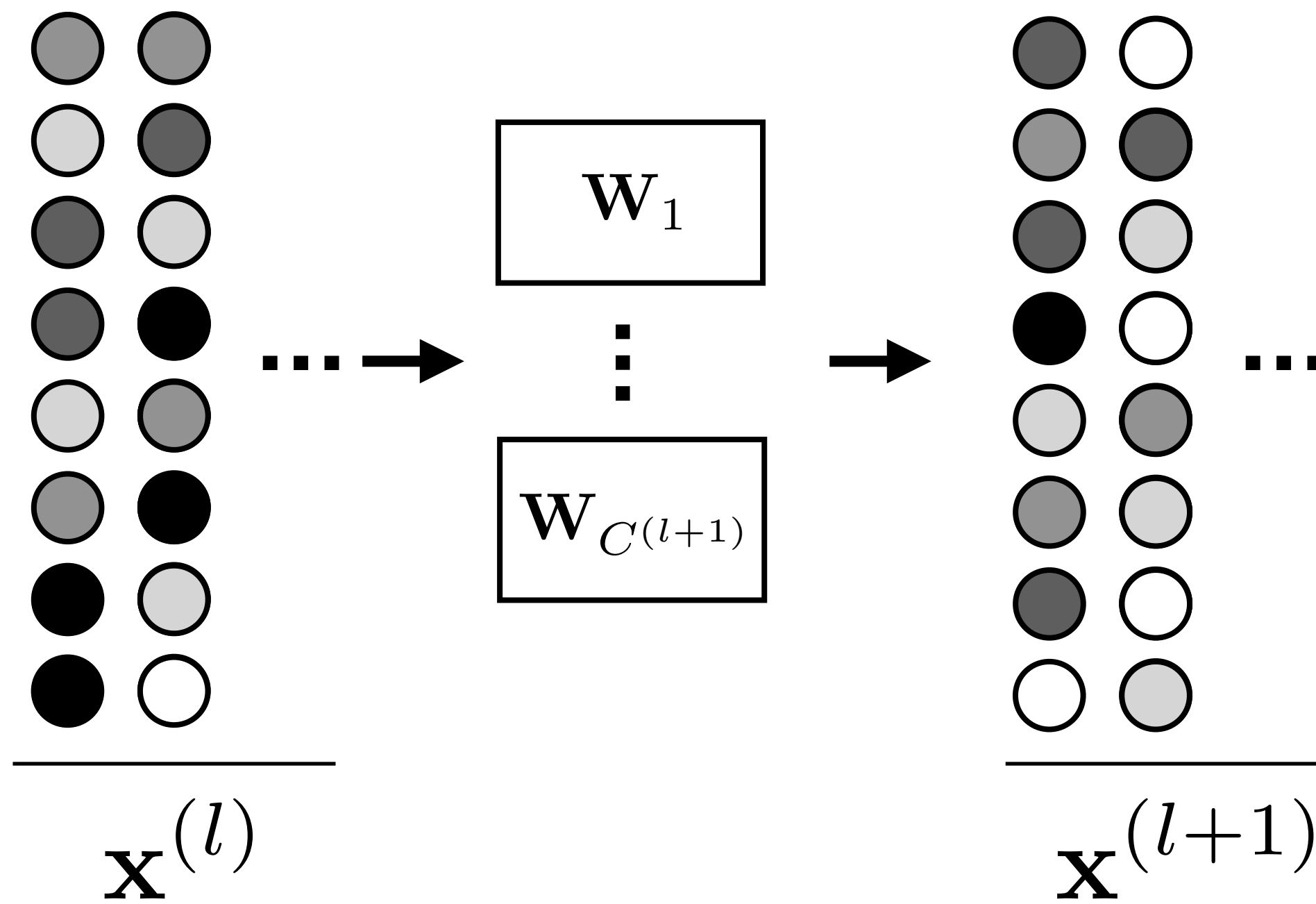
Conv layer



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

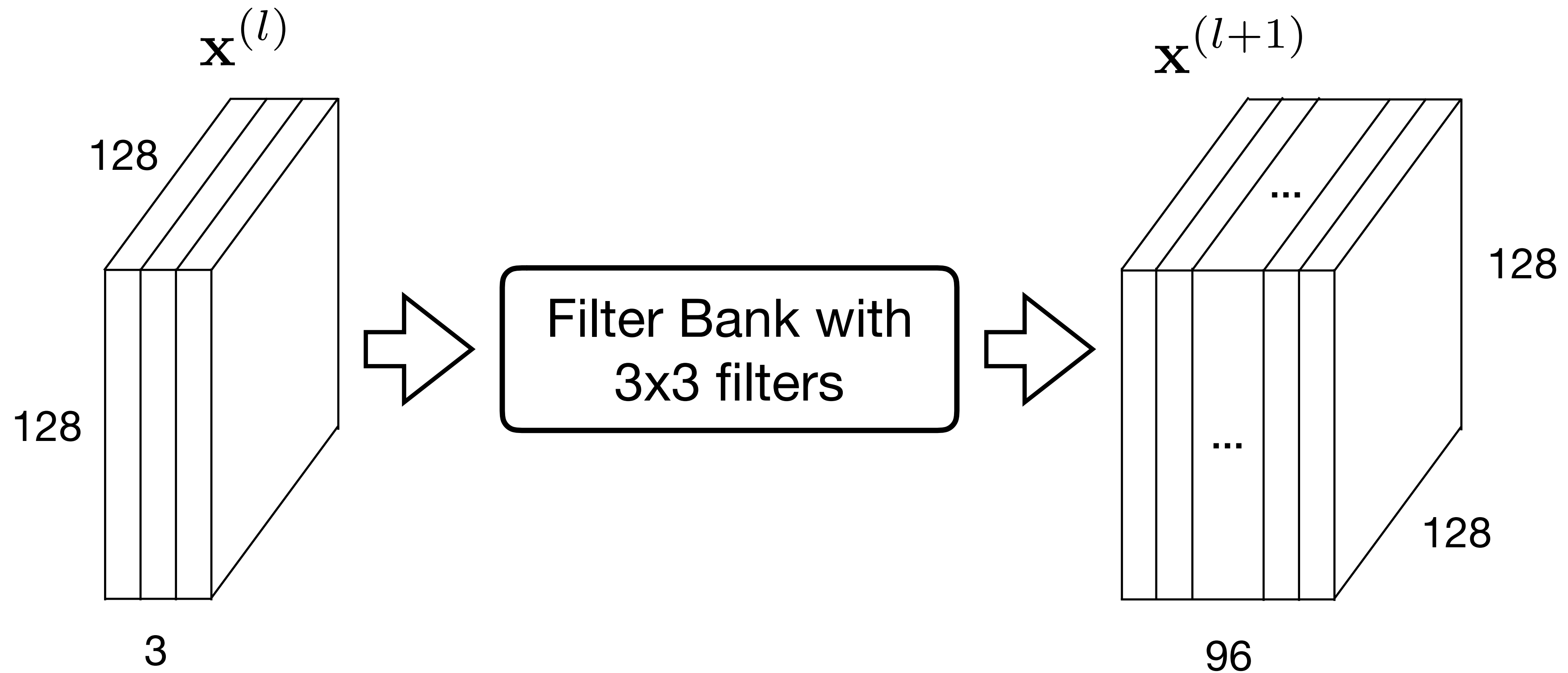
Multiple channels

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?

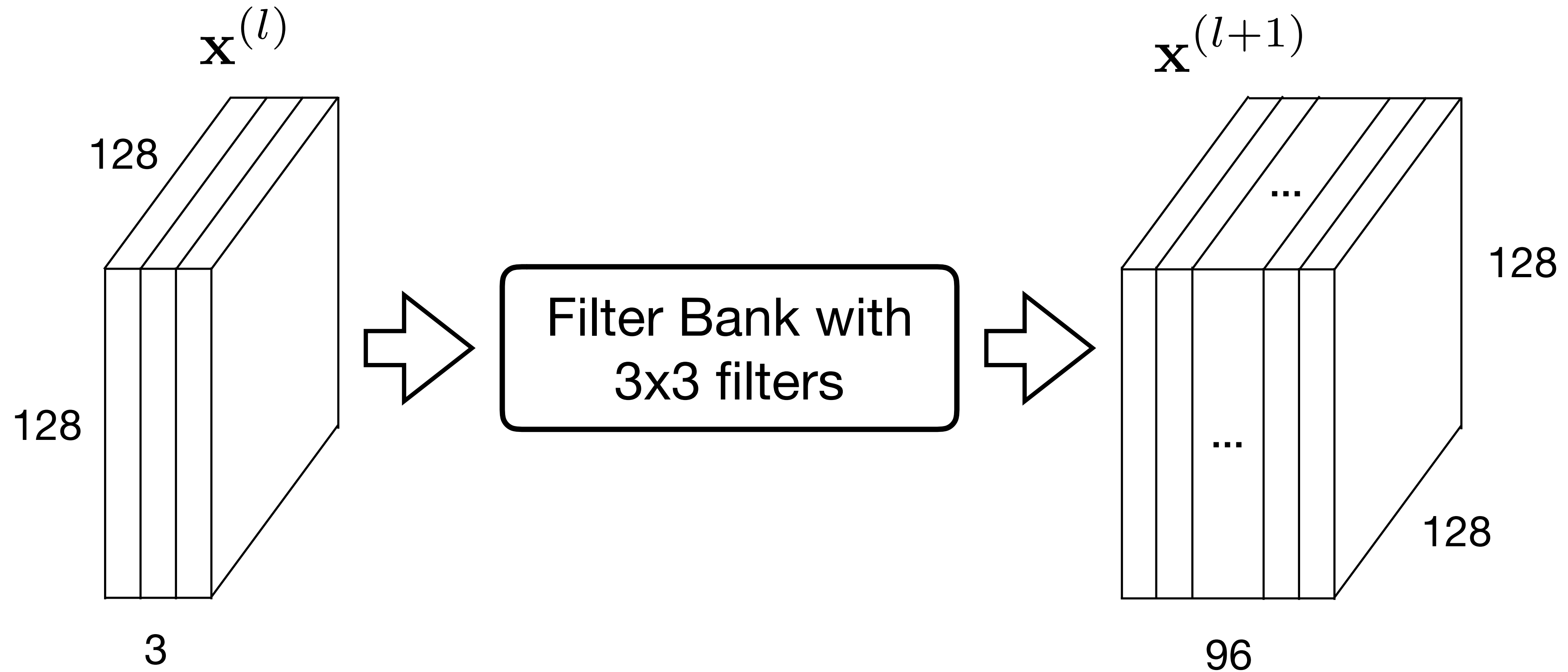
(a) 9

(b) 27

(c) 96

(d) 864

Multiple channels: Example



How many filters are in the bank?

- (a) 3 (b) 27 (c) 96 (d) can't say

Filter sizes

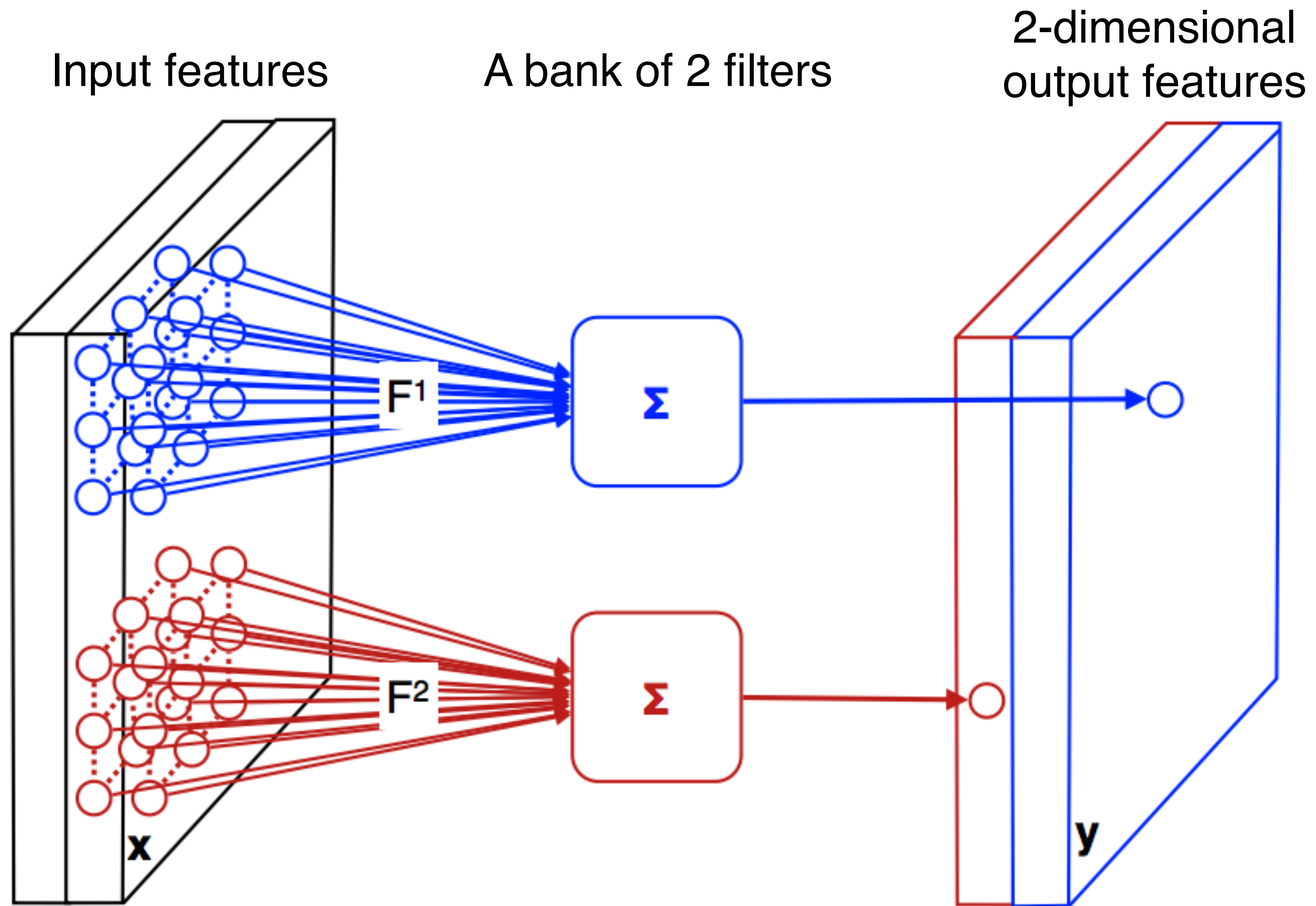
When mapping from

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

using an filter of spatial extent $M \times N$

Number of parameters per filter: $M \times N \times C^{(l)}$

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



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

[Figure from Andrea Vedaldi]

Image classification

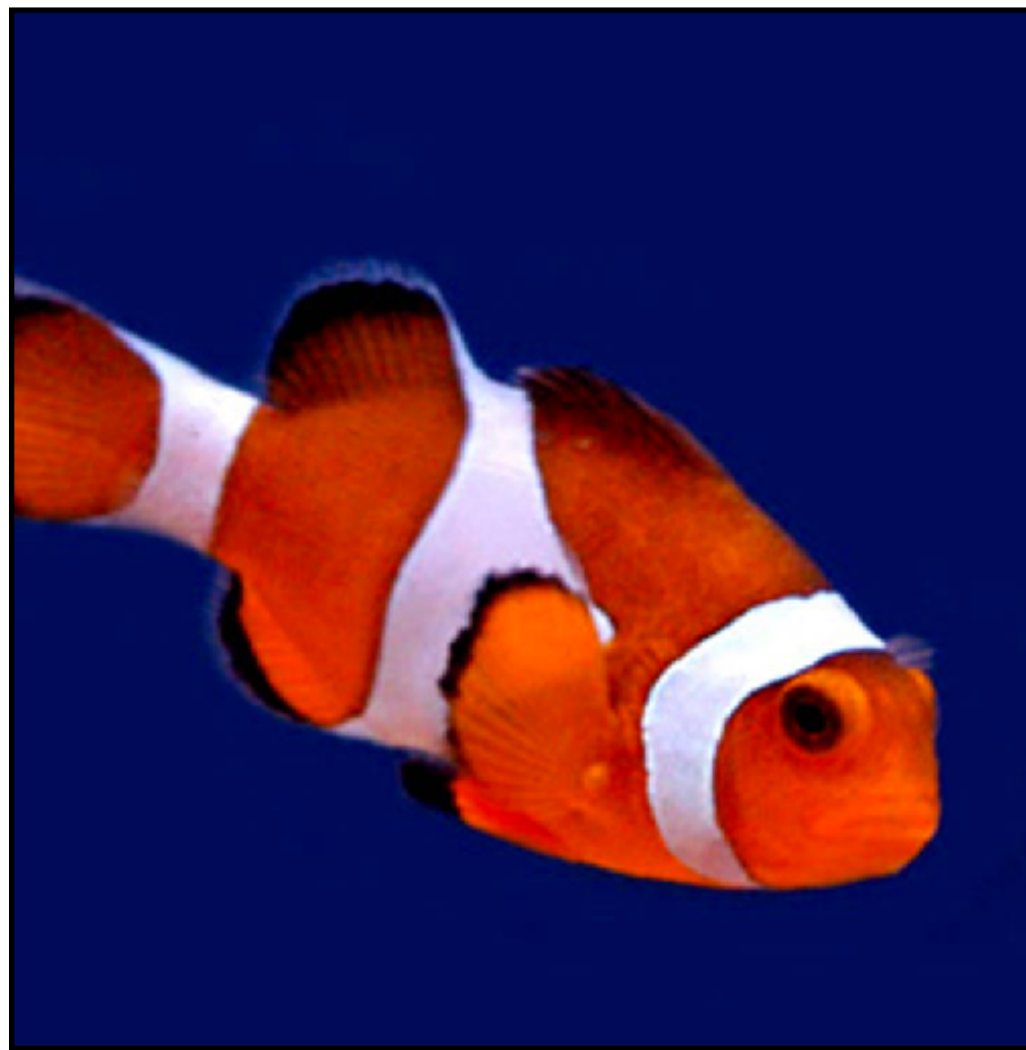


image x



"Fish"

label y

Image classification

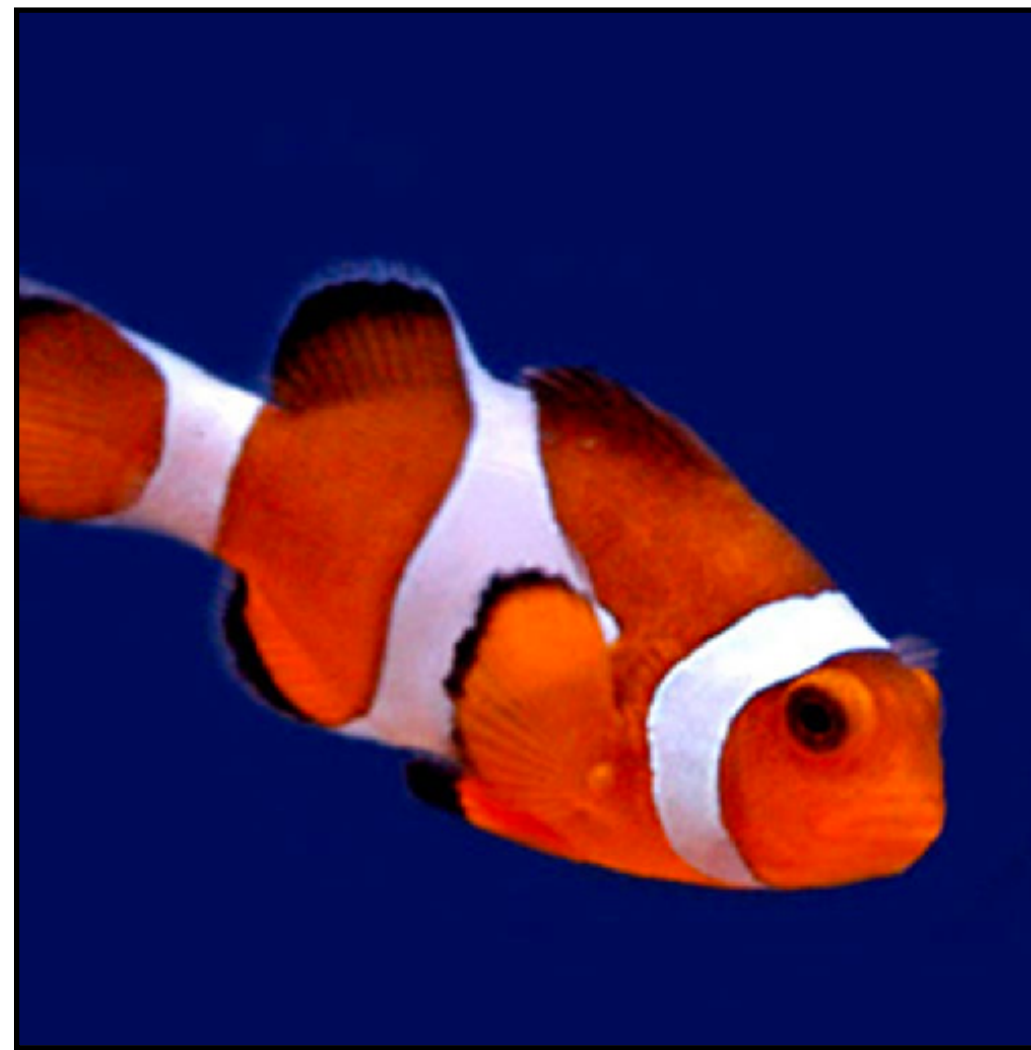
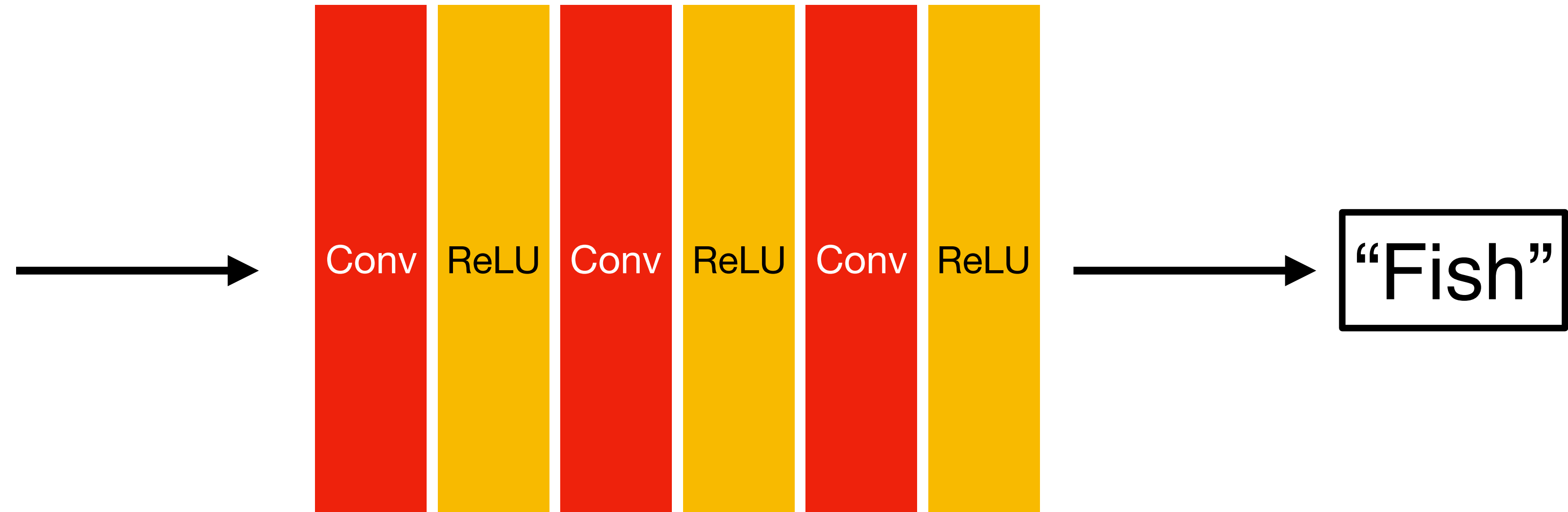
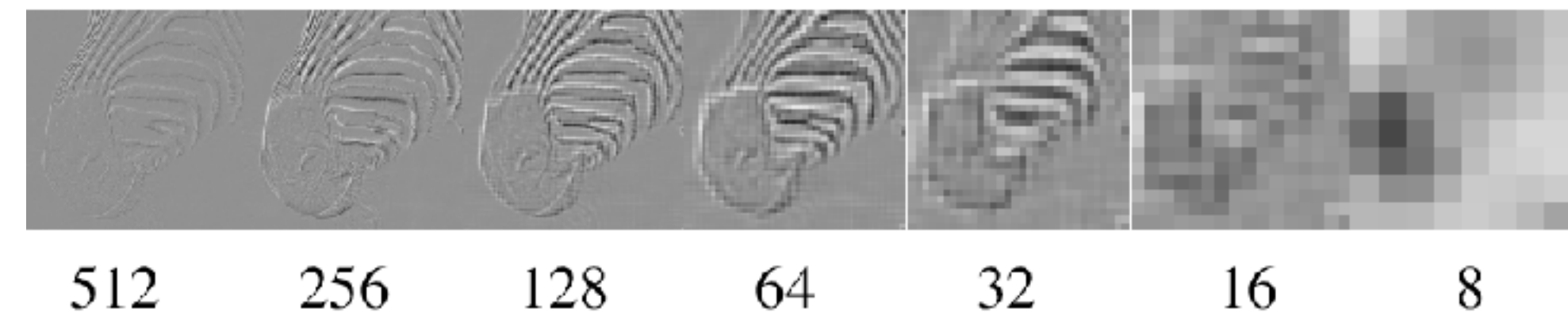
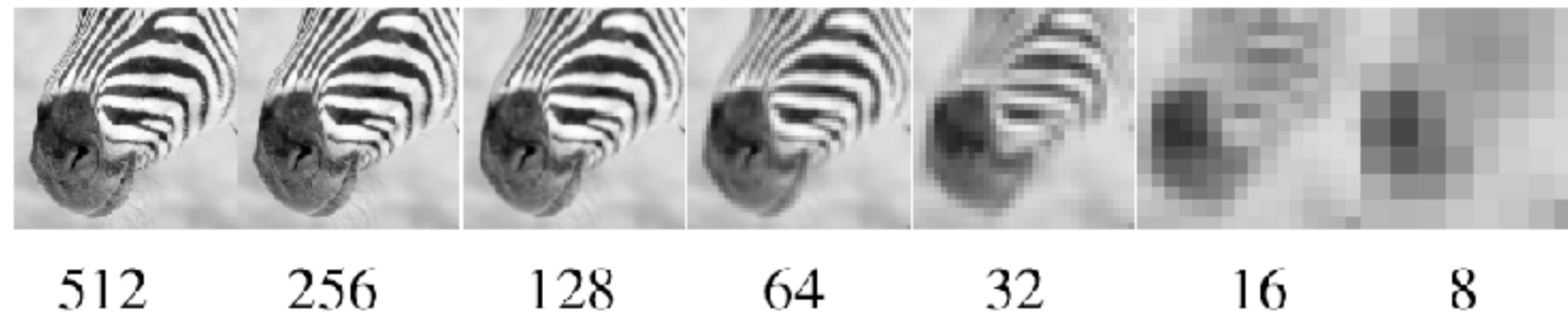


image x

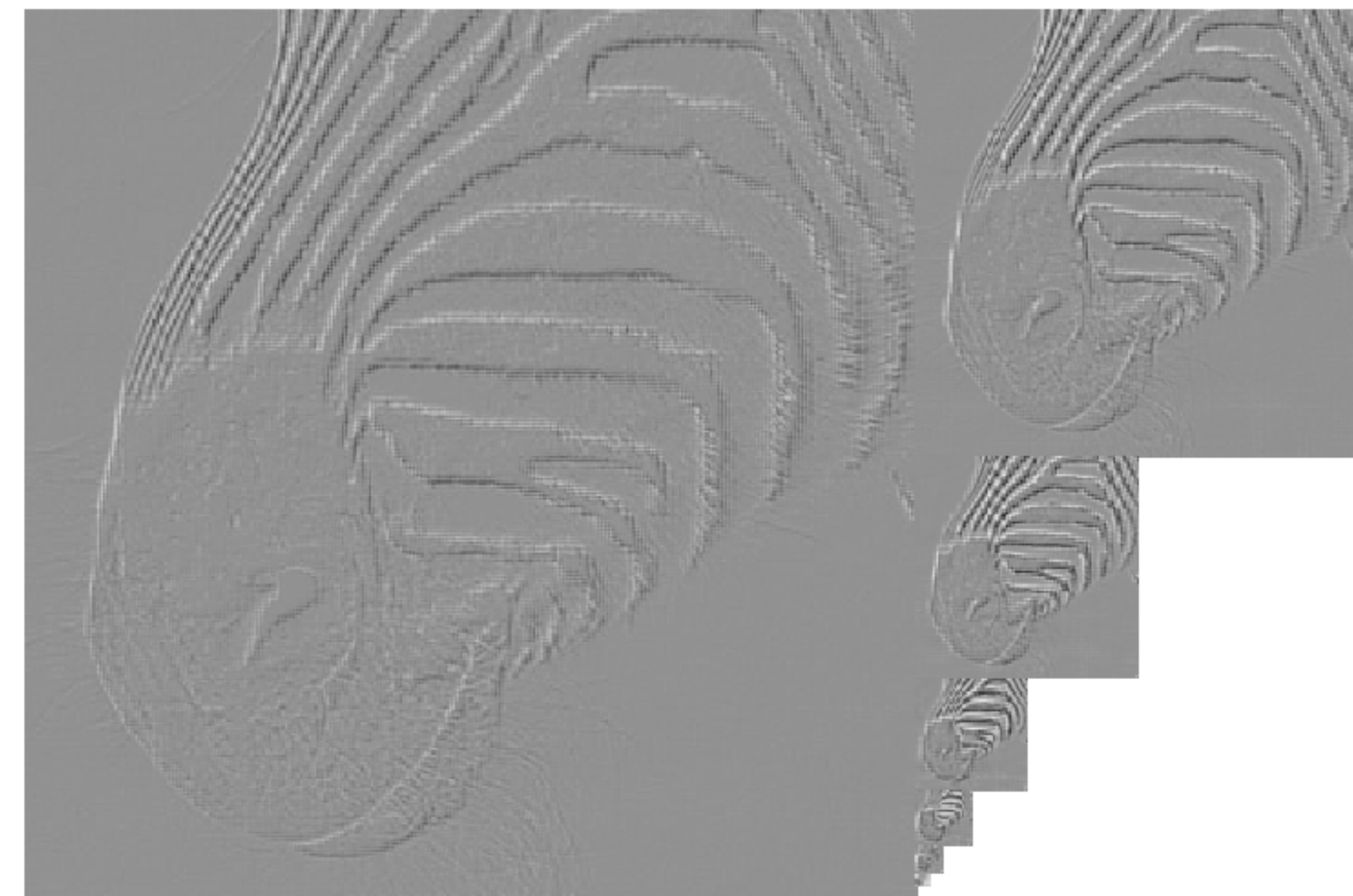


label y

Multiscale representations are great!



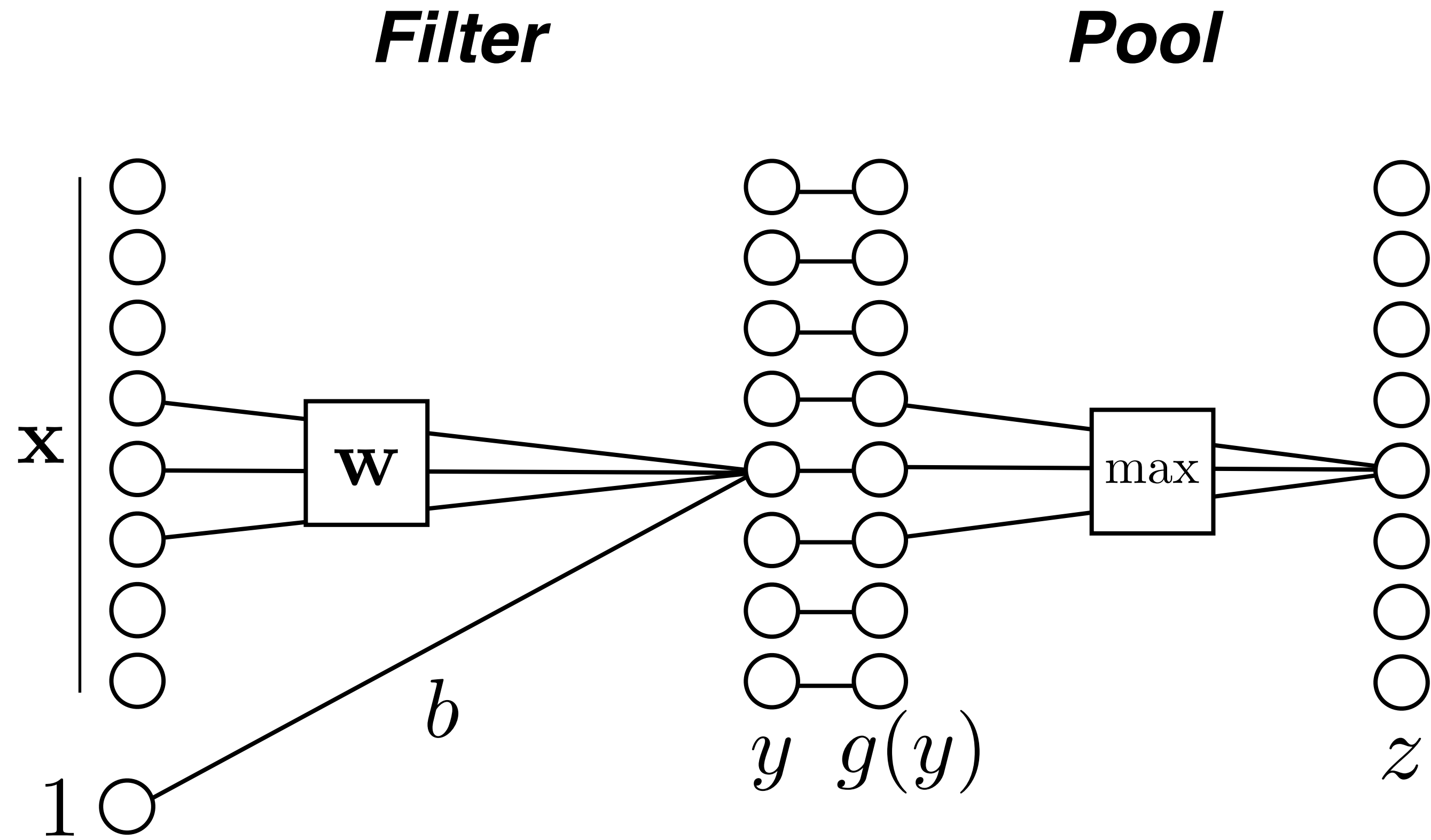
Gaussian Pyr



Laplacian Pyr

How can we use multi-scale modeling in Convnets?

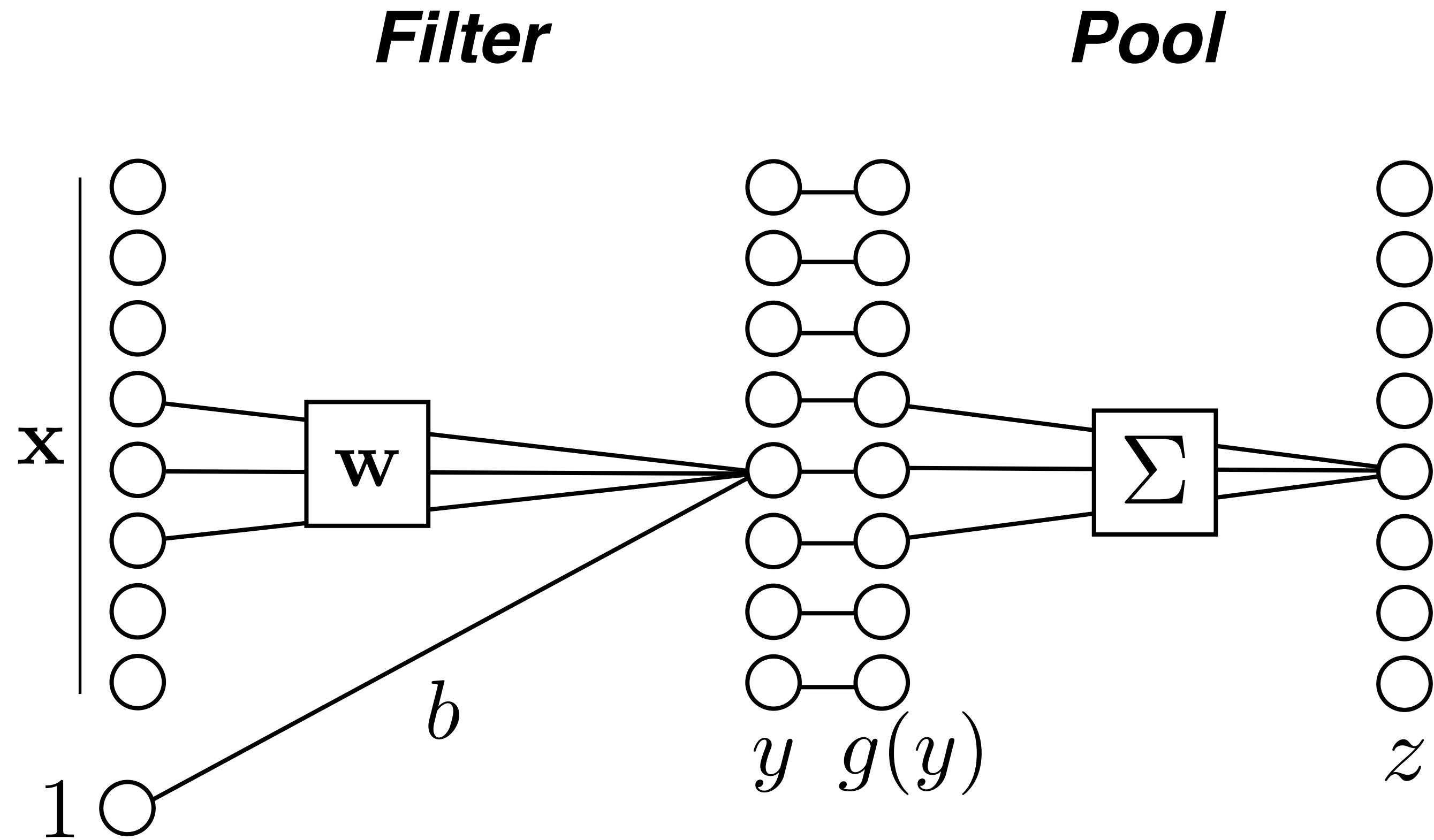
Pooling



Max pooling

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

Pooling



Max pooling

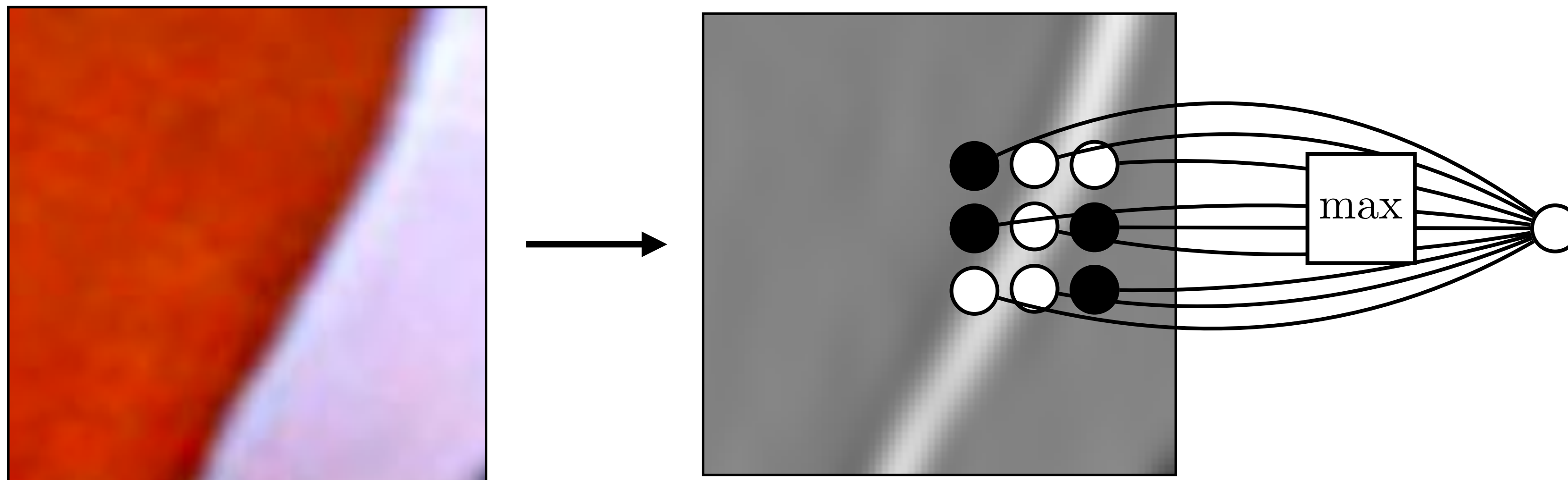
$$z_k = \max_{j \in \mathcal{N}(k)} g(y_j)$$

Mean pooling

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

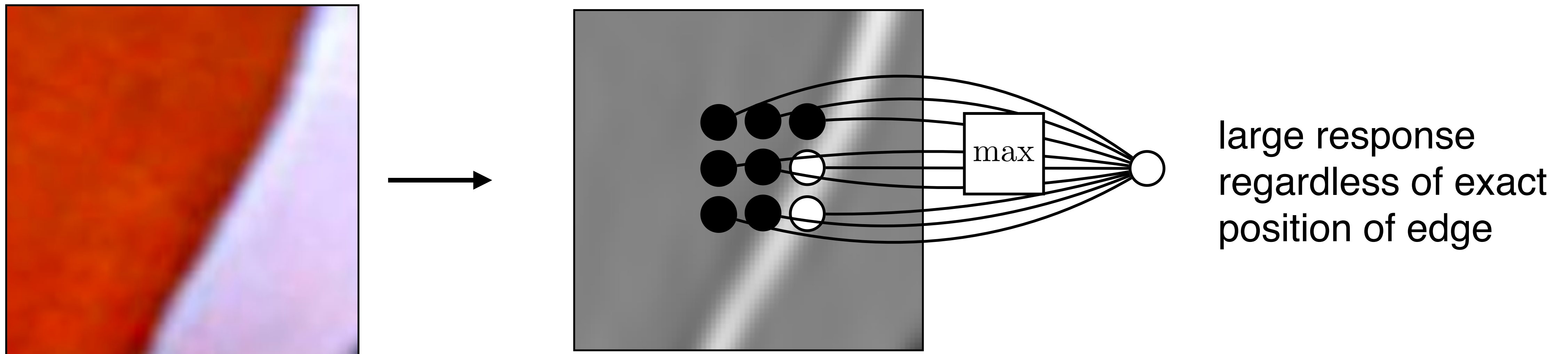
Pooling — Why?

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



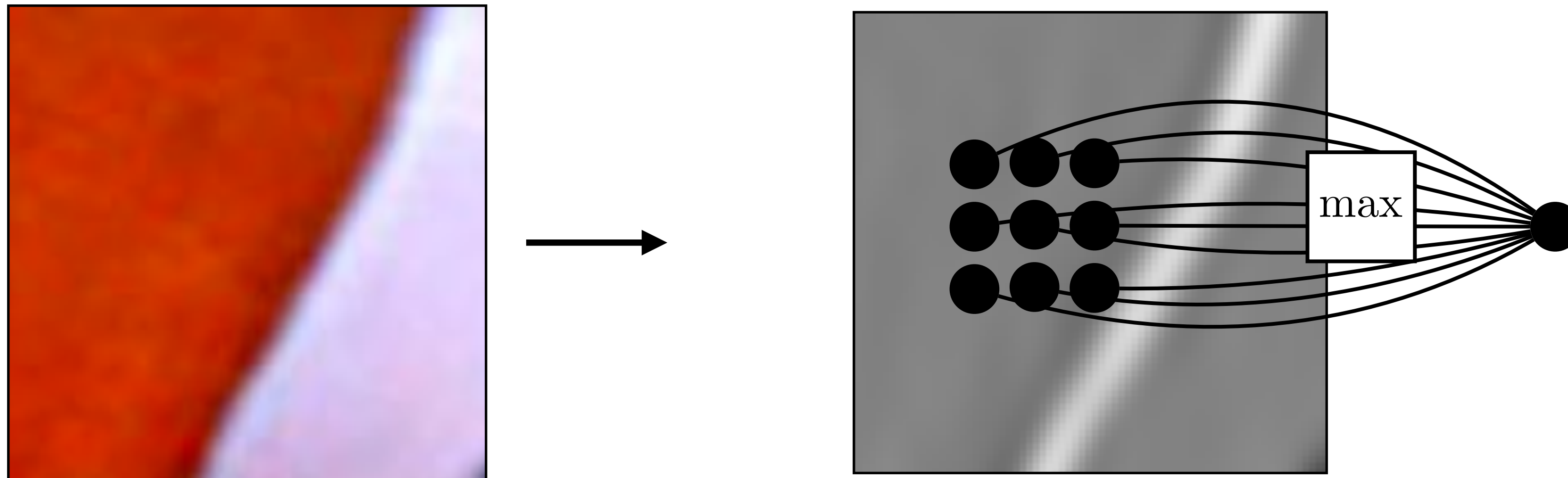
Pooling — Why?

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

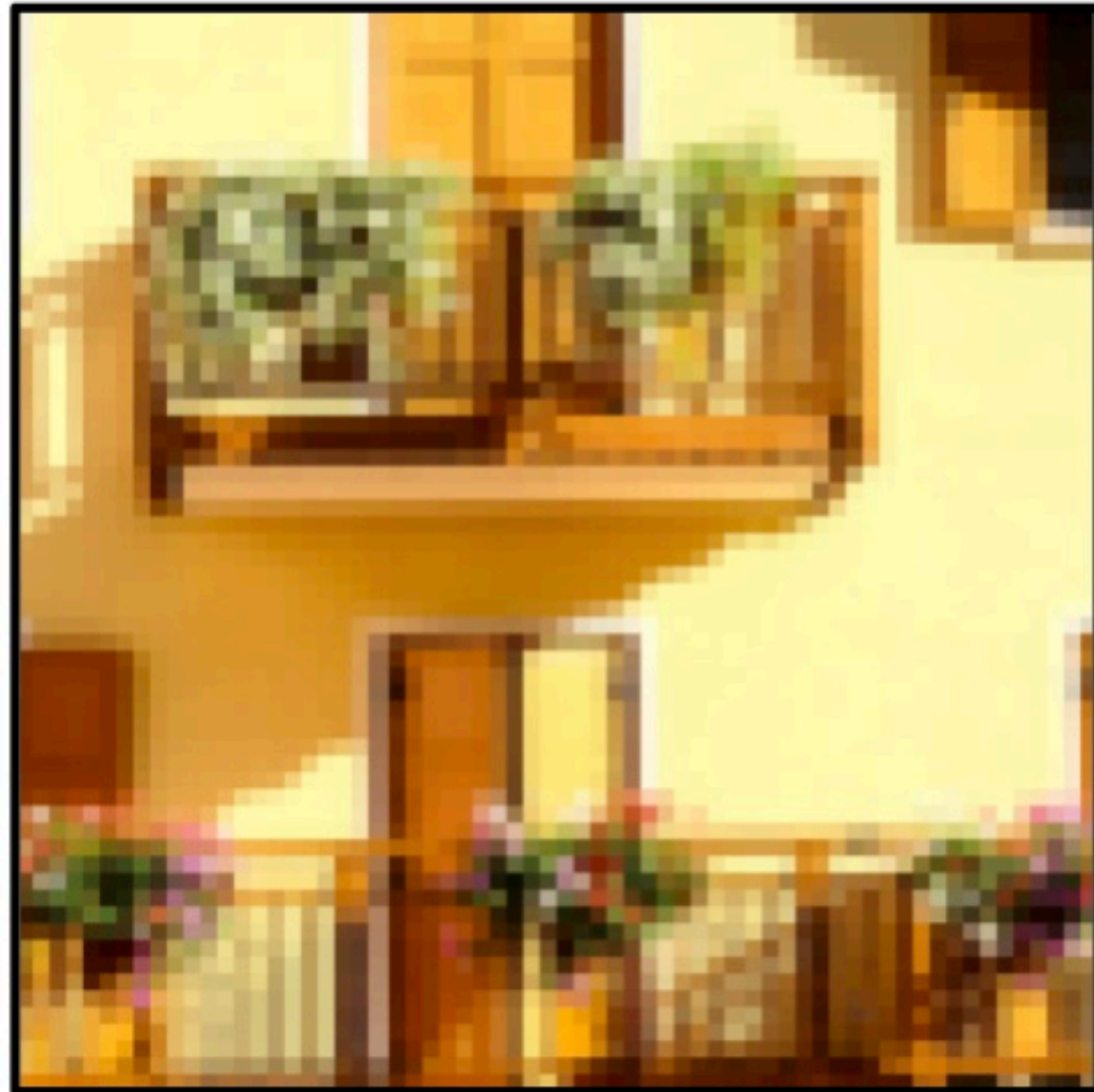


Pooling — Why?

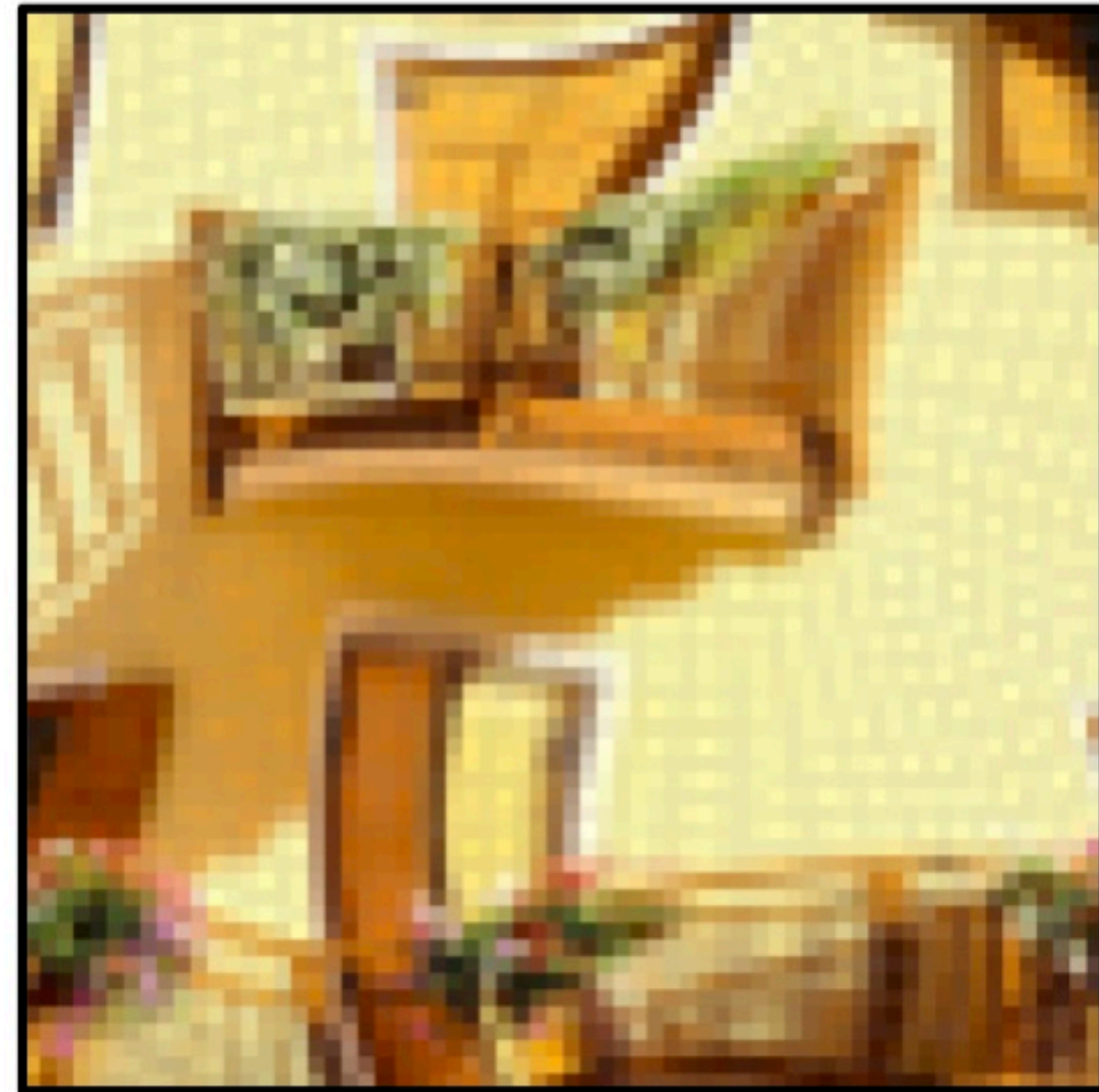
Pooling across spatial locations achieves stability w.r.t. small translations:



CNNs are stable w.r.t. diffeomorphisms



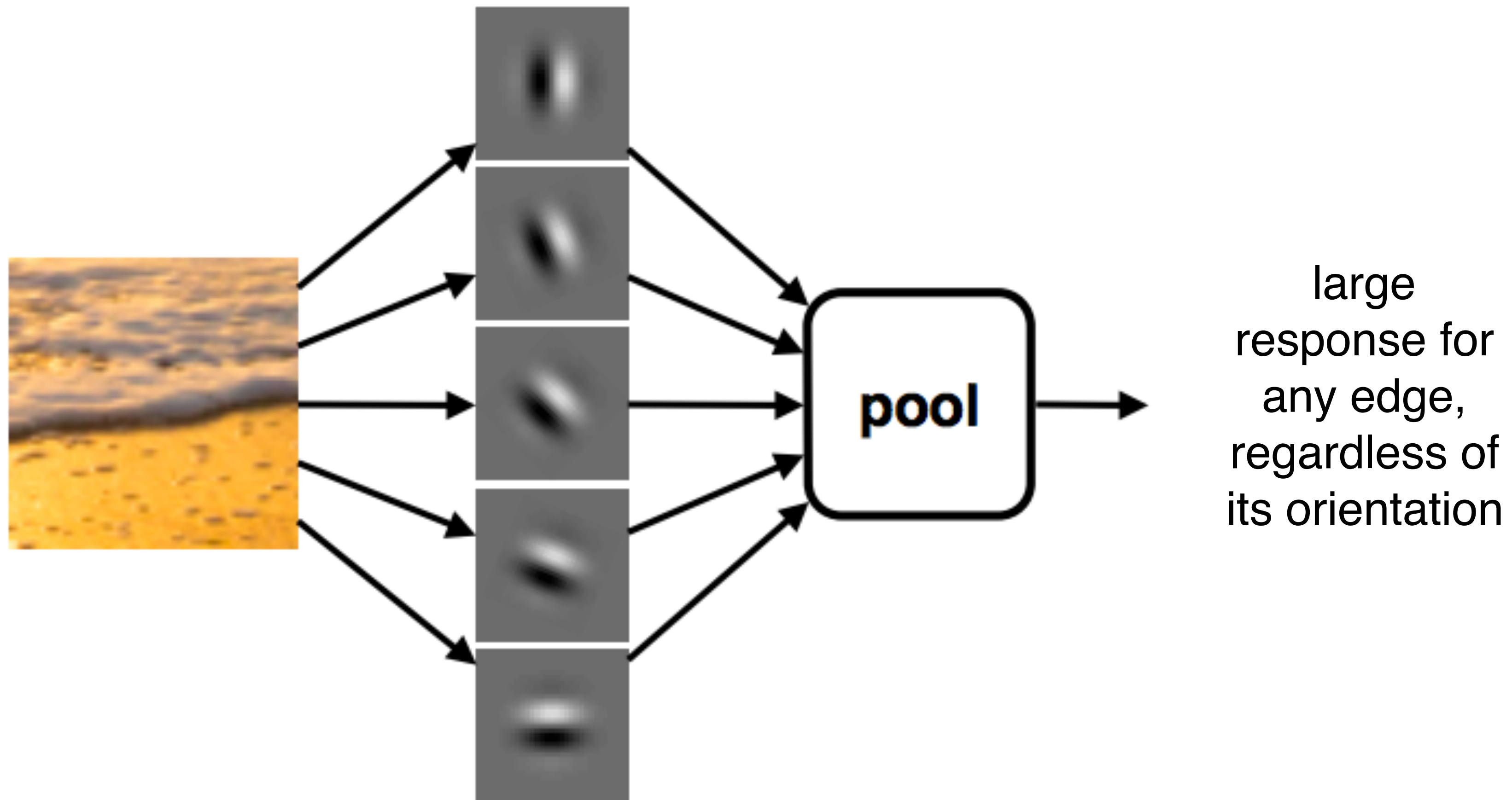
\approx



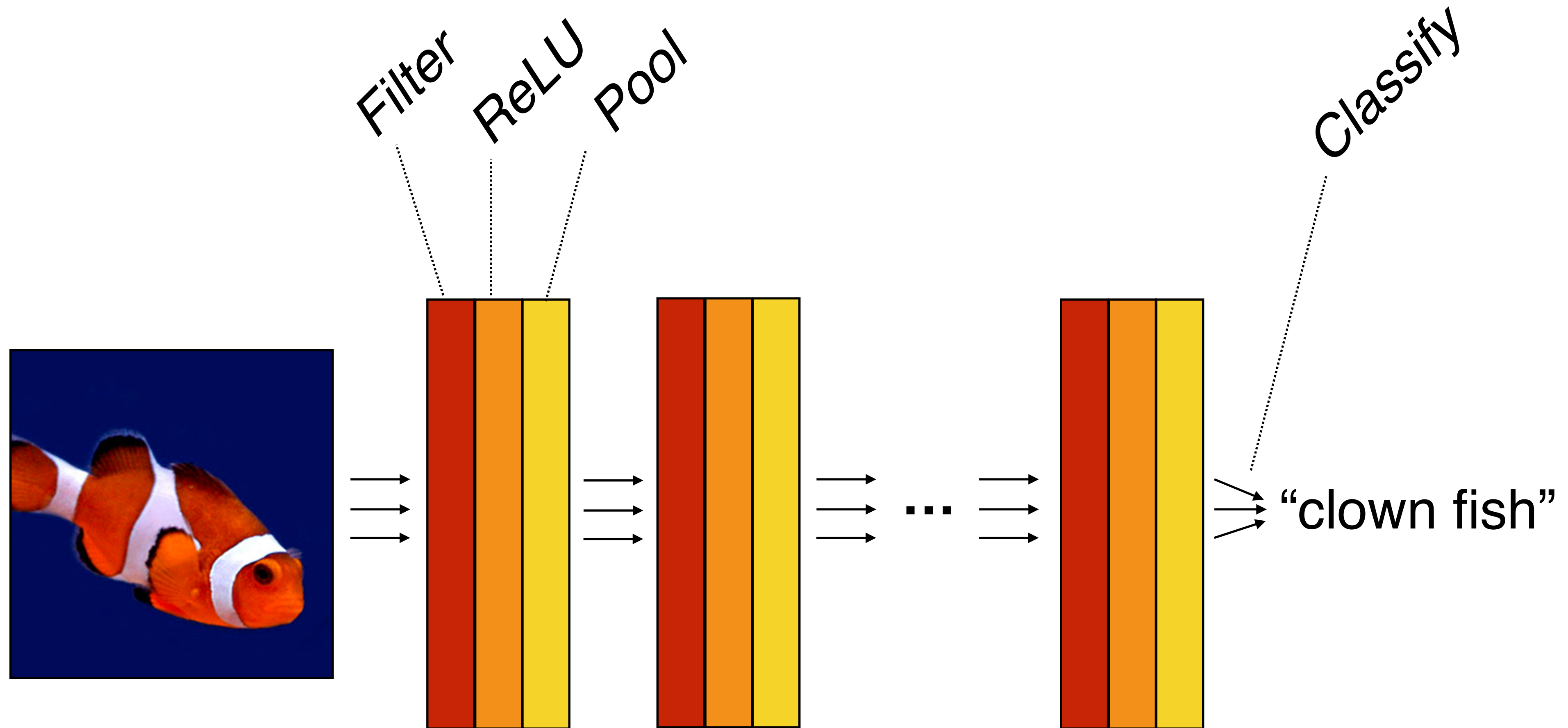
[“Unreasonable effectiveness of Deep Features as a Perceptual Metric”, Zhang et al. 2018]

Pooling — Why?

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

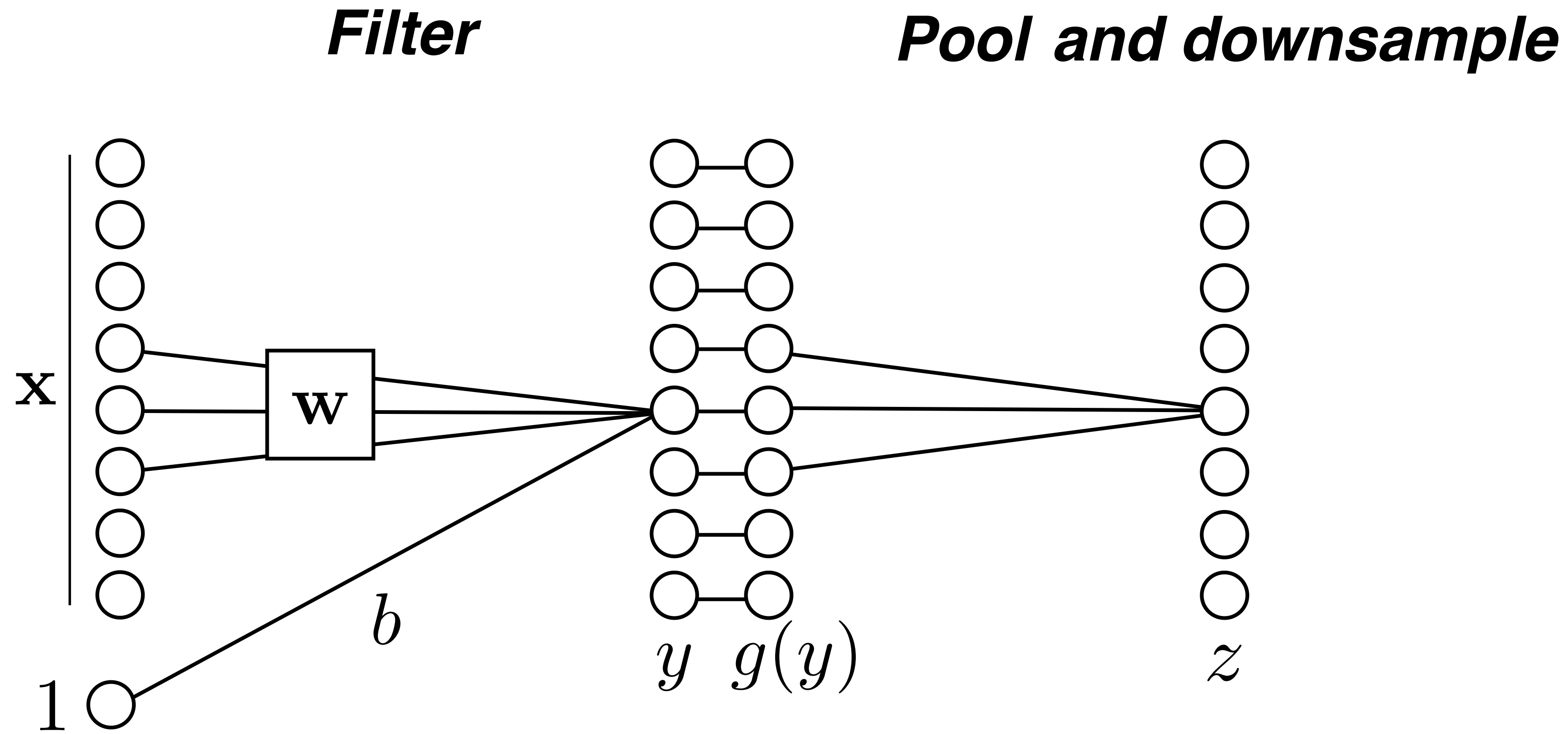


Computation in a neural net

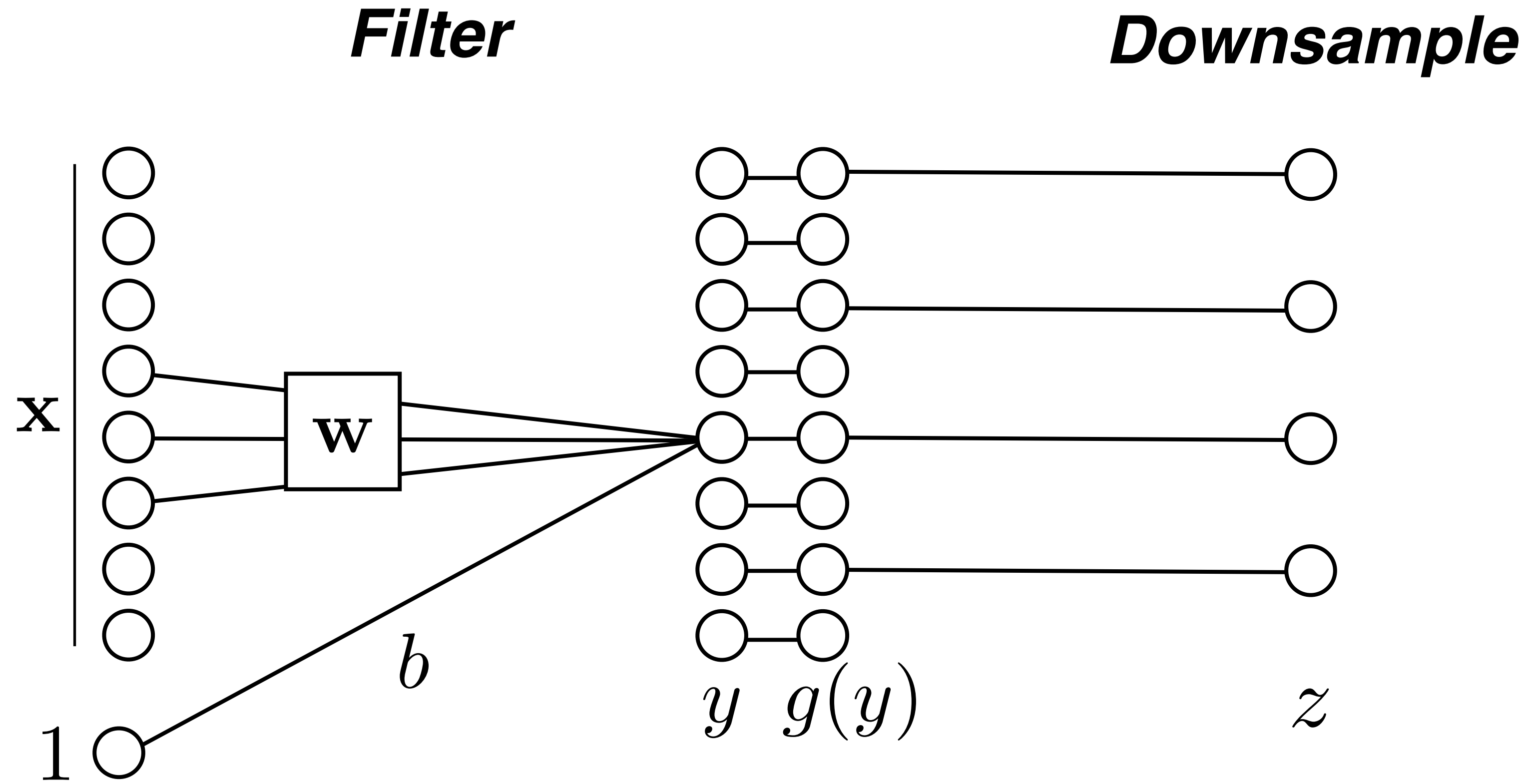


$$f(\mathbf{x}) = f_L(\dots f_2(f_1(\mathbf{x})))$$

Downsampling



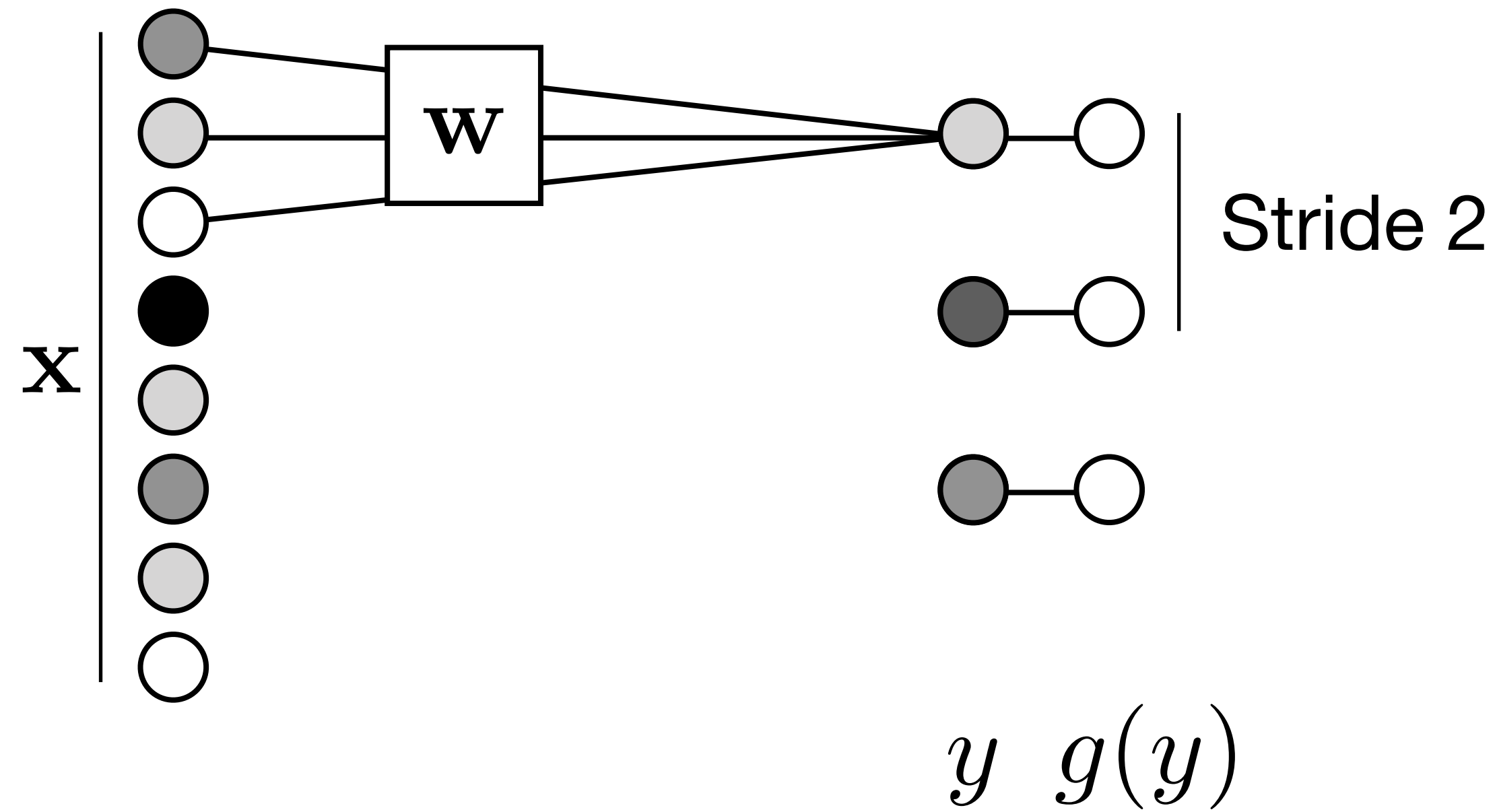
Downsampling



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

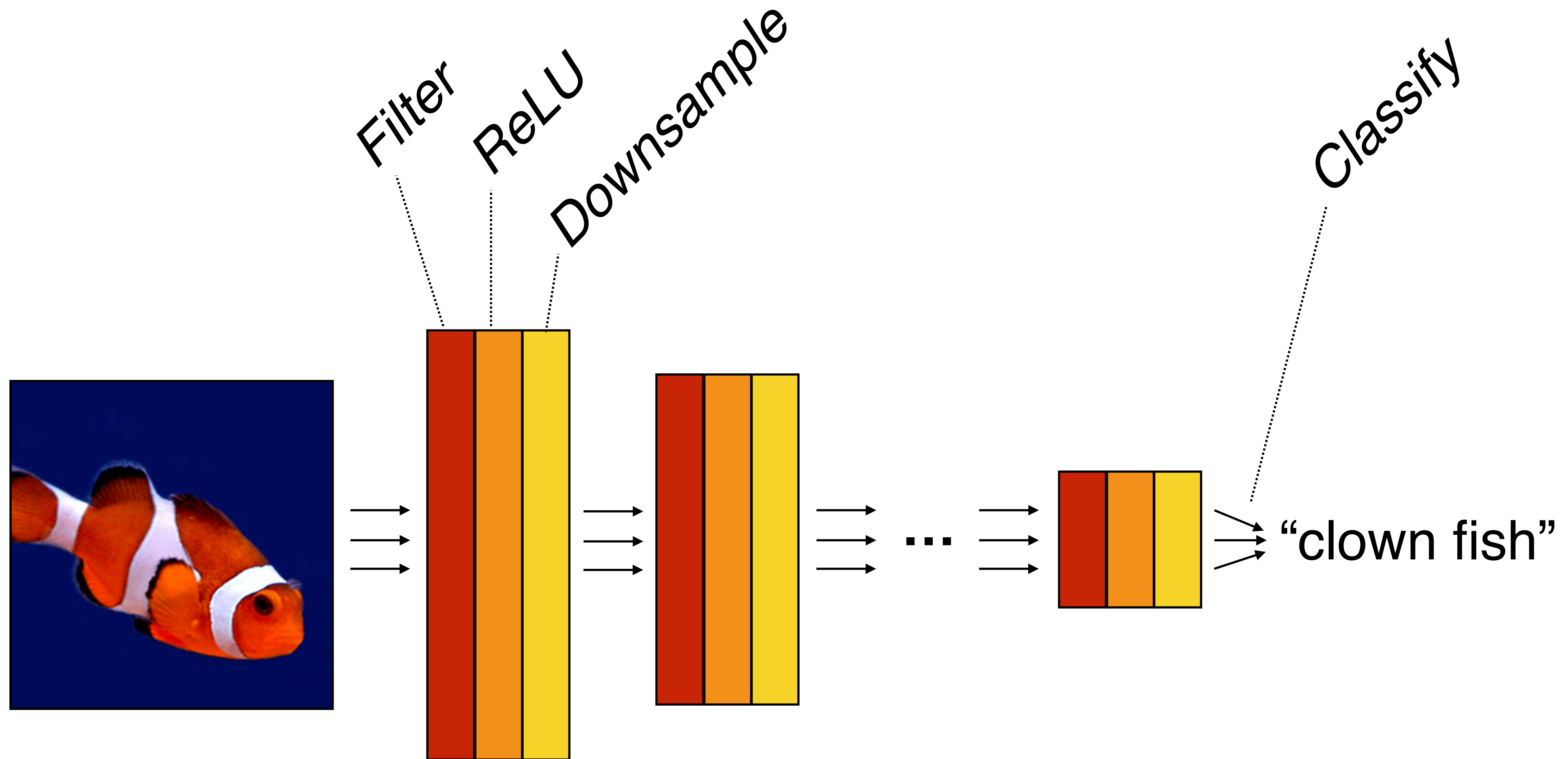
Strided operations

Conv layer



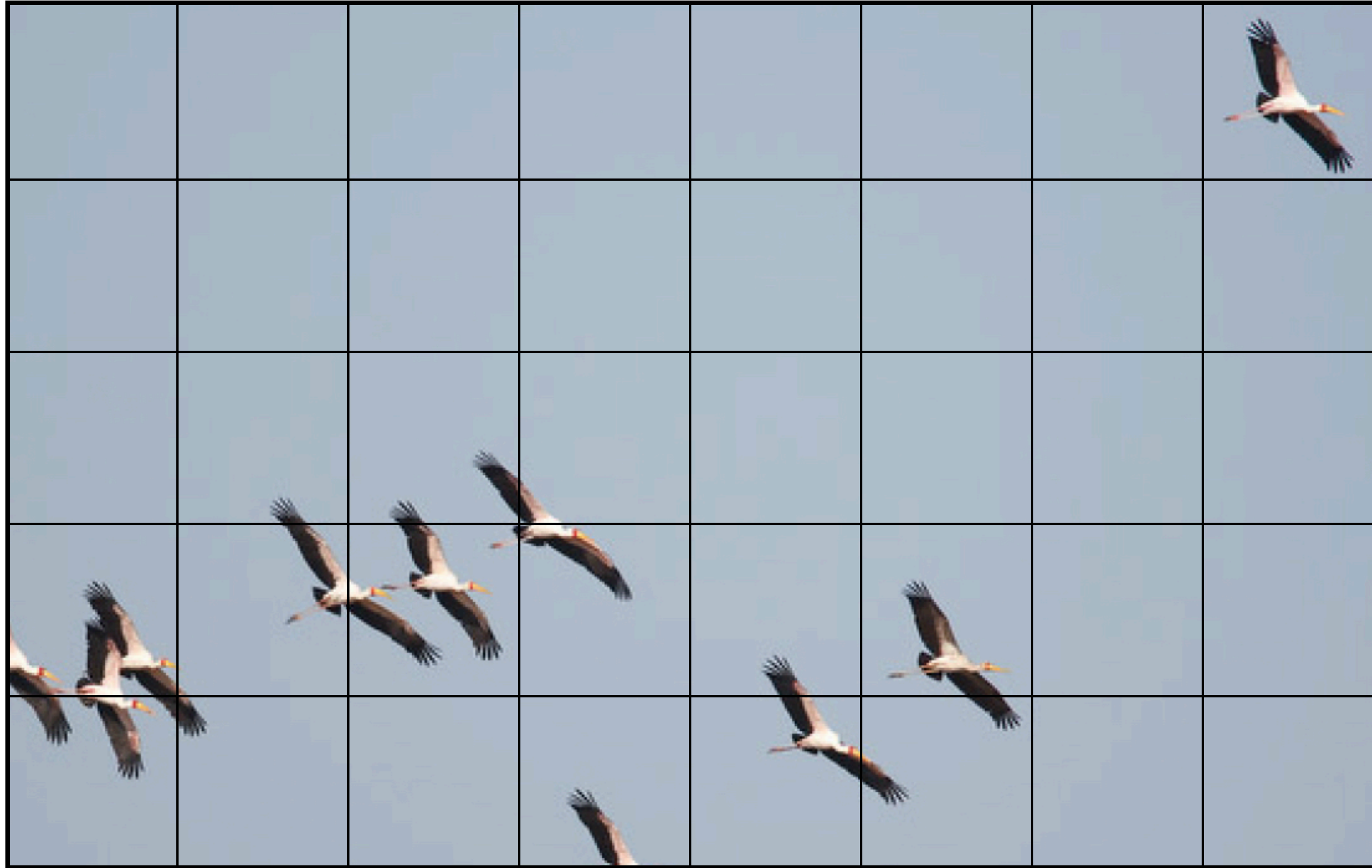
Strided operations combine a given operation (convolution or pooling) and downsampling into a single operation.

Computation in a neural net

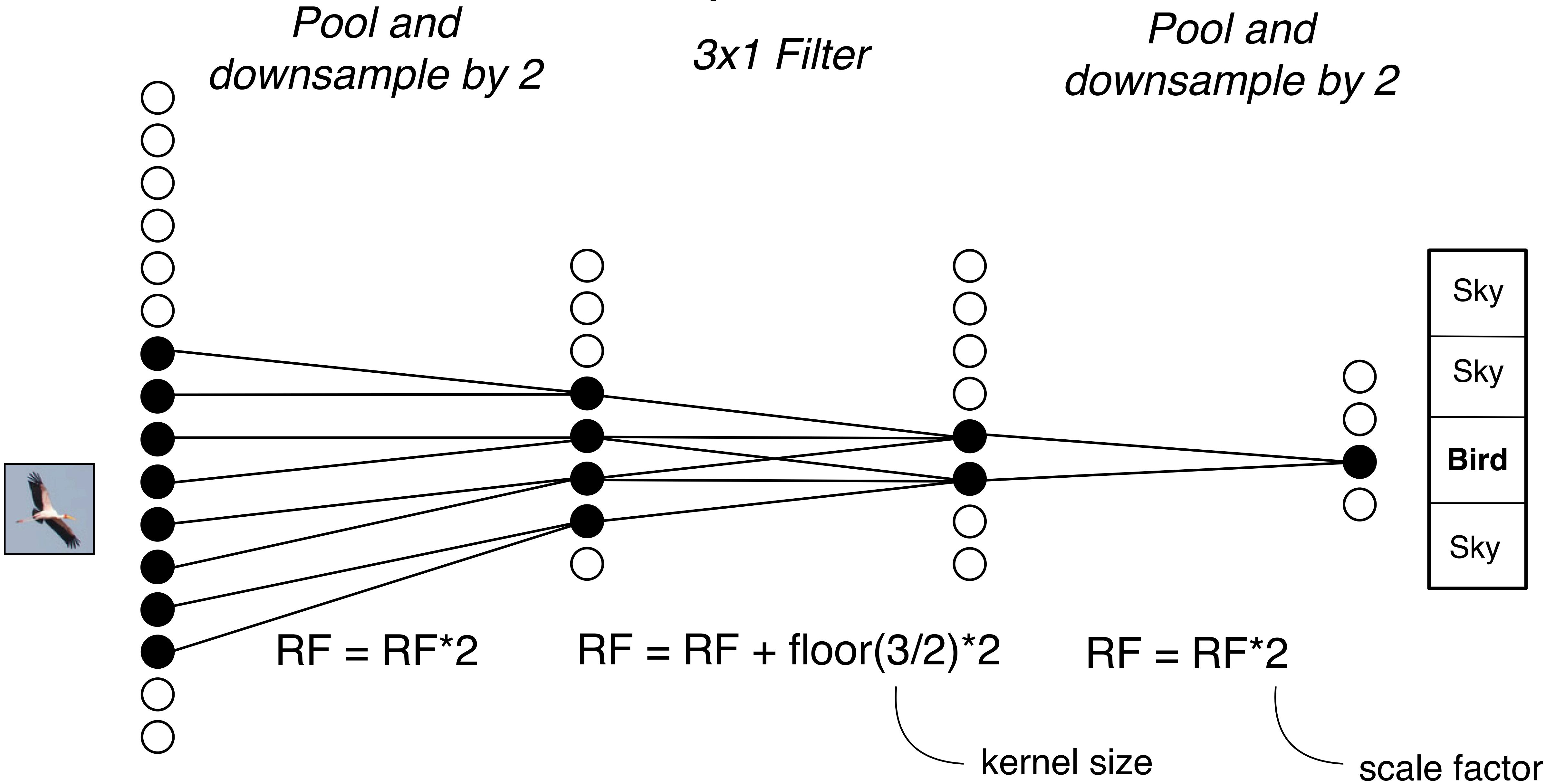


$$f(\mathbf{x}) = f_L(\dots f_2(f_1(\mathbf{x})))$$

Receptive fields



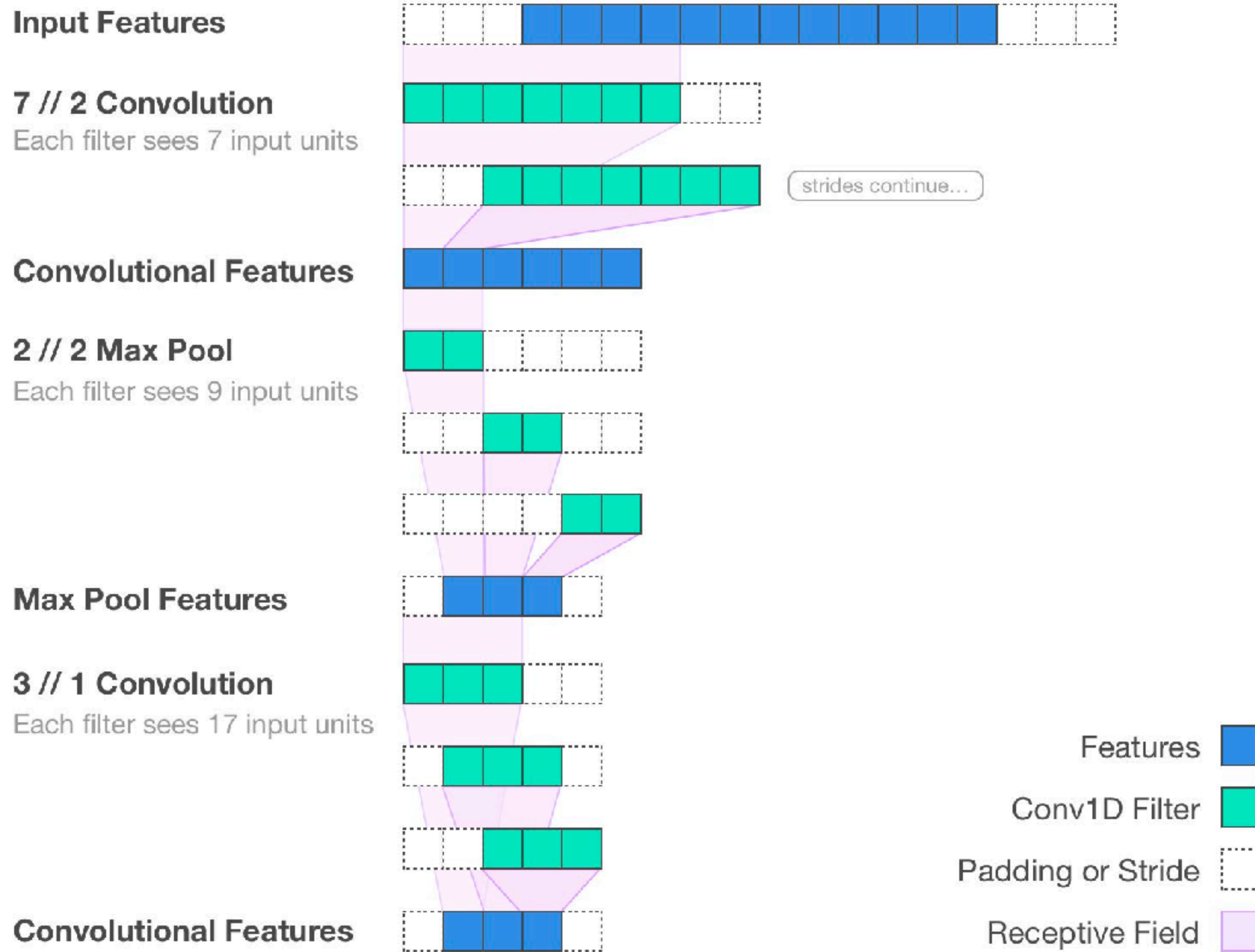
Receptive fields



Effective Receptive Field

Contributing input units to a convolutional filter.

@jimmfleming // fomoro.com

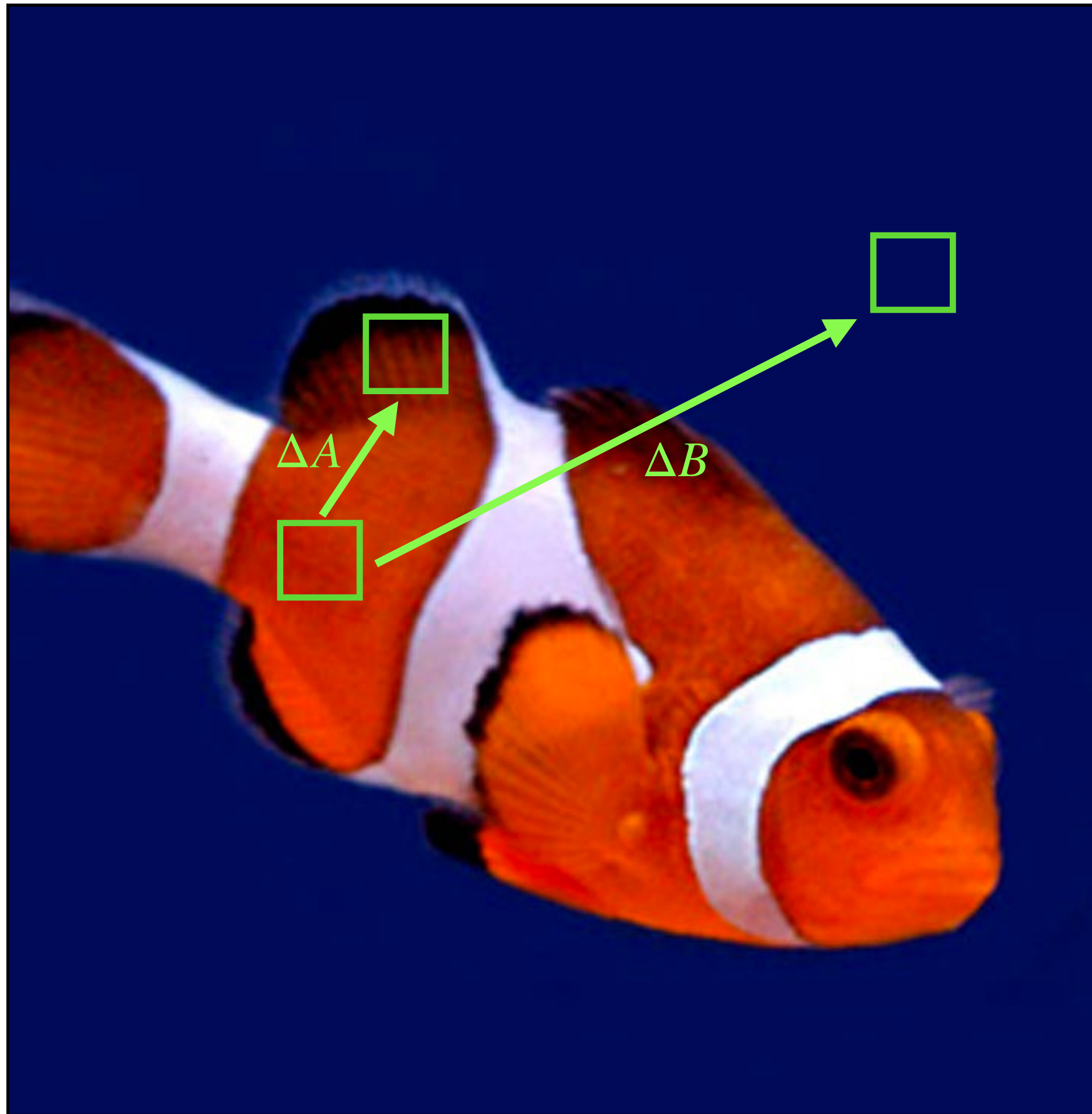


[<http://fomoro.com/tools/receptive-fields/index.html>]

Gradient / Backprop equations

... to be derived in the PSet (for Conv and Pool operations)

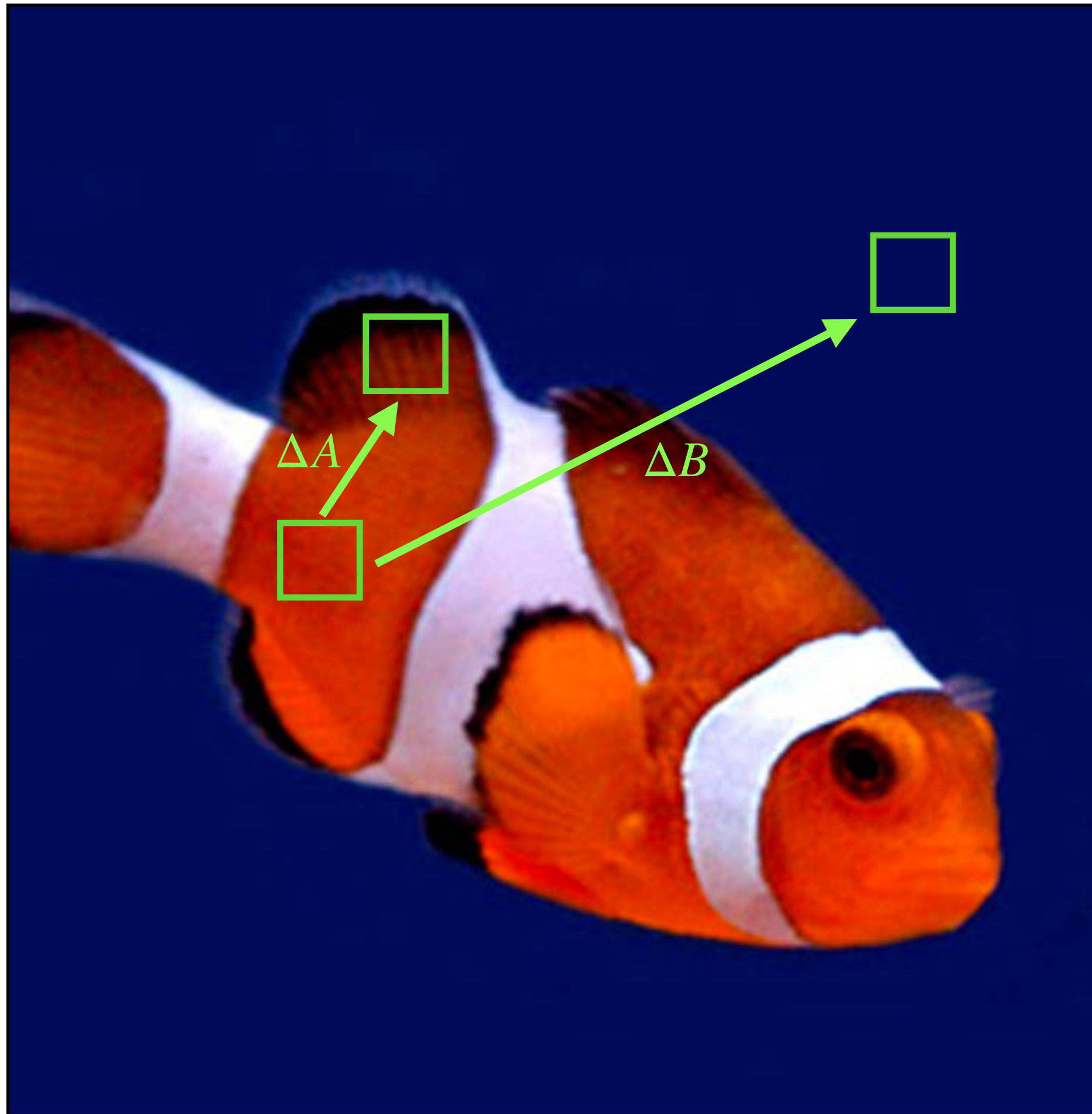
Local vs. global processing



Across all images, which is higher:

- (1) correlation between points with distance ΔA
- (2) correlation between points with distance ΔB
- (3) can't say

Local vs. global processing



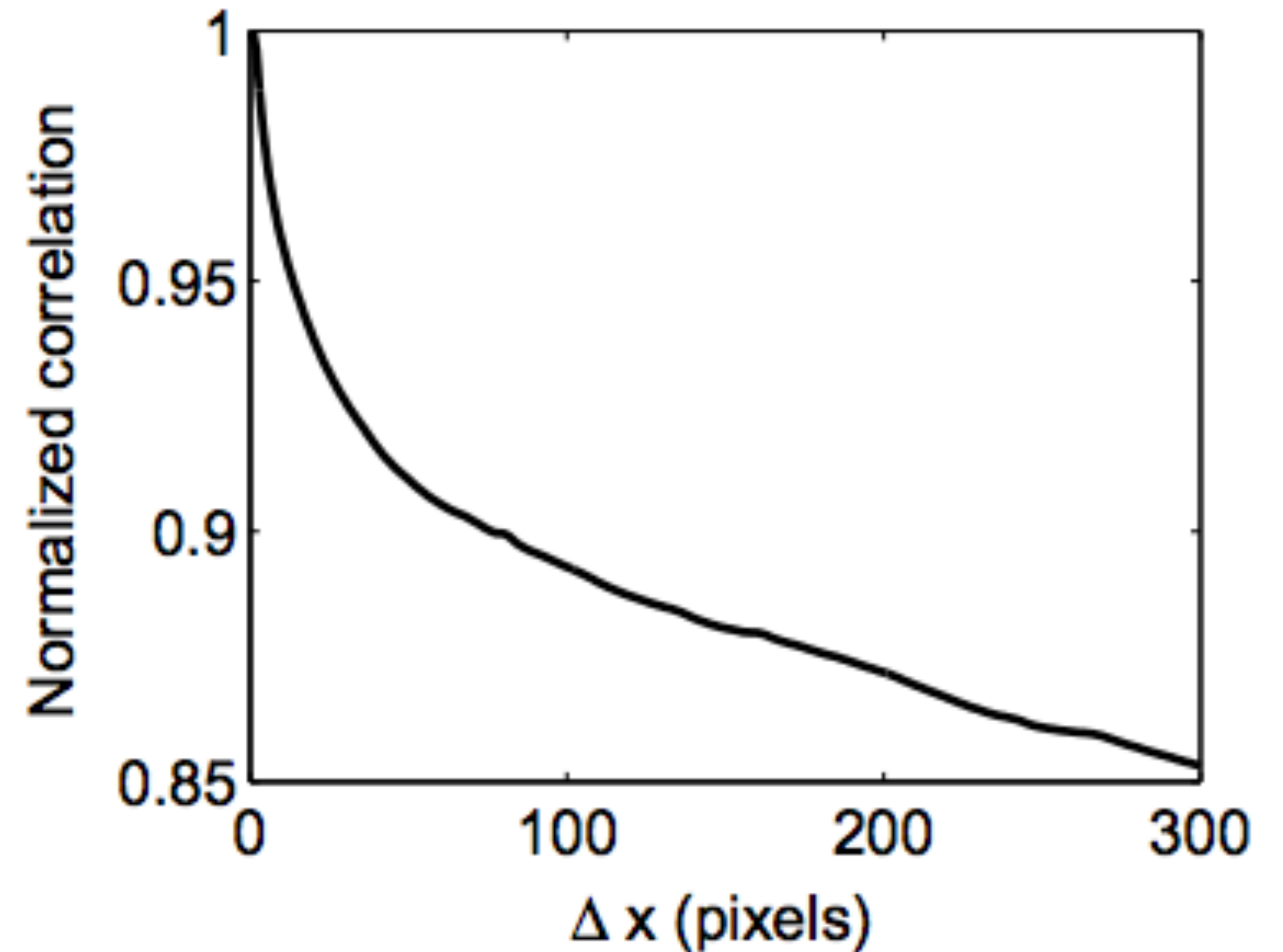
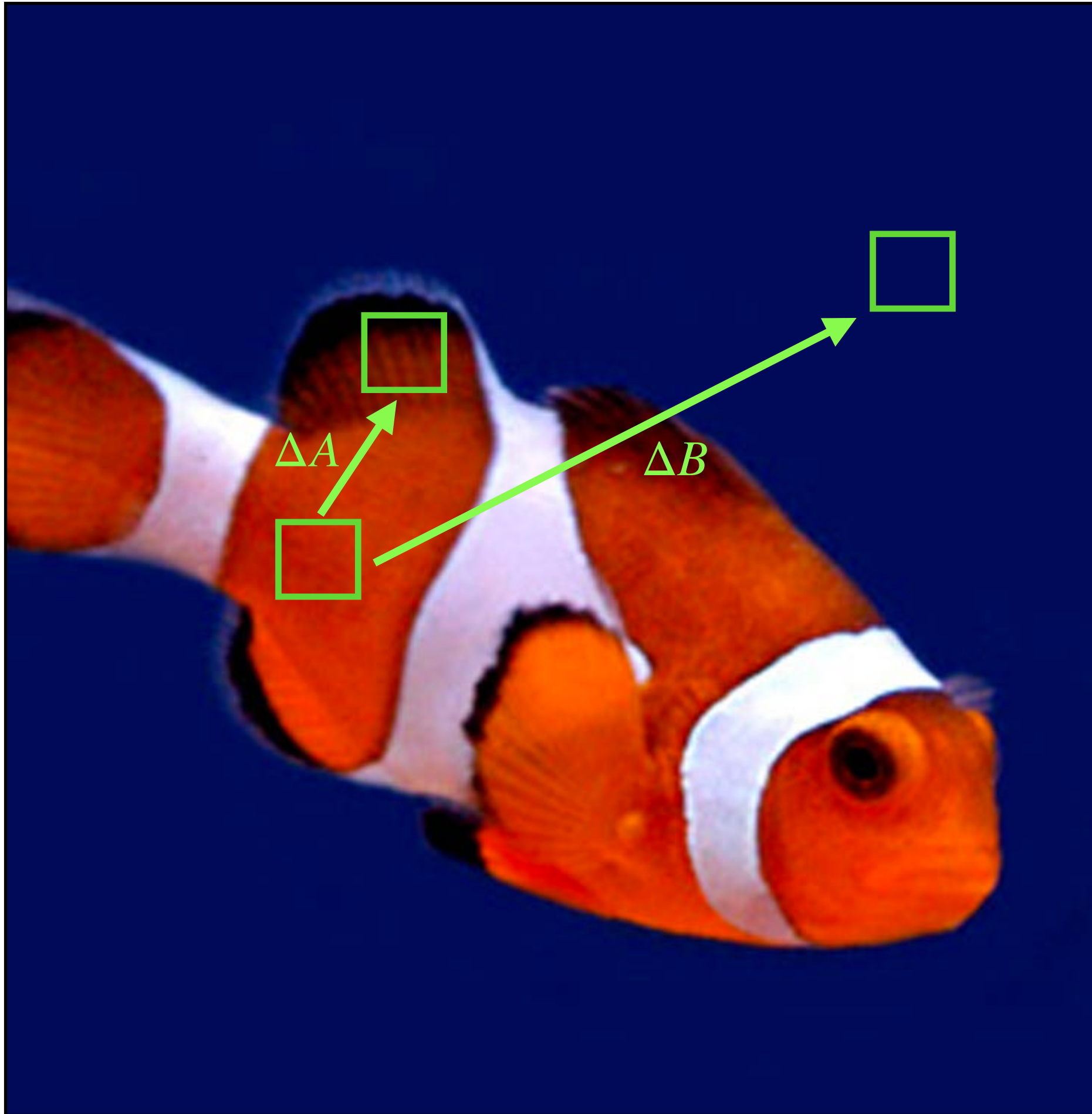
Across all images, which is higher:

(1) correlation between points with distance ΔA

(2) correlation between points with distance ΔB

(3) can't say

Local vs. global processing



[Simoncelli: *Statistical Modelling of Photographic Images*, 2005]

CNNs — Why?

Statistical dependences between pixels decay as a power law of distance between the pixels.

It is therefore often sufficient to model local dependences only. —> **Convolution**

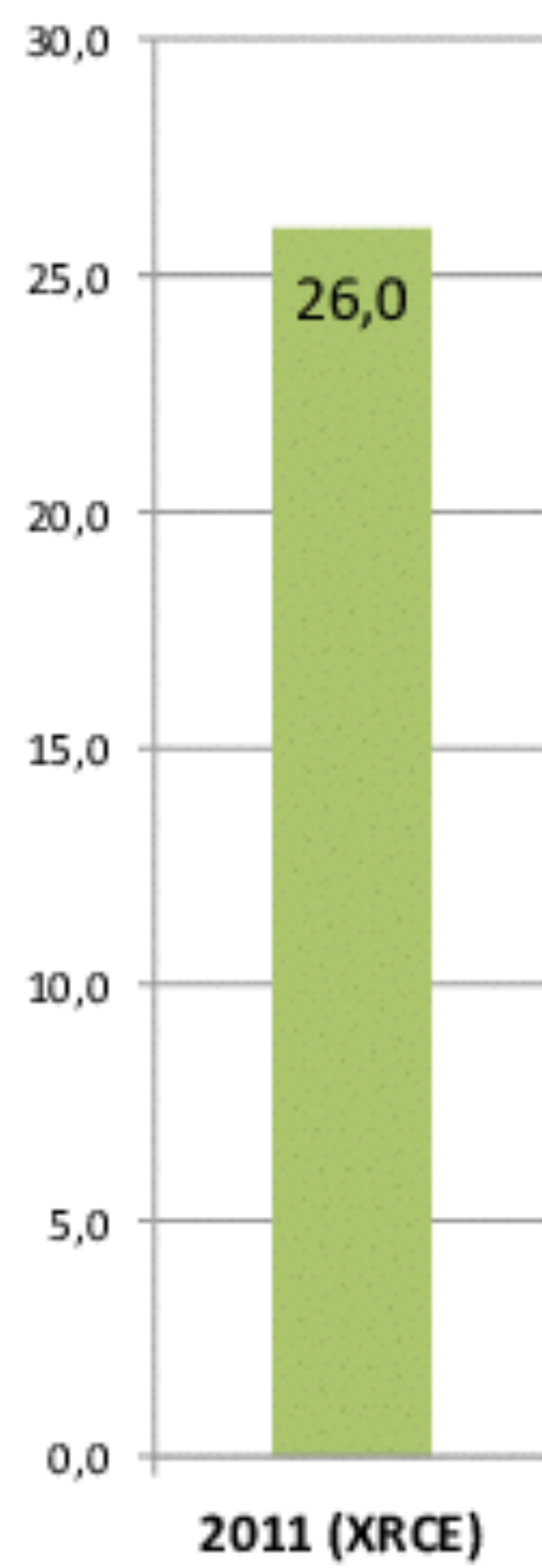
More generally, we should allocate parameters that model dependencies in proportion to the strength of those dependences. —> **Multiscale, hierarchical representations**

[For more discussion, see “Why does Deep and Cheap Learning Work So Well?”, Lin et al. 2017]

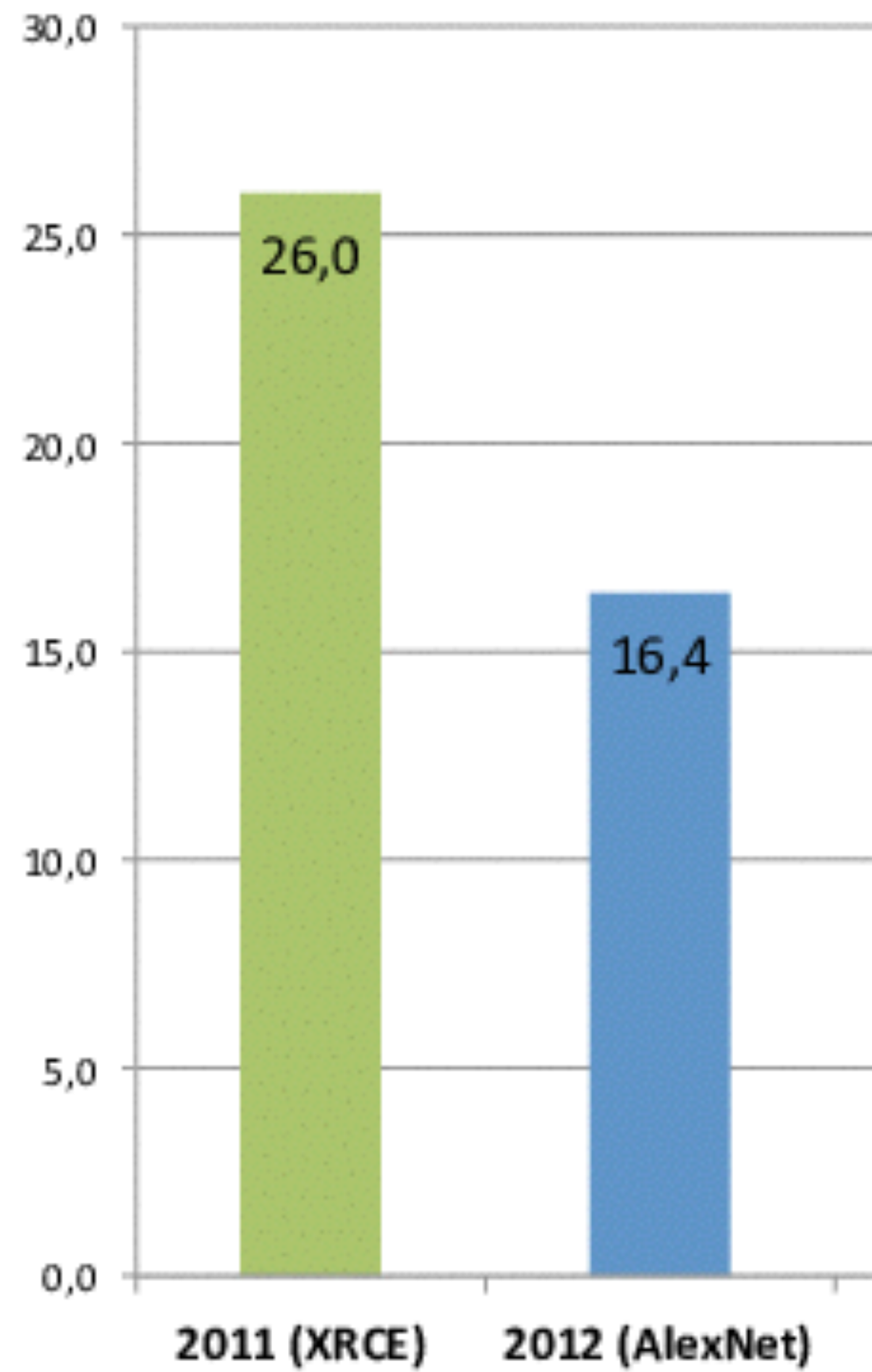
Some networks

... and what makes them work

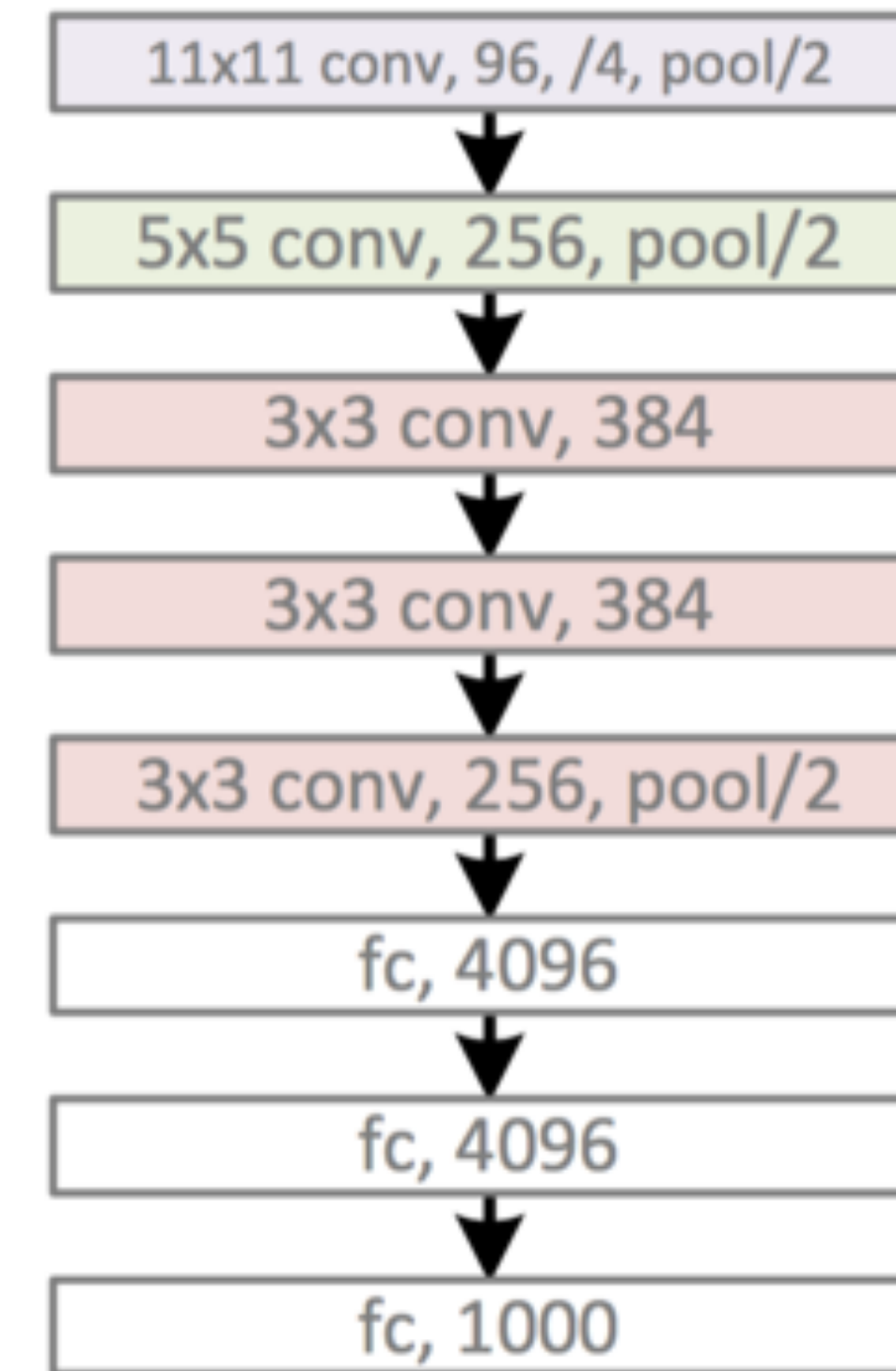
ImageNet Classification Error (Top 5)



ImageNet Classification Error (Top 5)

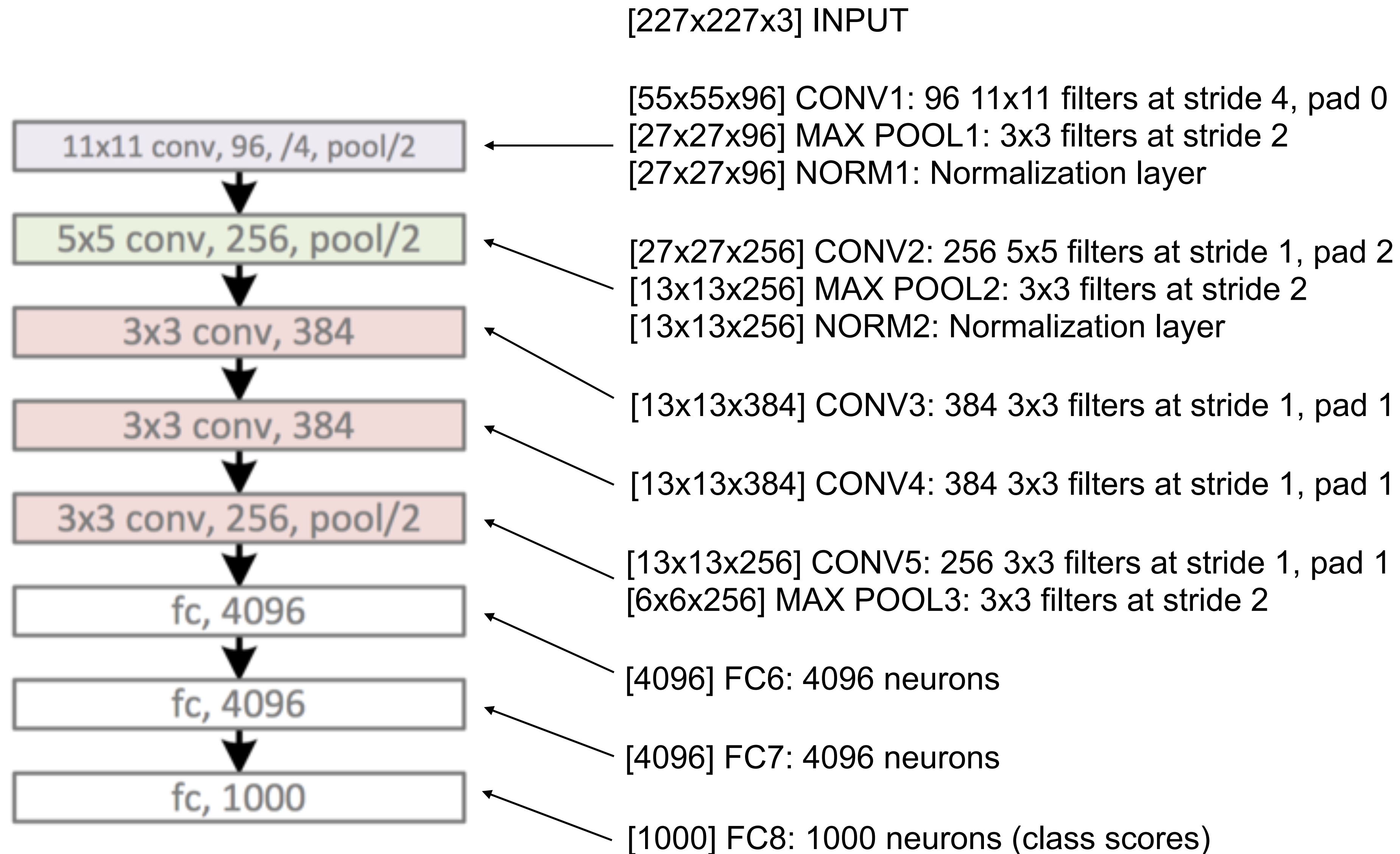


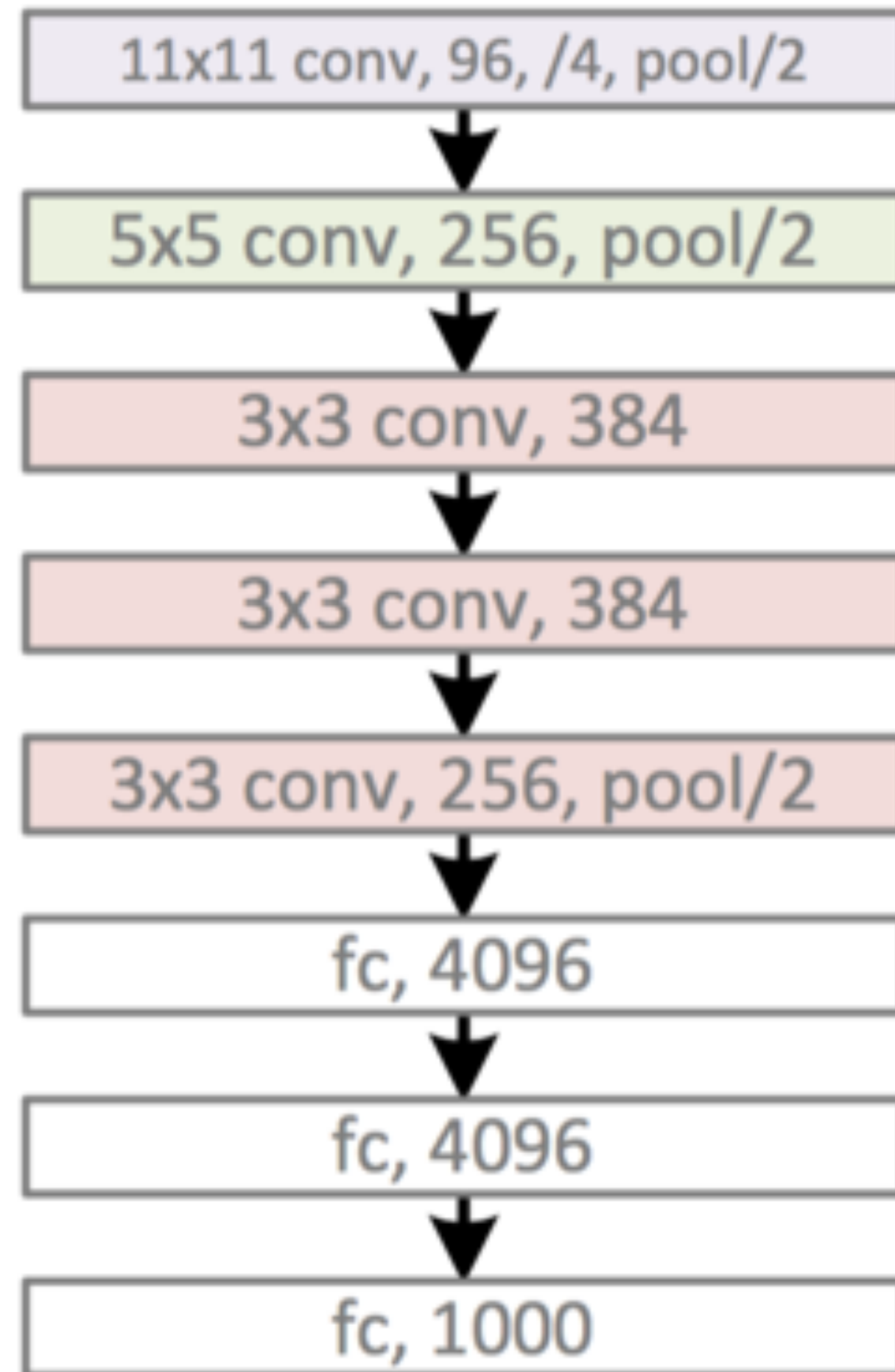
2012: AlexNet
5 conv. layers



Error: 16.4%

Alexnet — [Krizhevsky et al. NIPS 2012]

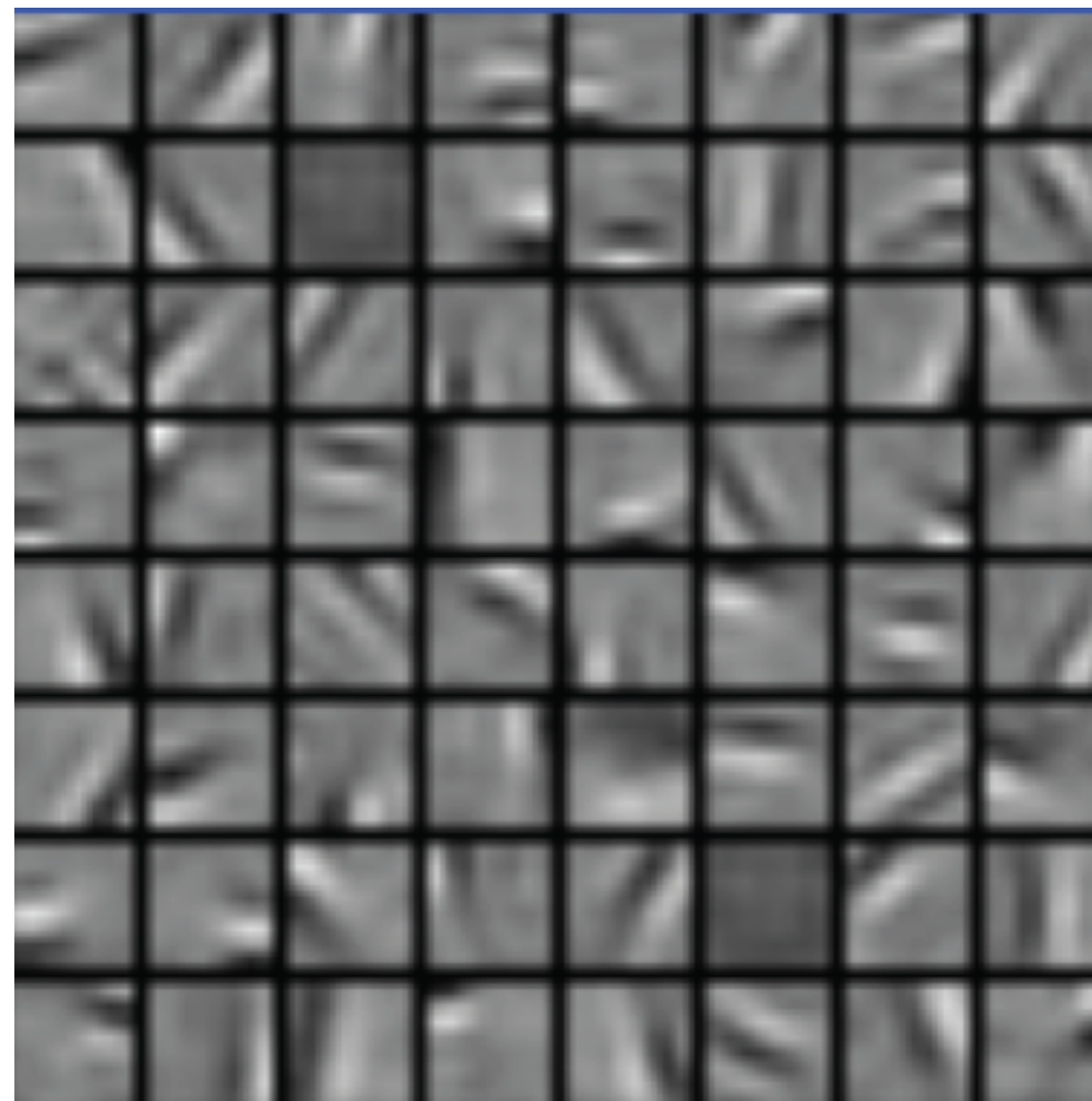




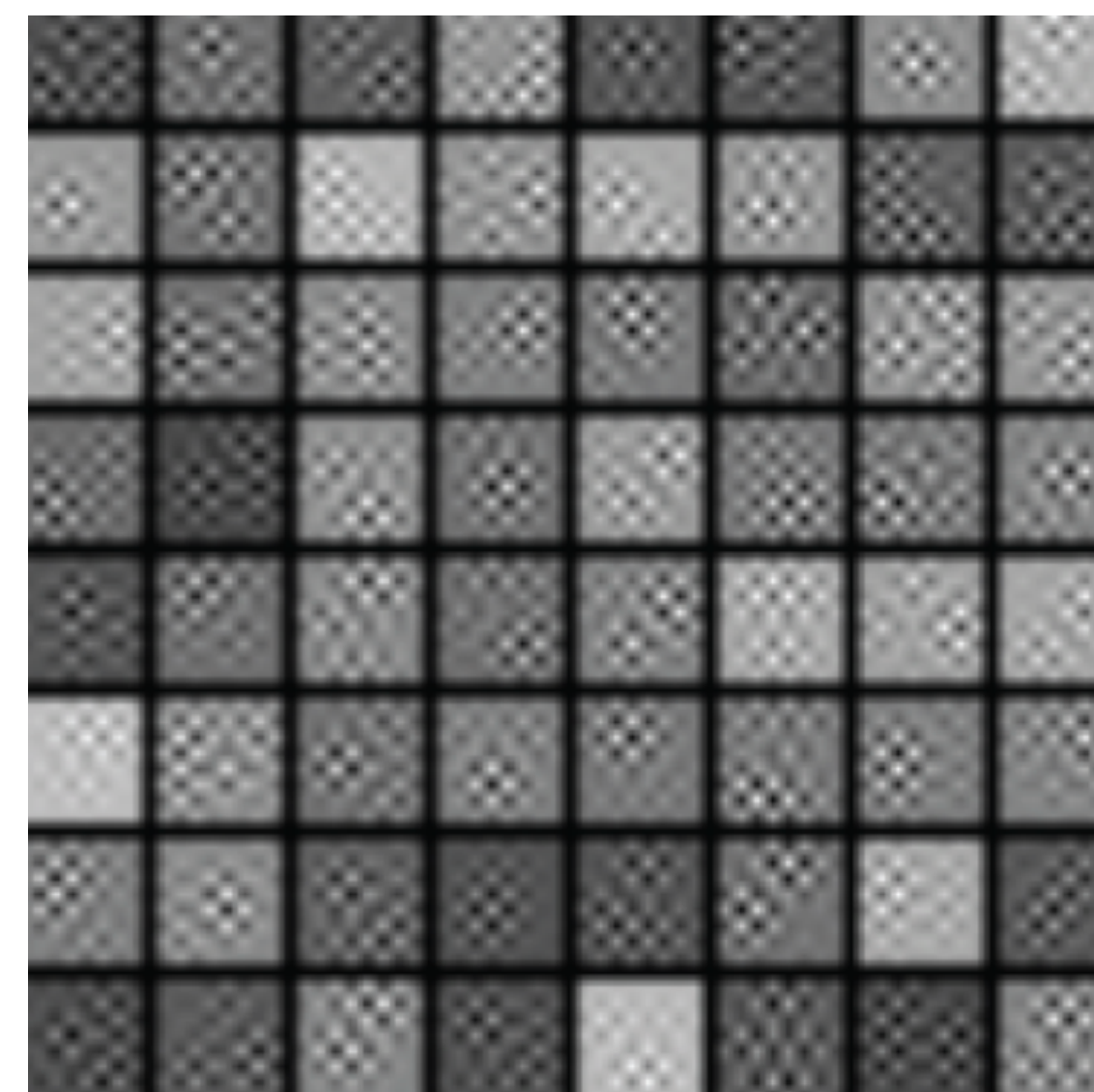
What filters are learned?

What filters are learned?

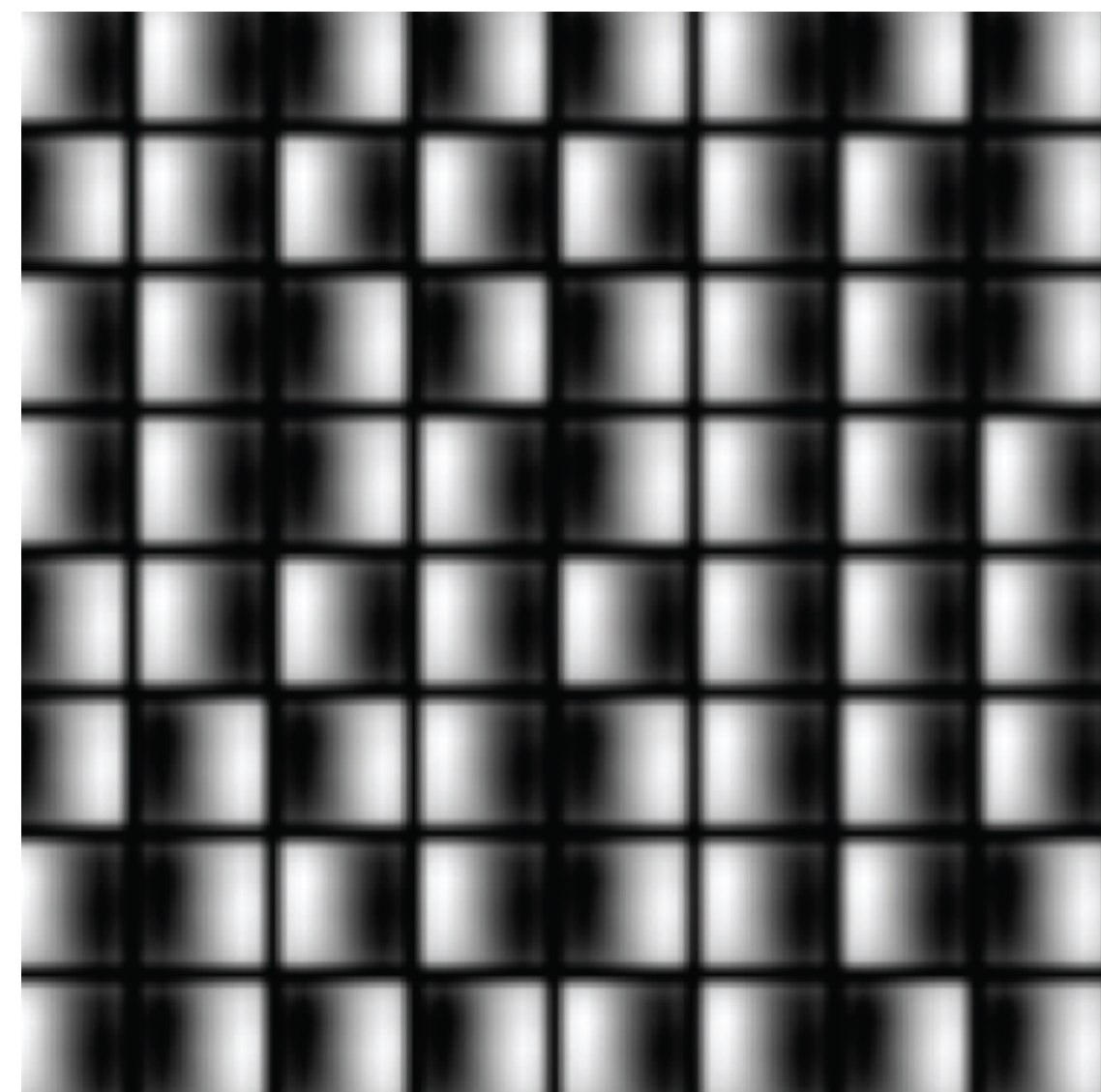
A



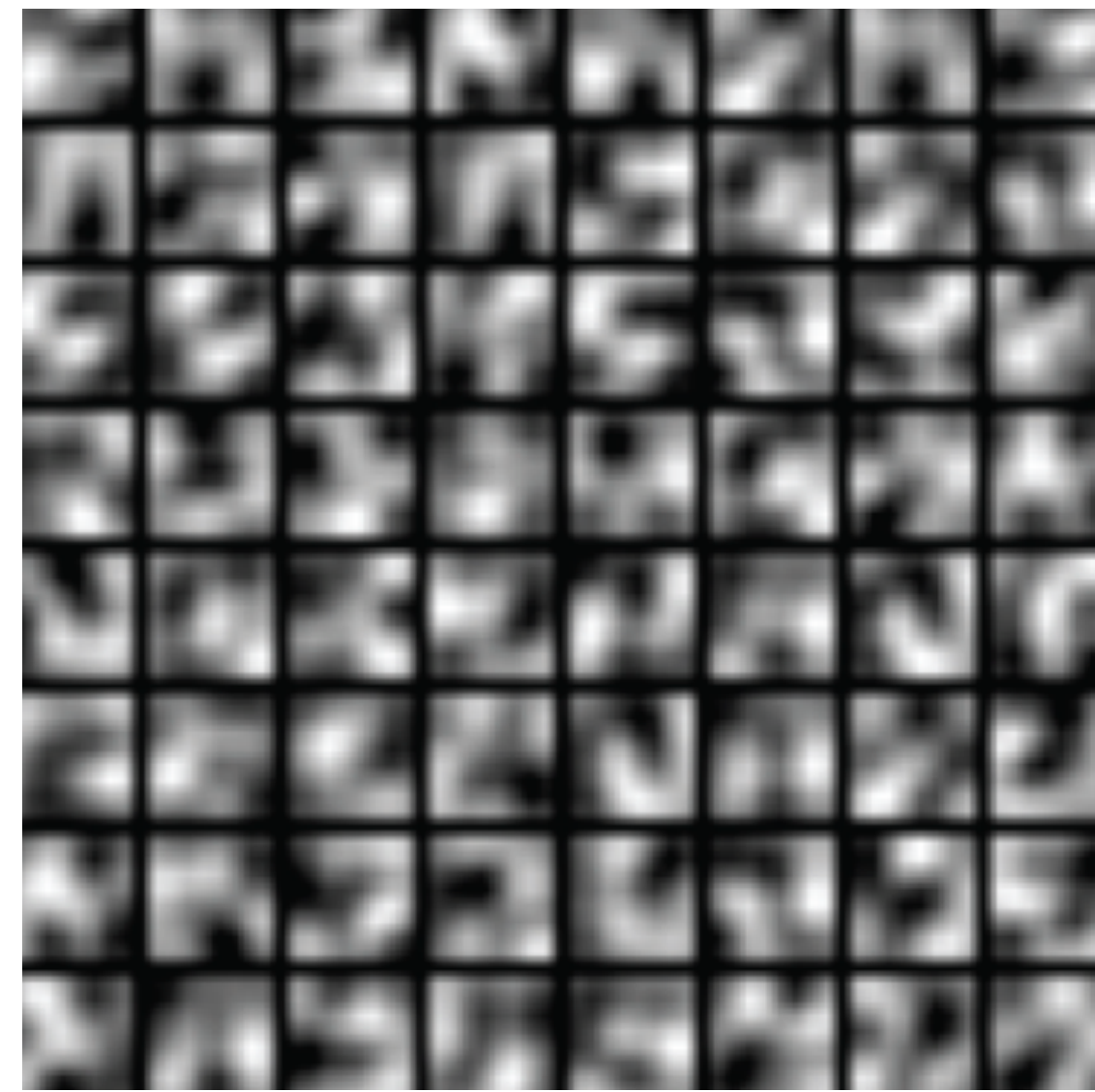
B



C



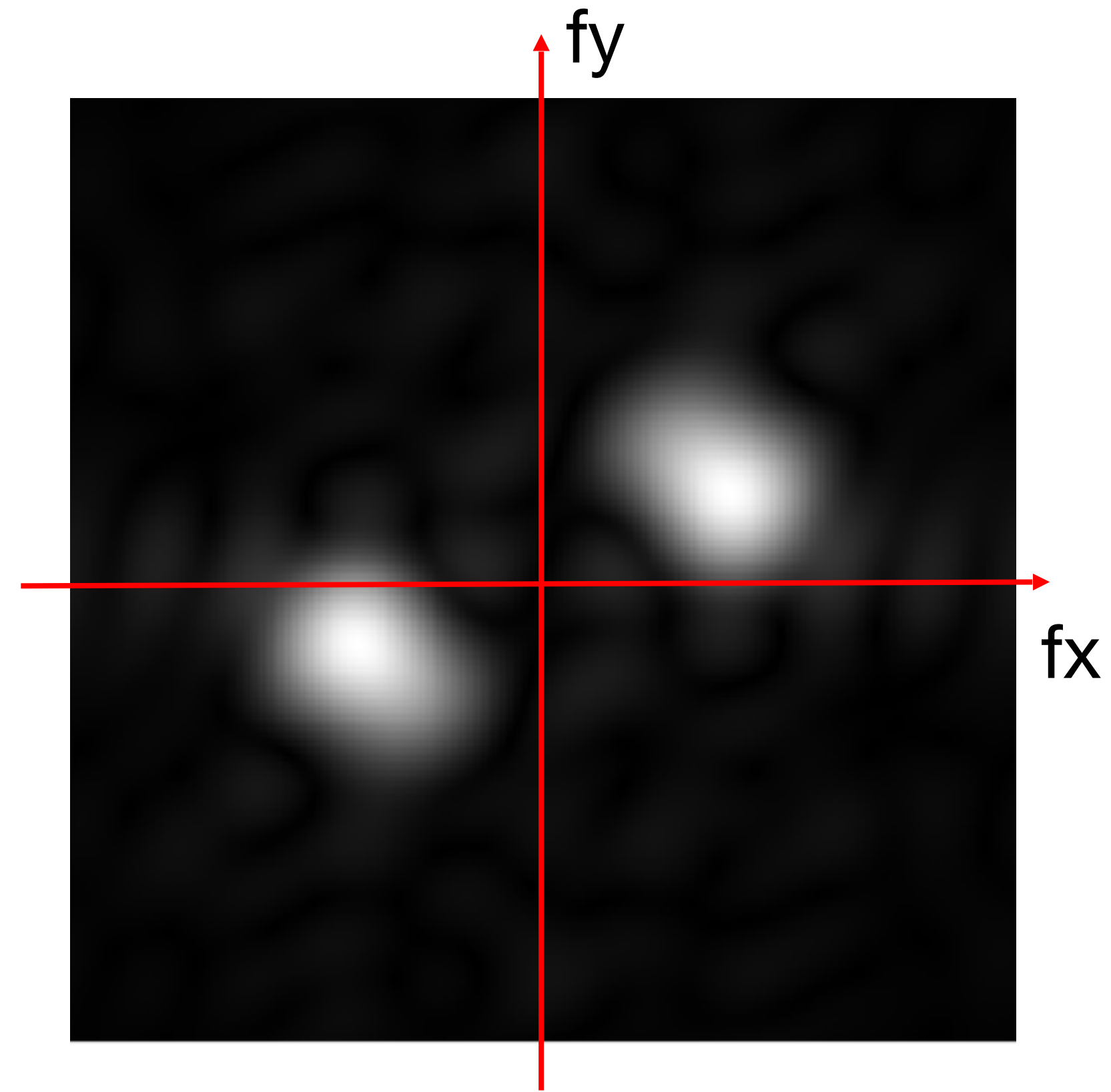
D



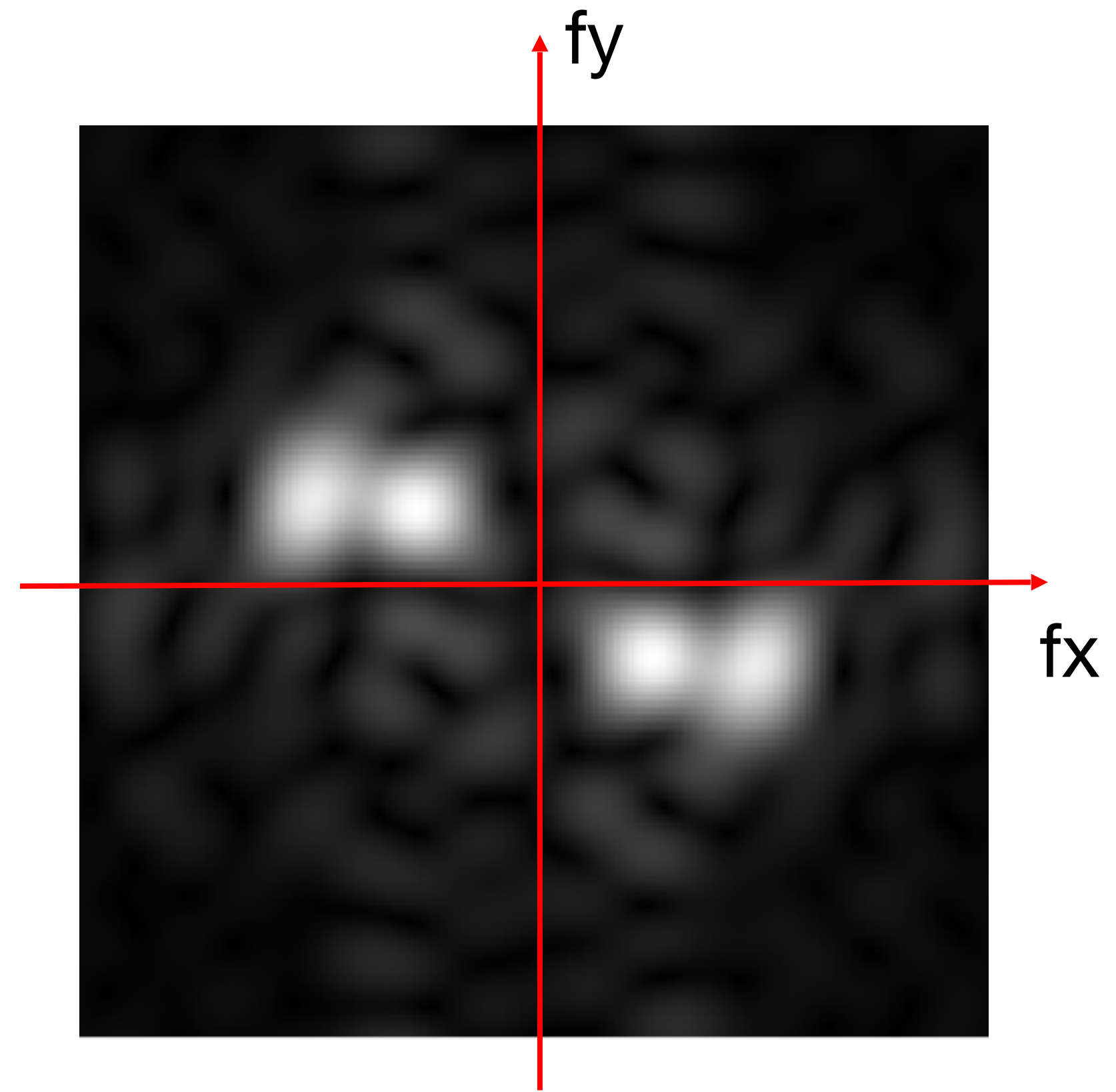
Get to know your units



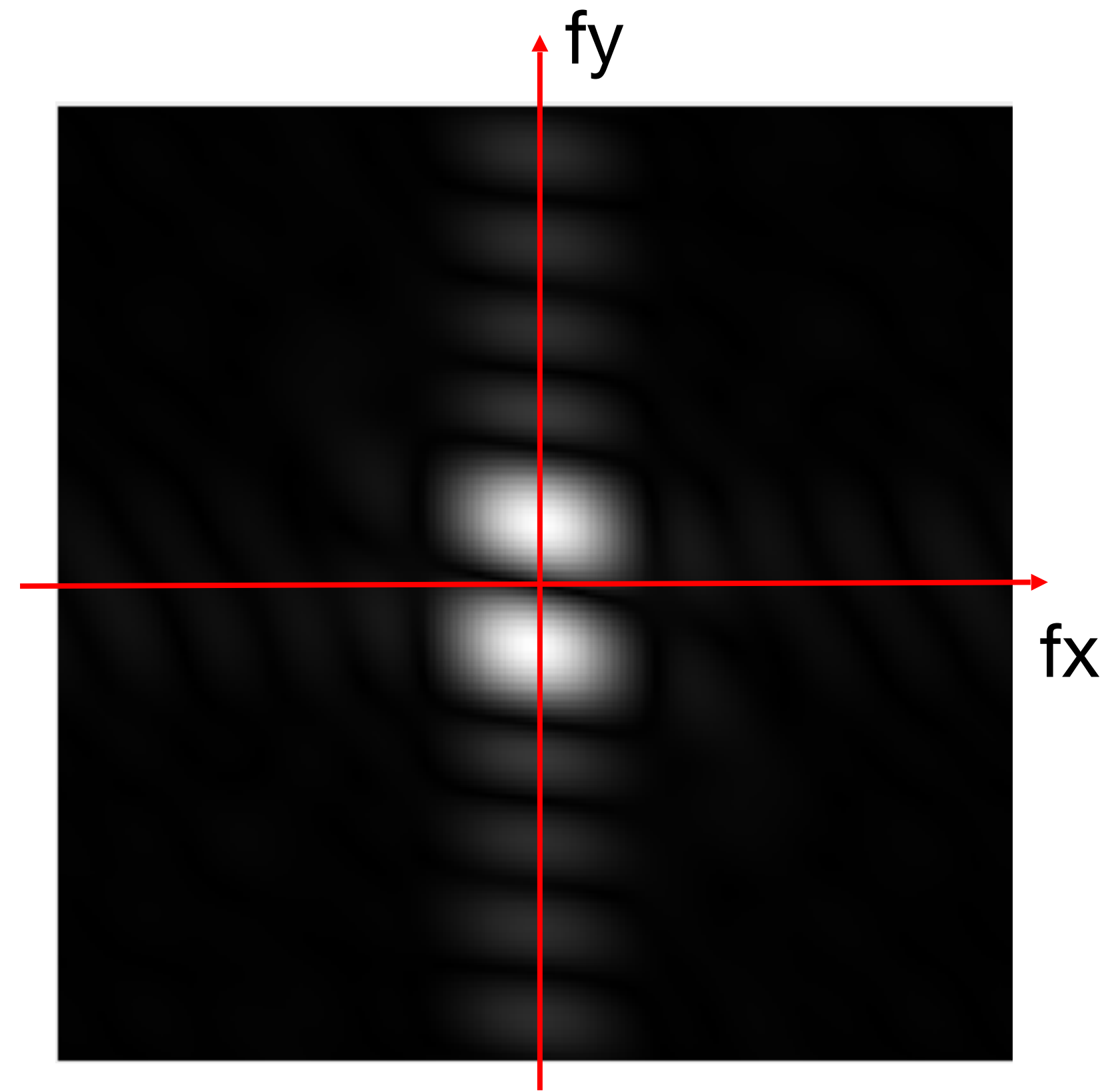
11x11 convolution kernel
(3 color channels)



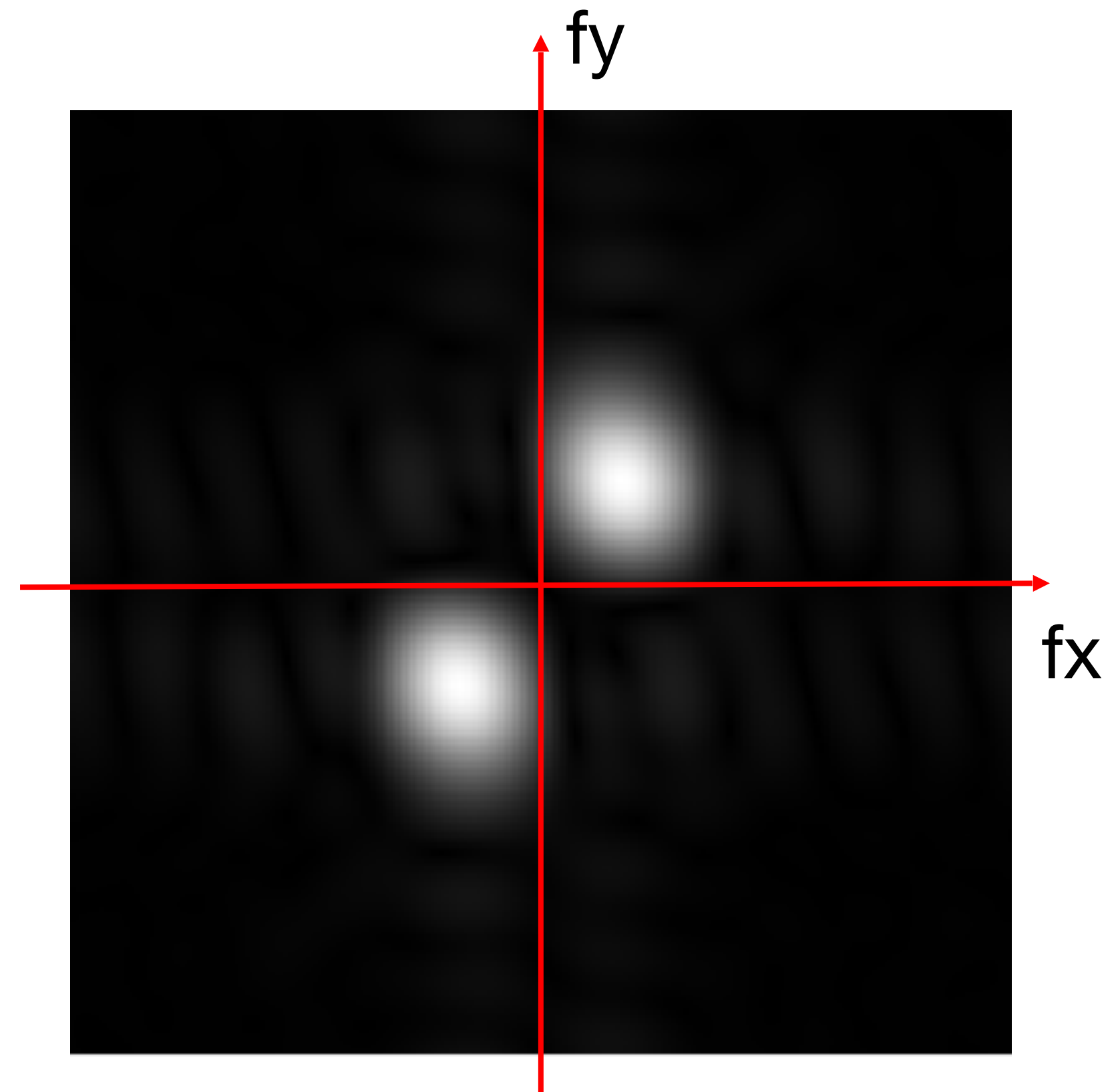
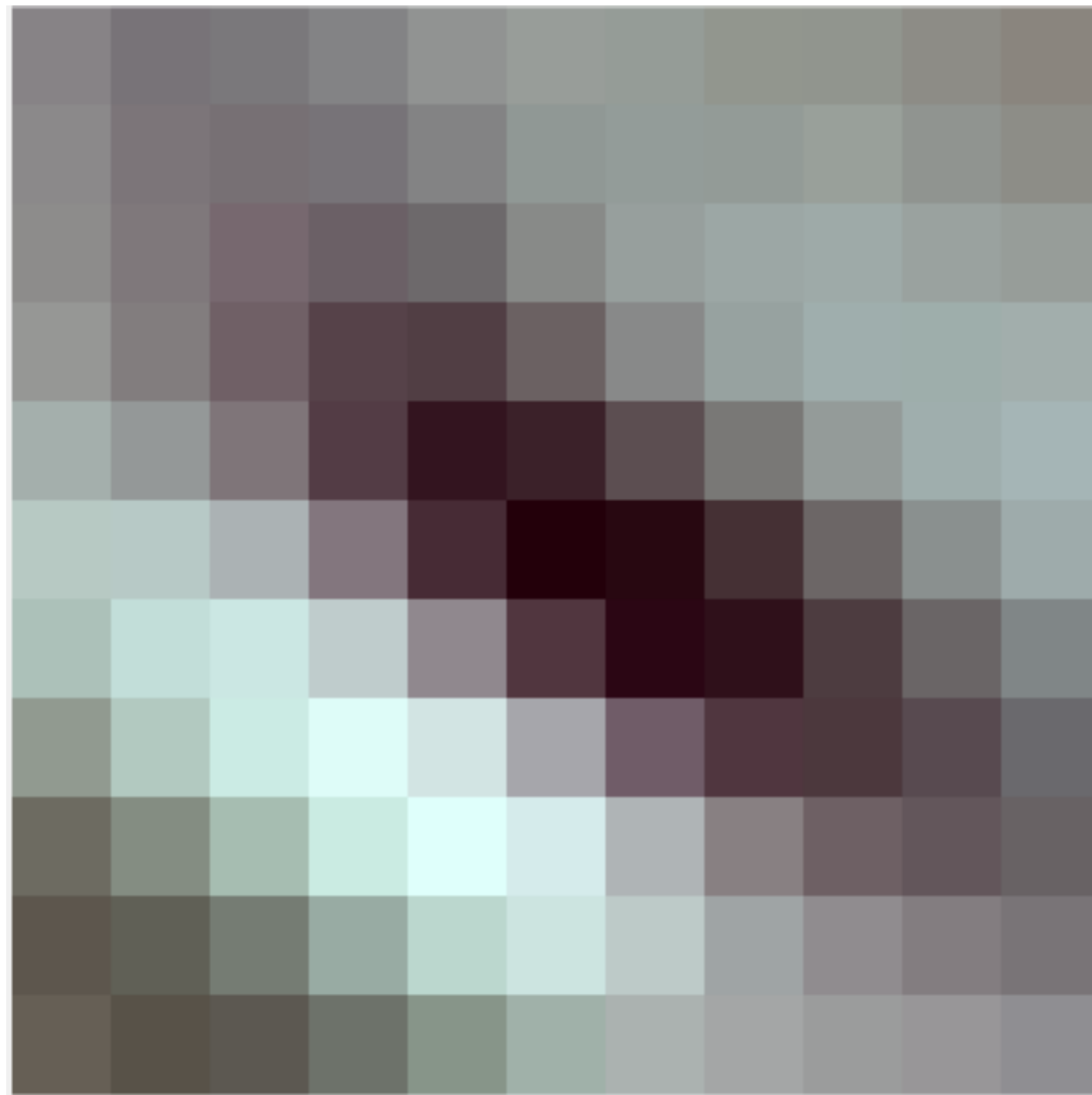
Get to know your units



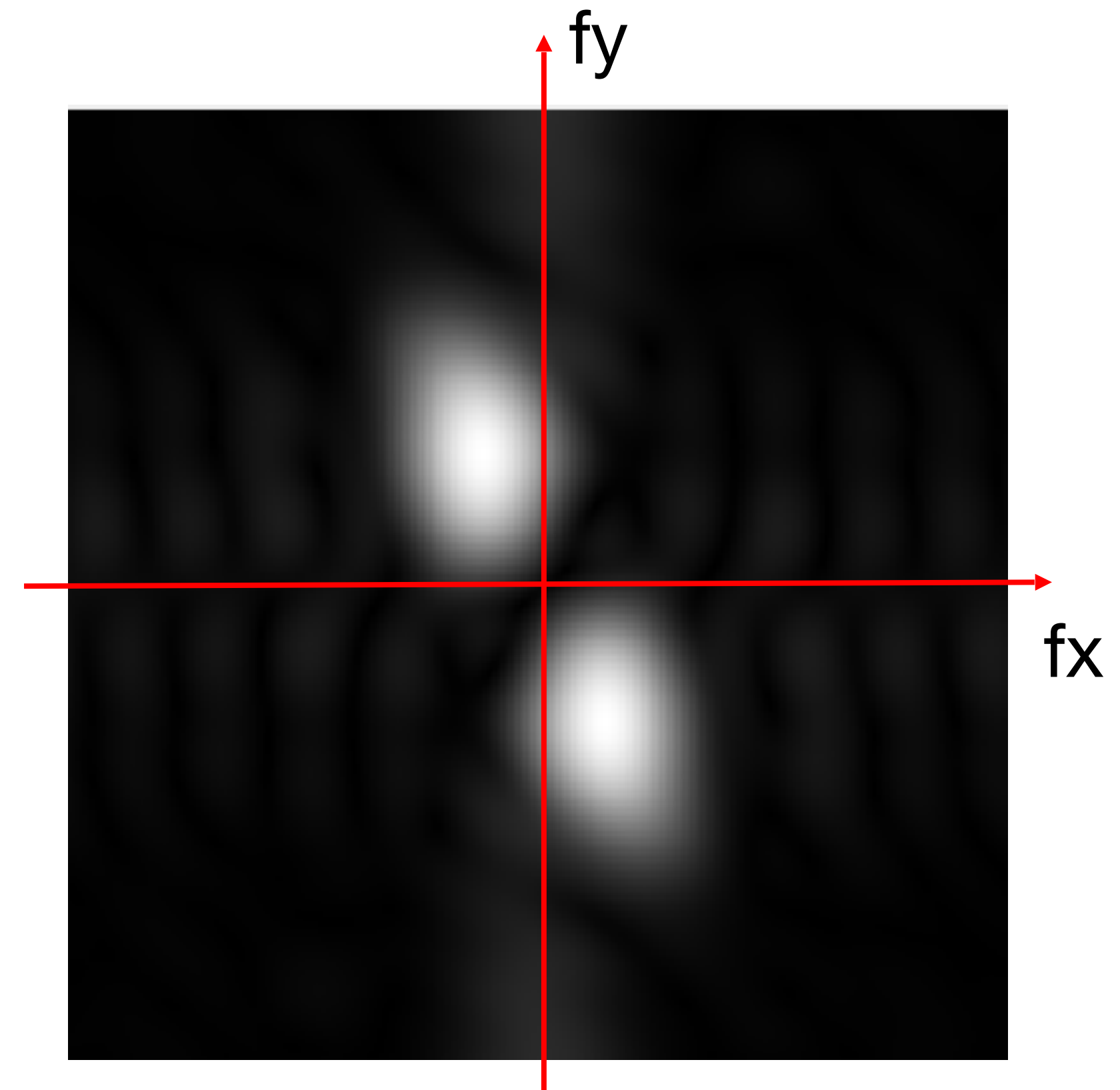
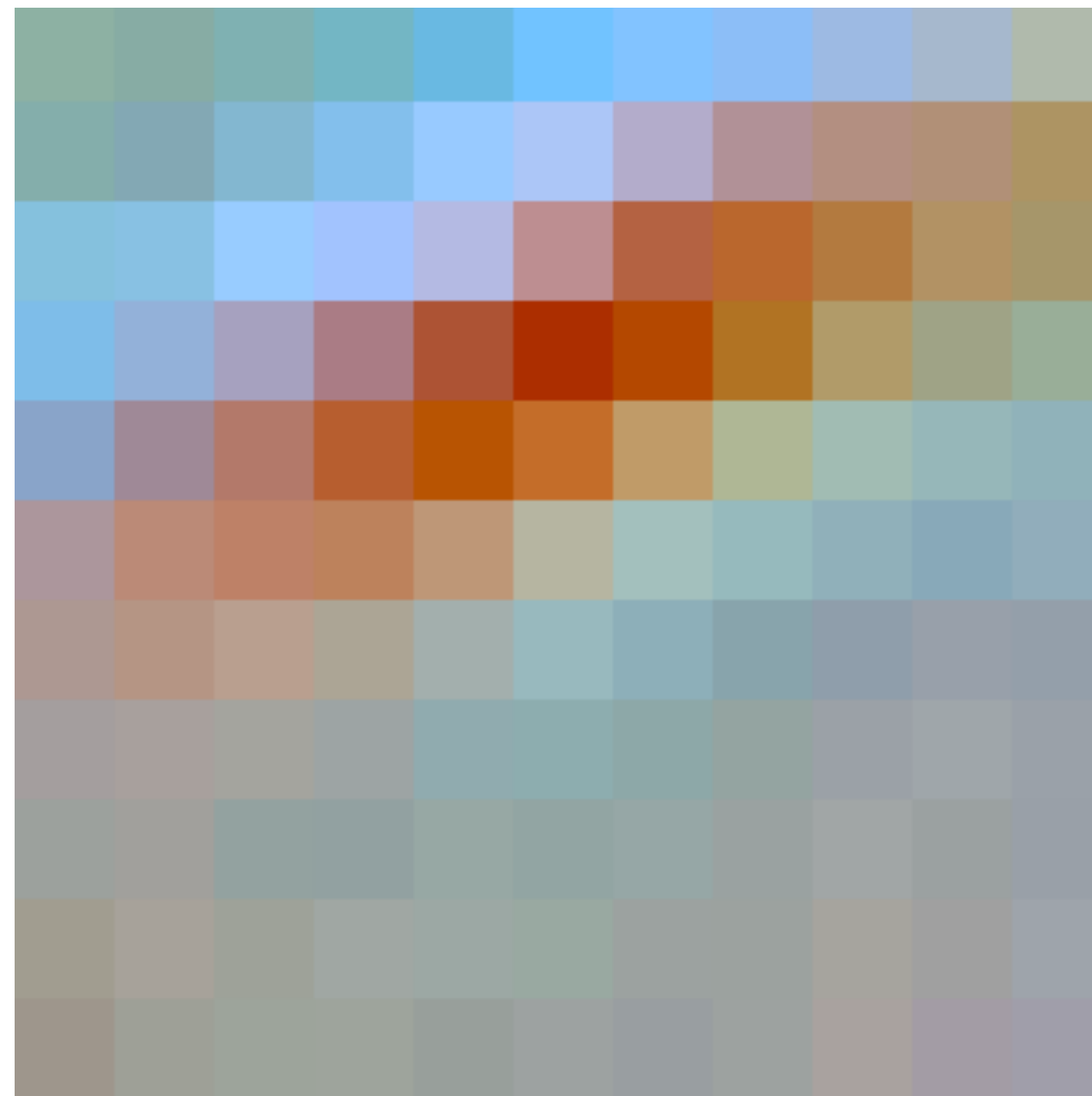
Get to know your units



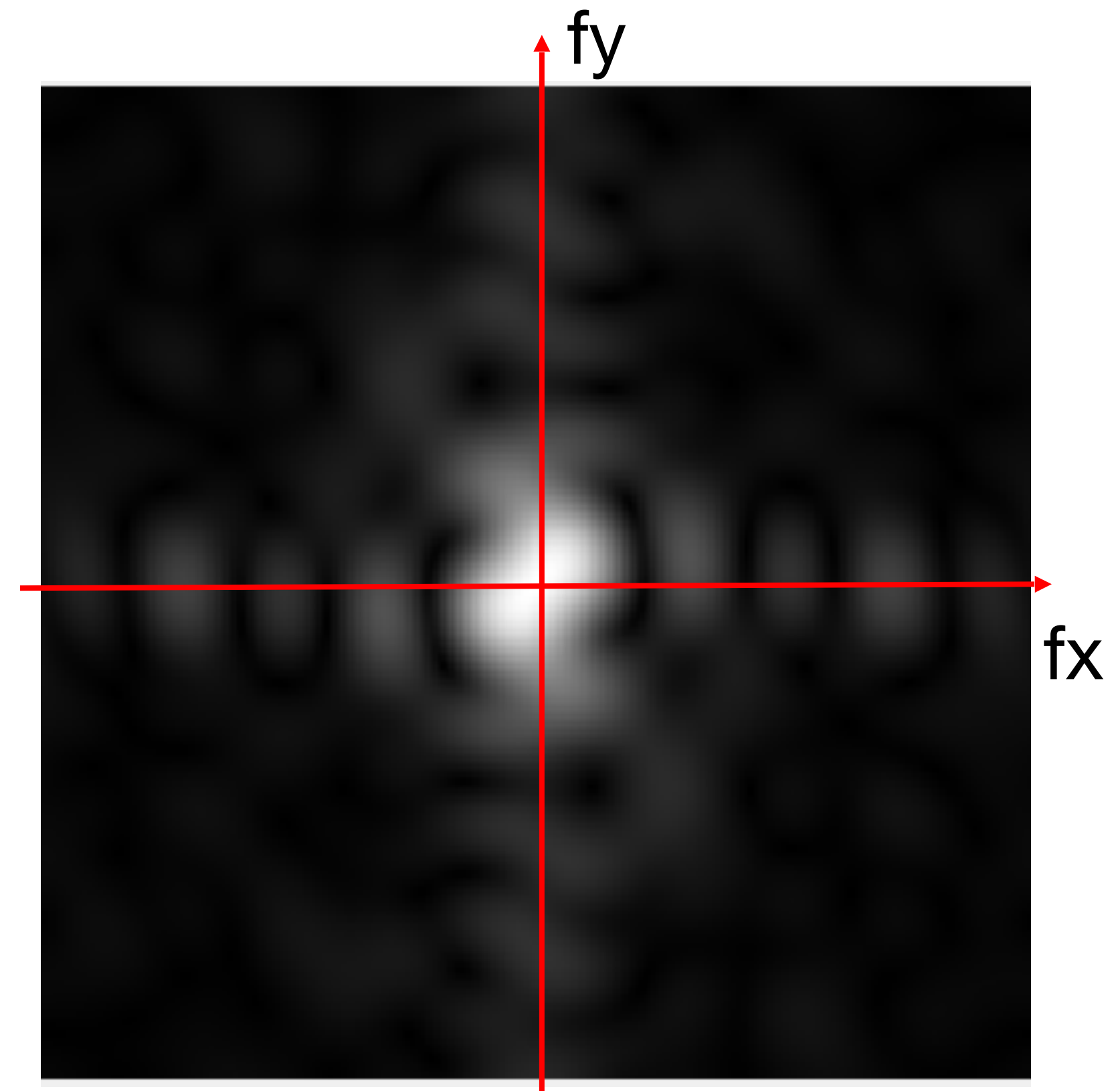
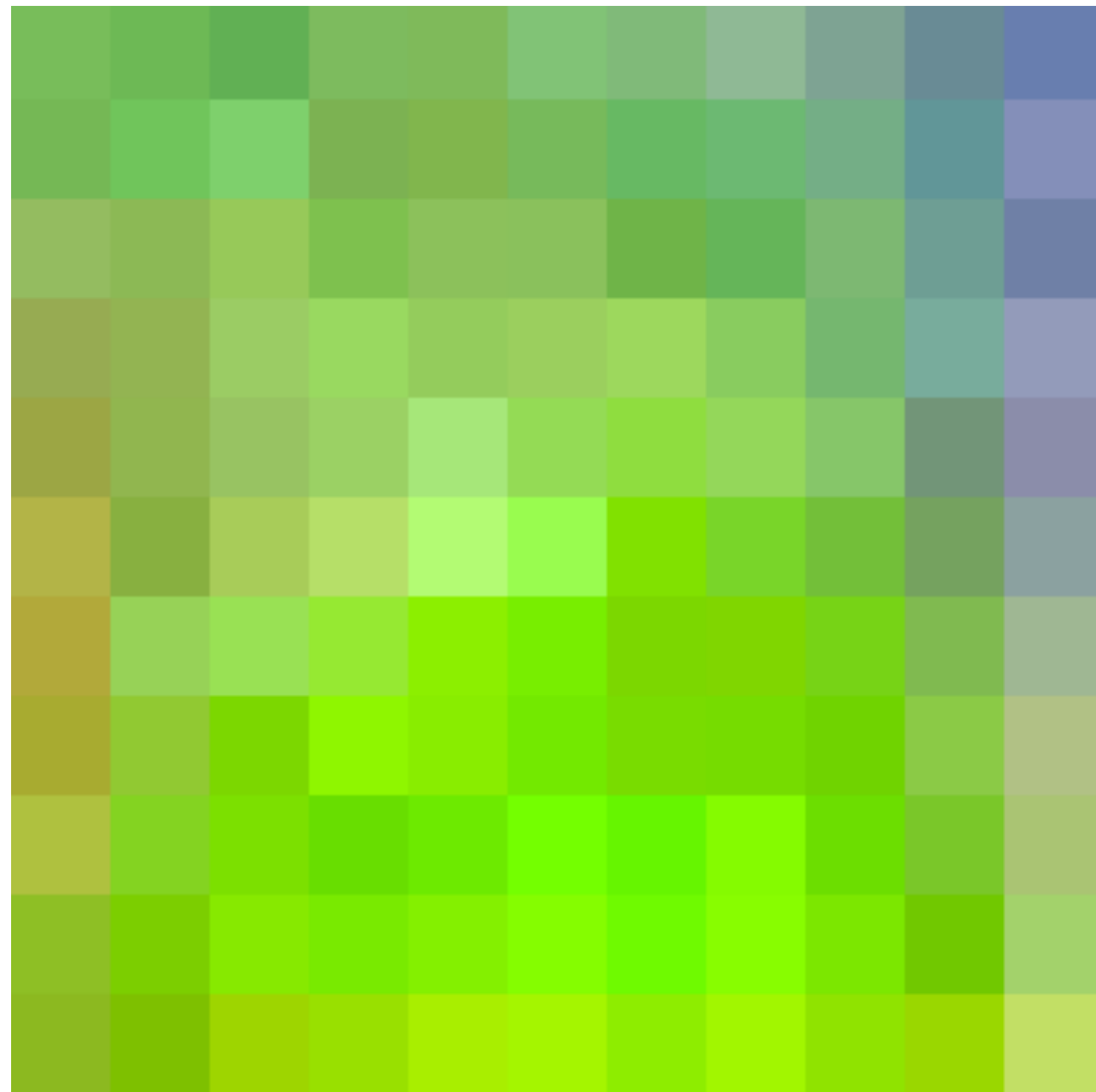
Get to know your units



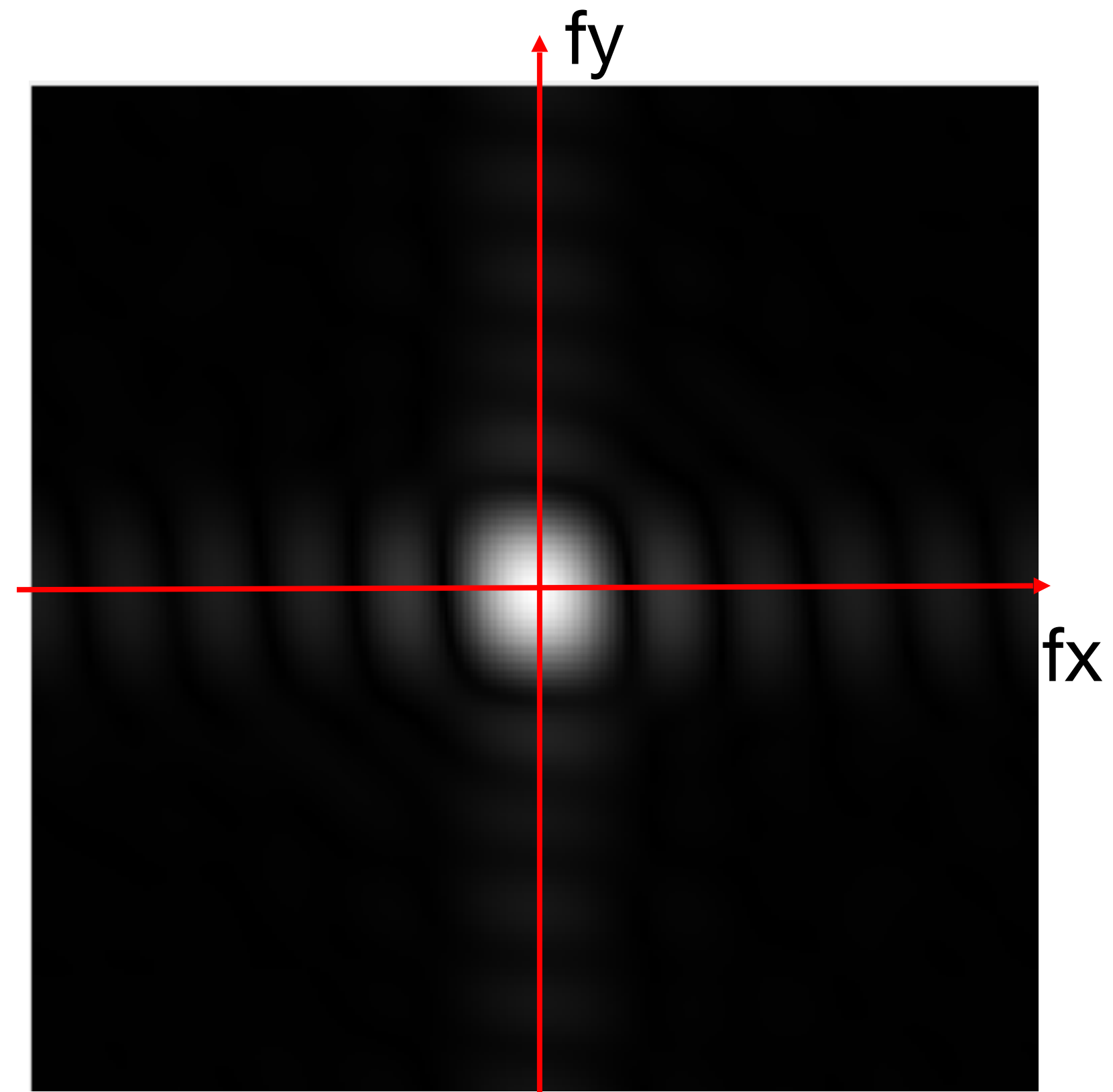
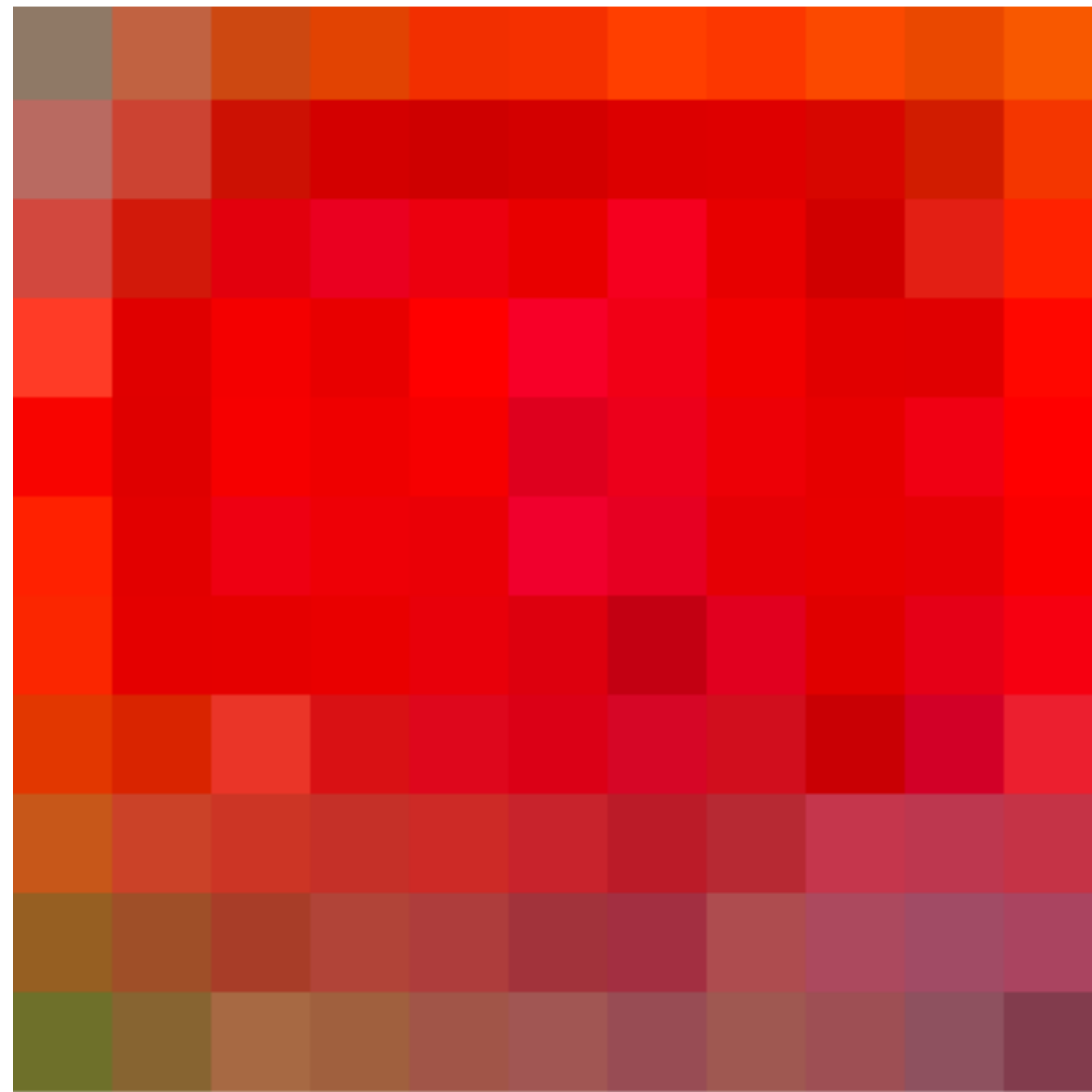
Get to know your units



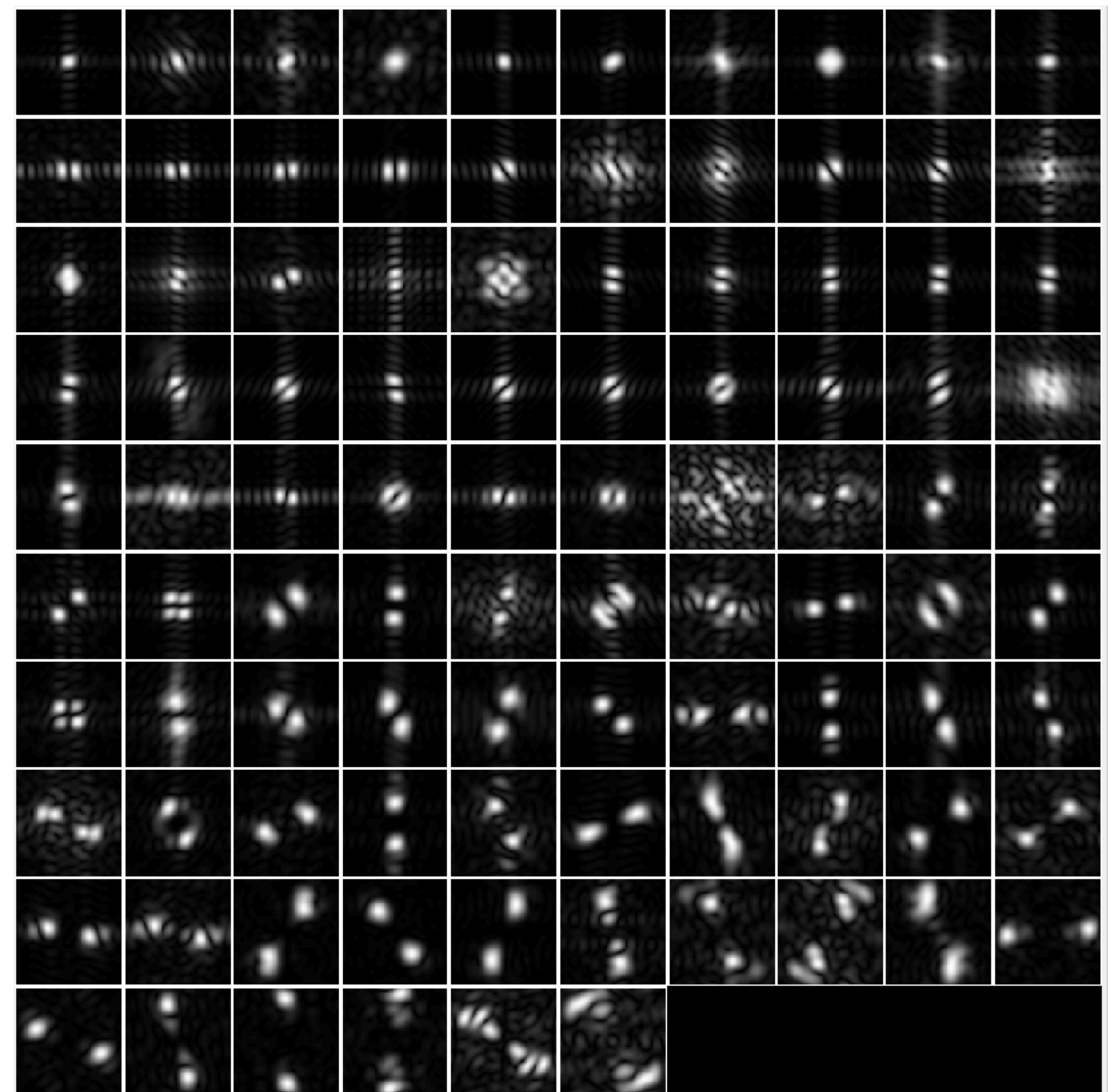
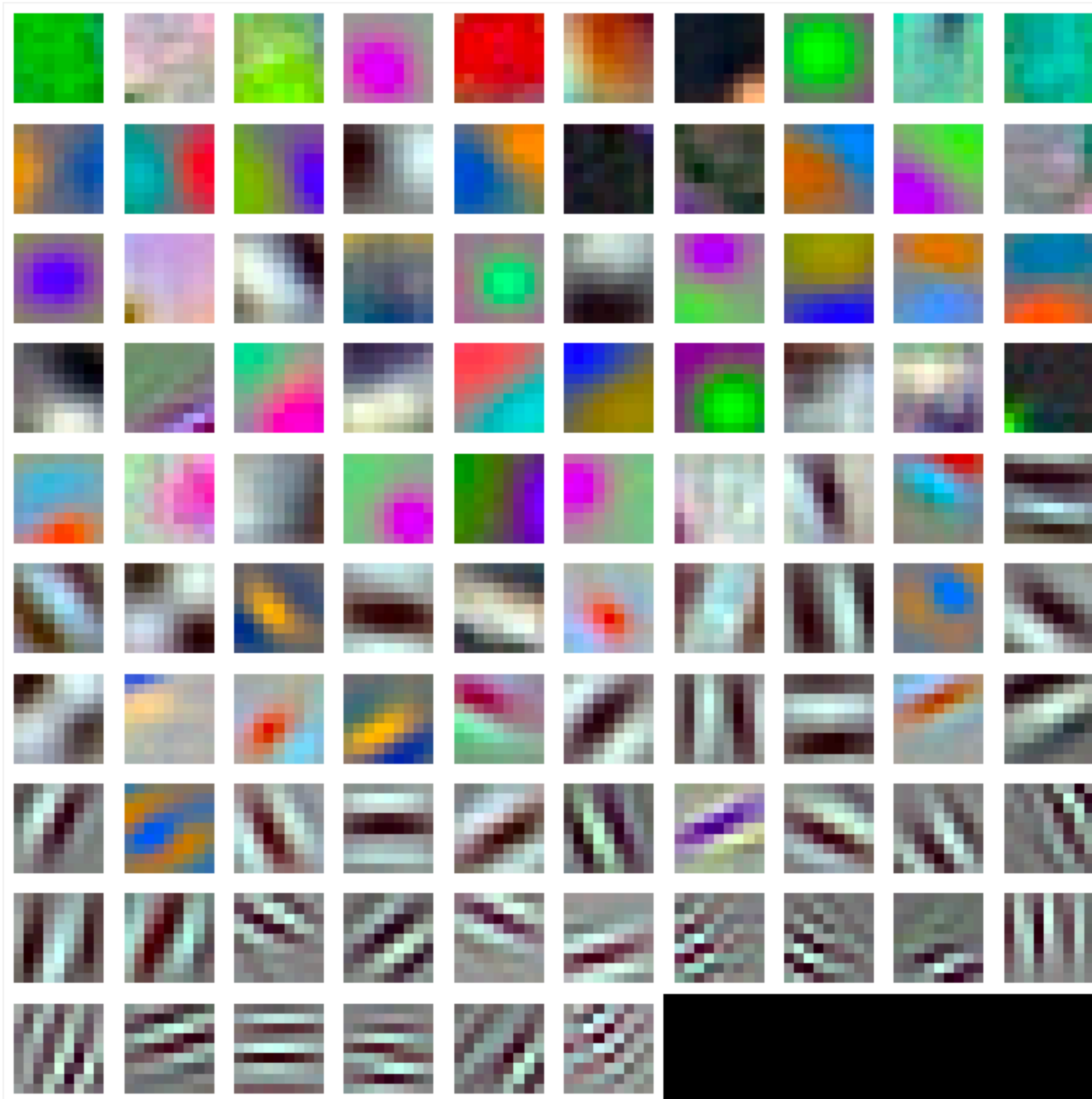
Get to know your units



Get to know your units

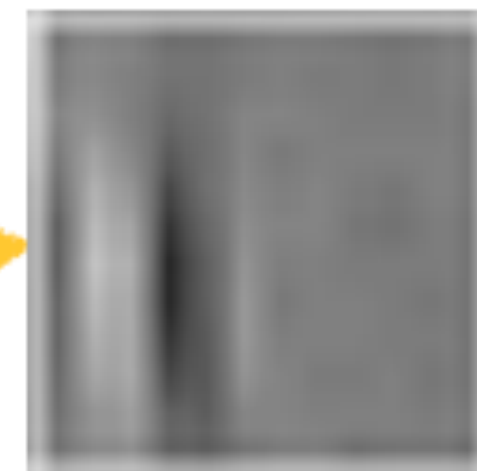
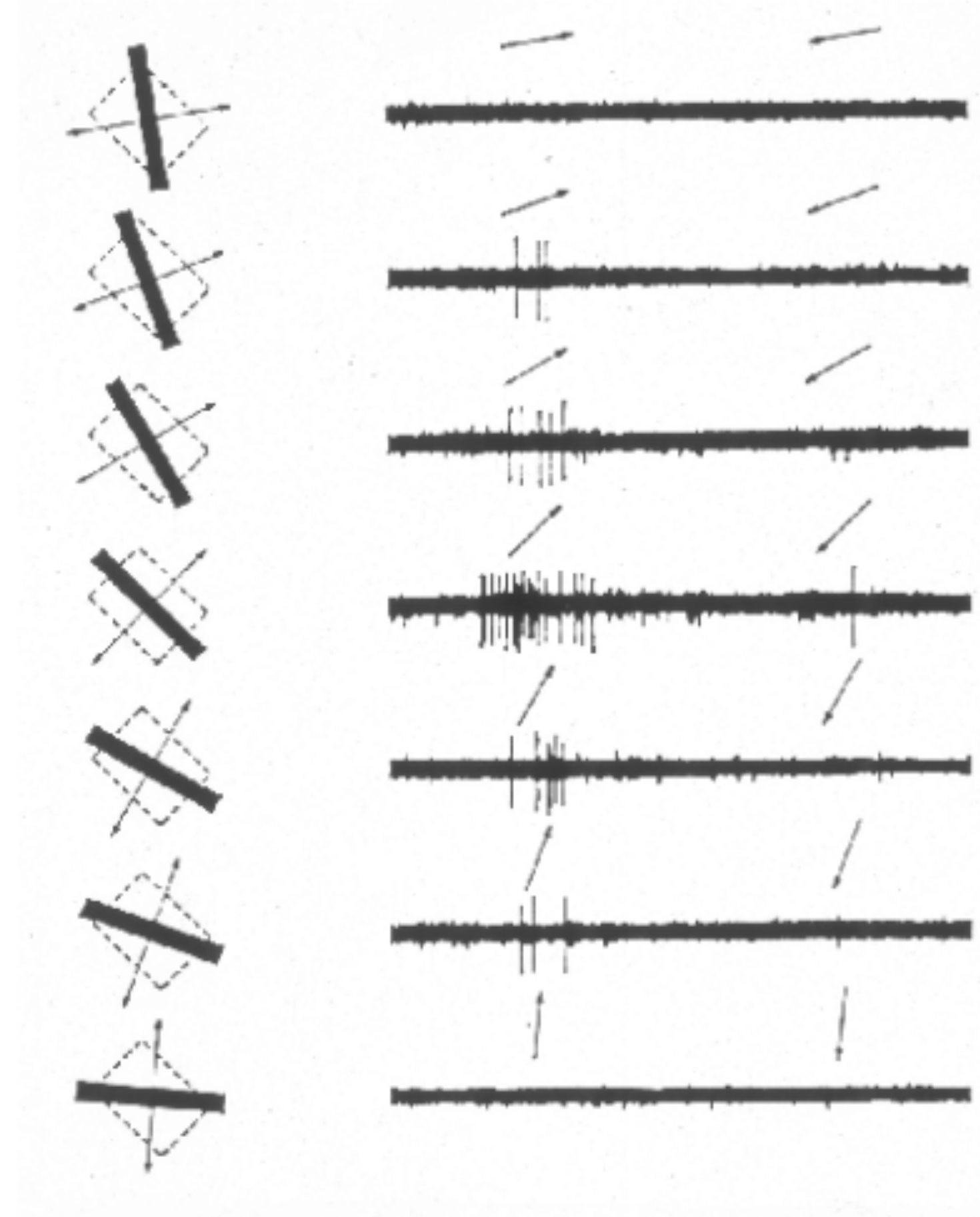
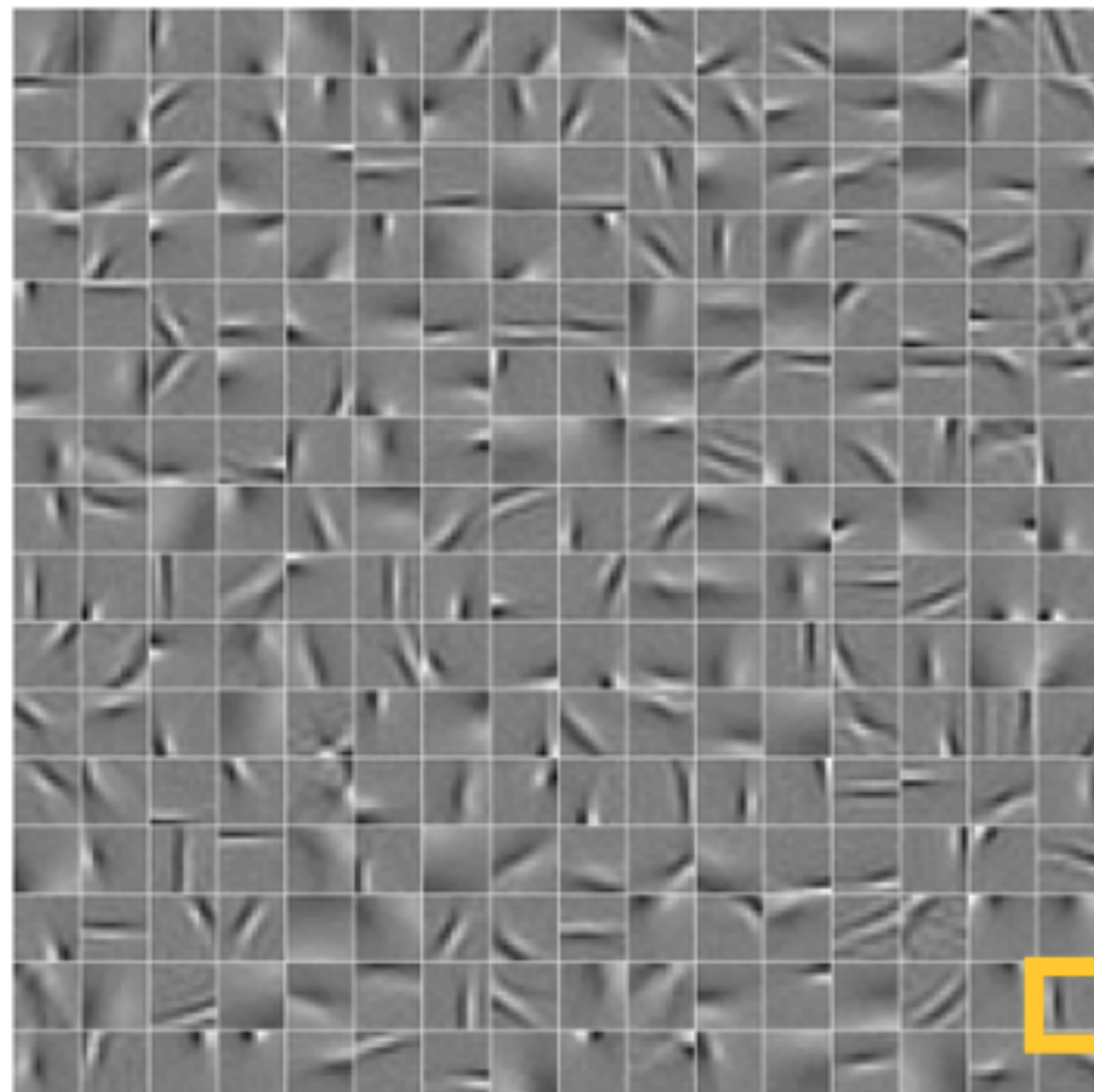
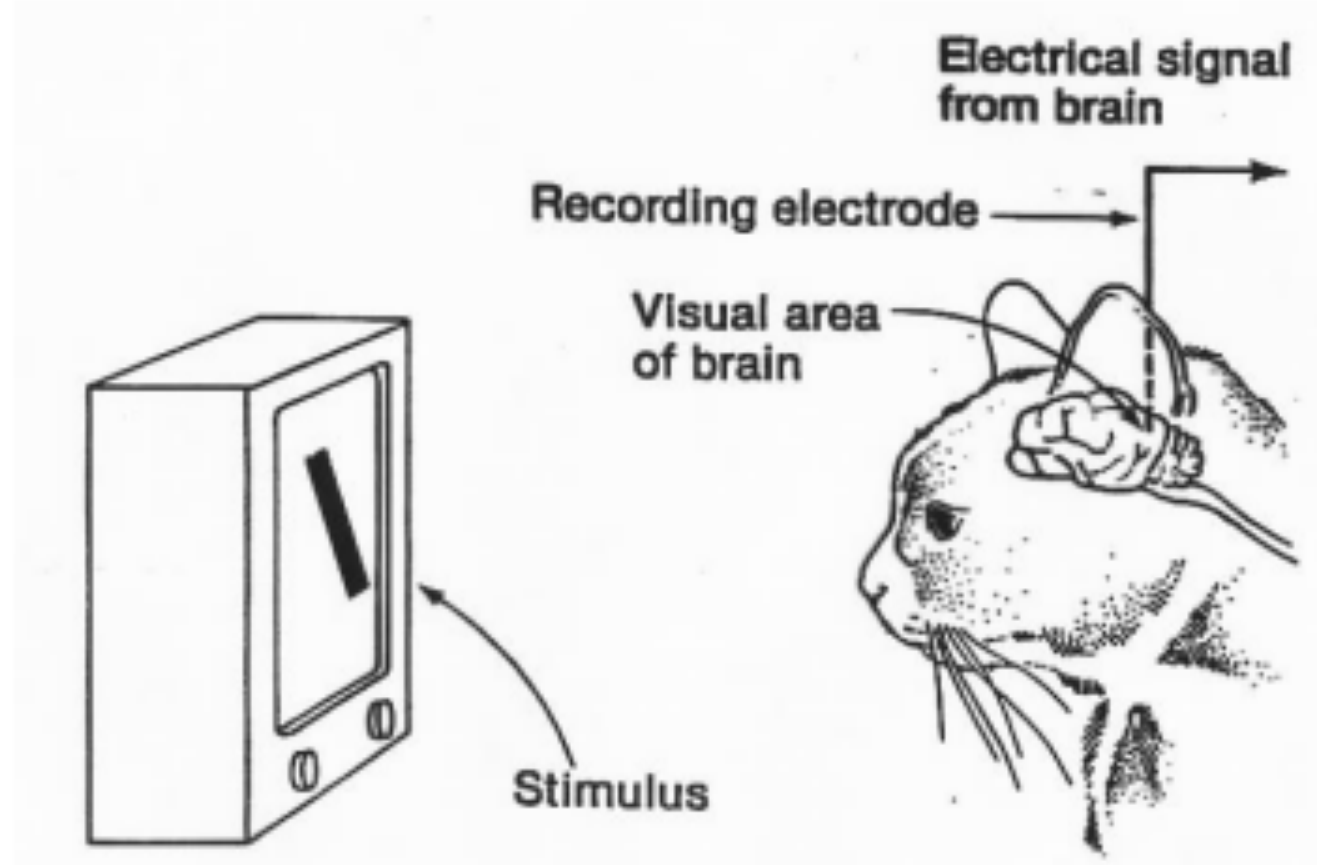


Get to know your units



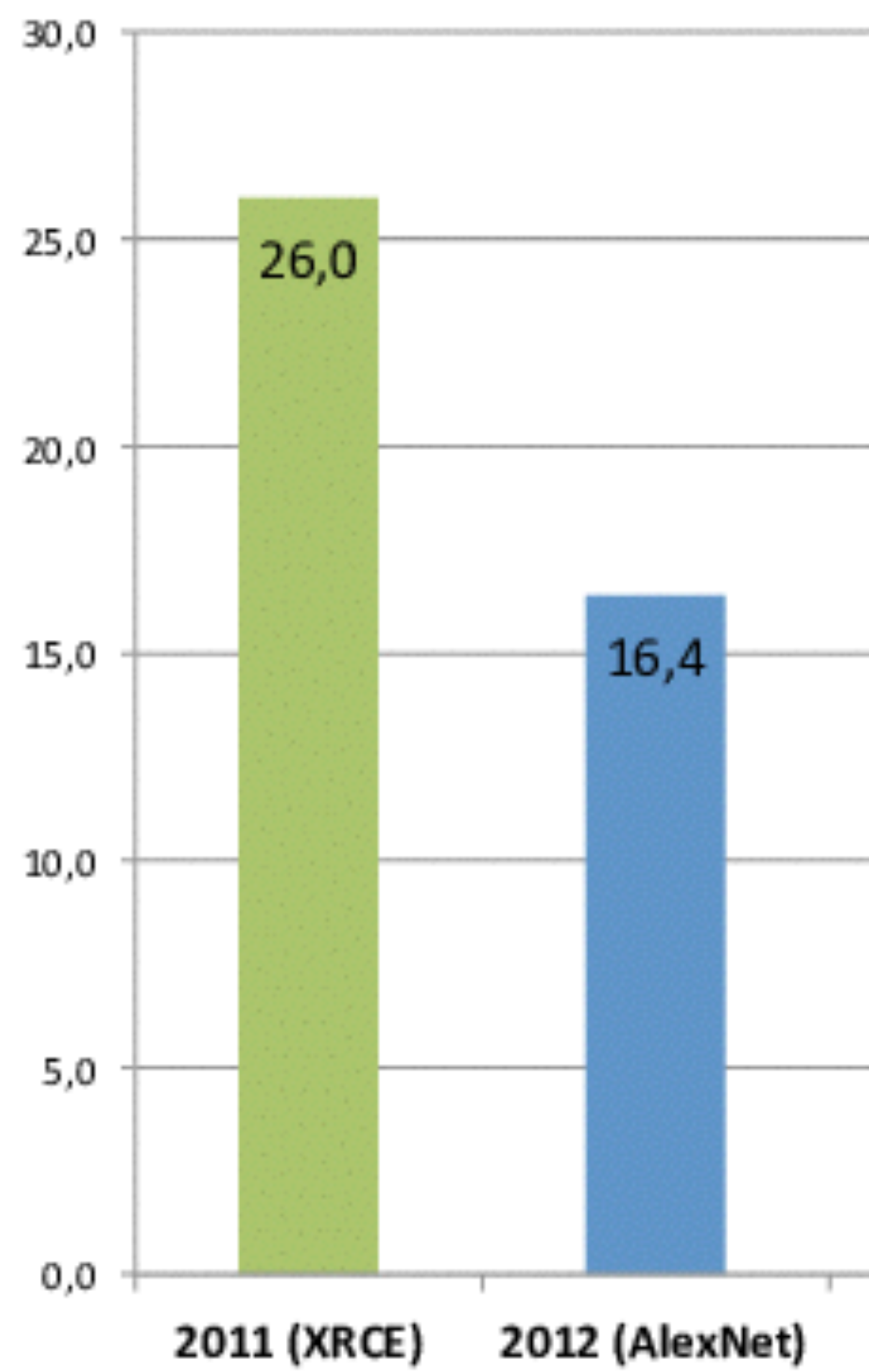
96 Units in conv1

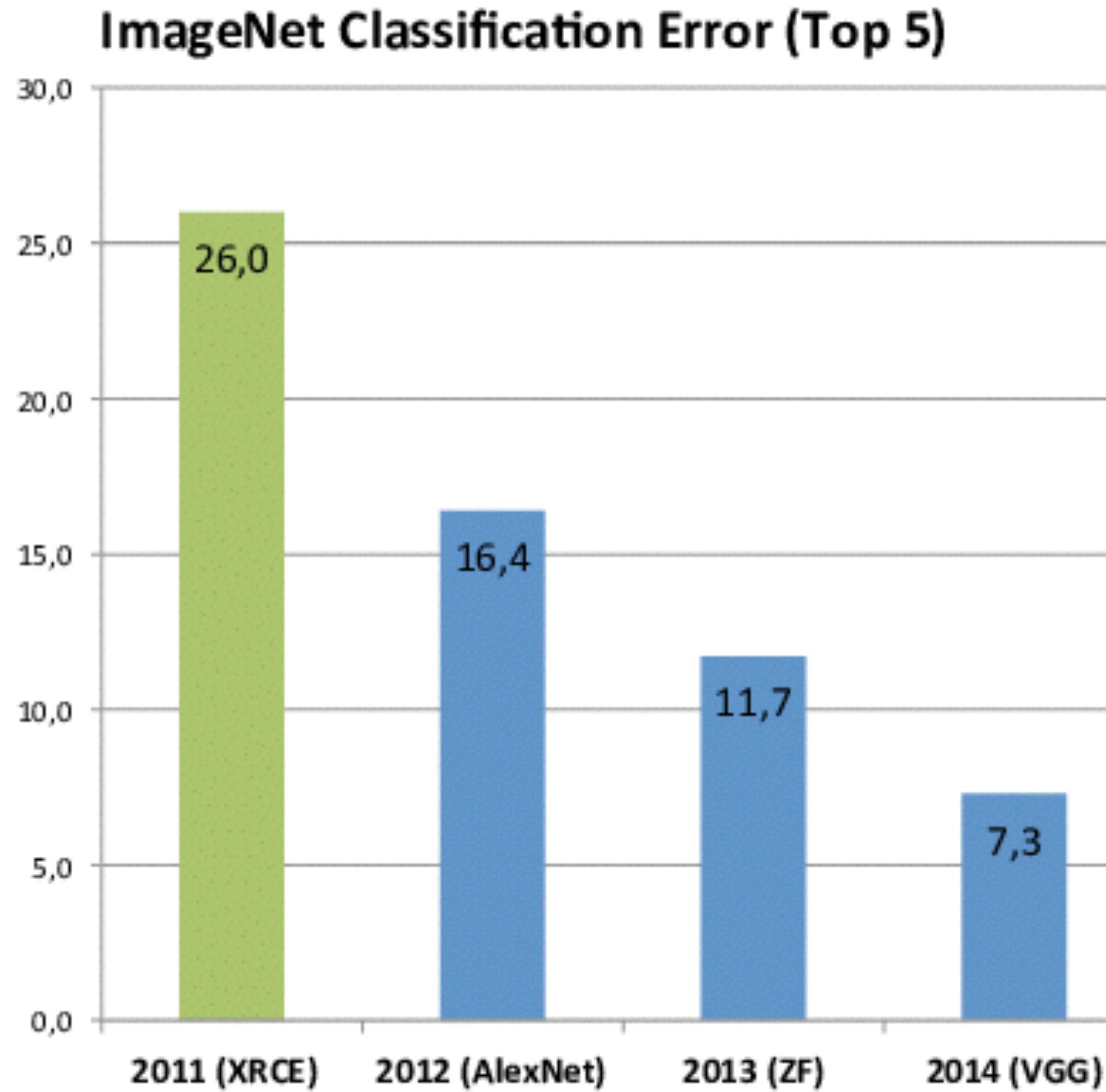
[Hubel and Wiesel 59]



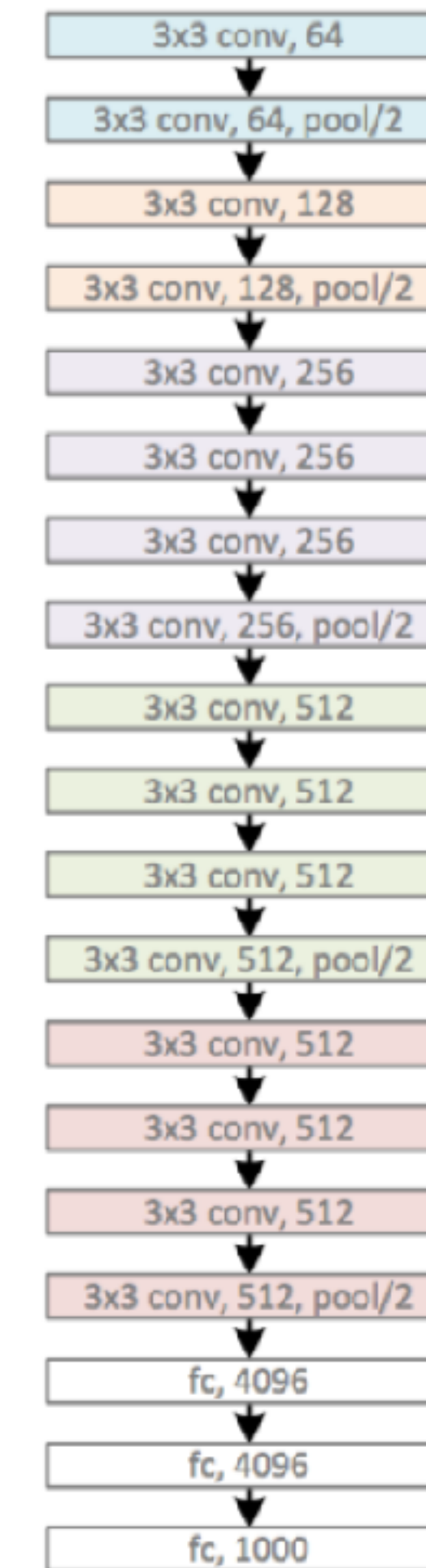
oriented filter

ImageNet Classification Error (Top 5)





2014: VGG
16 conv. layers

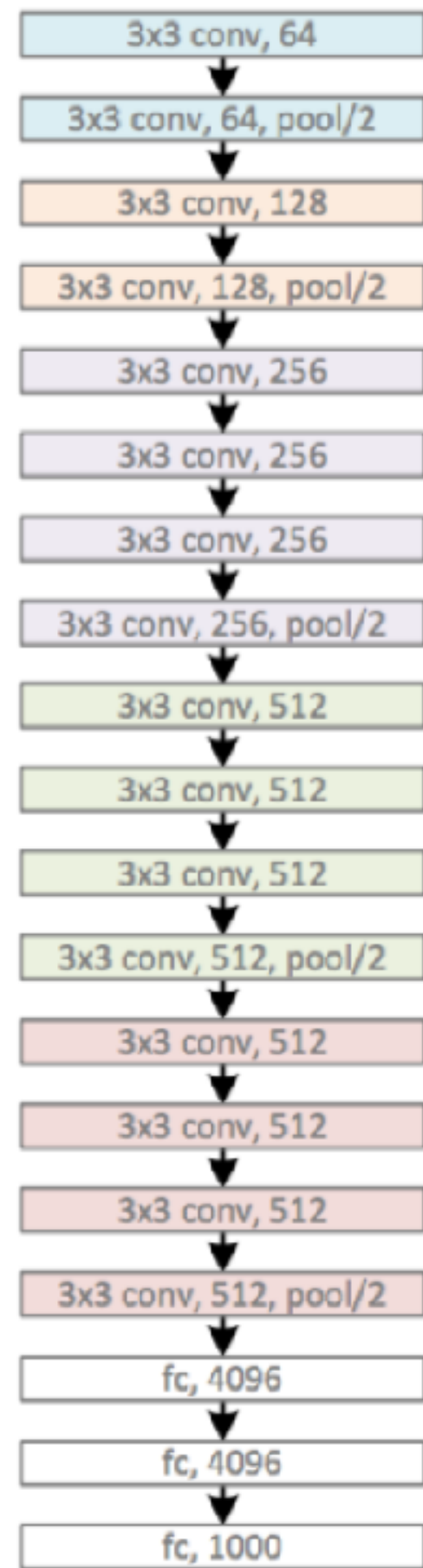


Error: 7.3%

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

VGG-Net [Simonyan & Zisserman, 2015]

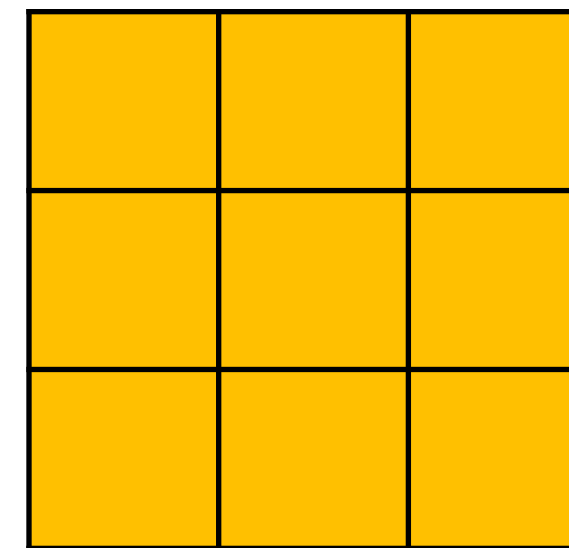
2014: VGG
16 conv. layers



Error: 7.3%

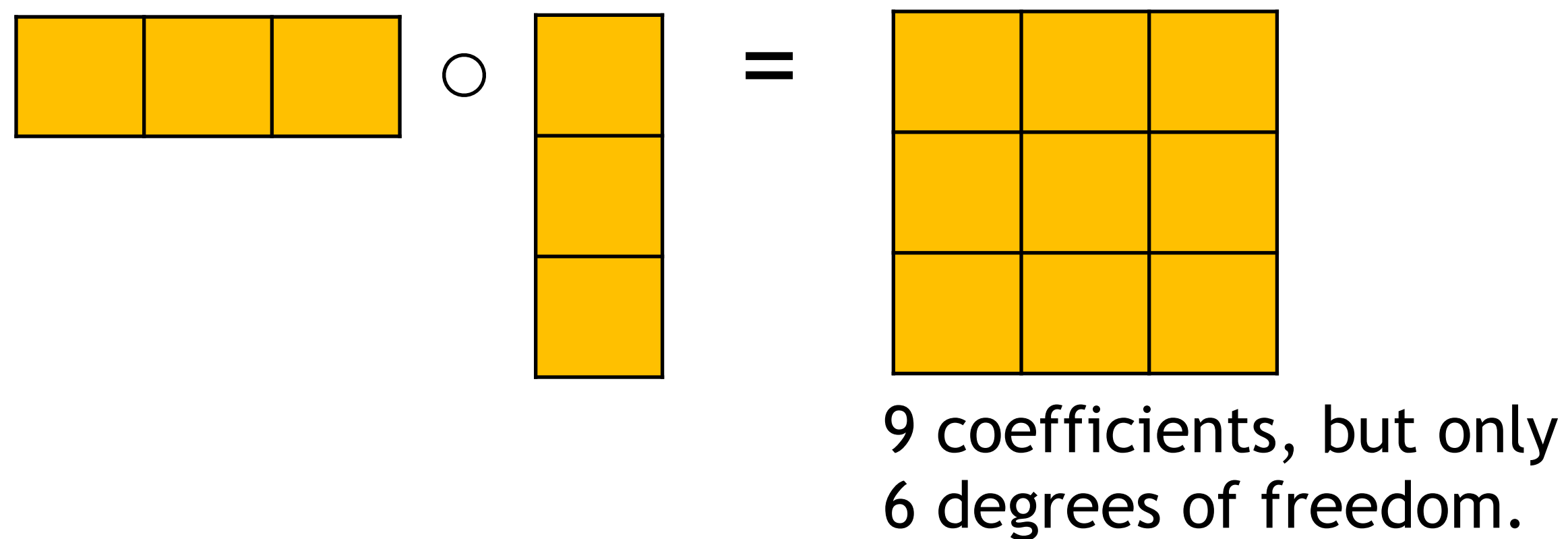
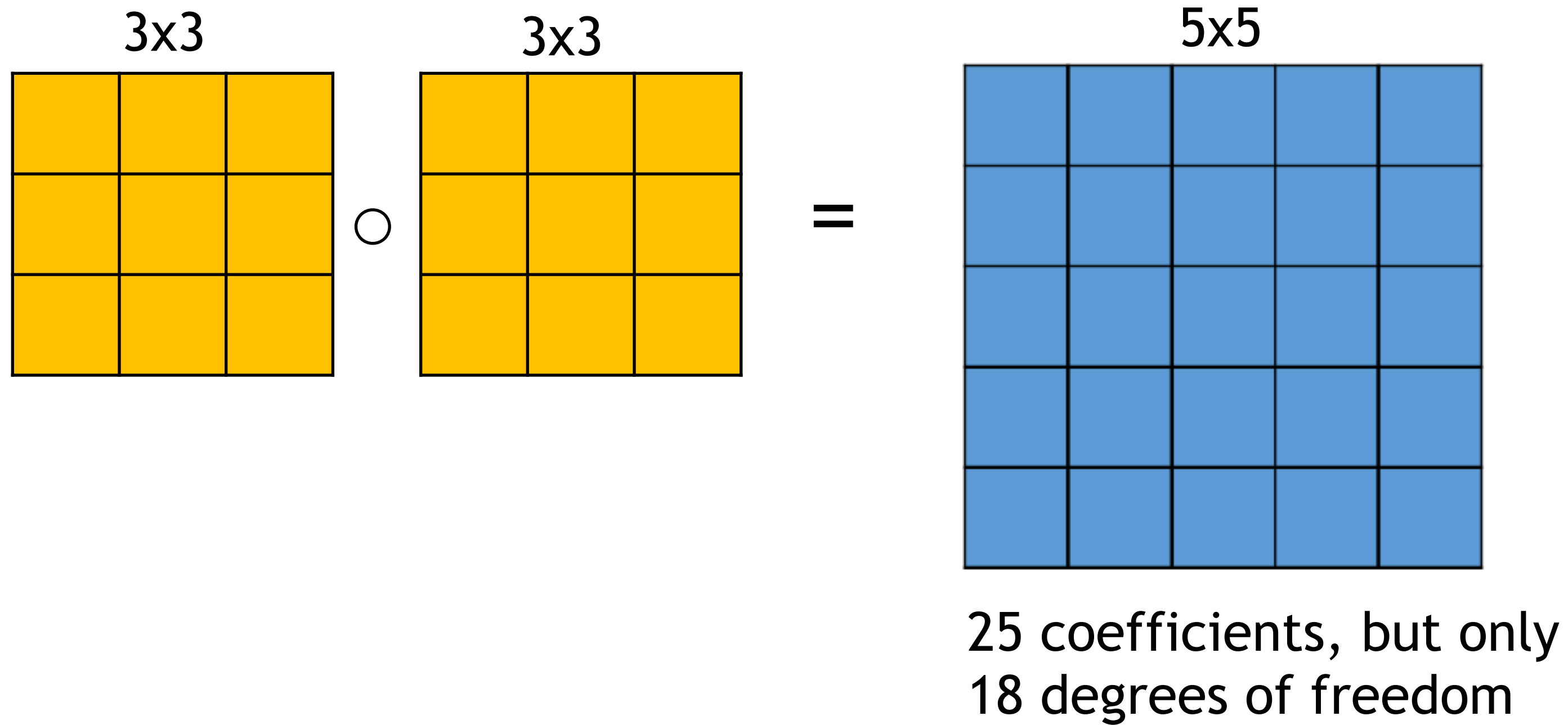
Main developments

- Small convolutional kernels: only 3x3



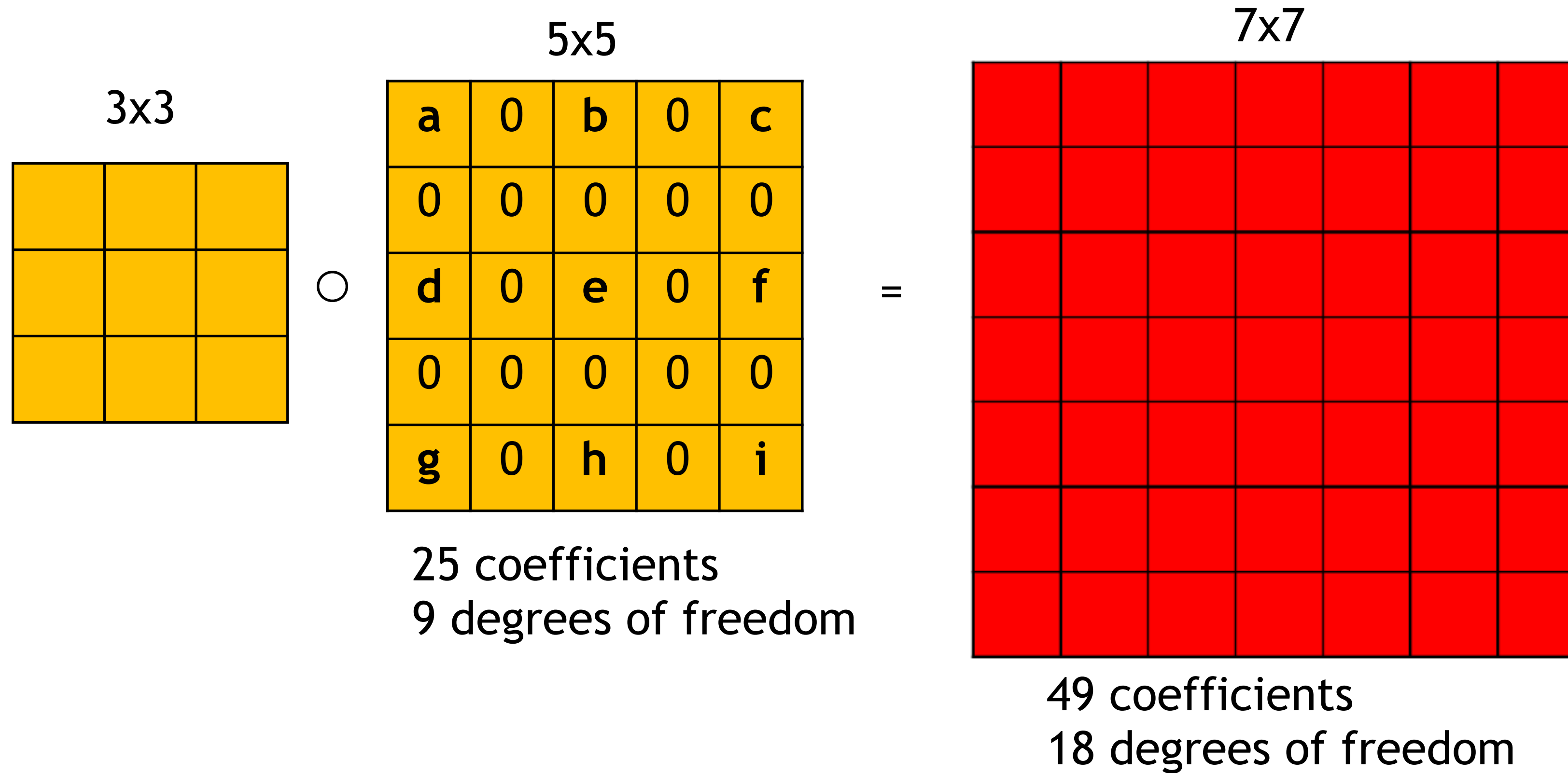
- Increased depth (5 -> 16/19 layers)

Chaining convolutions



Only separable filters... would this be enough?

Dilated convolutions



What is lost?

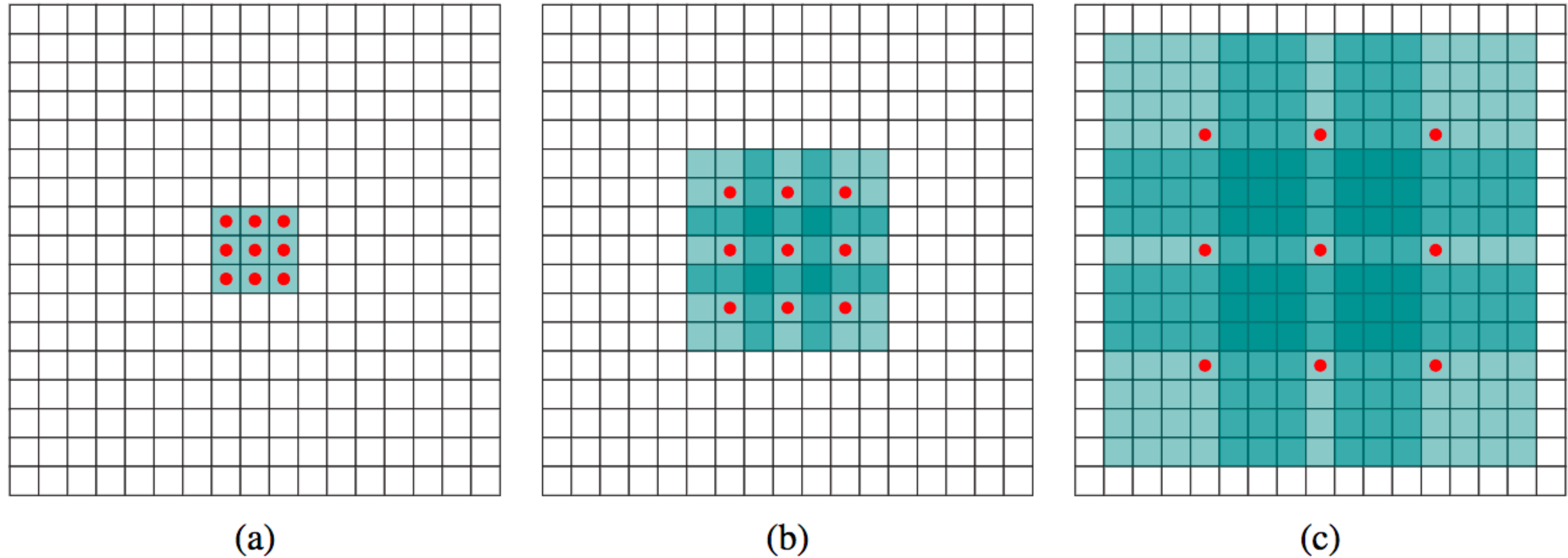
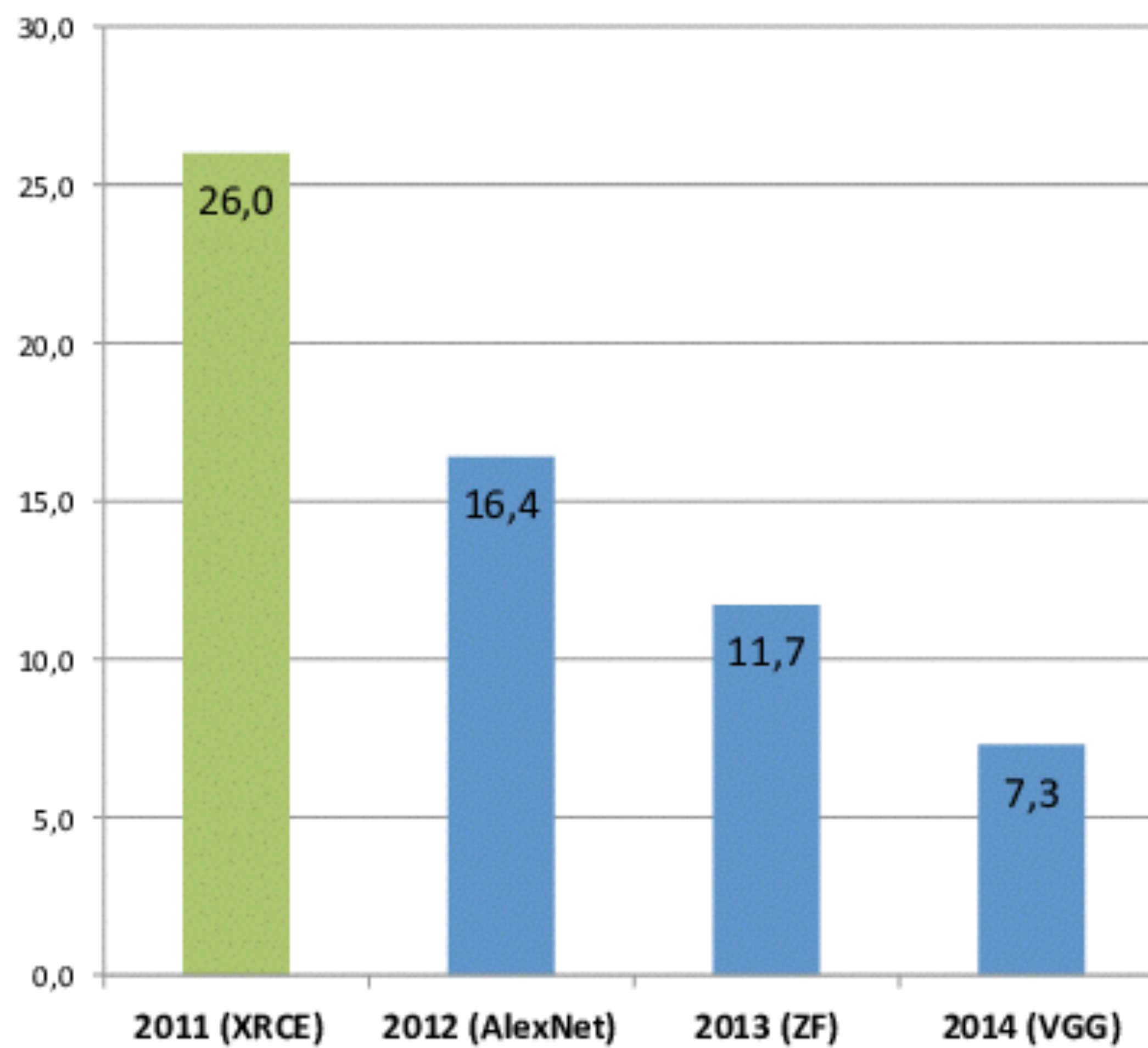
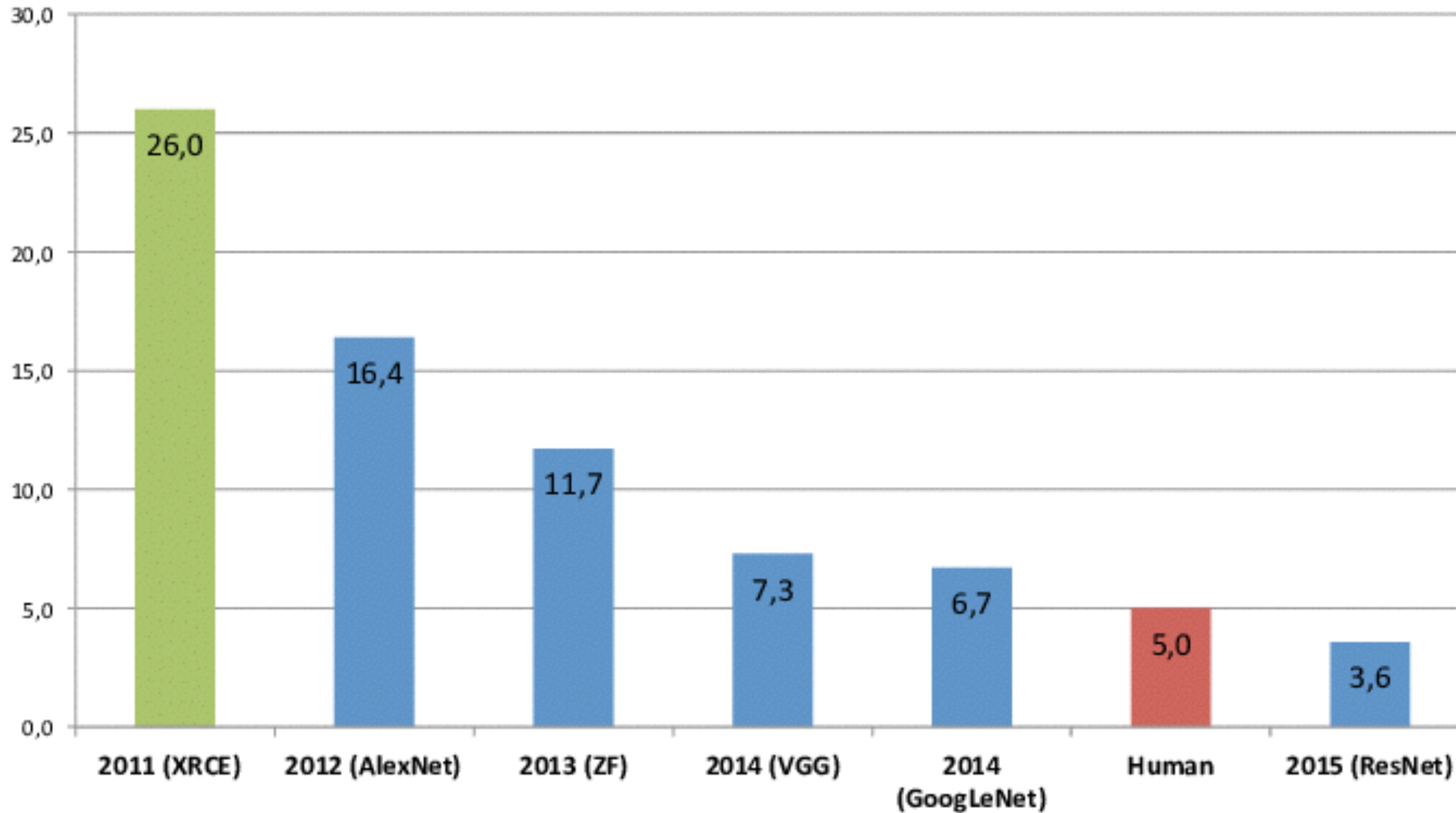


Figure 1: Systematic dilation supports exponential expansion of the receptive field without loss of resolution or coverage. (a) F_1 is produced from F_0 by a 1-dilated convolution; each element in F_1 has a receptive field of 3×3 . (b) F_2 is produced from F_1 by a 2-dilated convolution; each element in F_2 has a receptive field of 7×7 . (c) F_3 is produced from F_2 by a 4-dilated convolution; each element in F_3 has a receptive field of 15×15 . The number of parameters associated with each layer is identical. **The receptive field grows exponentially while the number of parameters grows linearly.**

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

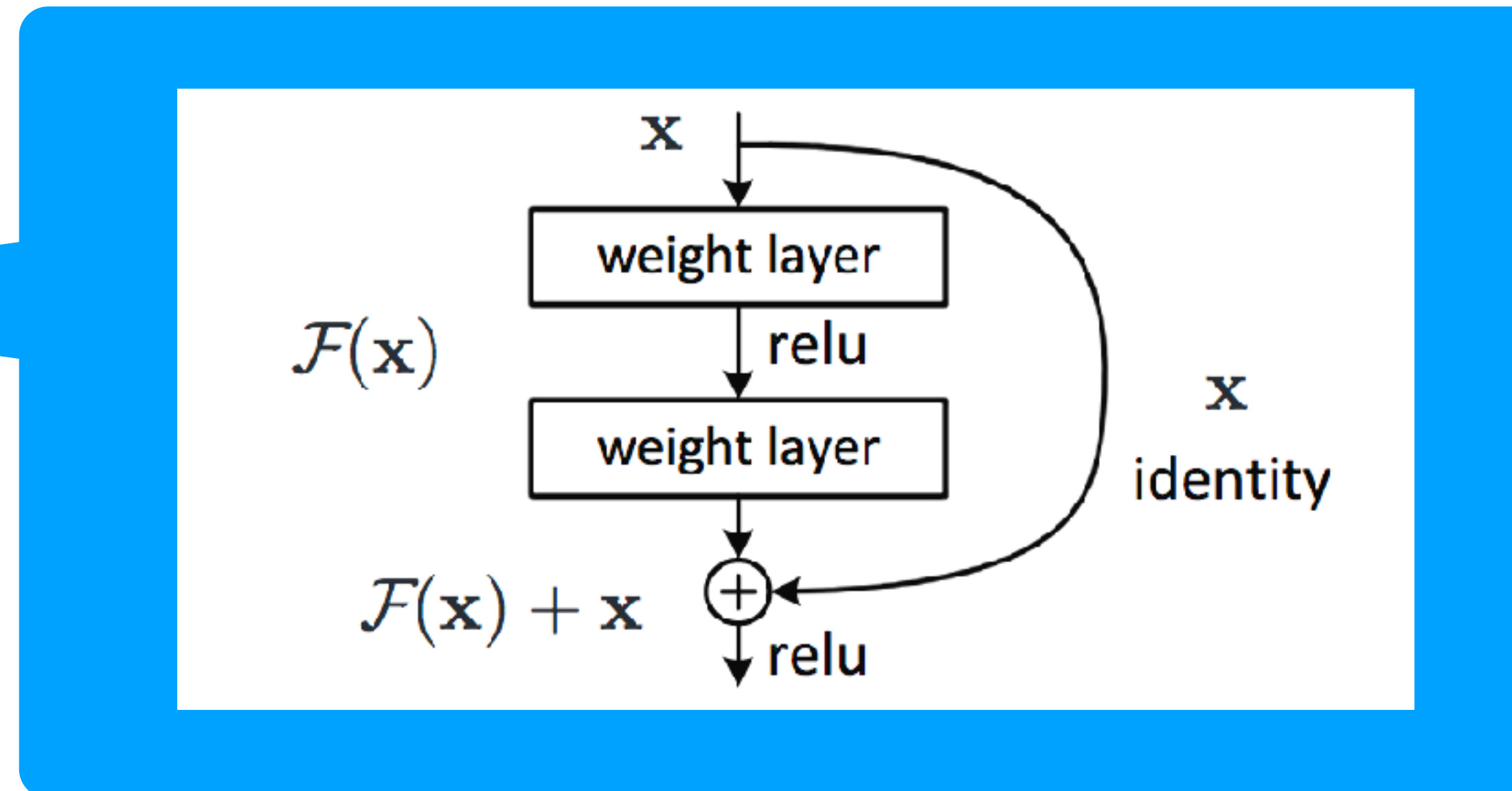
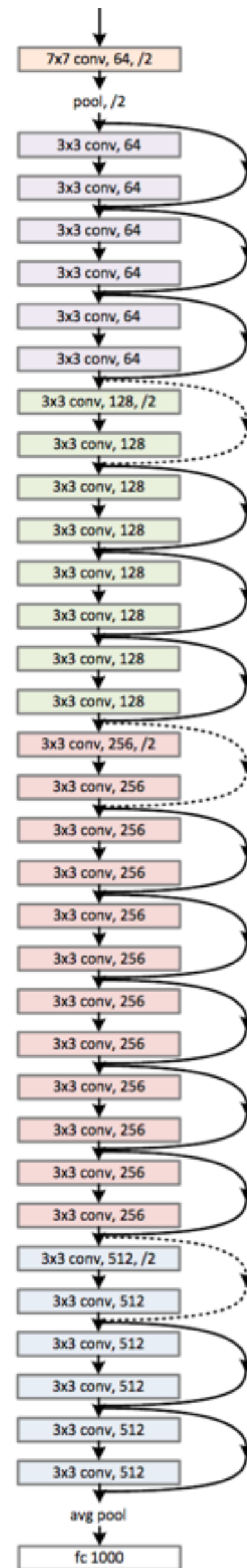
ResNet [He et al, 2016]

Main developments

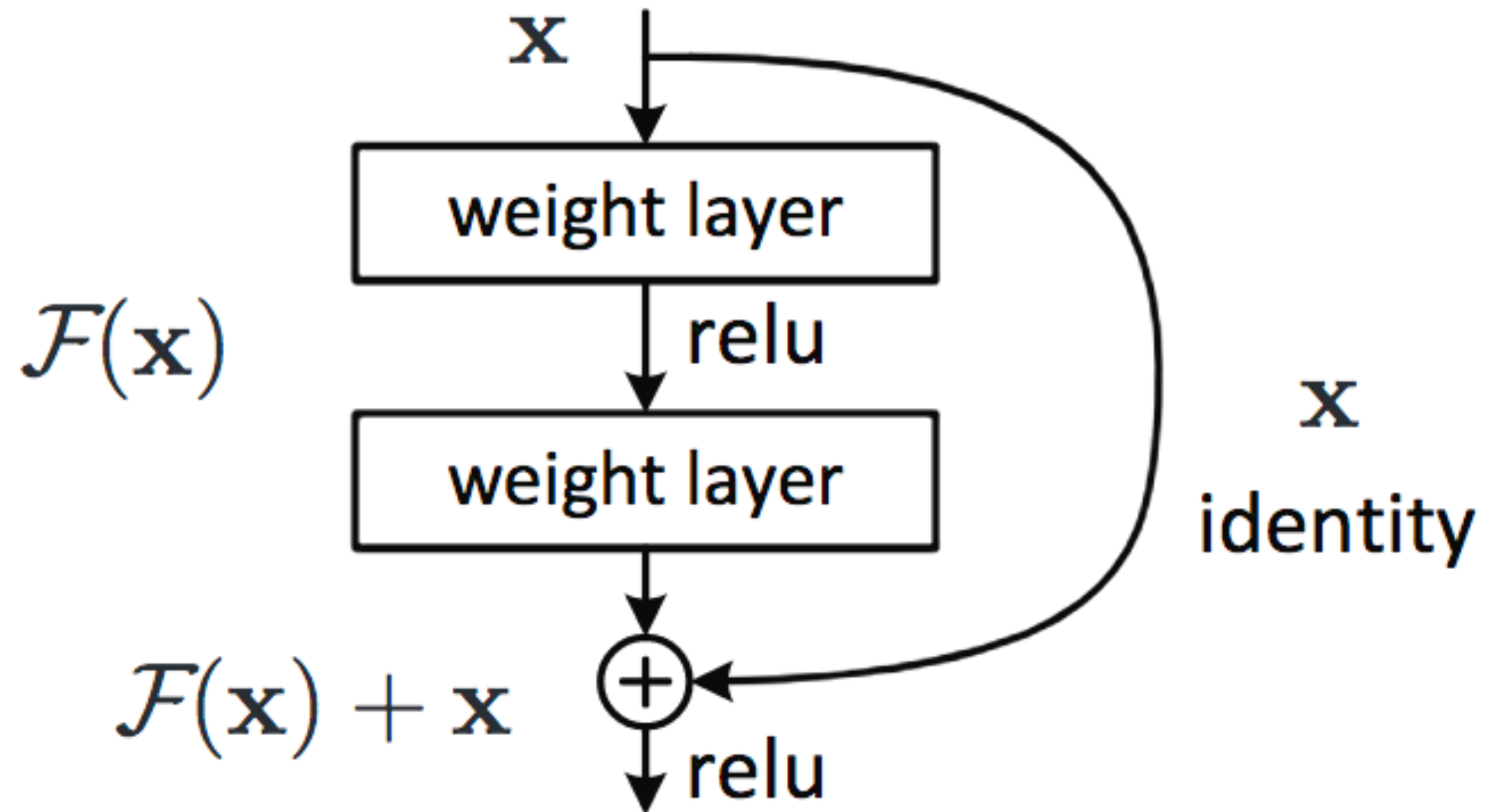
- Increased depth possible through residual blocks



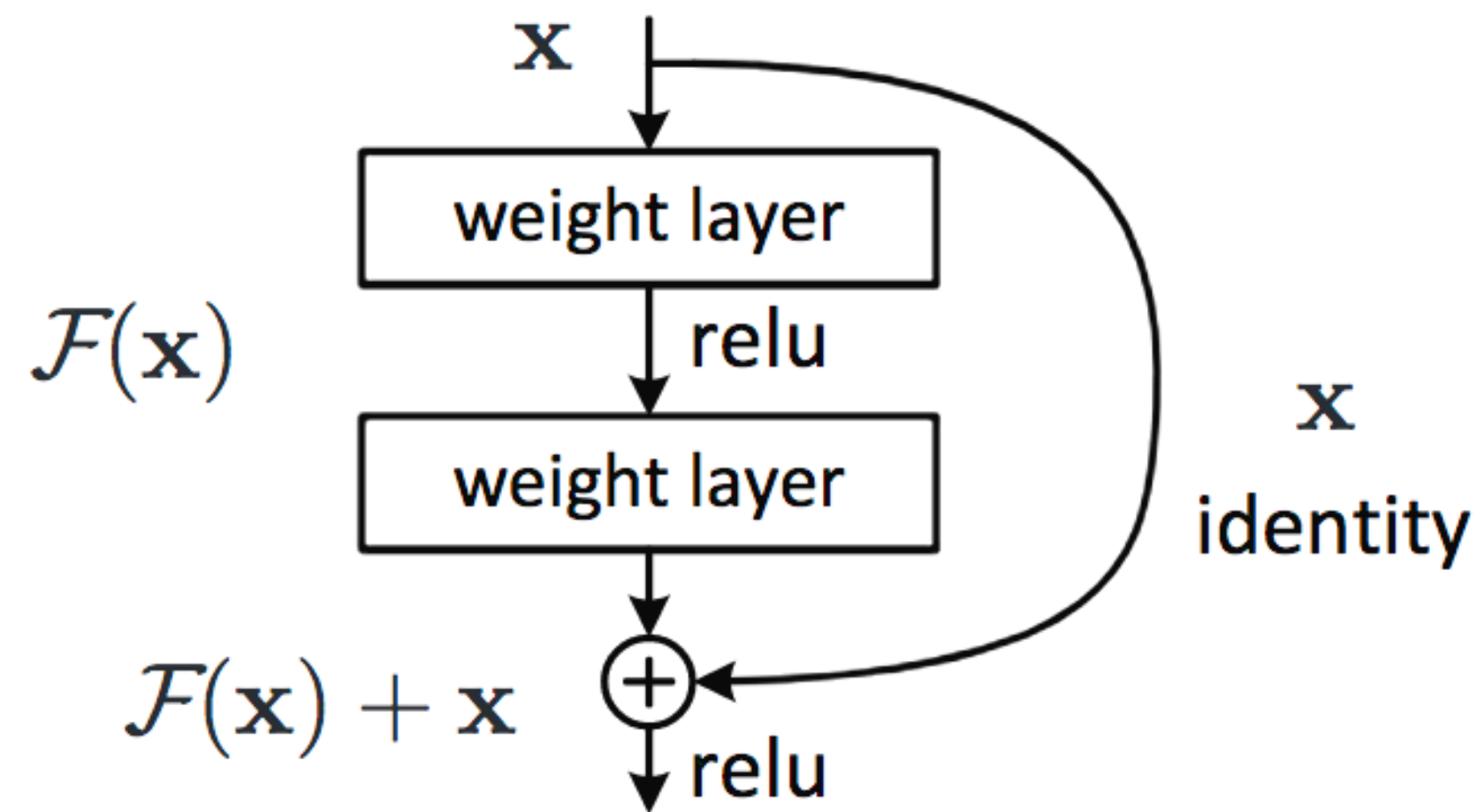
Error: 3.6%



Residual Blocks



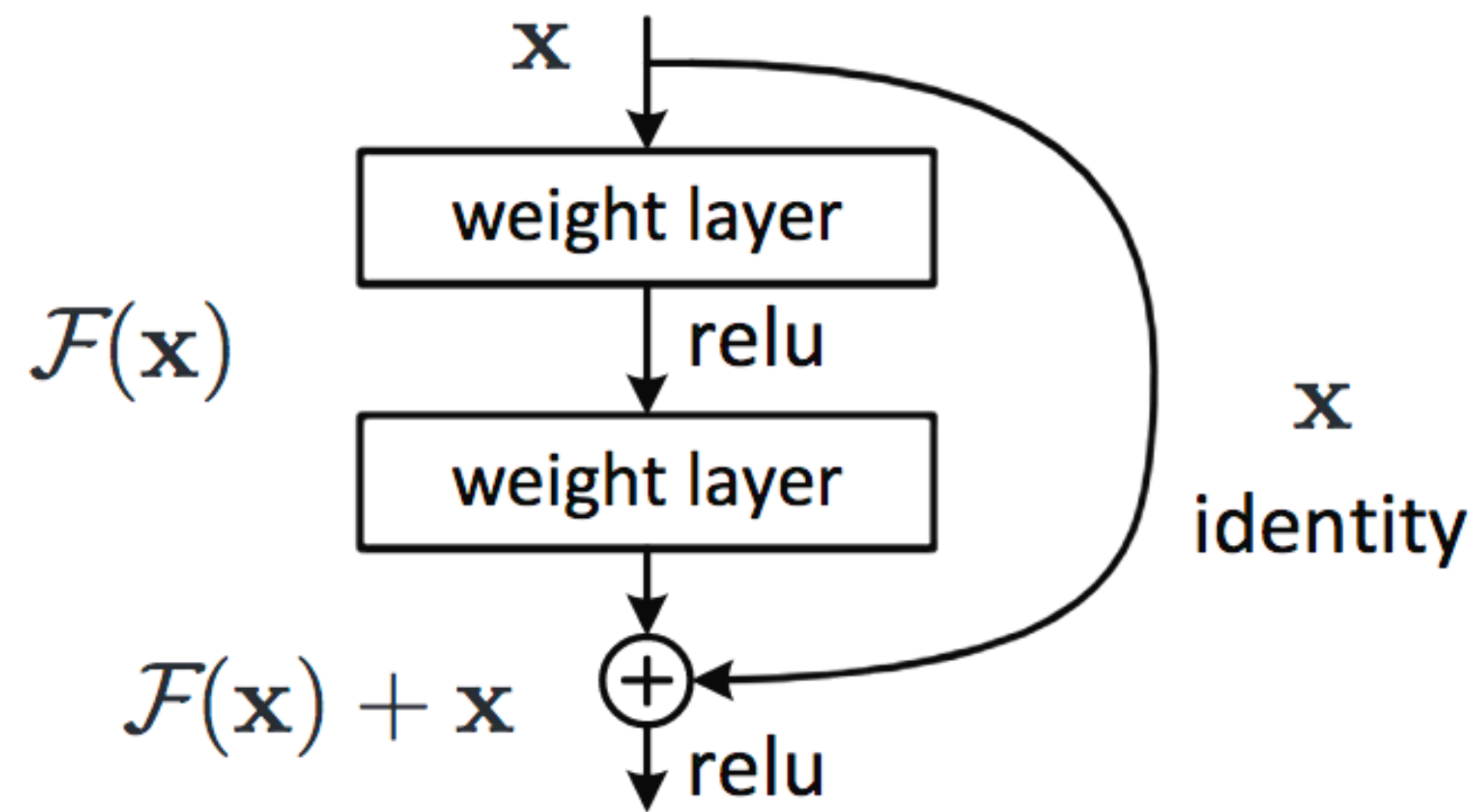
Residual Blocks



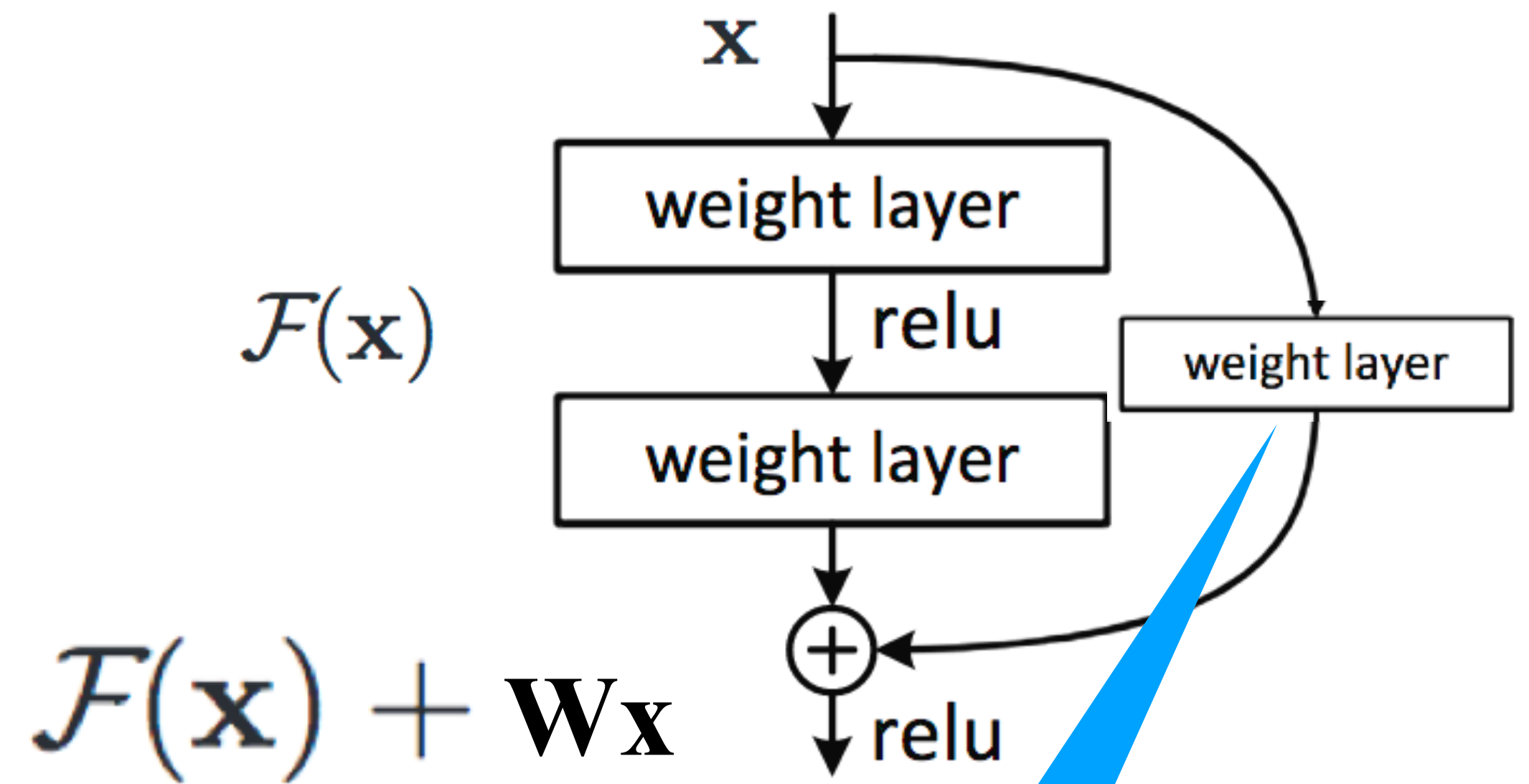
Why do they work?

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

If output has same size as input:



If output has a different size:



Projects into the right dimensionality:
 $\dim(\mathcal{F}(\mathbf{x})) = \dim(\mathbf{W}\mathbf{x})$

Some debugging advice

Other good things to know

- Check gradients numerically by finite differences
- Visualize hidden activations — should be uncorrelated and high variance



Good training: hidden units are sparse across samples and across features.

Other good things to know

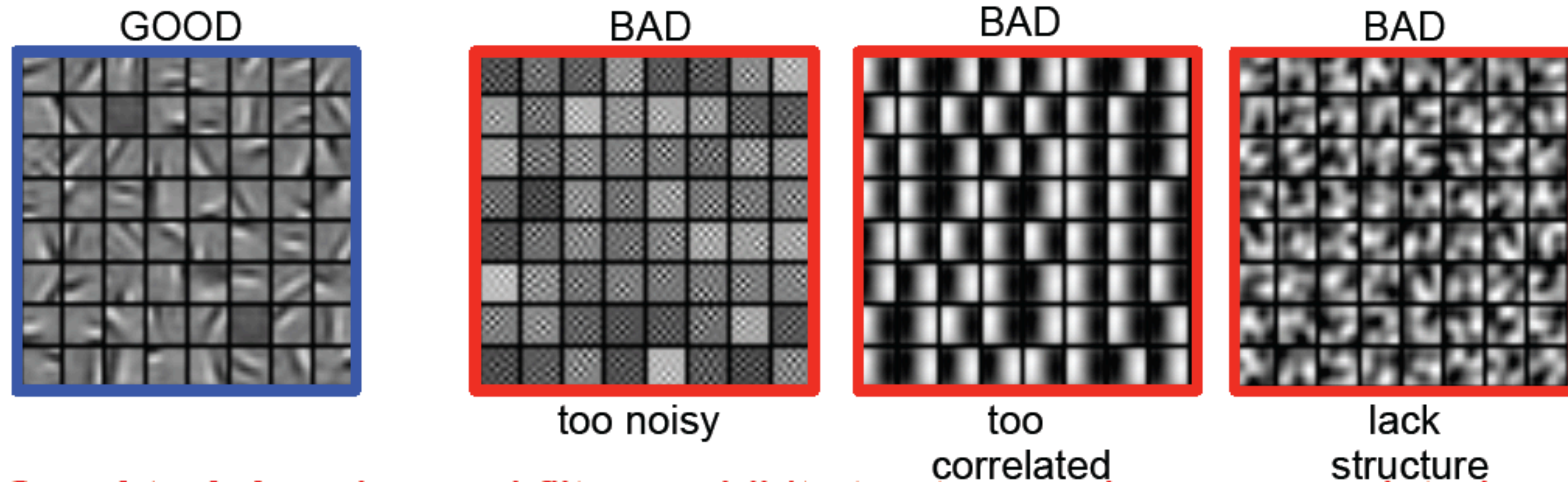
- Check gradients numerically by finite differences
- Visualize hidden activations — should be uncorrelated and high variance



Bad training: many hidden units ignore the input and/or exhibit strong correlations.

Other good things to know

- Check gradients numerically by finite differences
- Visualize hidden activations — should be uncorrelated and high variance
- Visualize filters



Good training: learned filters exhibit structure and are uncorrelated.

Next week:

Practical advice on training
and debugging networks