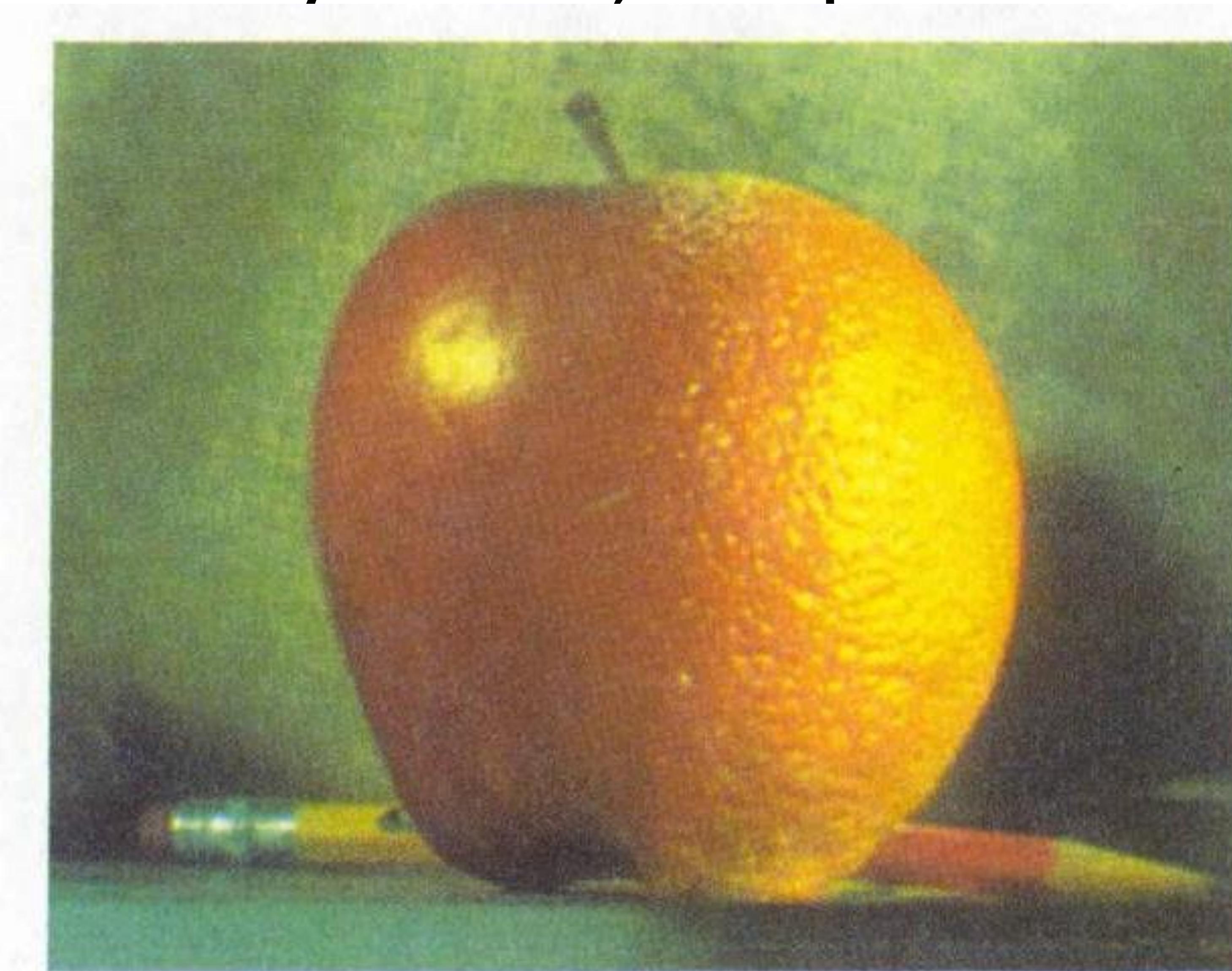


Pyramids, Deep Nets

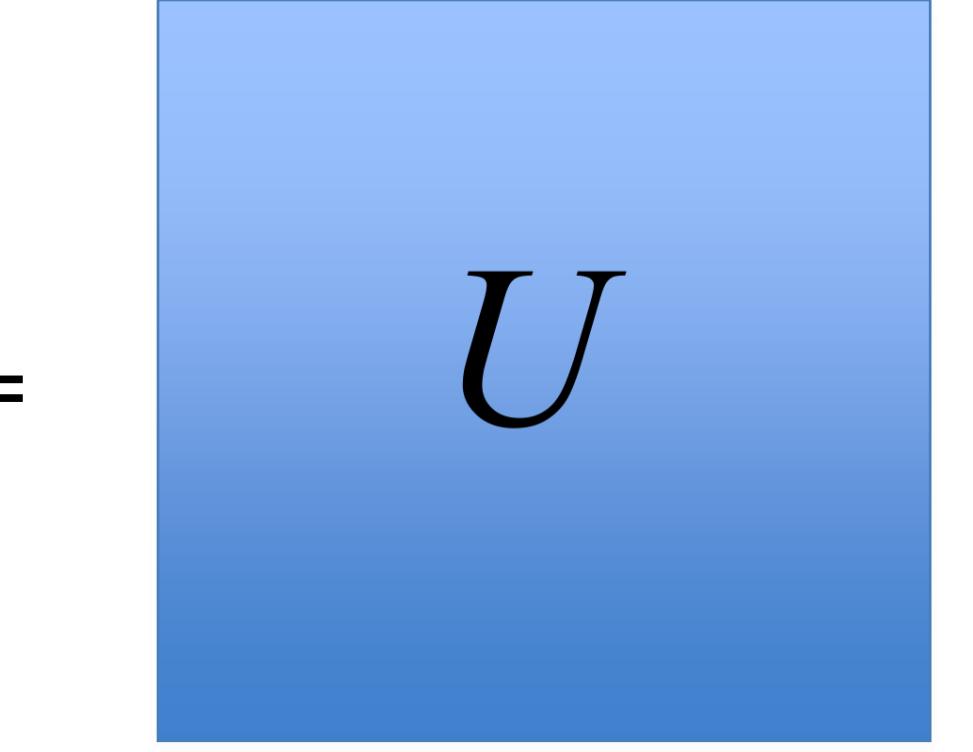


Linear image transformations

In analyzing images, it's often useful to make a change of basis.

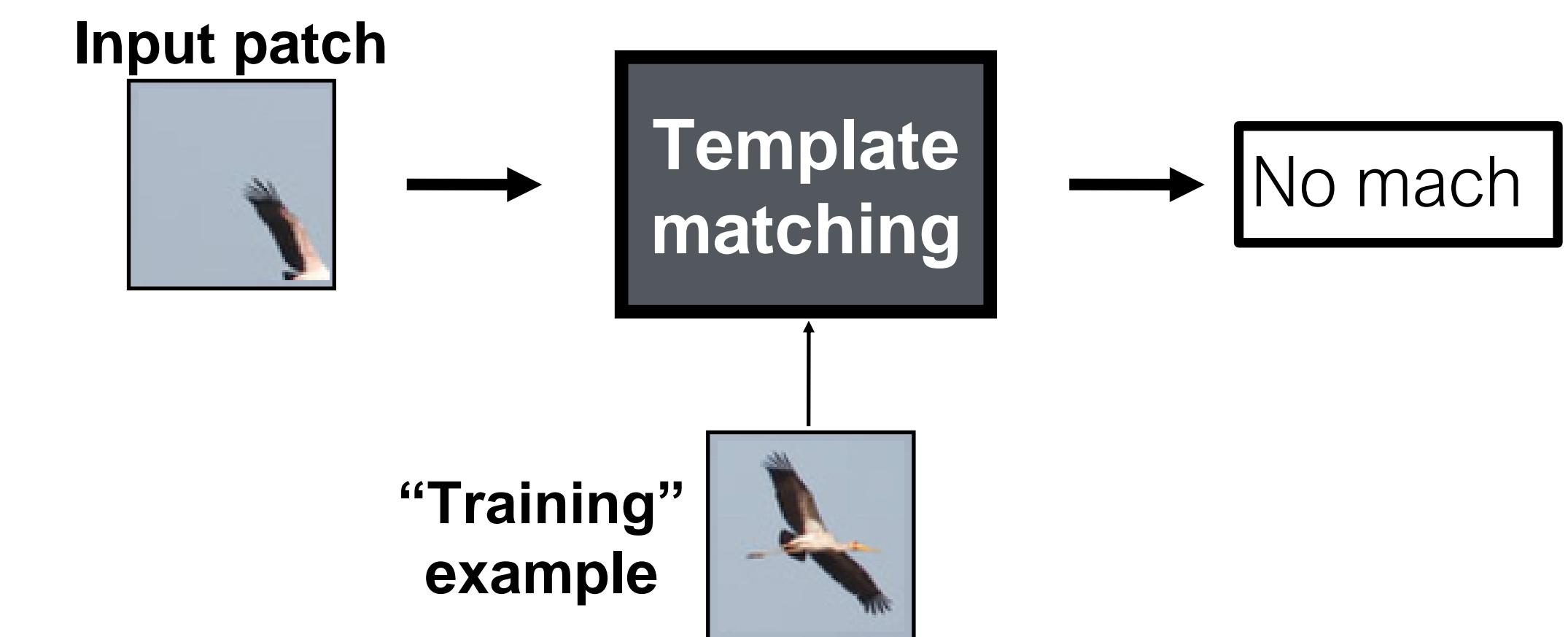
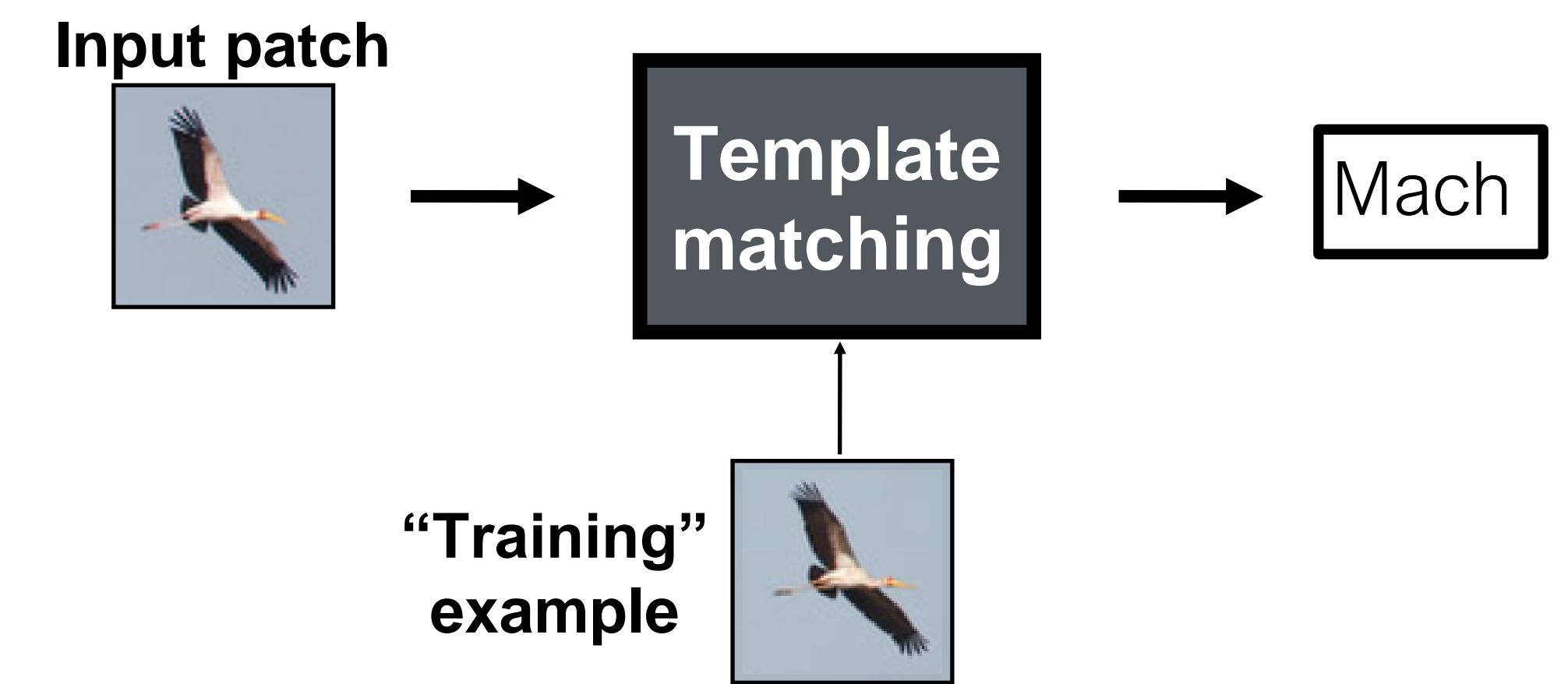
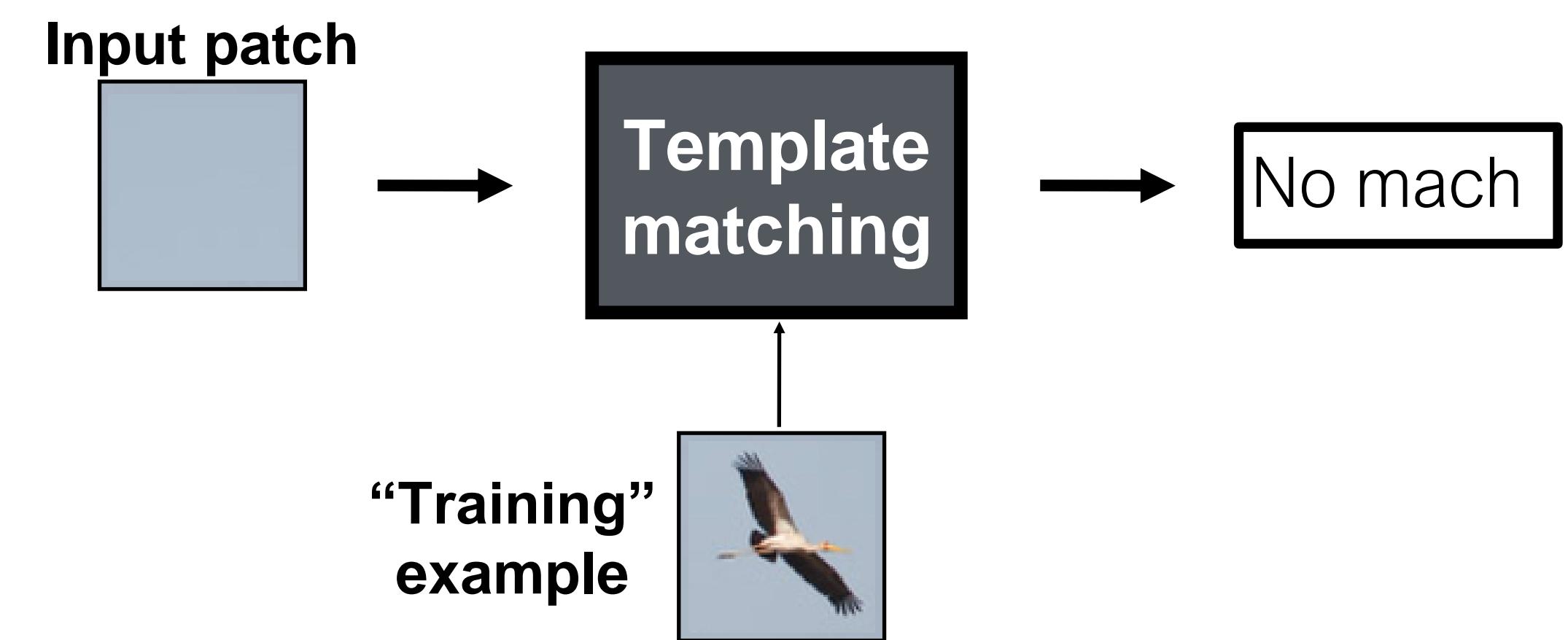
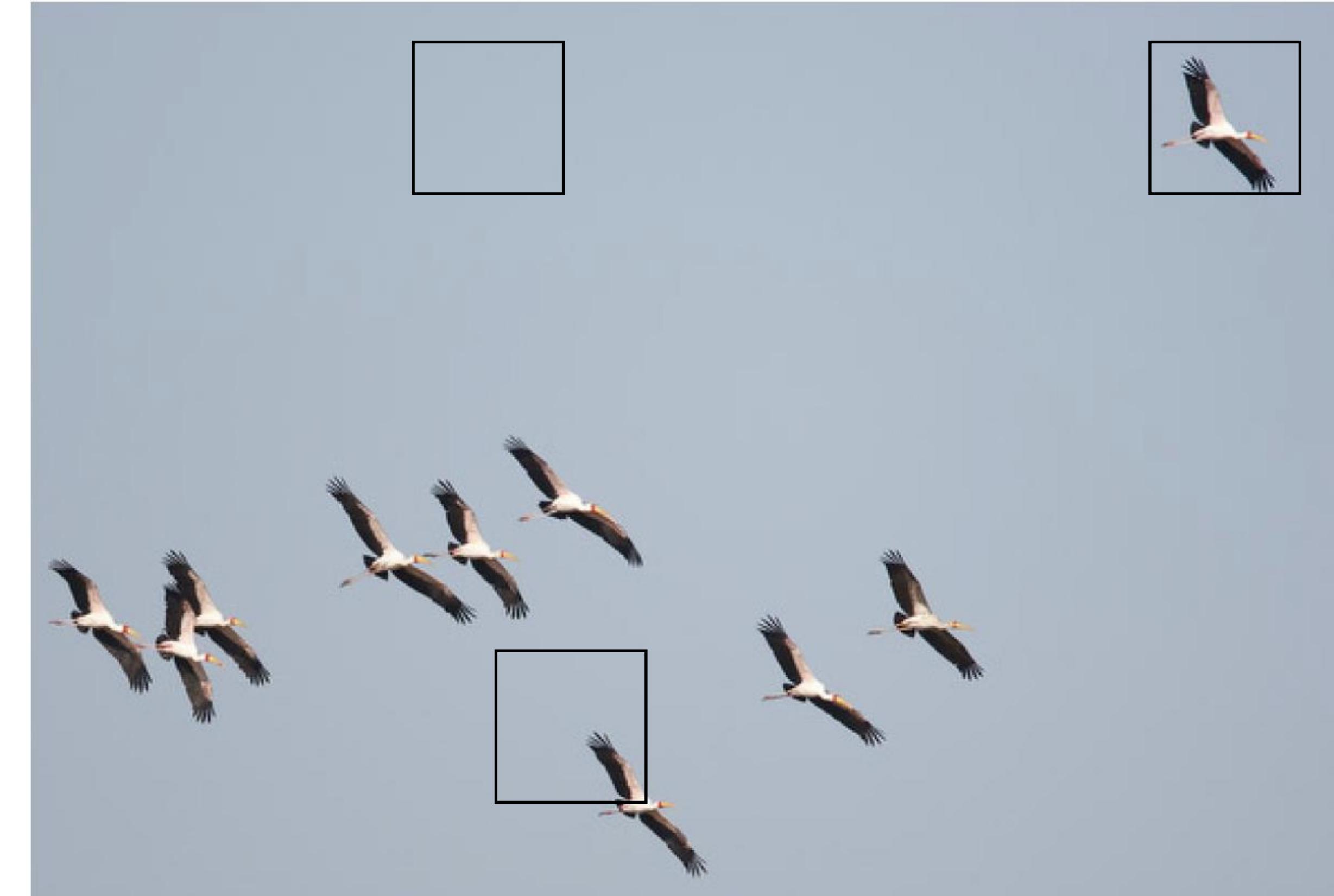
$$\vec{\bar{F}} = \vec{U} \vec{f}$$

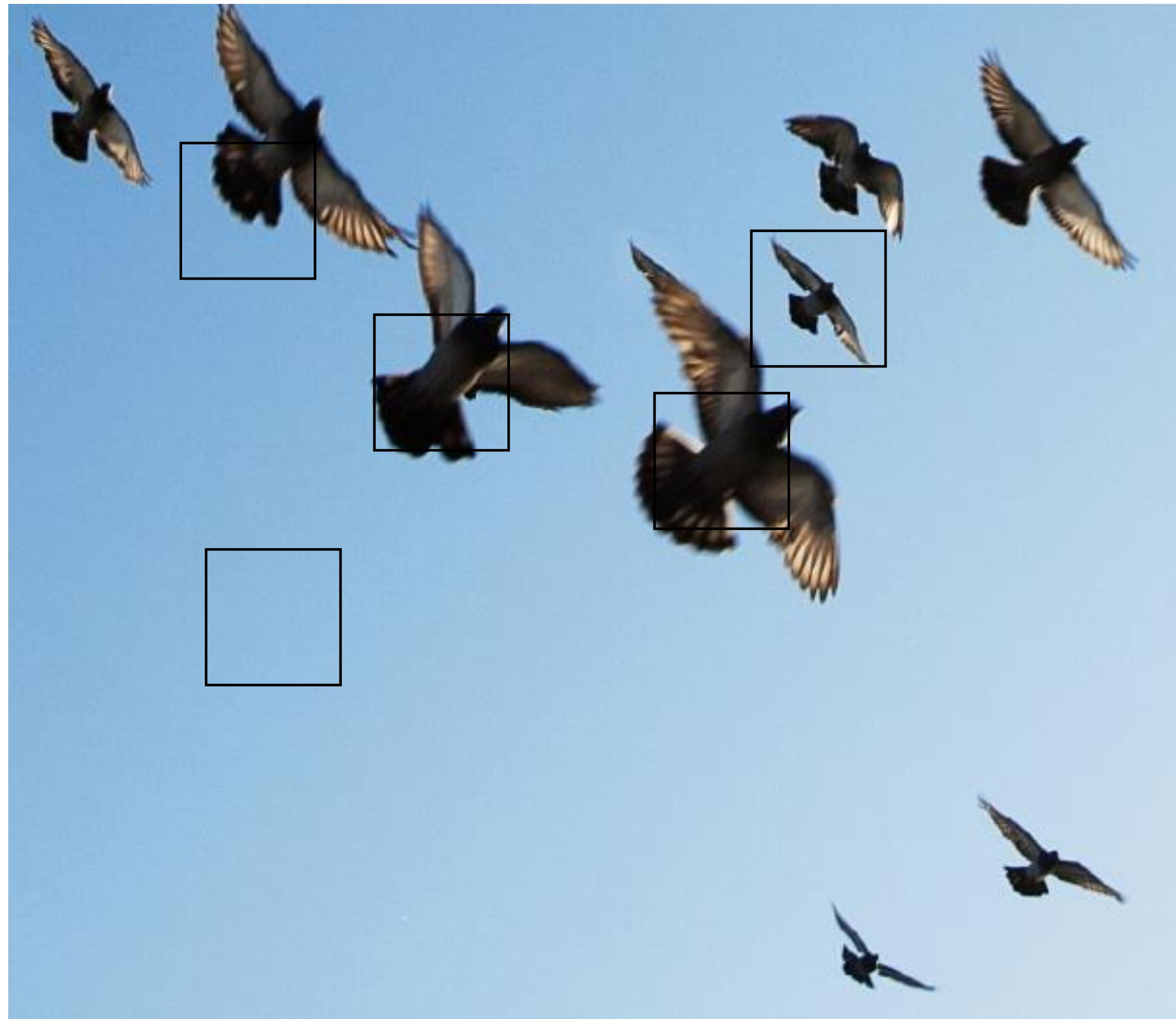
Transformed image → $\vec{\bar{F}}$ = $\vec{U} \vec{f}$ ← Vectorized image

=  Fourier transform, or
Wavelet transform, or
Steerable pyramid transform

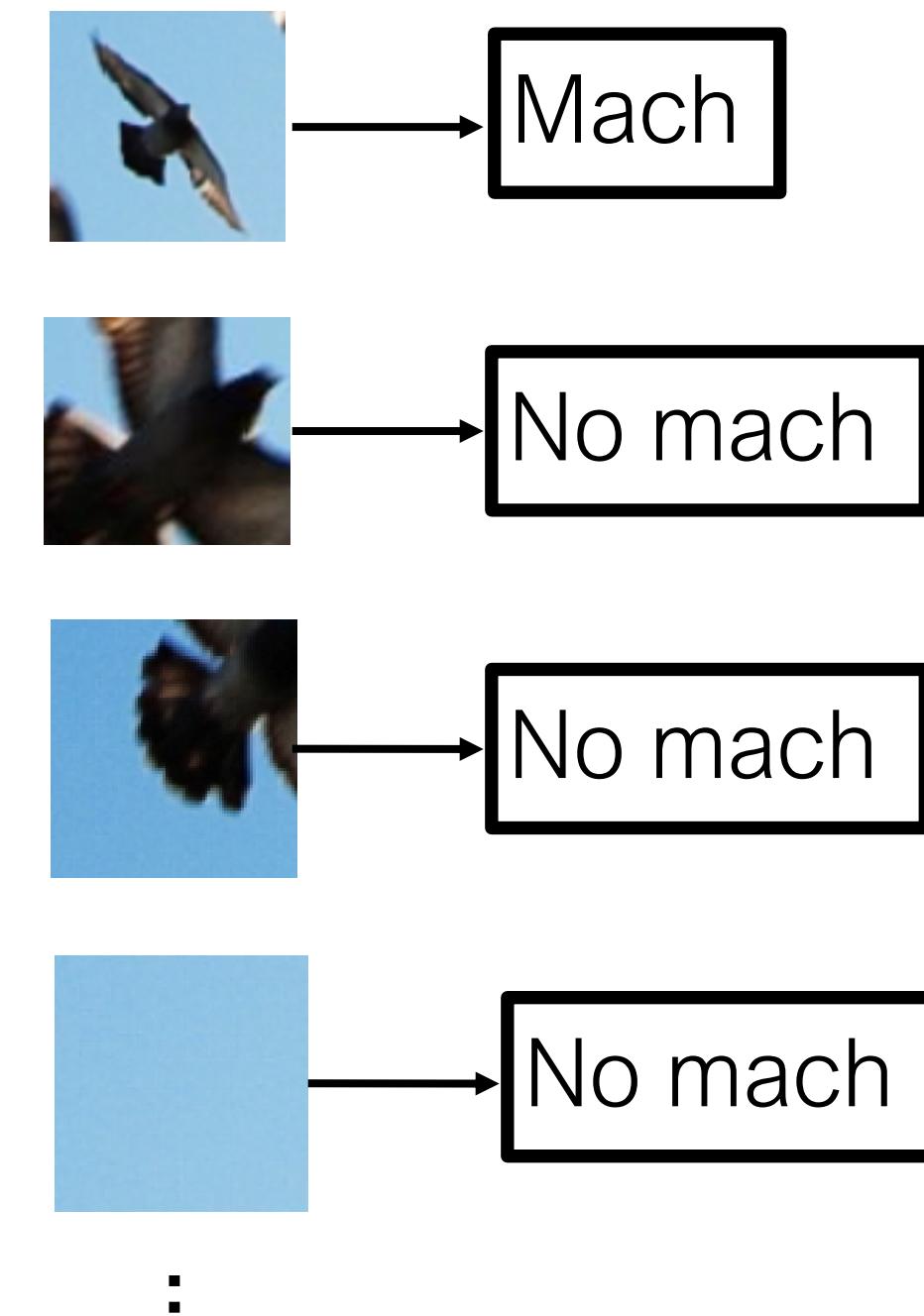


We need translation invariance

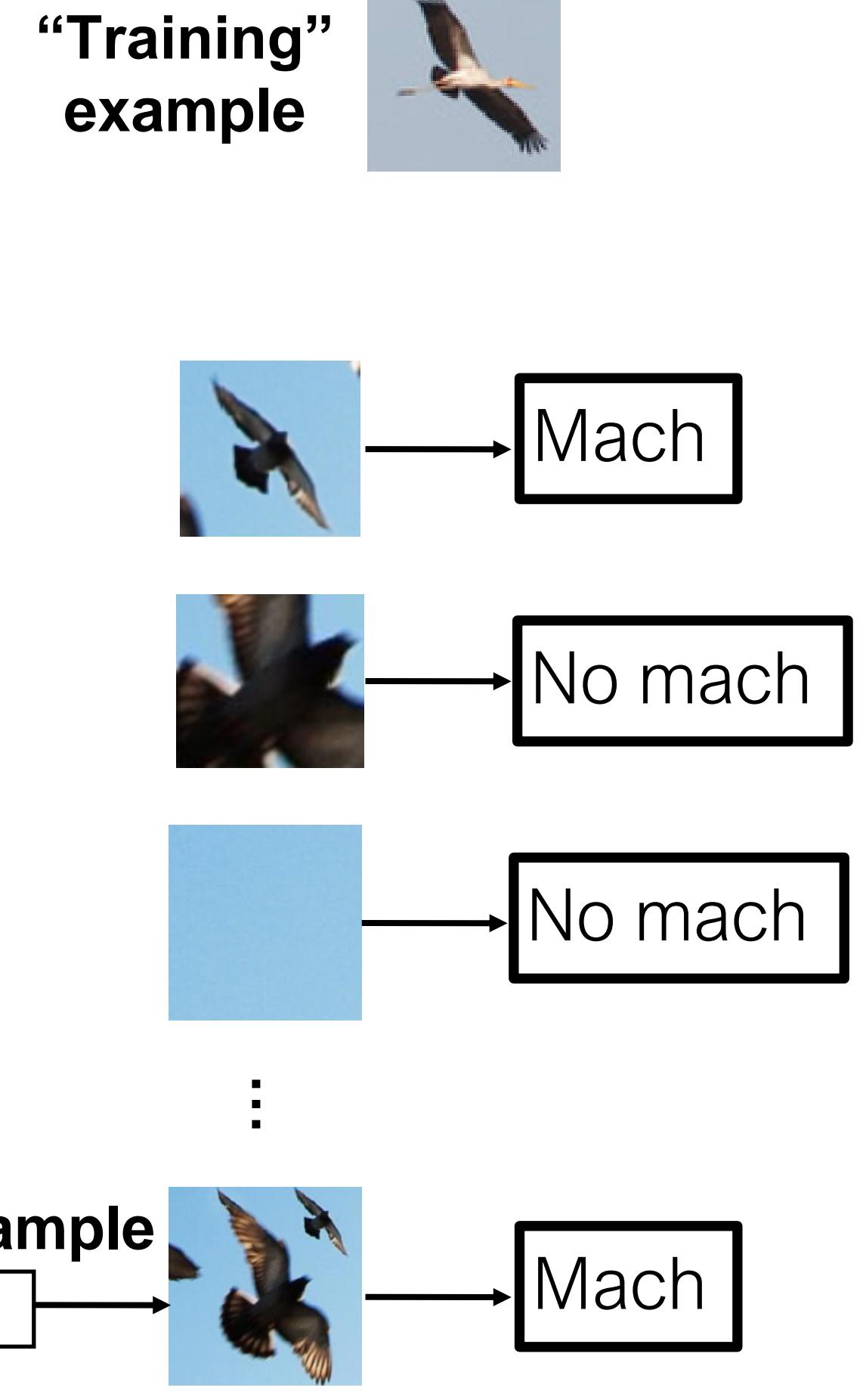
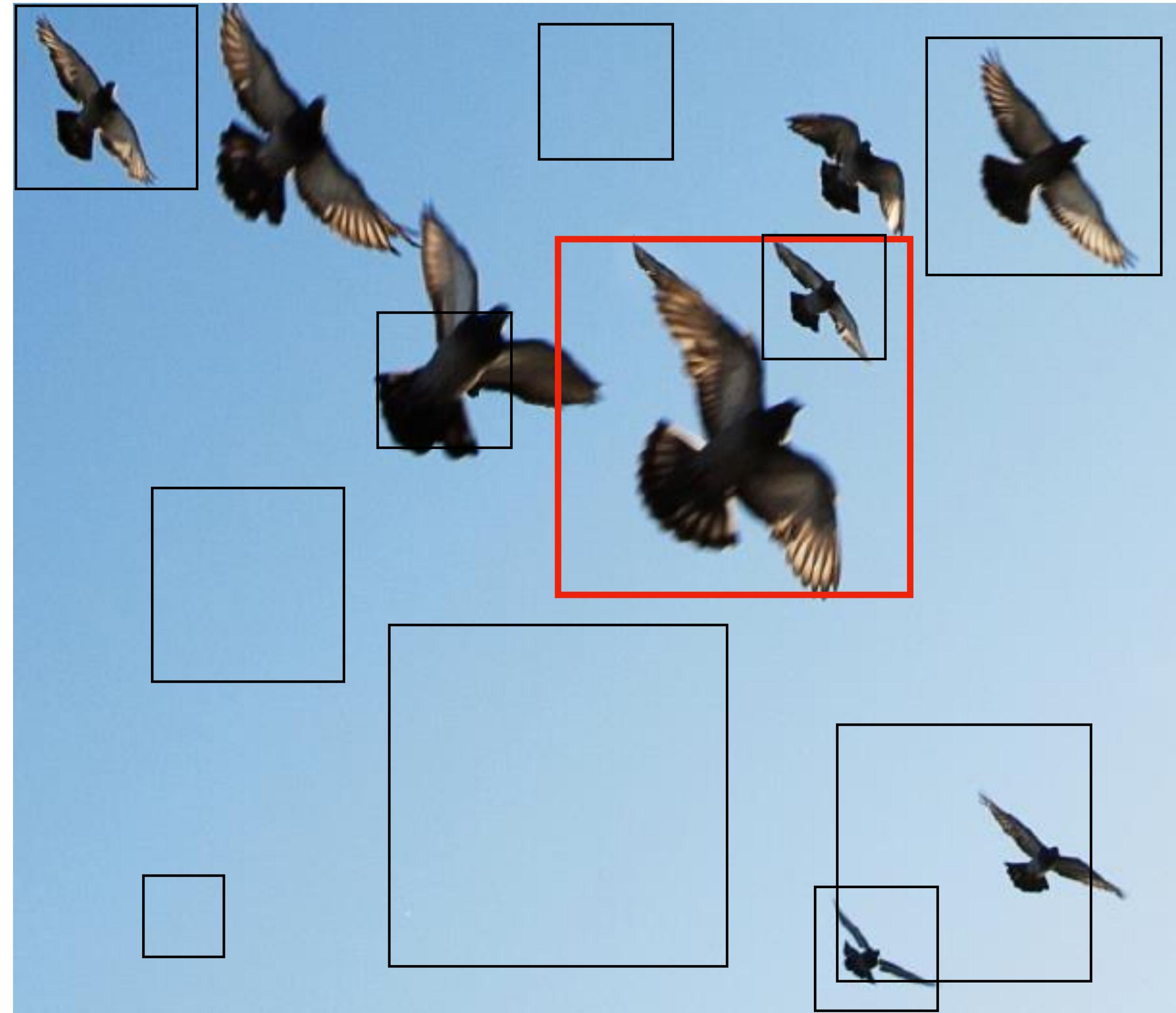




“Training” example



We need translation and scale invariance



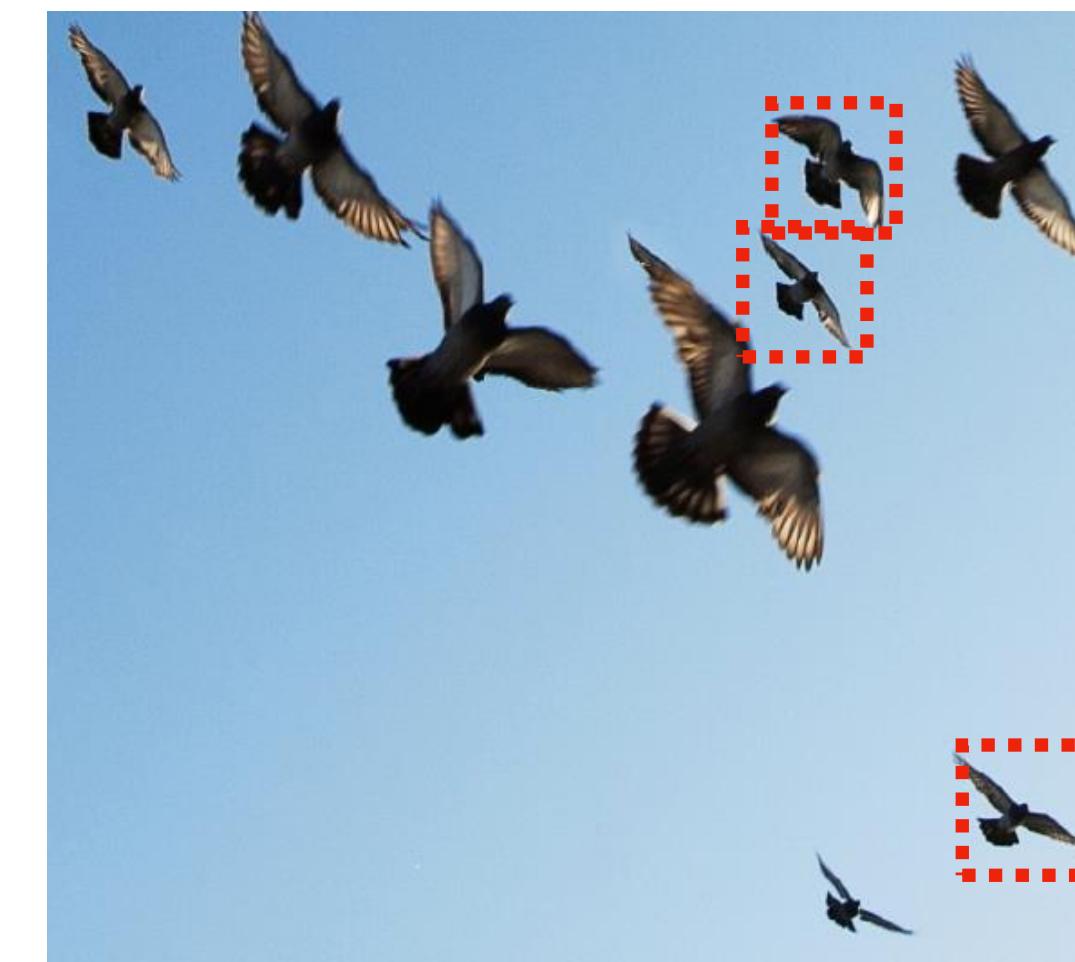
Consider all patch location and sizes

Multi-Scale Pyramids

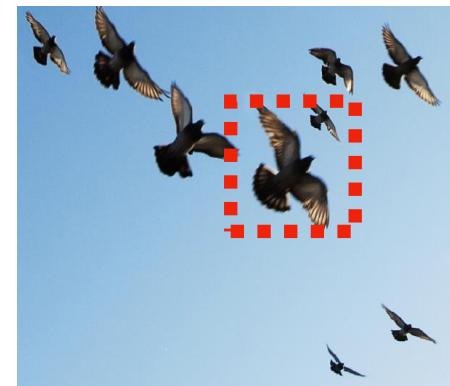
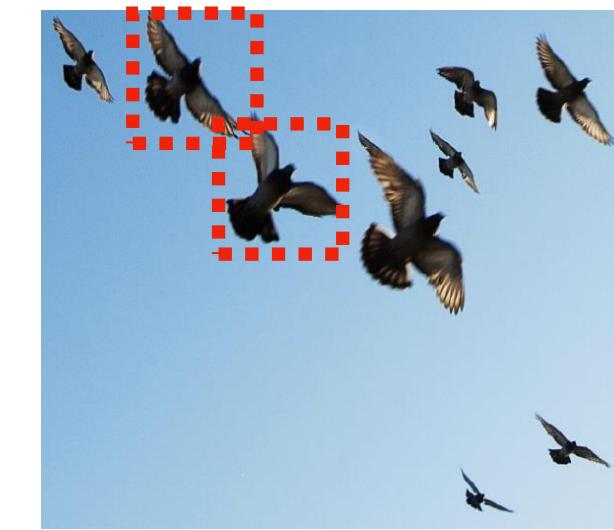
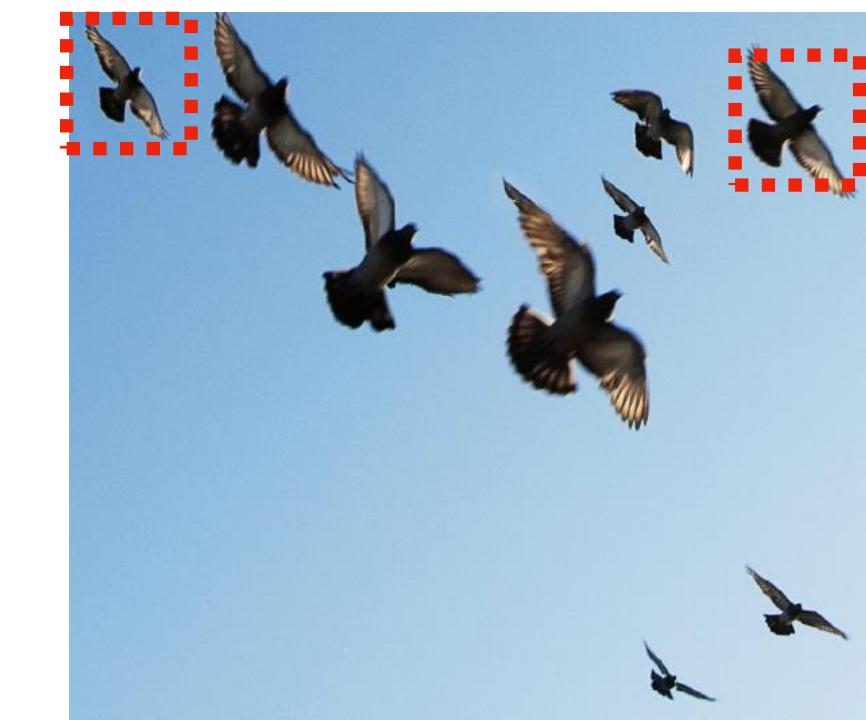


Template

Multi-Scale Pyramids



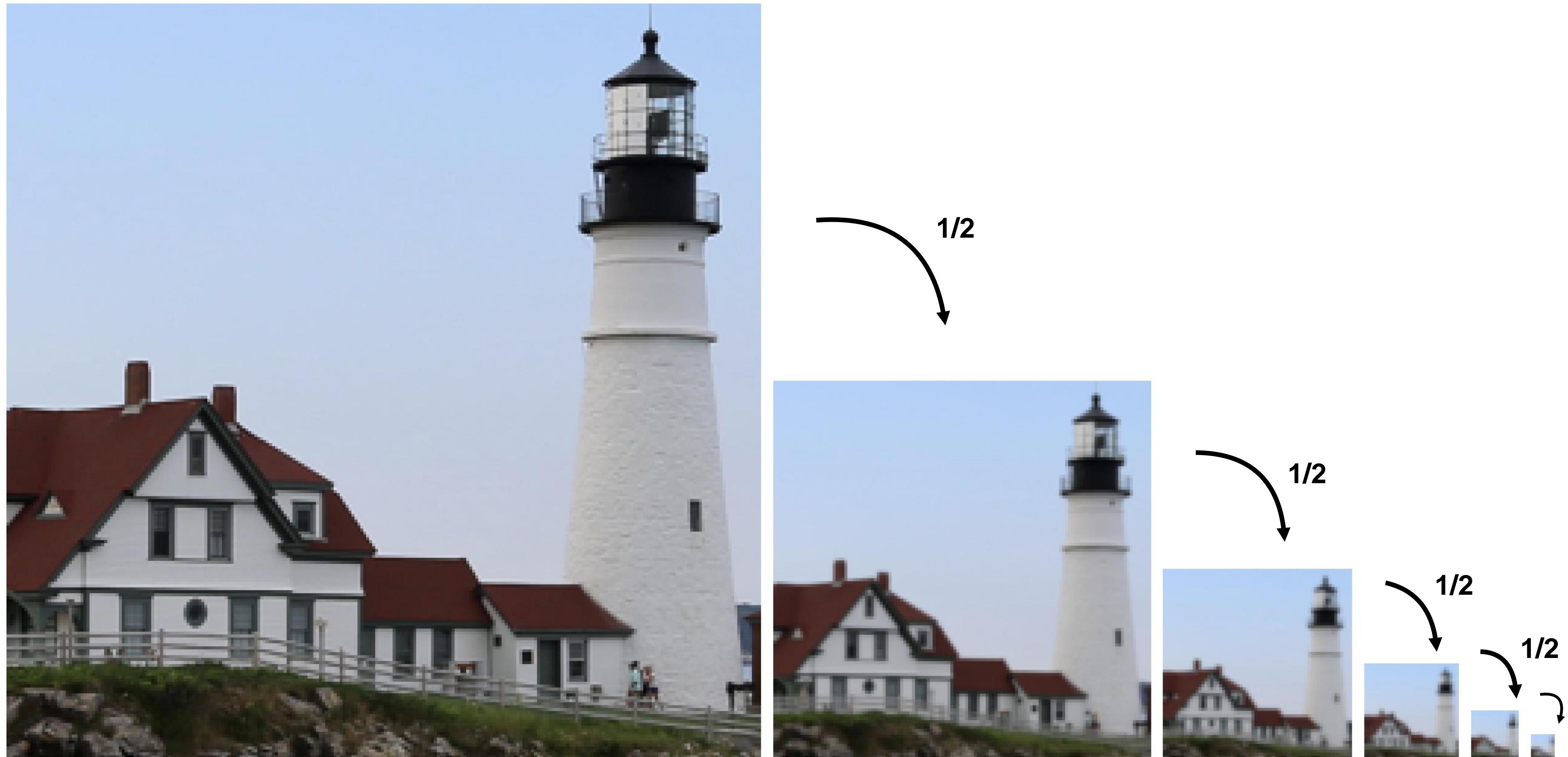
Multiscale image pyramid



Template

A multiscale image pyramid provides an alternative image representation to achieve translation and scale invariance

Gaussian Pyramid

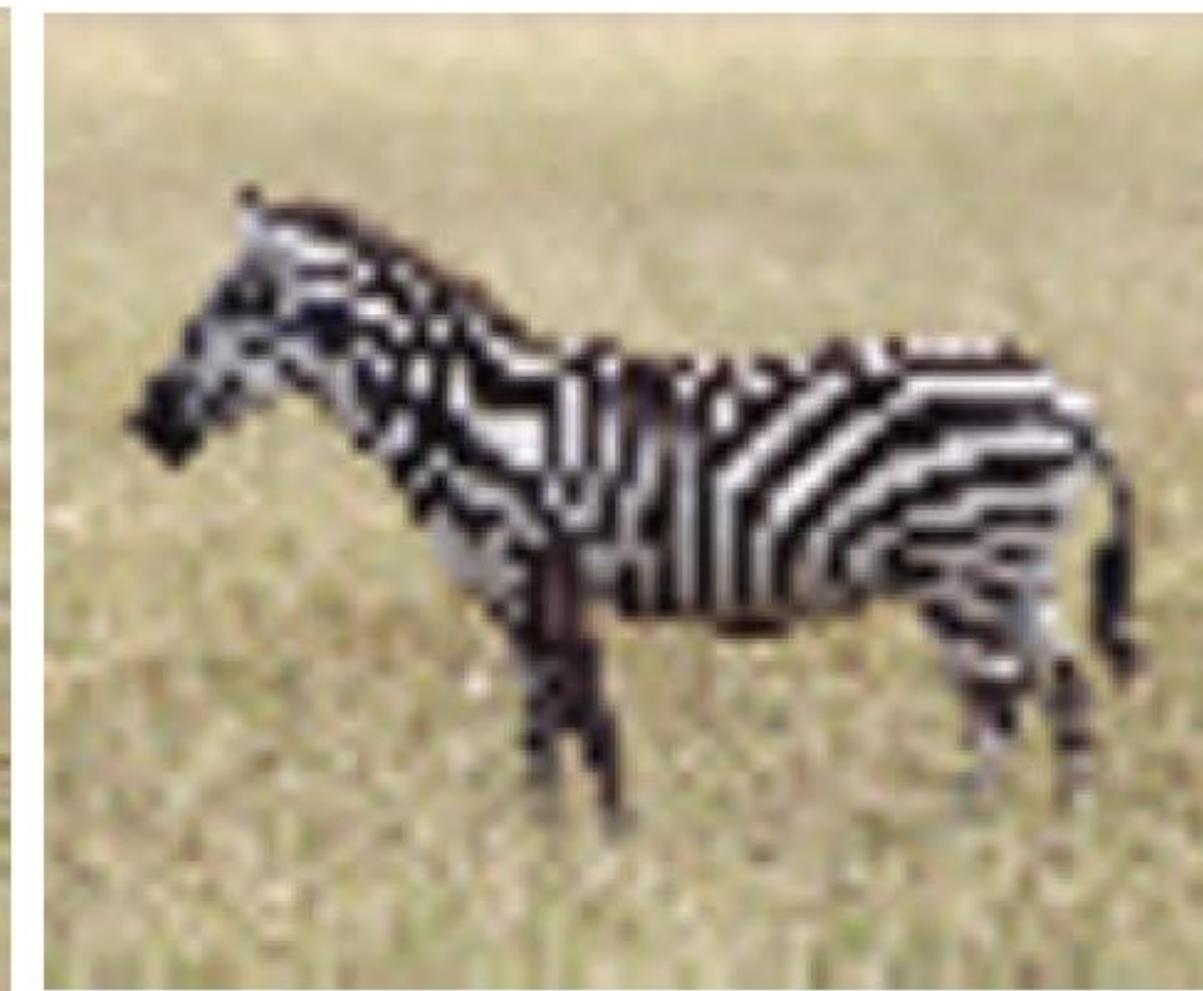


Subsampling and aliasing

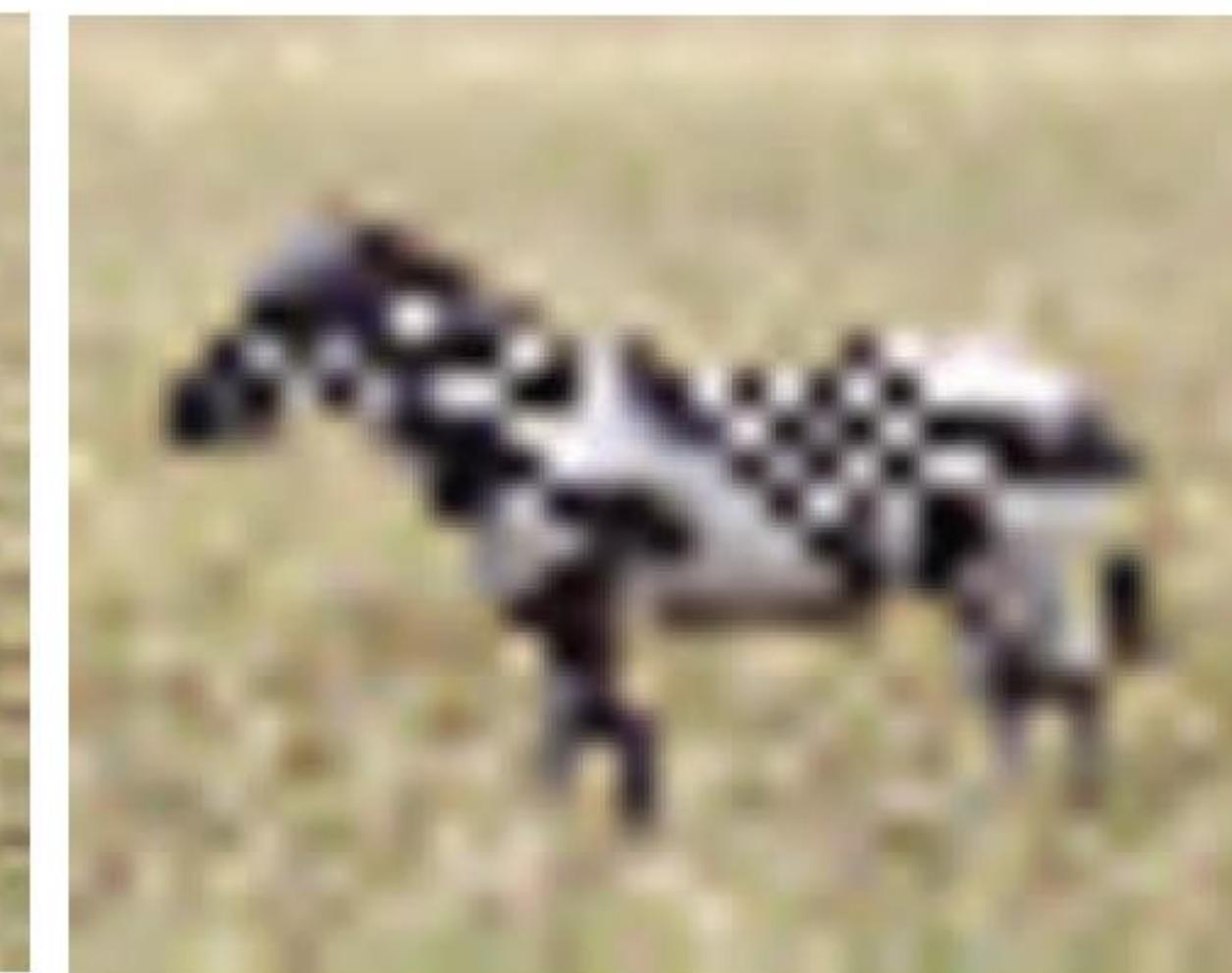
103×128



52×64



26×32



1/2

1/2

The Gaussian pyramid

For each level

1. Blur input image with a Gaussian filter

$$[1, 4, 6, 4, 1]$$



$$\bigcirc [1, 4, 6, 4, 1] \bigcirc \xrightarrow{6 \rightarrow} [1 \\ 4 \\ 6 \\ 4 \\ 1]$$



(The Gaussian filter is approximated with a binomial filter)

The Gaussian pyramid

For each level

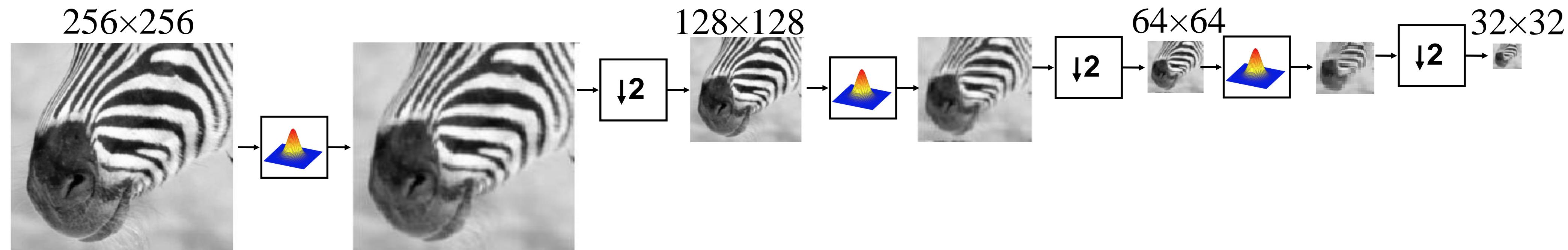
1. Blur input image with a Gaussian filter
2. Downsample image



$$\bigcirc [1, 4, 6, 4, 1] \bigcirc \begin{matrix} 6 \\ \longrightarrow \\ 4 \\ 1 \end{matrix}$$



The Gaussian pyramid



The Gaussian pyramid

512×512



(original image)

256×256



128×128



64×64



32×32

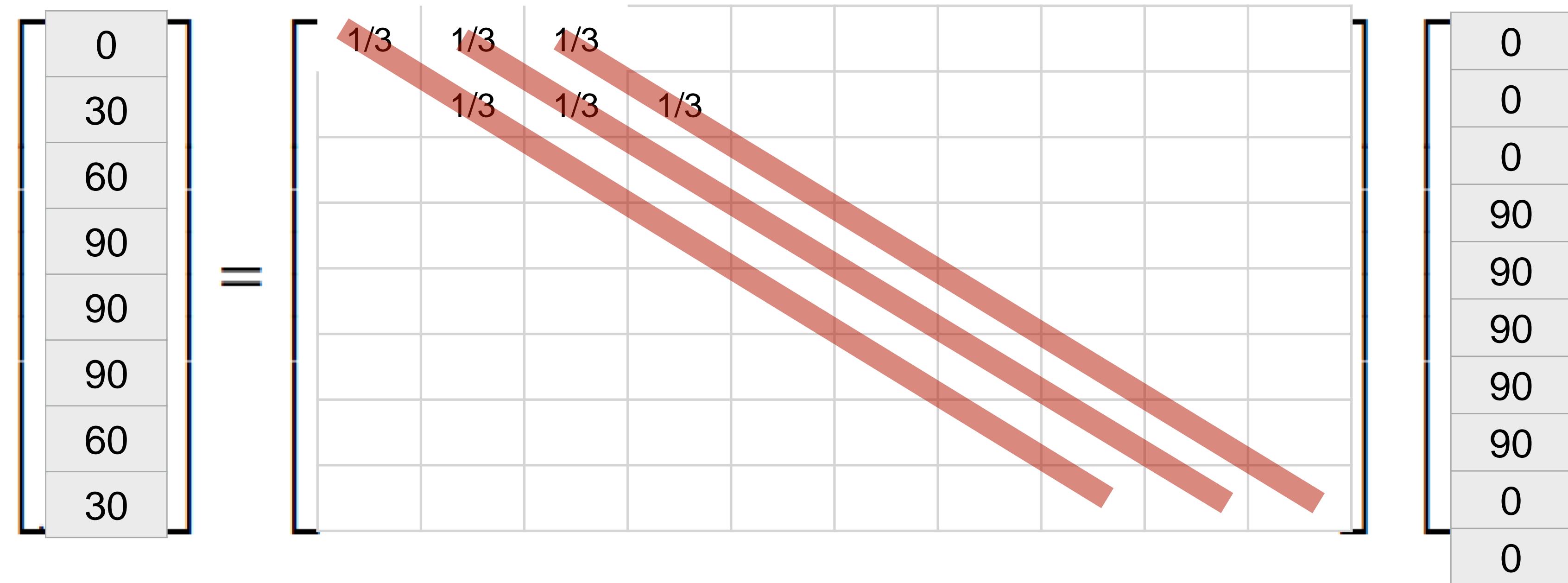


Convolution as matrix multiplication

In the 1D case, it helps to make explicit the structure of the matrix:

$$\begin{bmatrix} 0 & 0 & 0 & 90 & 90 & 90 & 90 & 90 & 0 & 0 \end{bmatrix} \circ \begin{bmatrix} 1/3 & 1/3 & 1/3 \end{bmatrix} = \begin{bmatrix} & 0 & 30 & 60 & 90 & 90 & 60 & 30 & \end{bmatrix}$$

In the 1D case, it helps to make explicit the structure of the matrix:



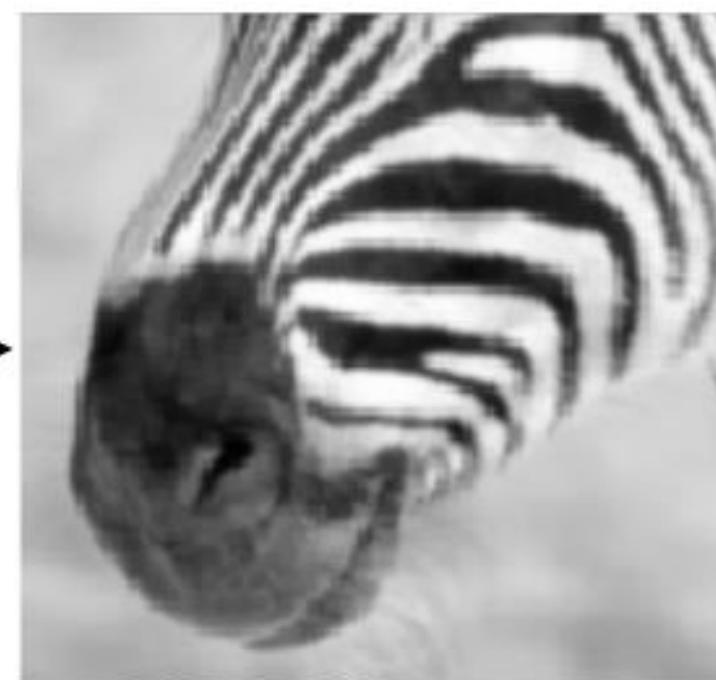
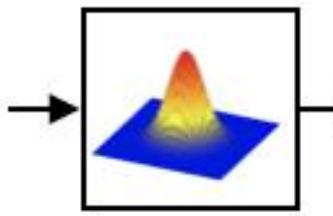
The Gaussian pyramid

$$\frac{1}{16} \begin{bmatrix} 6 & 4 & 1 & 0 & 0 & 0 & 0 & 0 \\ 4 & 6 & 4 & 1 & 0 & 0 & 0 & 0 \\ 1 & 4 & 6 & 4 & 1 & 0 & 0 & 0 \\ 0 & 1 & 4 & 6 & 4 & 1 & 0 & 0 \\ 0 & 0 & 1 & 4 & 6 & 4 & 1 & 0 \\ 0 & 0 & 0 & 1 & 4 & 6 & 4 & 1 \\ 0 & 0 & 0 & 0 & 1 & 4 & 6 & 4 \\ 0 & 0 & 0 & 0 & 0 & 1 & 4 & 6 \end{bmatrix}$$

256×256

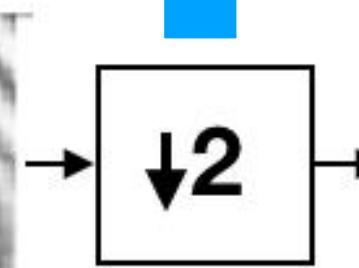


$[1, 4, 6, 4, 1]$



$$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$

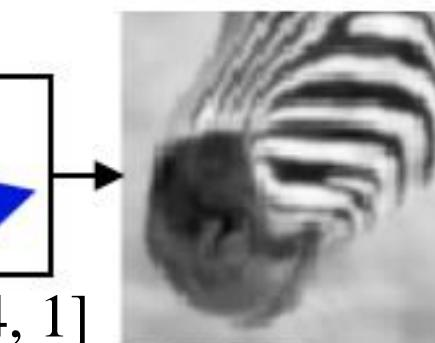
128×128



g_1



$[1, 4, 6, 4, 1]$



64×64

g_2

g_0

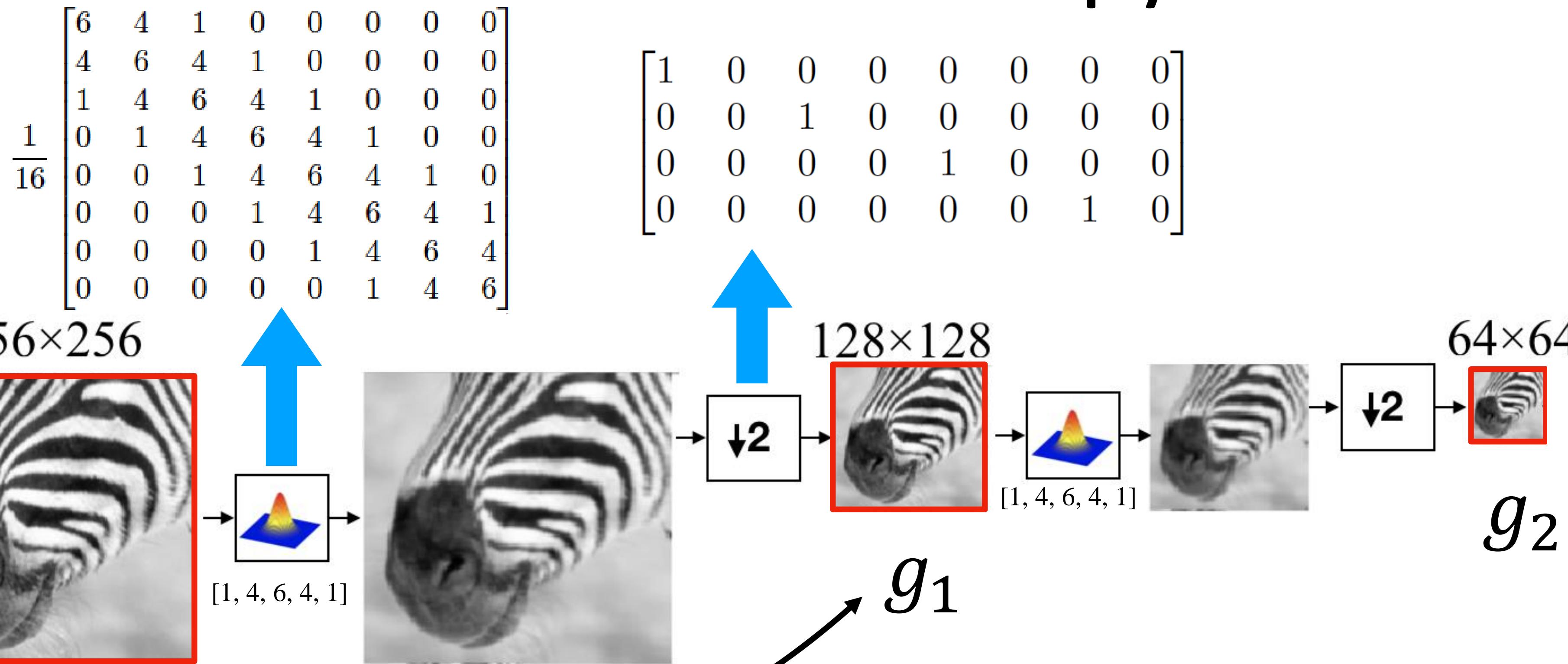
$g_1 = G_0 g_0$

$$G_0 = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix} \frac{1}{16}$$

$$\begin{bmatrix} 6 & 4 & 1 & 0 & 0 & 0 & 0 & 0 \\ 4 & 6 & 4 & 1 & 0 & 0 & 0 & 0 \\ 1 & 4 & 6 & 4 & 1 & 0 & 0 & 0 \\ 0 & 1 & 4 & 6 & 4 & 1 & 0 & 0 \\ 0 & 0 & 1 & 4 & 6 & 4 & 1 & 0 \\ 0 & 0 & 0 & 1 & 4 & 6 & 4 & 1 \\ 0 & 0 & 0 & 0 & 1 & 4 & 6 & 4 \\ 0 & 0 & 0 & 0 & 0 & 1 & 4 & 6 \end{bmatrix}$$

(The arrays shown here are for 1D signals)

The Gaussian pyramid

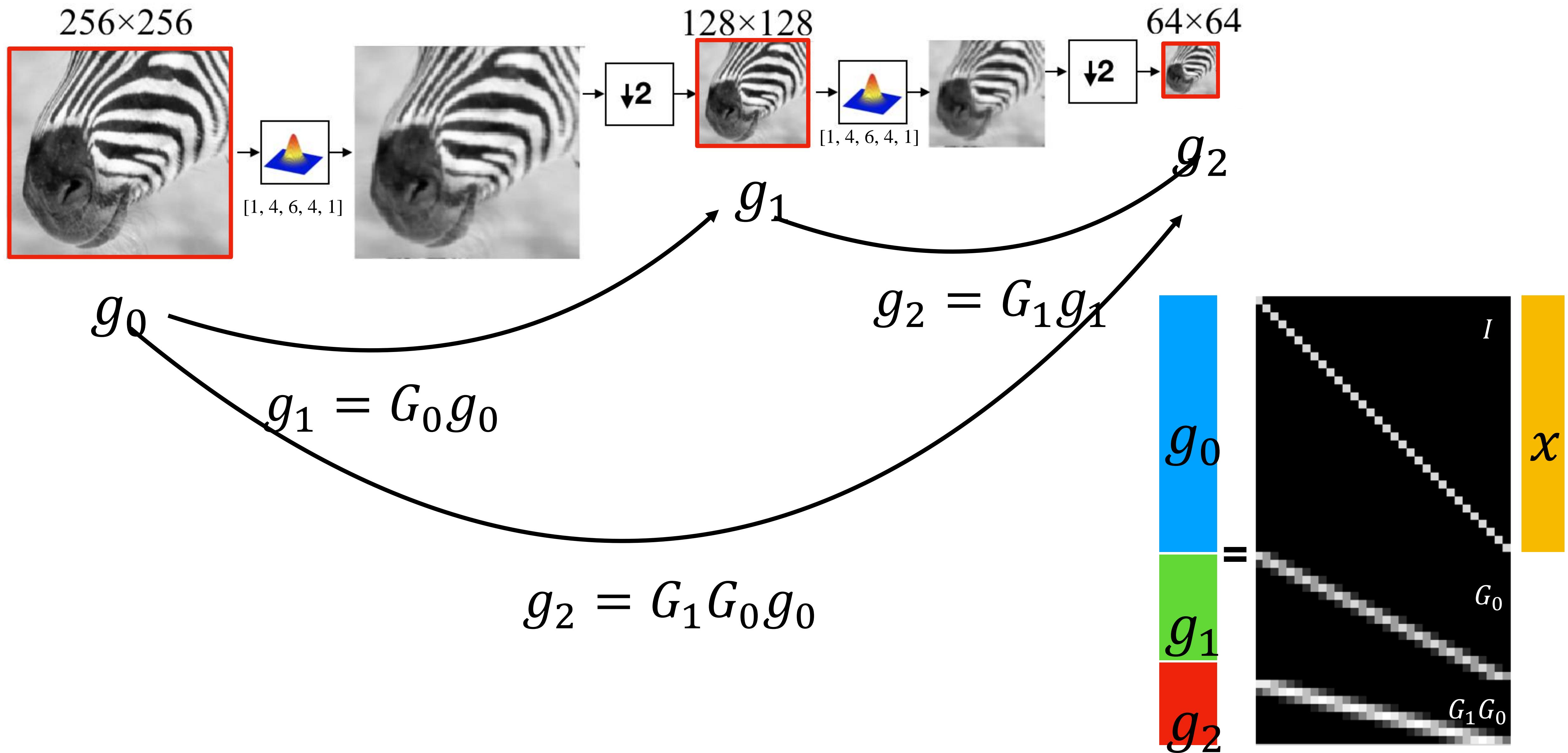


g_0

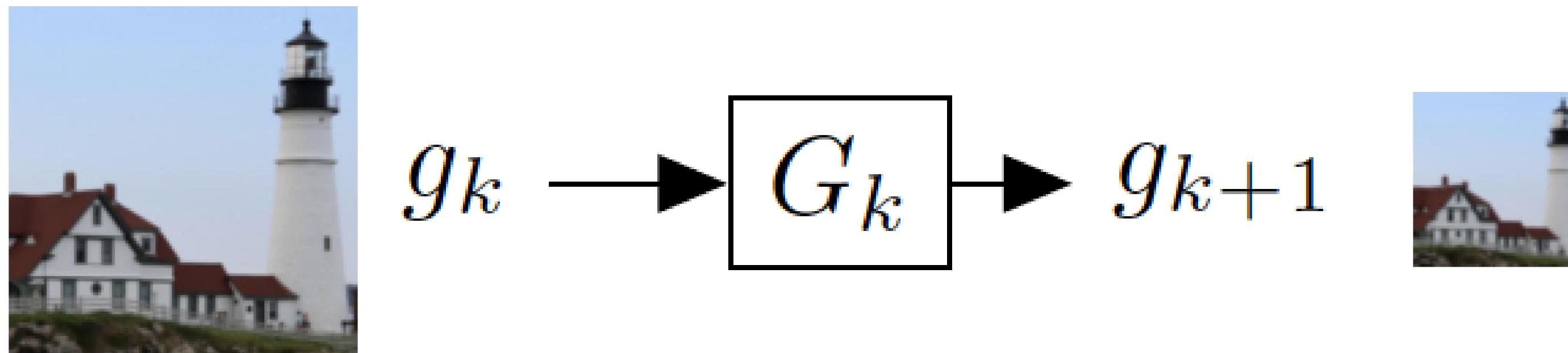
$$g_1 = G_0 g_0$$

$$G_0 = \frac{1}{16} \begin{bmatrix} 6 & 4 & 1 & 0 & 0 & 0 & 0 \\ 1 & 4 & 6 & 4 & 1 & 0 & 0 \\ 0 & 0 & 1 & 4 & 6 & 4 & 0 \\ 0 & 0 & 0 & 1 & 4 & 6 & 4 \\ 0 & 0 & 0 & 0 & 1 & 4 & 6 \\ 0 & 0 & 0 & 0 & 0 & 1 & 4 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

The Gaussian pyramid



The Gaussian pyramid



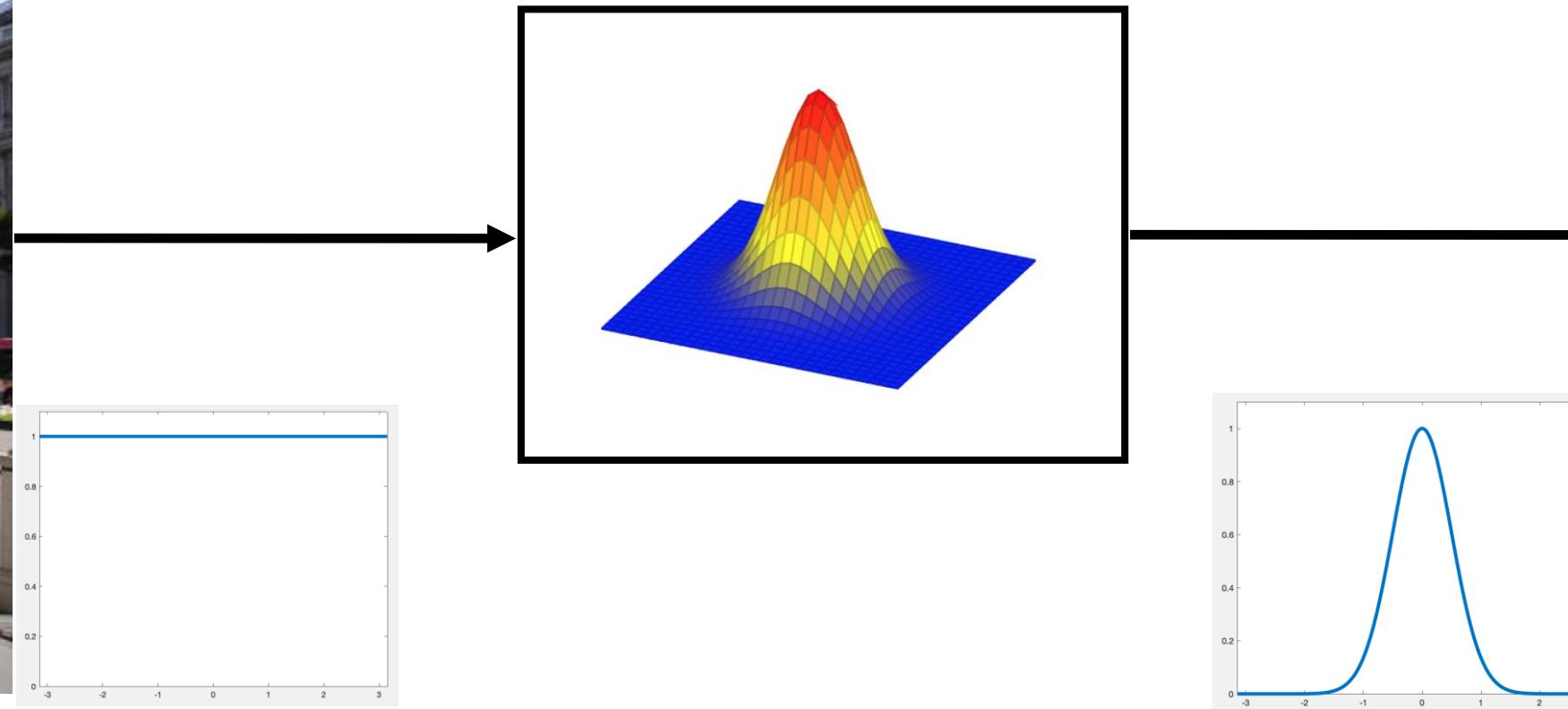
For each level

1. Blur input image with a Gaussian filter
2. Downsample image

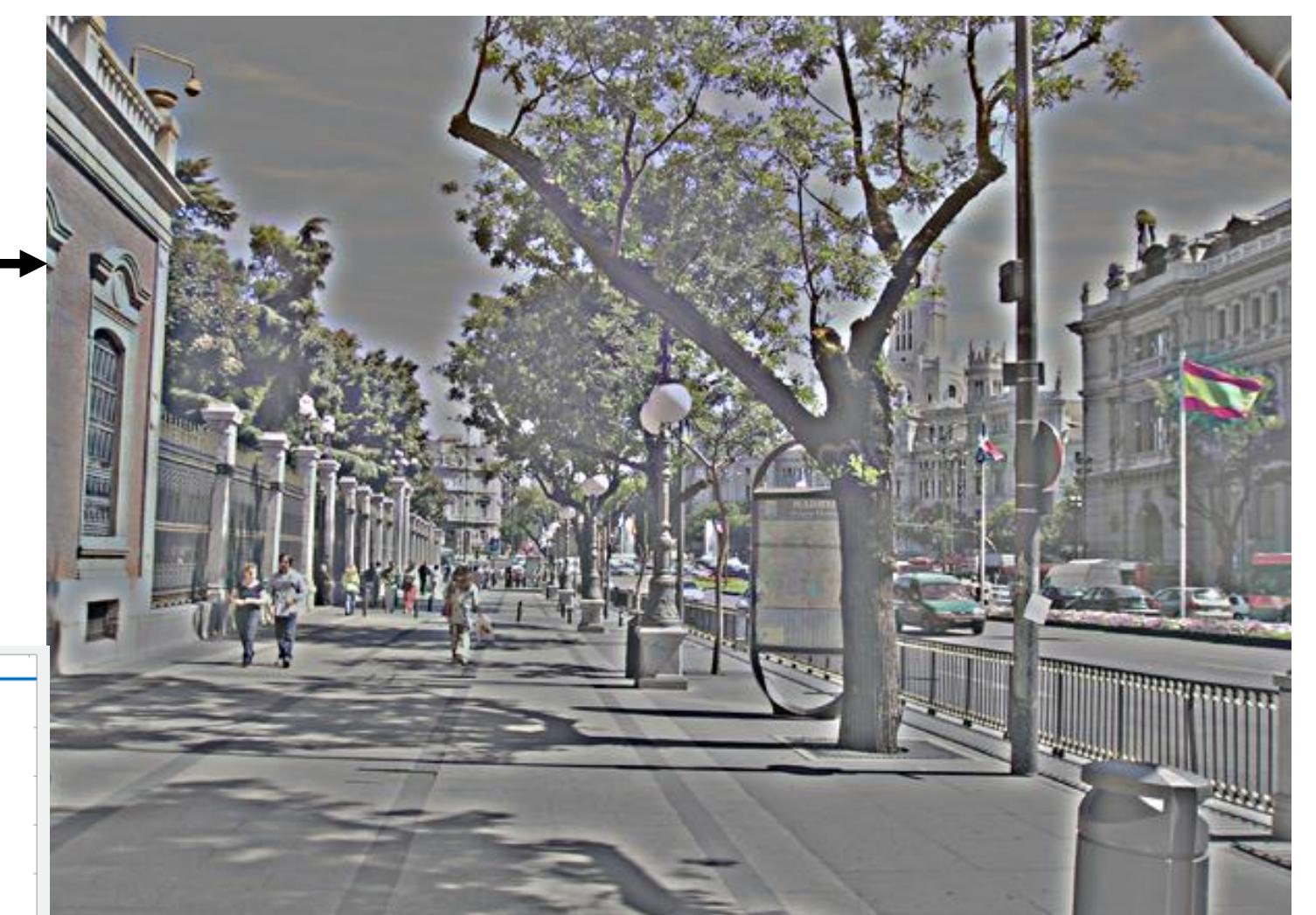
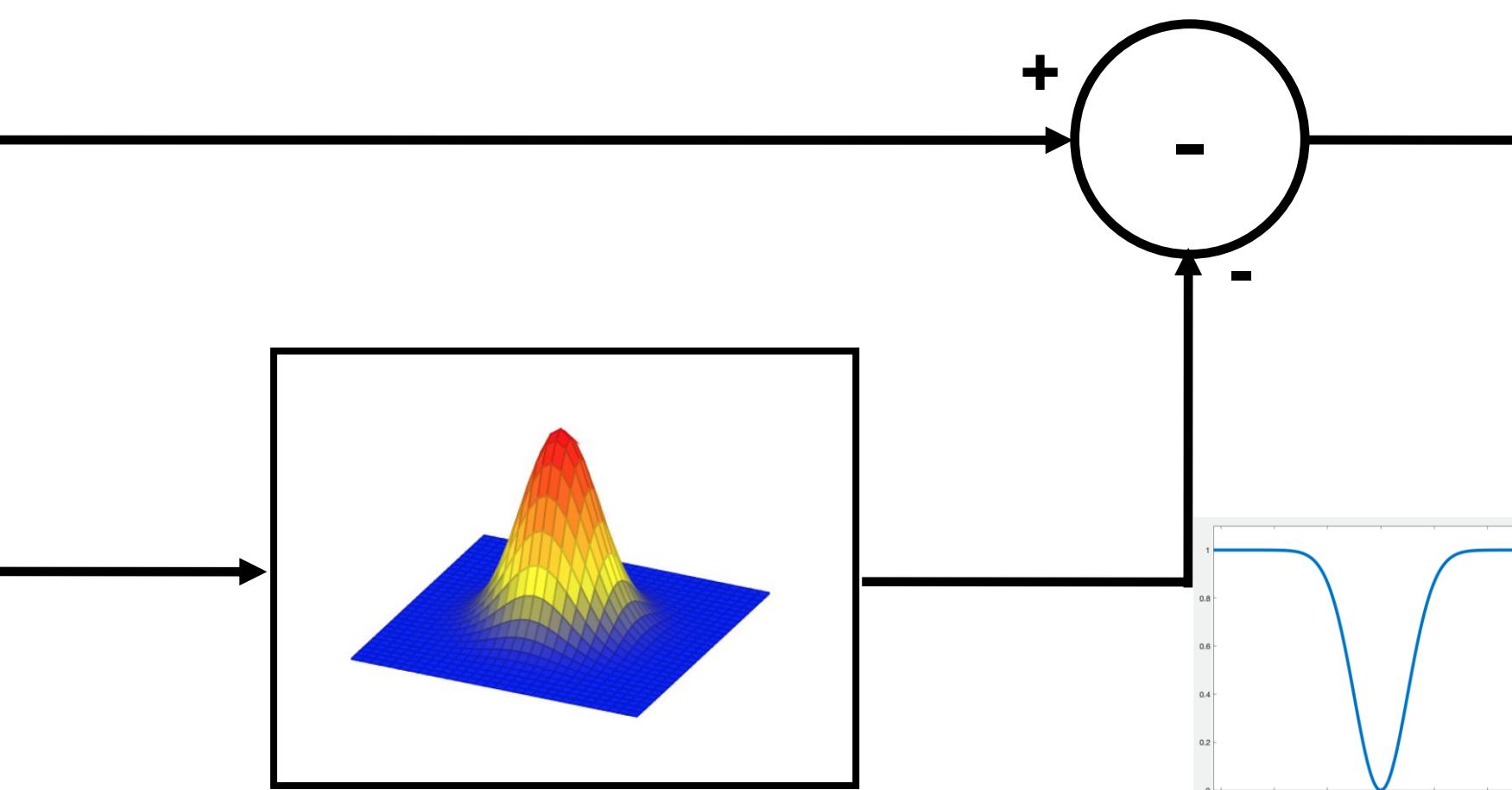
What about the opposite of blurring?



Gaussian filter

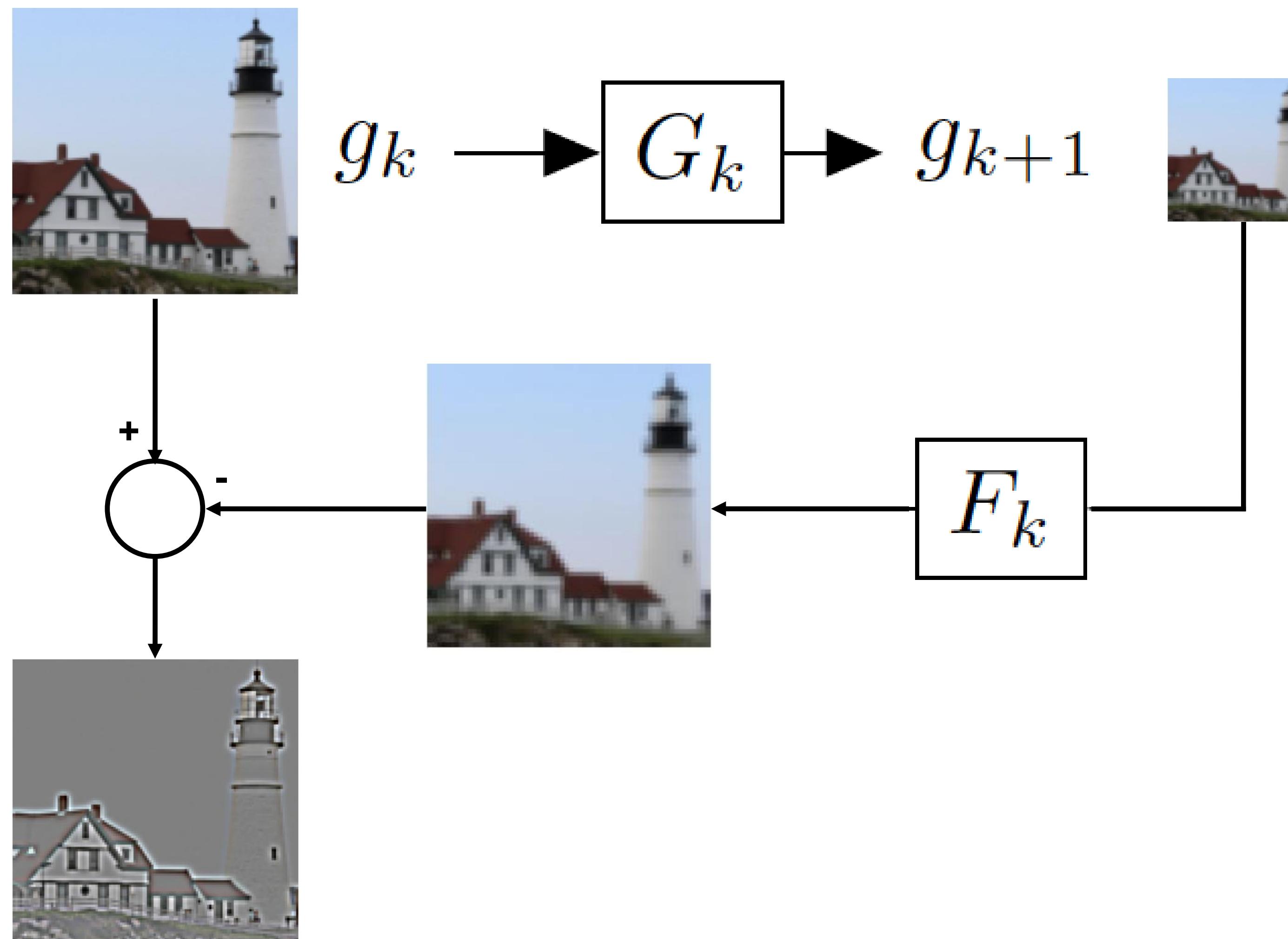


Laplacian filter

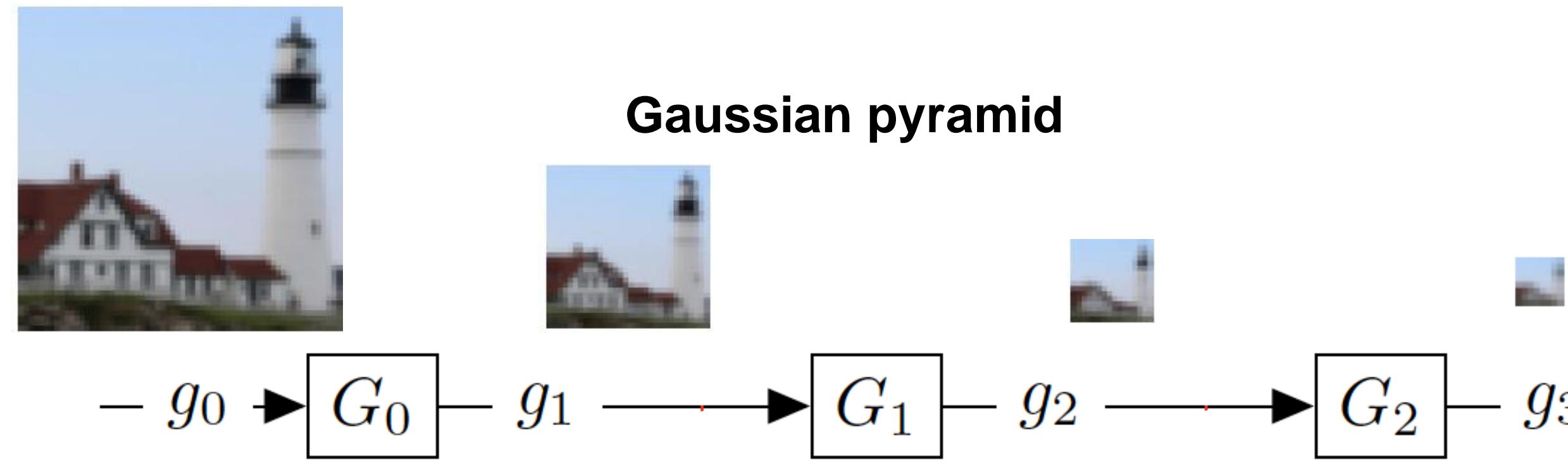


The Laplacian Pyramid

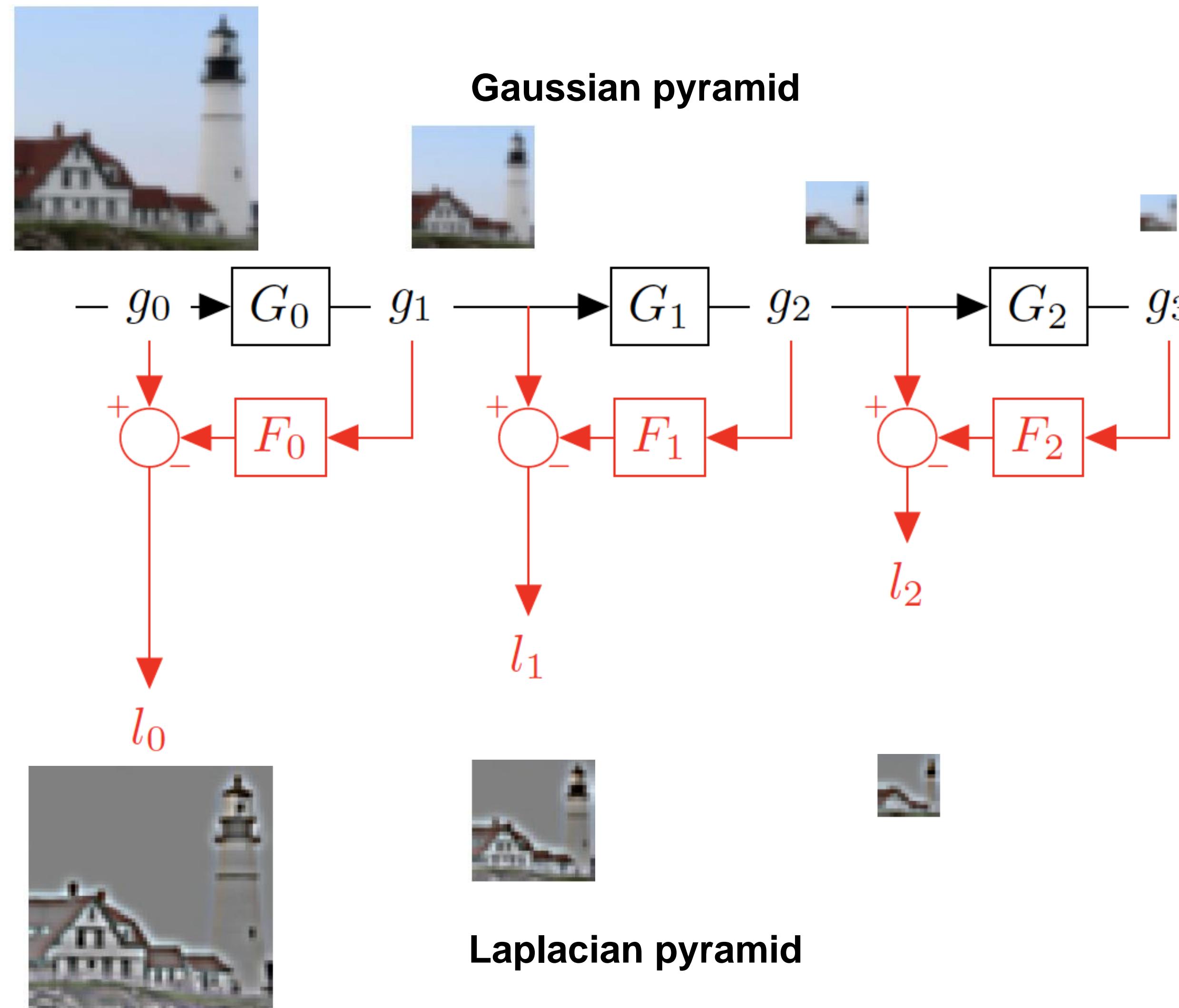
Compute the difference between upsampled Gaussian pyramid level $k+1$ and Gaussian pyramid level k .



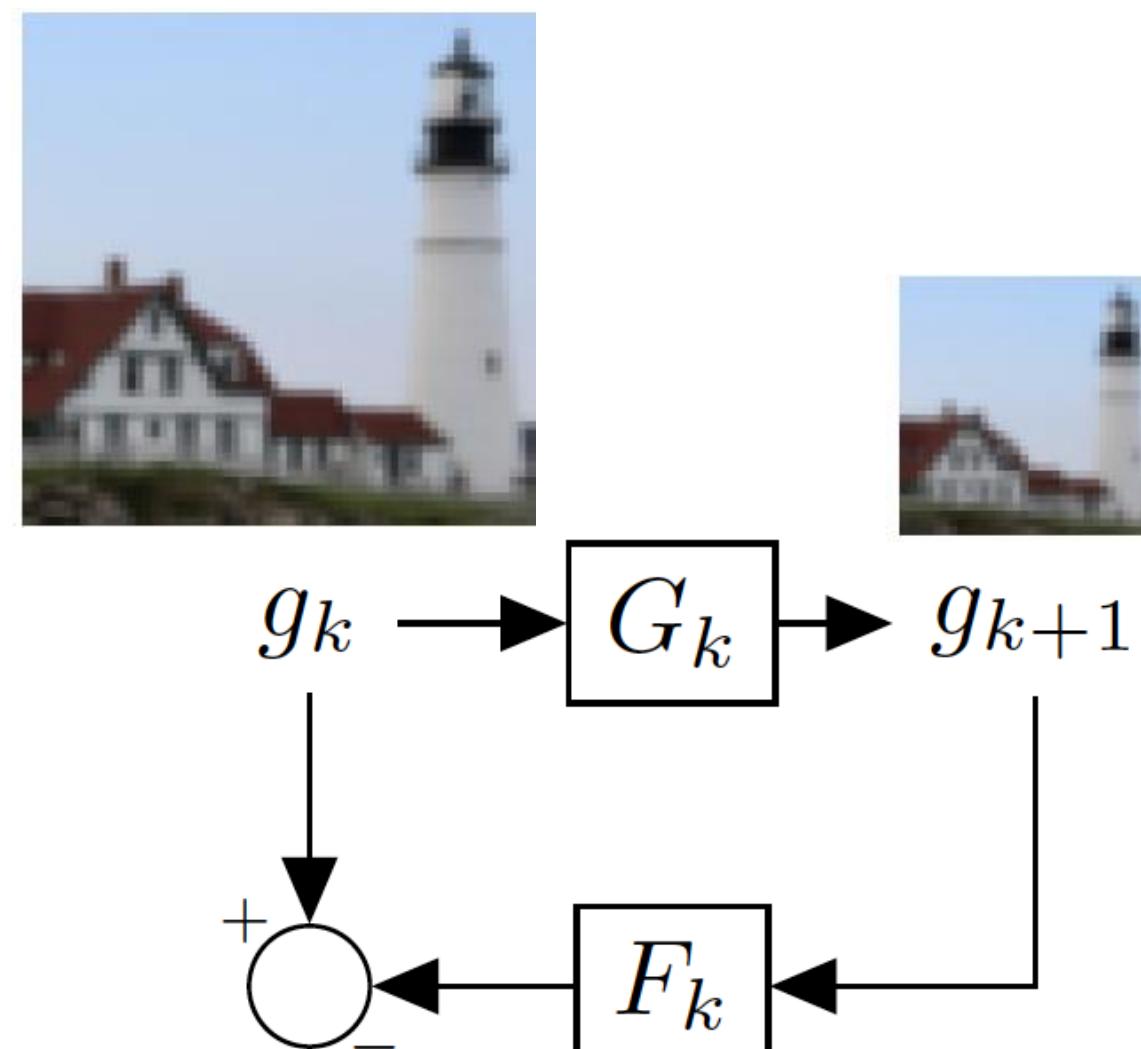
The Laplacian Pyramid



The Laplacian Pyramid



The Laplacian Pyramid



Blurring and downsampling:

$$G_0 = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix} \frac{1}{16}$$

(Downsampling by 2)

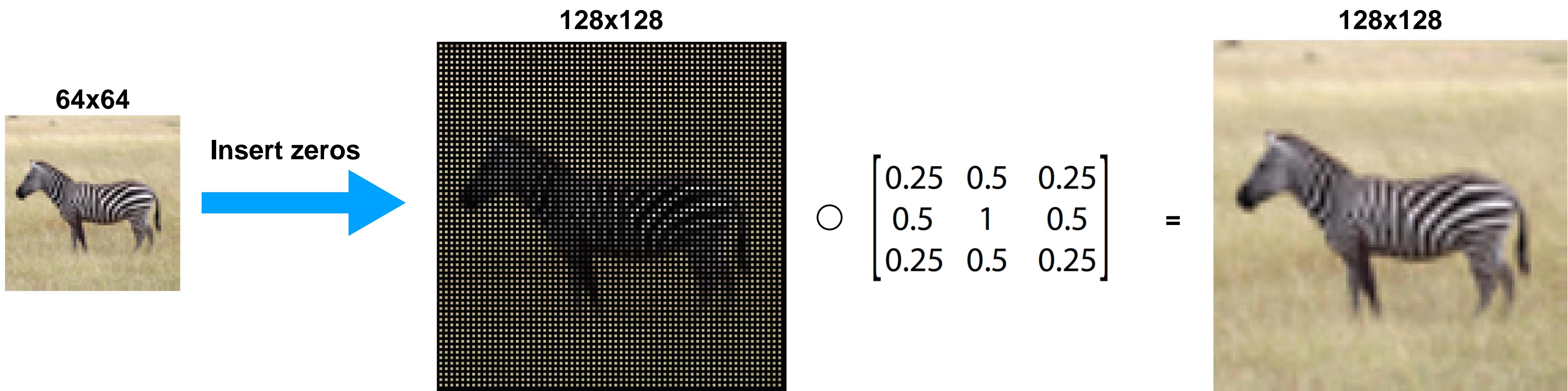
Upsampling and blurring:

$$F_0 =$$

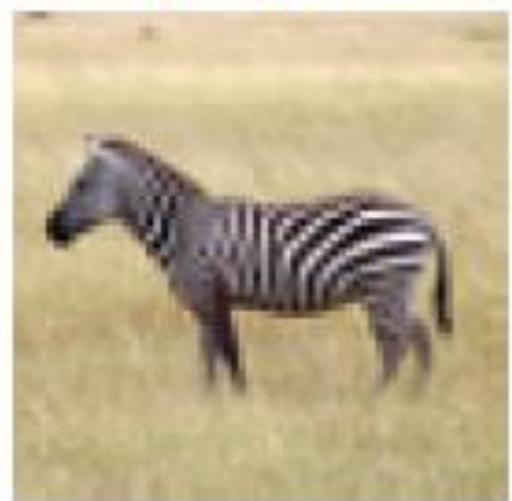
$$\begin{bmatrix} 6 & 4 & 1 & 0 & 0 & 0 & 0 & 0 \\ 4 & 6 & 4 & 1 & 0 & 0 & 0 & 0 \\ 1 & 4 & 6 & 4 & 1 & 0 & 0 & 0 \\ 0 & 1 & 4 & 6 & 4 & 1 & 0 & 0 \\ 0 & 0 & 1 & 4 & 6 & 4 & 1 & 0 \\ 0 & 0 & 0 & 1 & 4 & 6 & 4 & 1 \\ 0 & 0 & 0 & 0 & 1 & 4 & 6 & 4 \\ 0 & 0 & 0 & 0 & 0 & 1 & 4 & 6 \end{bmatrix}$$

(blur)

Upsampling



64x64



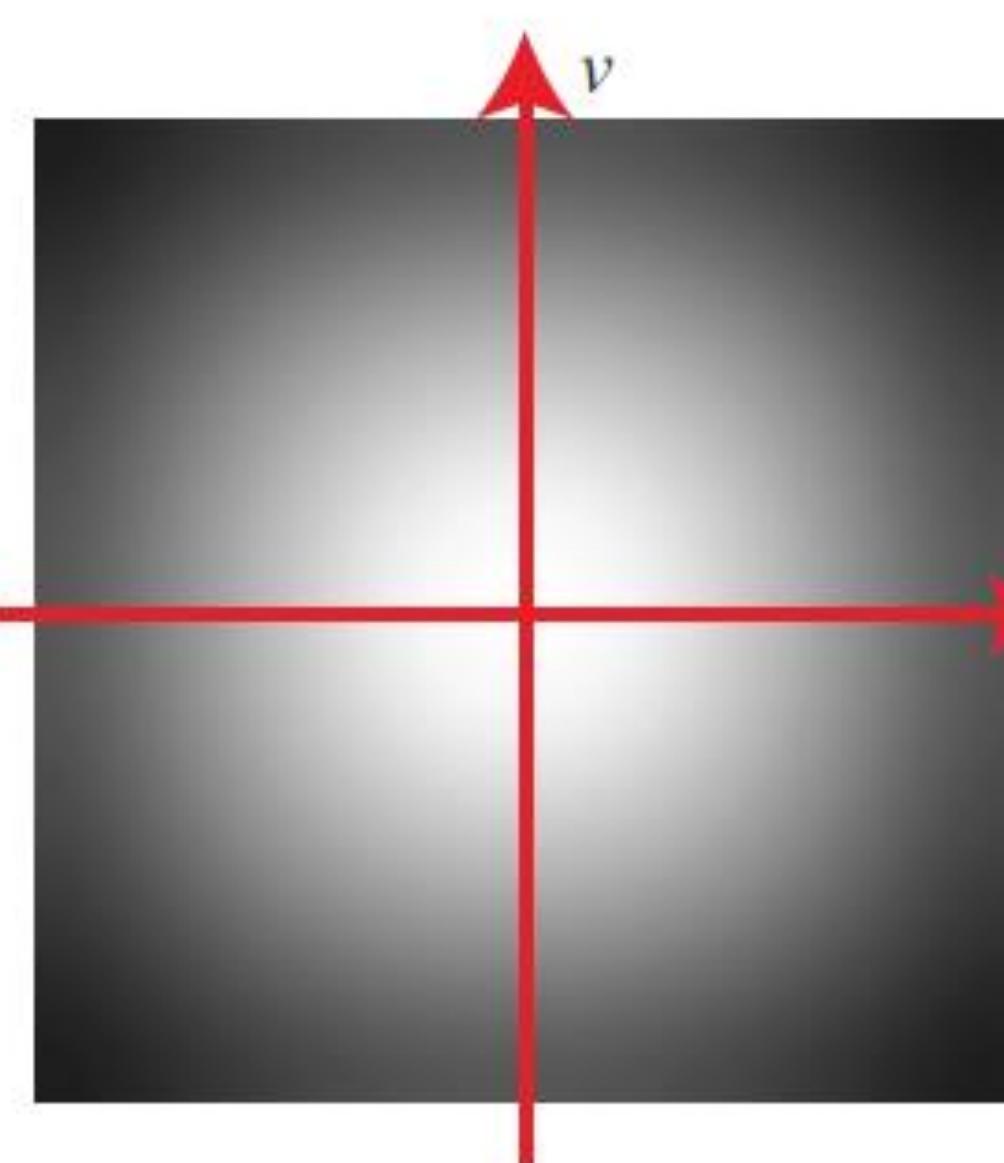
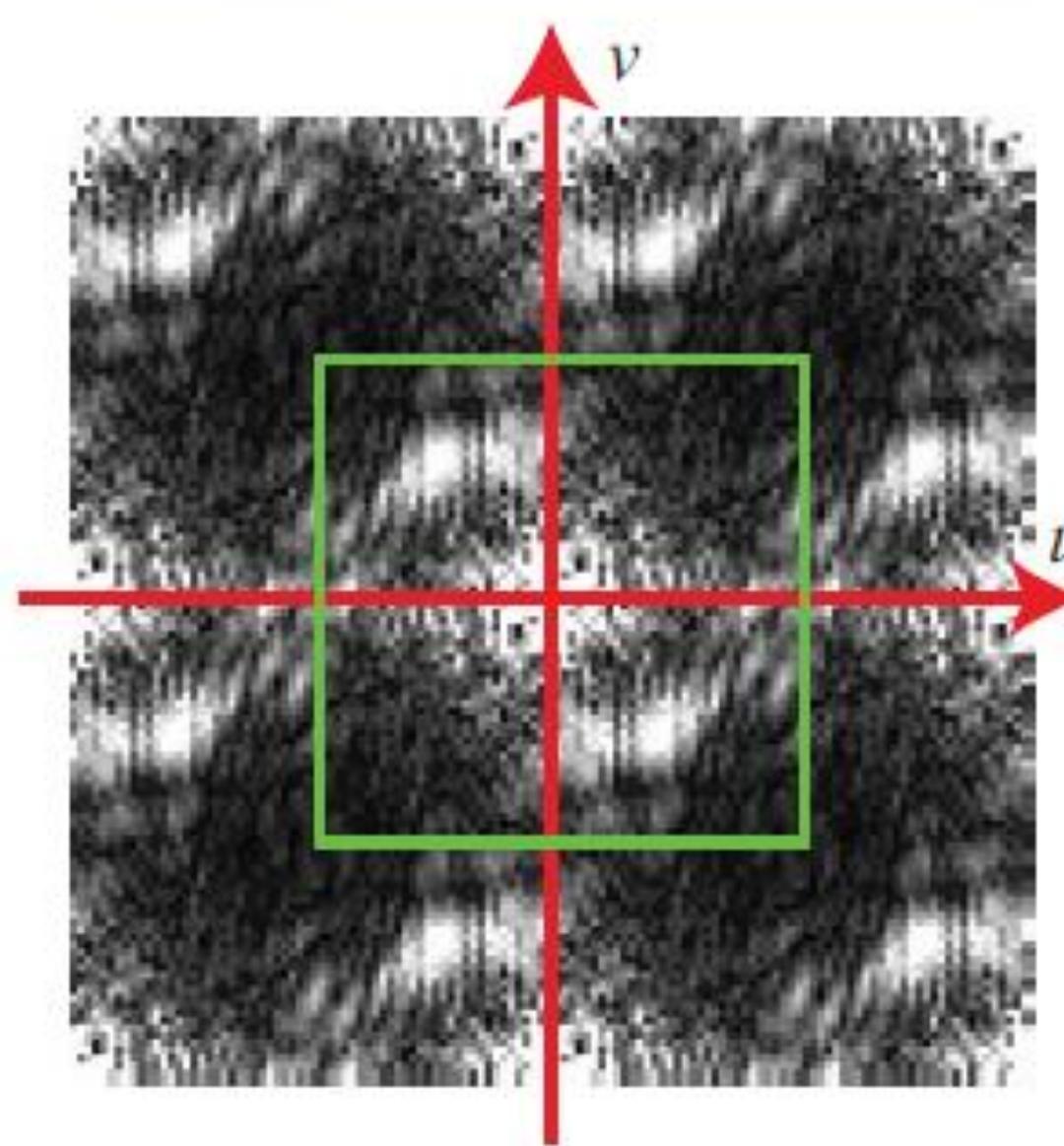
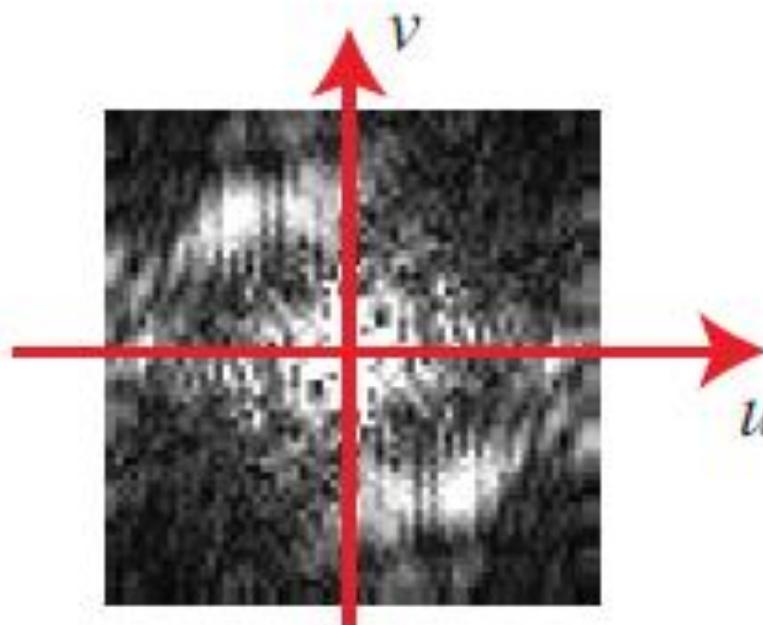
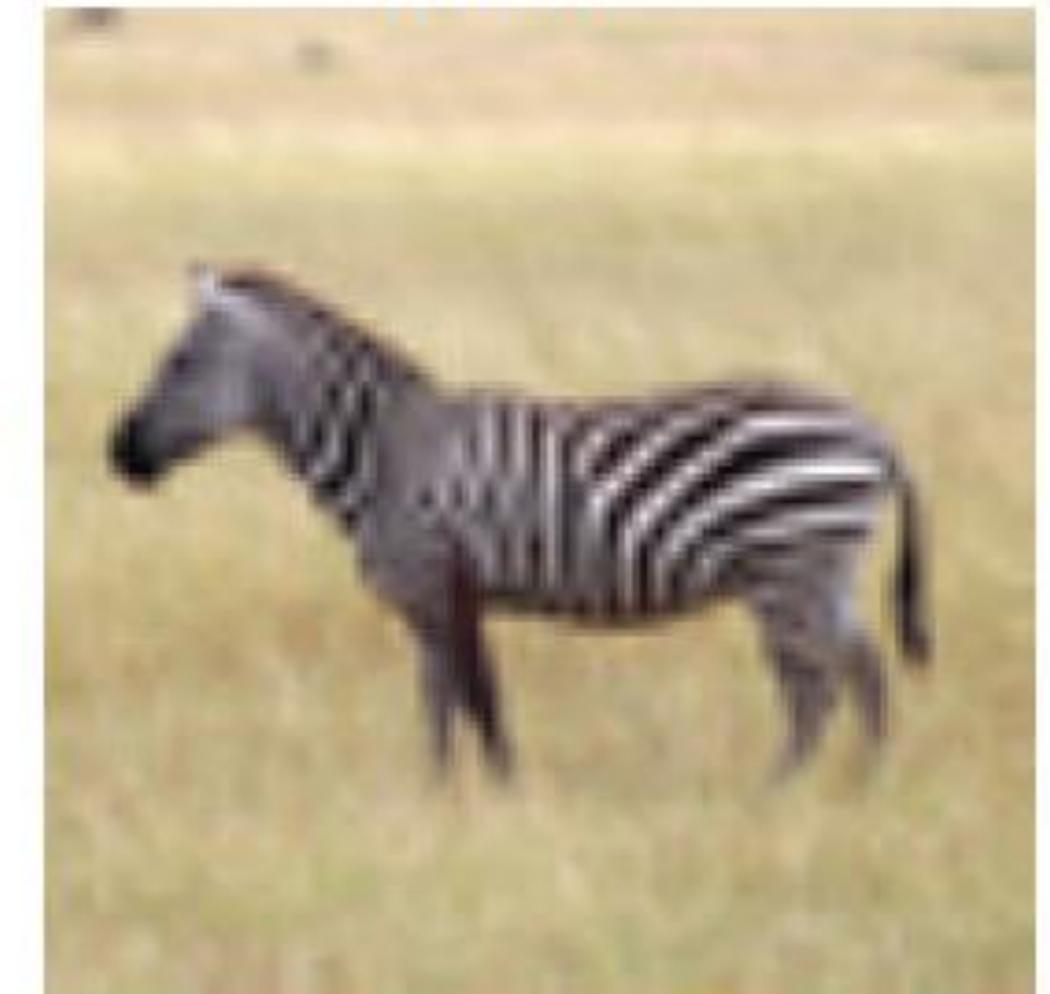
128x128



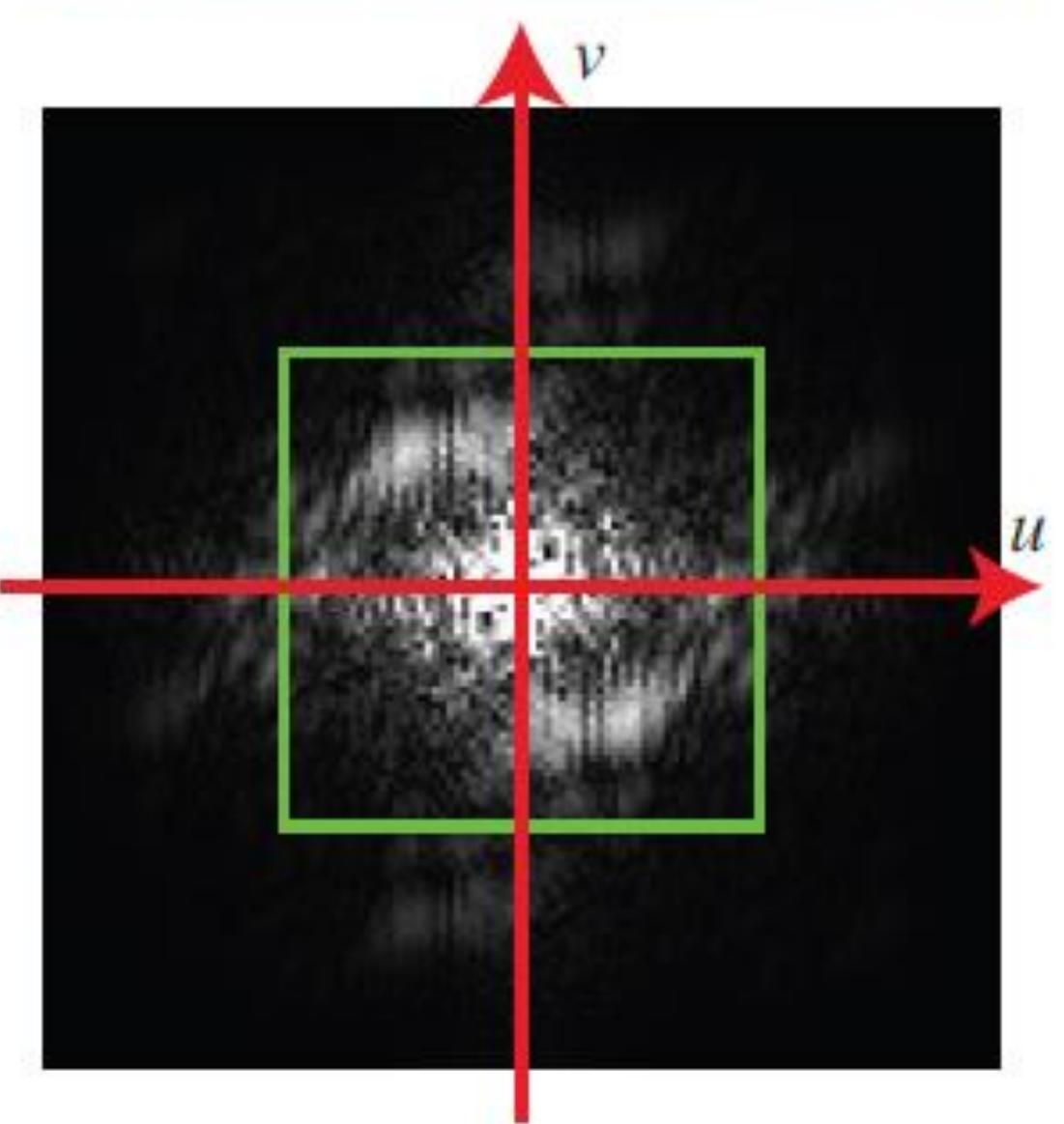
$$\begin{bmatrix} 0.25 & 0.5 & 0.25 \\ 0.5 & 1 & 0.5 \\ 0.25 & 0.5 & 0.25 \end{bmatrix}$$

=

128x128



=



Upsampling

1	0	1	0	1
0	0	0	0	0
1	0	1	0	1
0	0	0	0	0
1	0	1	0	1

$$\circ \begin{bmatrix} 0.25 & 0.5 & 0.25 \\ 0.5 & 1 & 0.5 \\ 0.25 & 0.5 & 0.25 \end{bmatrix} = ?$$

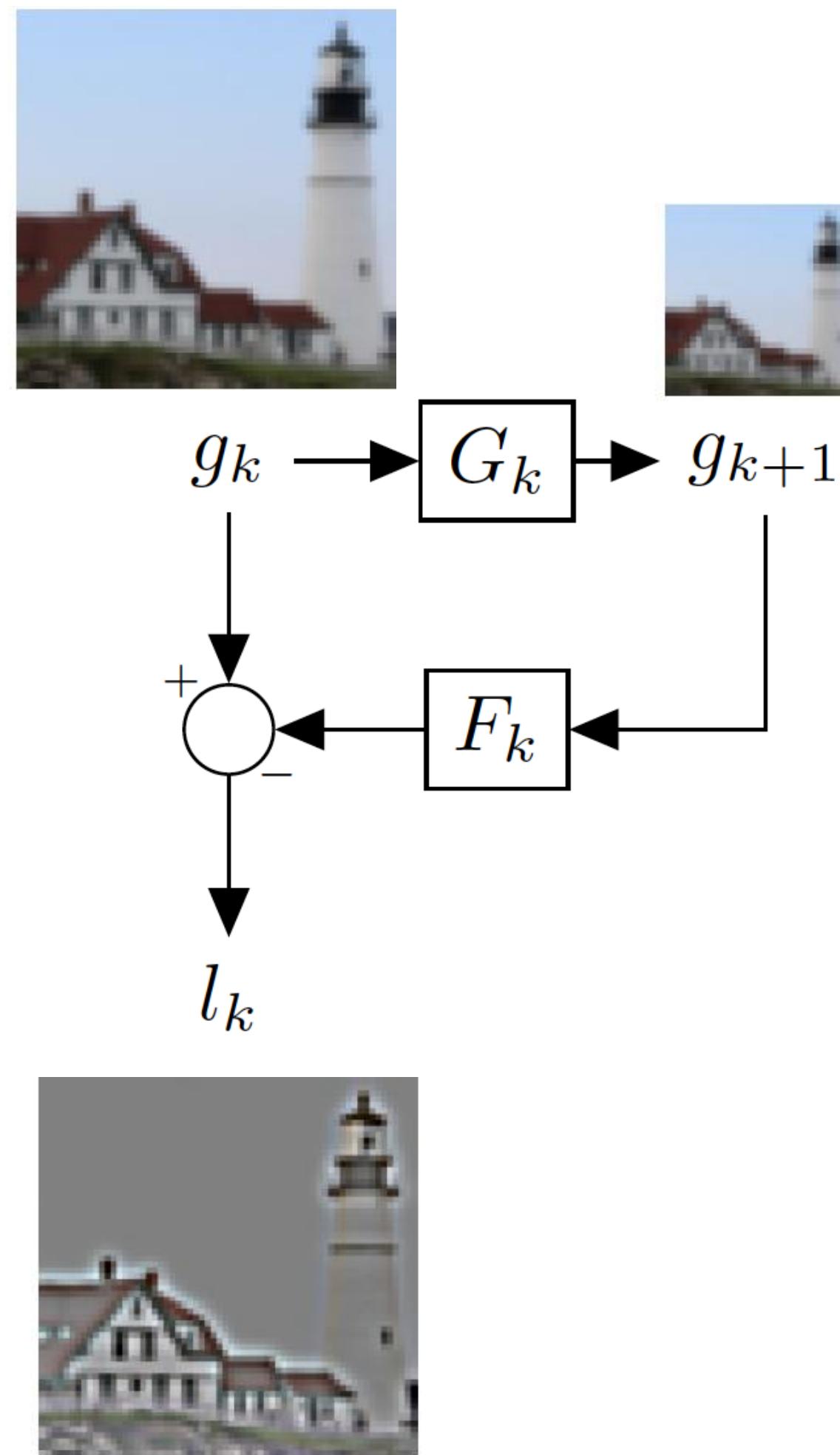
Upsampling

1	0	1	0	1
0	0	0	0	0
1	0	1	0	1
0	0	0	0	0
1	0	1	0	1

$$\odot \begin{bmatrix} 0.25 & 0.5 & 0.25 \\ 0.5 & 1 & 0.5 \\ 0.25 & 0.5 & 0.25 \end{bmatrix} =$$

1	1	1	1	1
1	1	1	1	1
1	1	1	1	1
1	1	1	1	1
1	1	1	1	1

The Laplacian Pyramid



Blurring and downsampling:

$$G_0 = \frac{1}{16} \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix} \quad \text{(Downsampling by 2)}$$

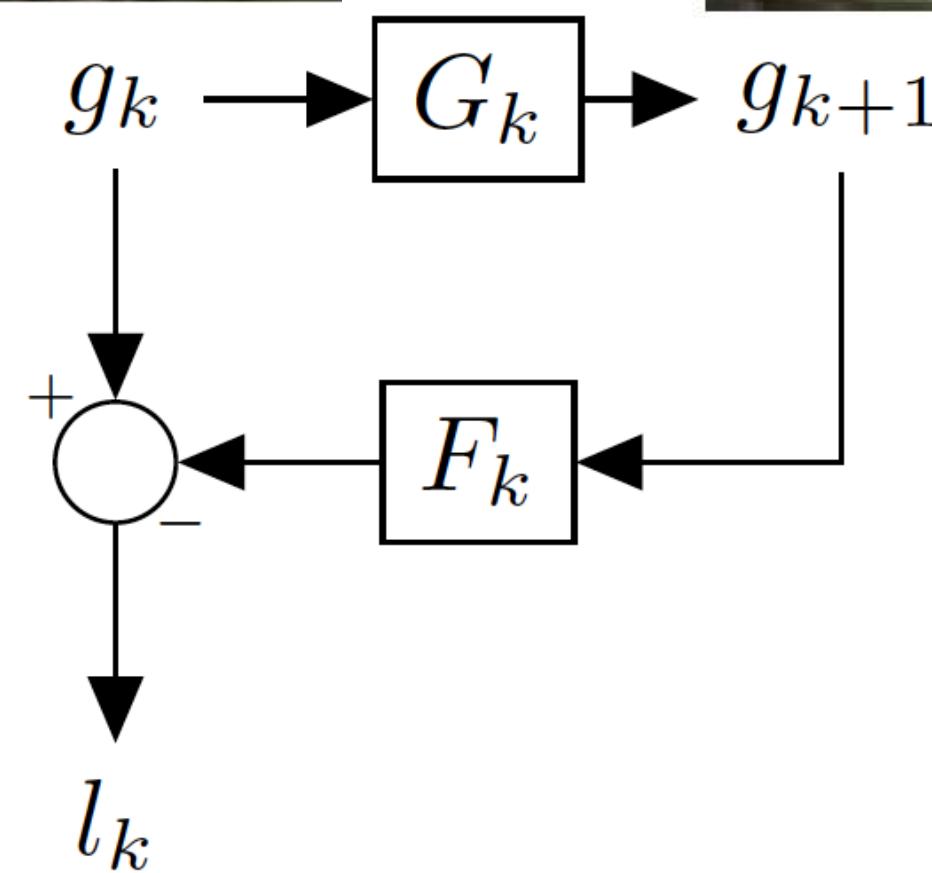
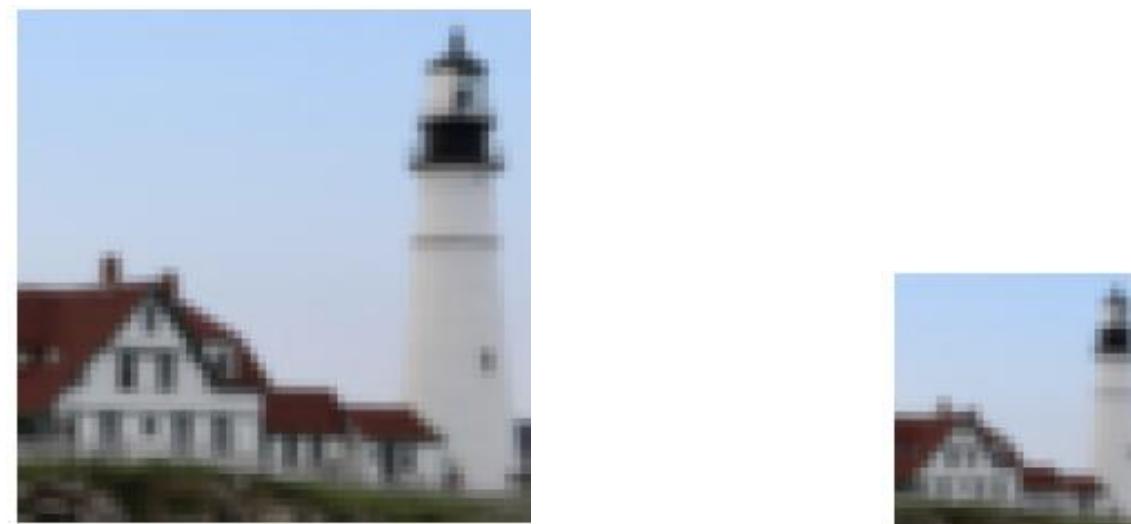
$$\begin{bmatrix} 6 & 4 & 1 & 0 & 0 & 0 & 0 & 0 \\ 4 & 6 & 4 & 1 & 0 & 0 & 0 & 0 \\ 1 & 4 & 6 & 4 & 1 & 0 & 0 & 0 \\ 0 & 1 & 4 & 6 & 4 & 1 & 0 & 0 \\ 0 & 0 & 1 & 4 & 6 & 4 & 1 & 0 \\ 0 & 0 & 0 & 1 & 4 & 6 & 4 & 1 \\ 0 & 0 & 0 & 0 & 1 & 4 & 6 & 4 \\ 0 & 0 & 0 & 0 & 0 & 1 & 4 & 6 \end{bmatrix} \quad \text{(blur)}$$

Upsampling and blurring:

$$F_0 = \frac{1}{8} \begin{bmatrix} 6 & 4 & 1 & 0 & 0 & 0 & 0 & 0 \\ 4 & 6 & 4 & 1 & 0 & 0 & 0 & 0 \\ 1 & 4 & 6 & 4 & 1 & 0 & 0 & 0 \\ 0 & 1 & 4 & 6 & 4 & 1 & 0 & 0 \\ 0 & 0 & 1 & 4 & 6 & 4 & 1 & 0 \\ 0 & 0 & 0 & 1 & 4 & 6 & 4 & 1 \\ 0 & 0 & 0 & 0 & 1 & 4 & 6 & 4 \\ 0 & 0 & 0 & 0 & 0 & 1 & 4 & 6 \end{bmatrix} \quad \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix} \quad \begin{bmatrix} \text{(blur)} & & & \\ & \text{(Upsampling by 2)} & & \end{bmatrix}$$

$$l_0 = (I_0 - F_0 G_0) g_0$$

The Laplacian Pyramid



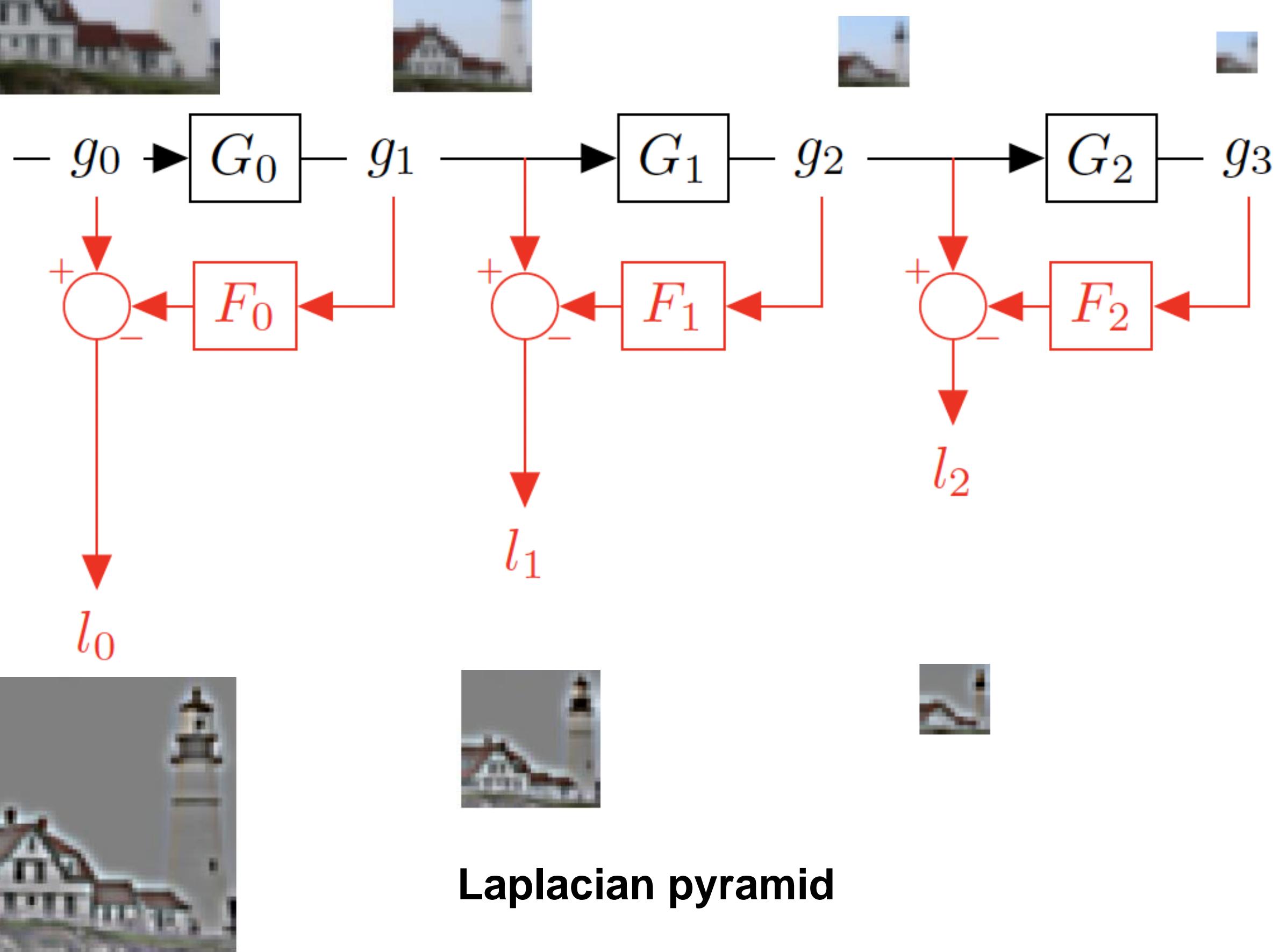
$$l_0 = (I_0 - F_0 G_0) g_0$$

$$= \frac{1}{256} \begin{bmatrix} 182 & -56 & -24 & -8 & -2 & 0 & 0 & 0 \\ -56 & 192 & -56 & -32 & -8 & 0 & 0 & 0 \\ -24 & -56 & 180 & -56 & -24 & -8 & -2 & 0 \\ -8 & -32 & -56 & 192 & -56 & -32 & -8 & 0 \\ -2 & -8 & -24 & -56 & 180 & -56 & -24 & -8 \\ 0 & 0 & -8 & -32 & -56 & 192 & -56 & -32 \\ 0 & 0 & -2 & -8 & -24 & -56 & 182 & -48 \\ 0 & 0 & 0 & 0 & -8 & -32 & -48 & 224 \end{bmatrix} x$$

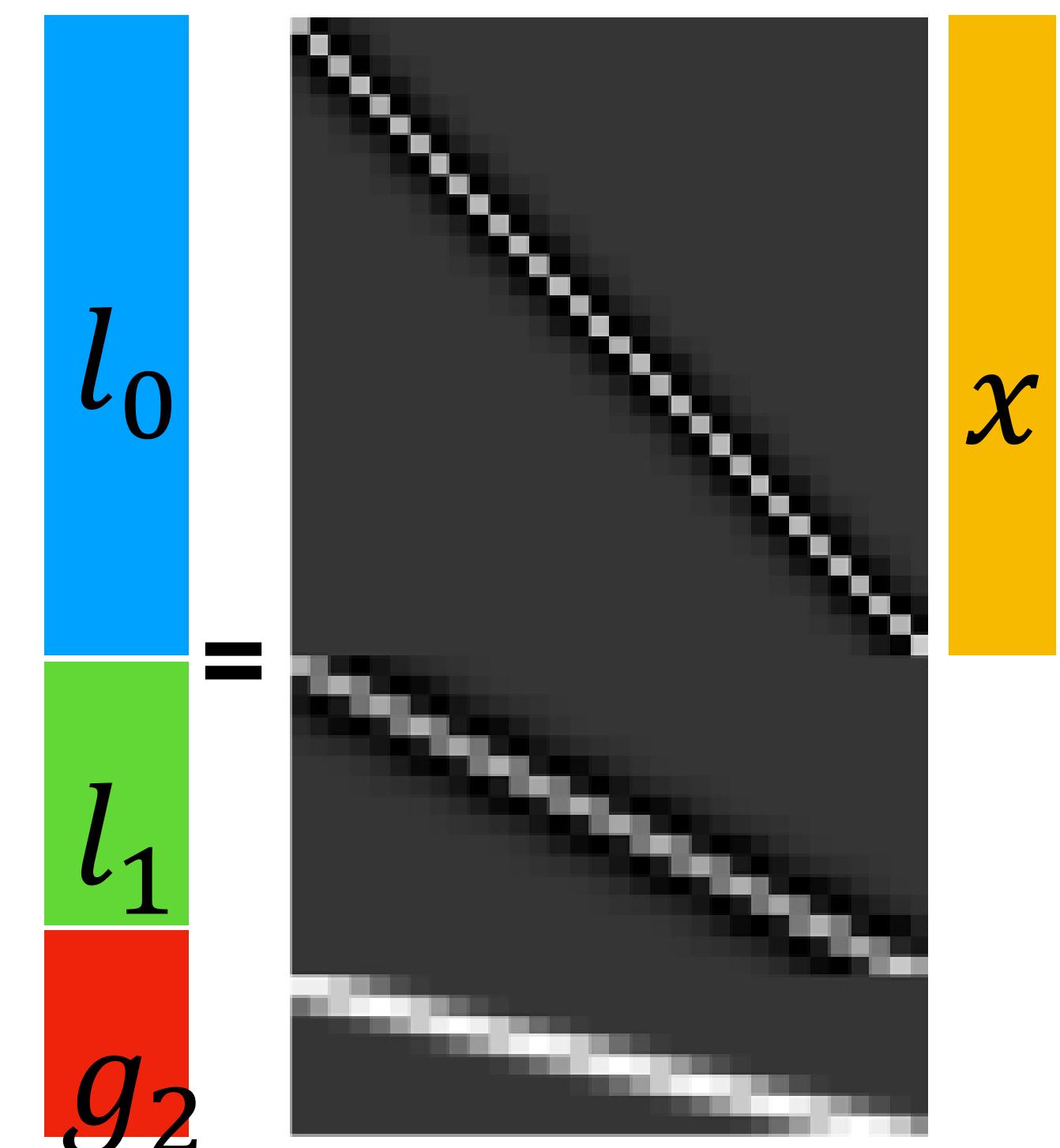
The Laplacian Pyramid



Gaussian pyramid



Laplacian pyramid



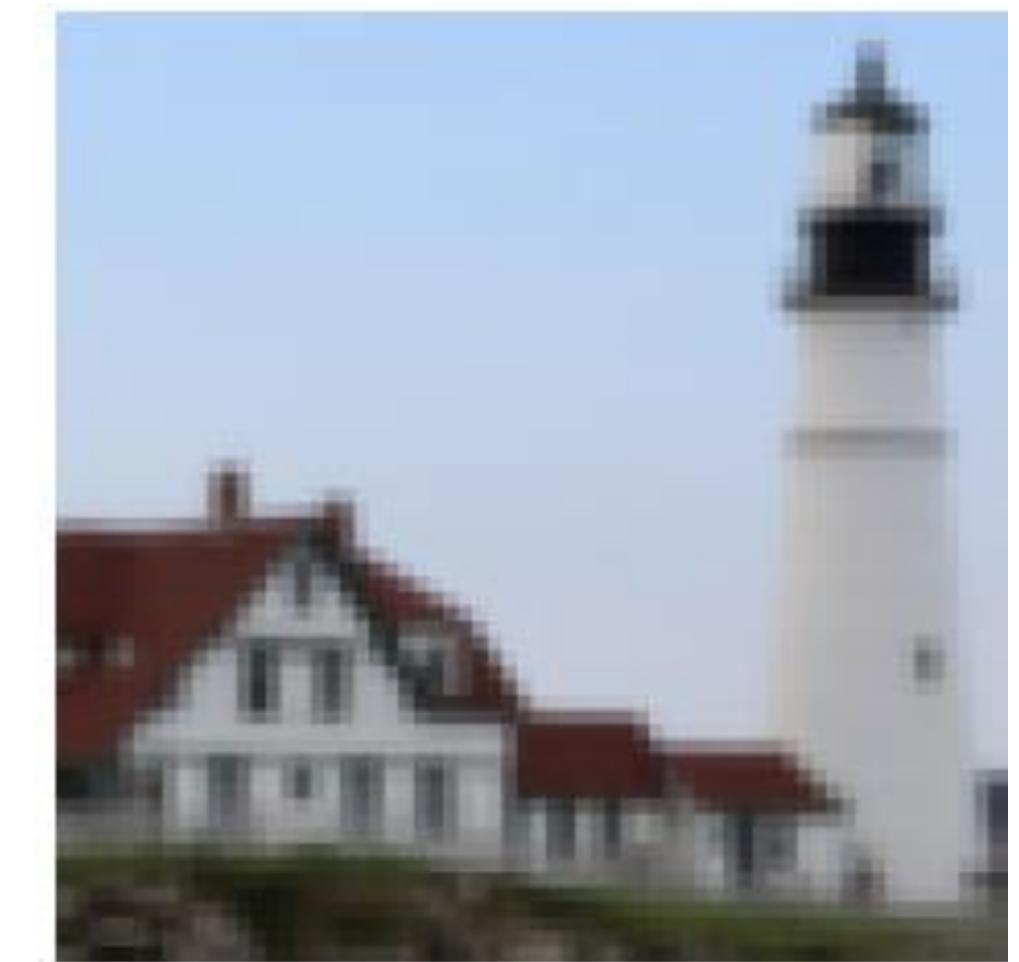
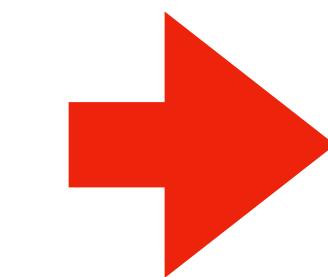
The Laplacian Pyramid



Laplacian pyramid



Gaussian
residual

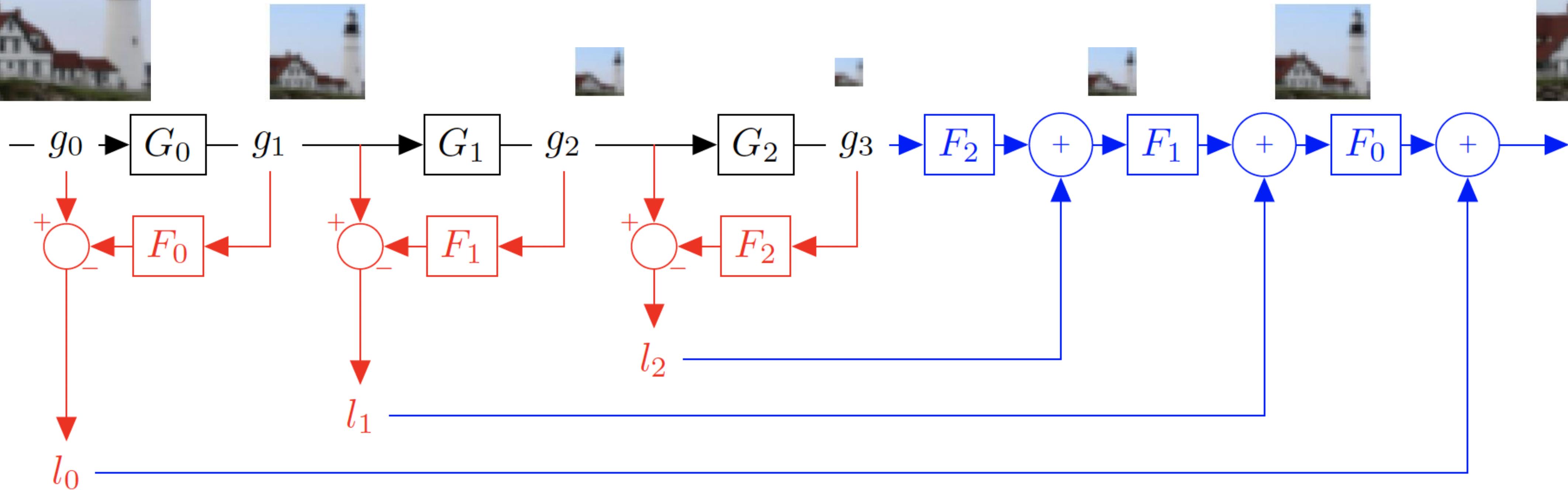


Can we invert the
Laplacian Pyramid?

The Laplacian Pyramid



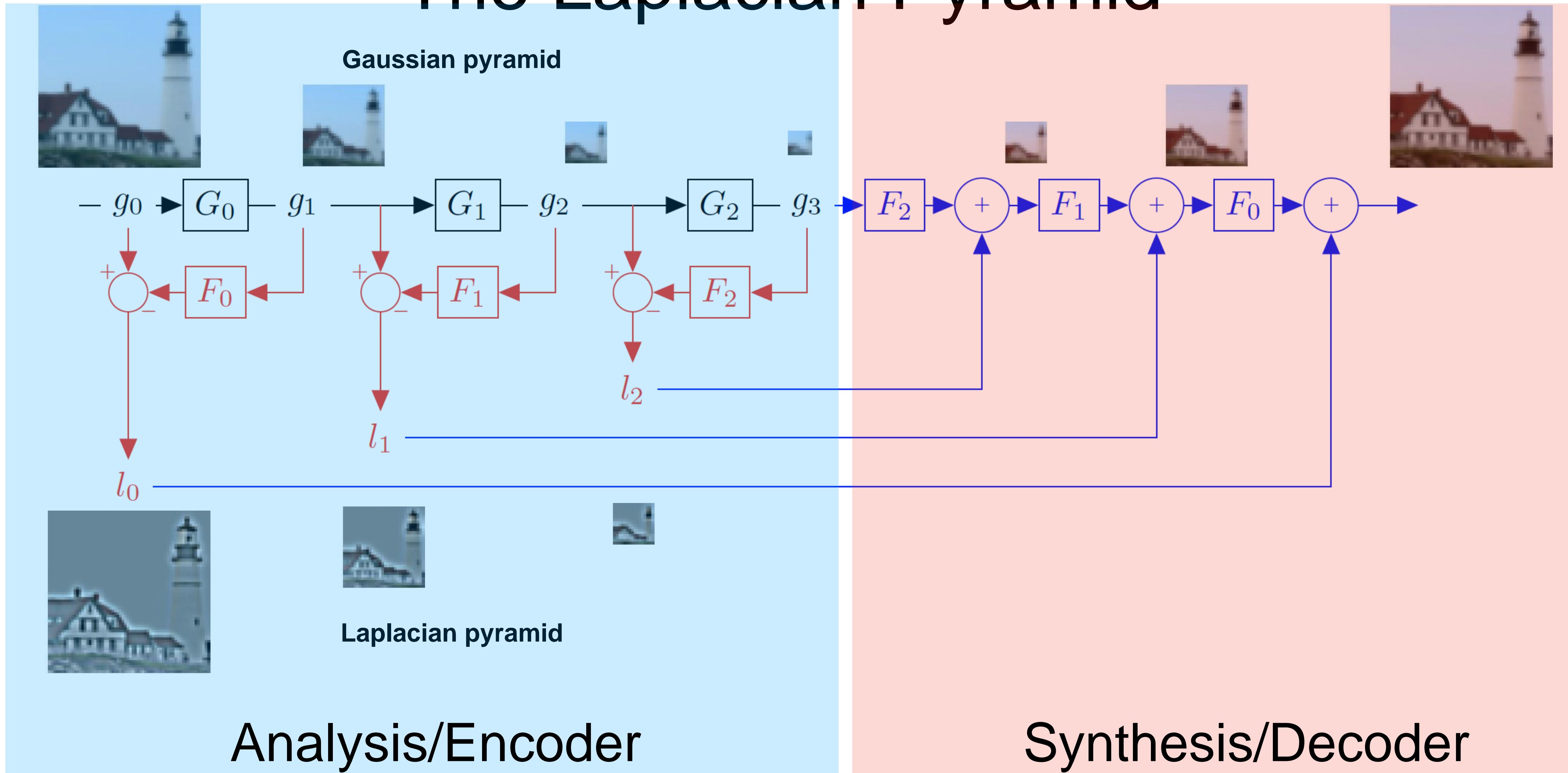
Gaussian pyramid



Laplacian pyramid



The Laplacian Pyramid



The Laplacian Pyramid

Laplacian pyramid



Gaussian residual



Laplacian pyramid



Gaussian residual

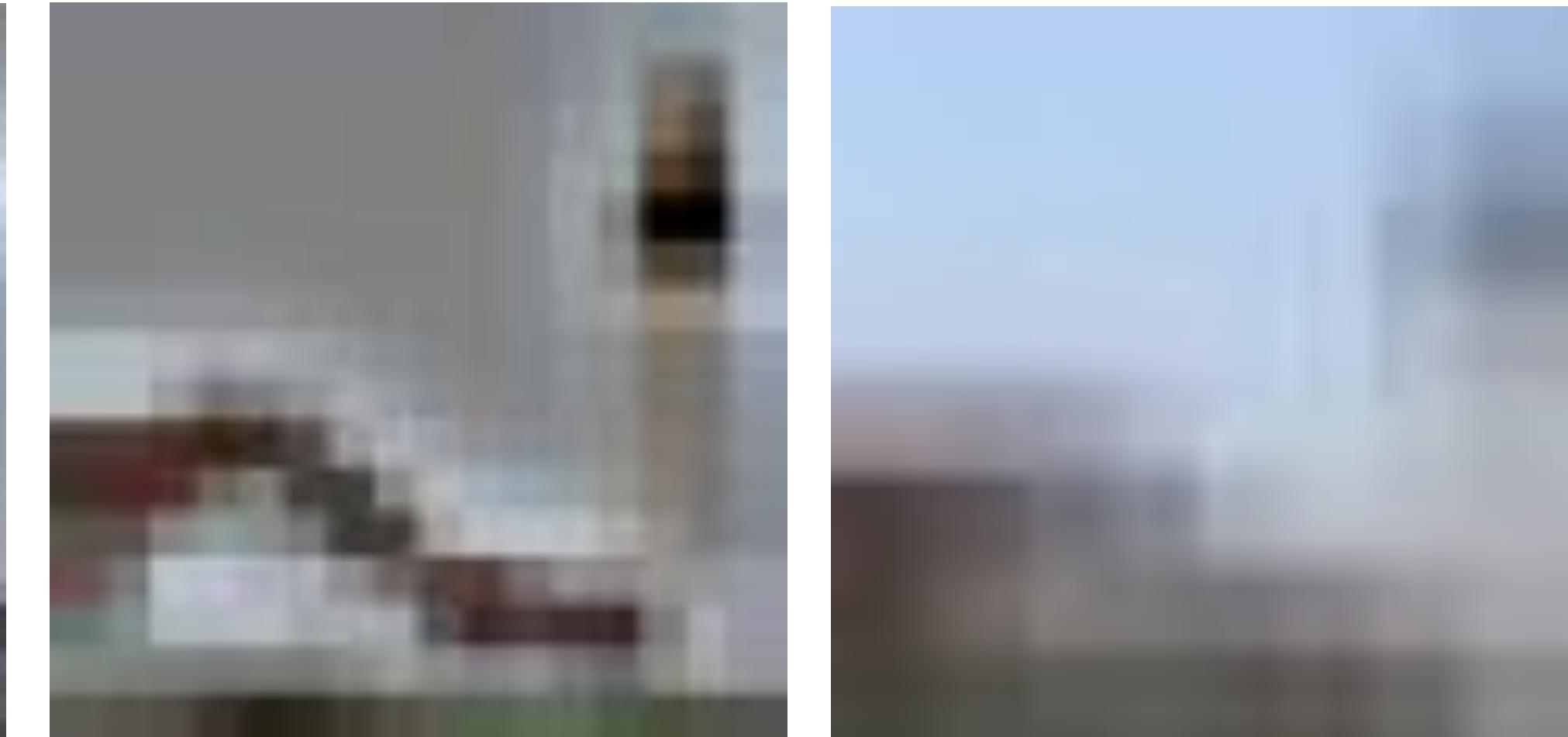


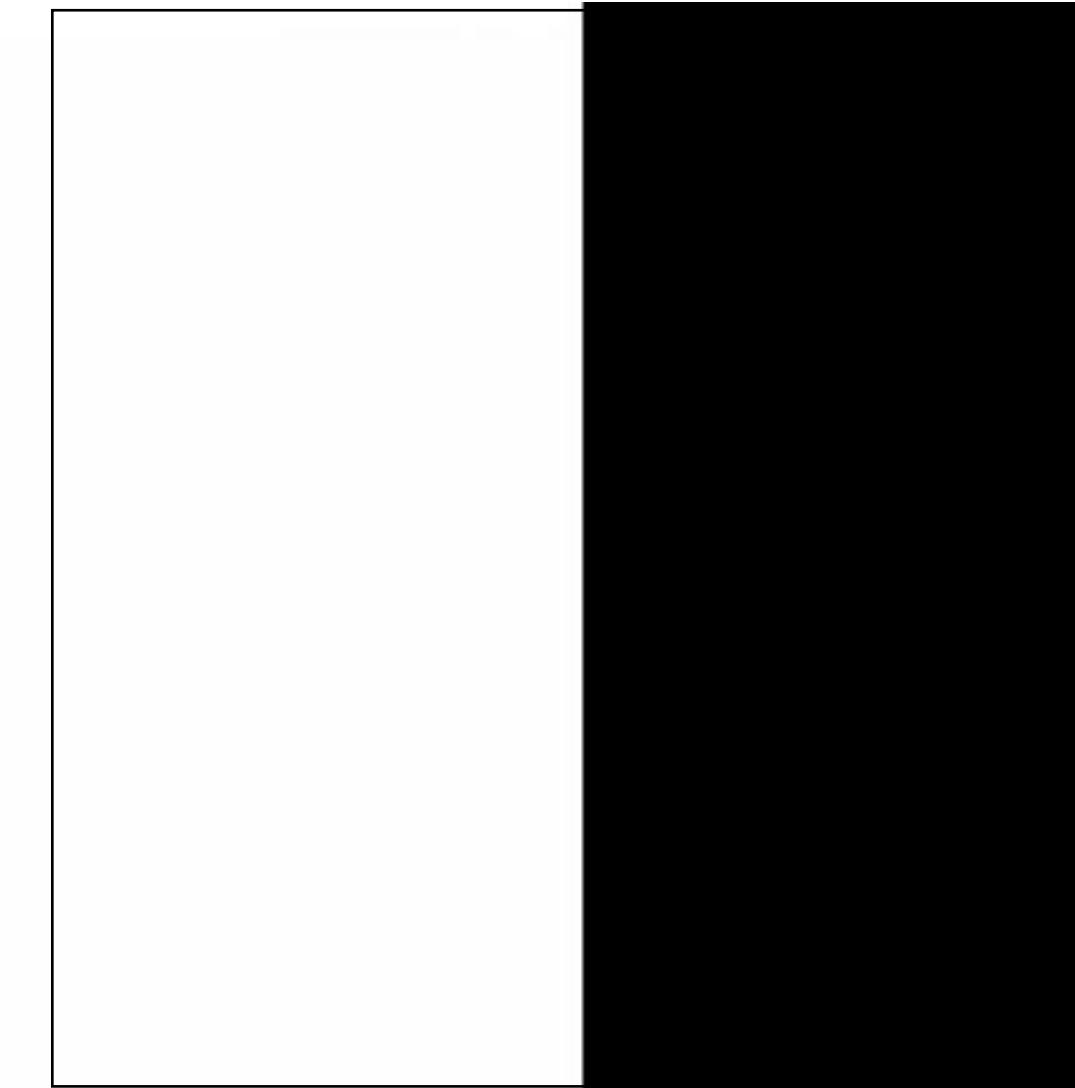
Image Blending



Image Blending

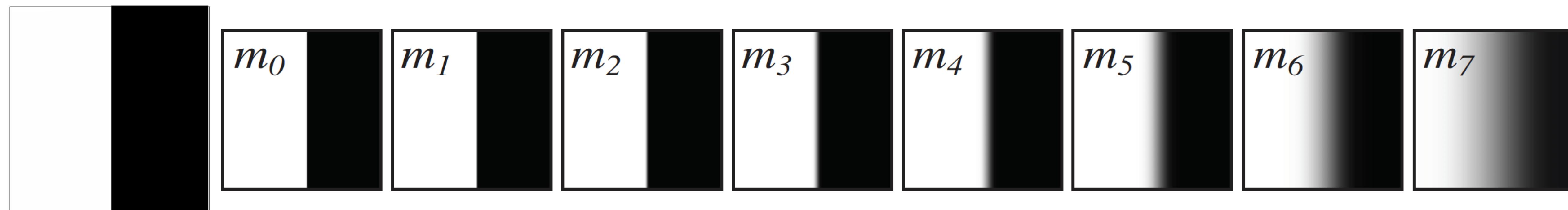
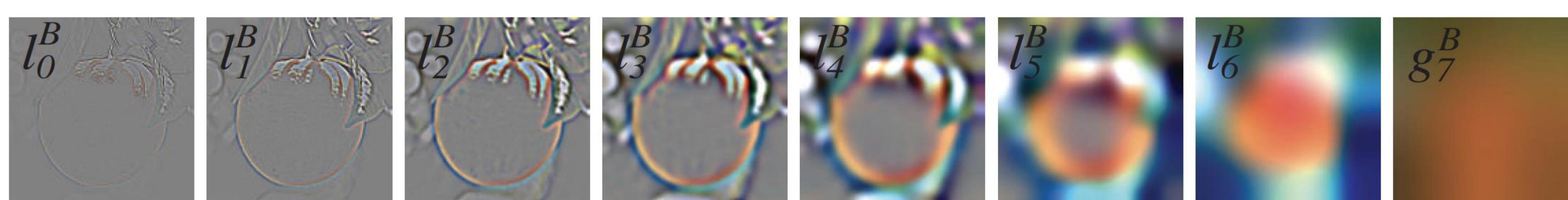
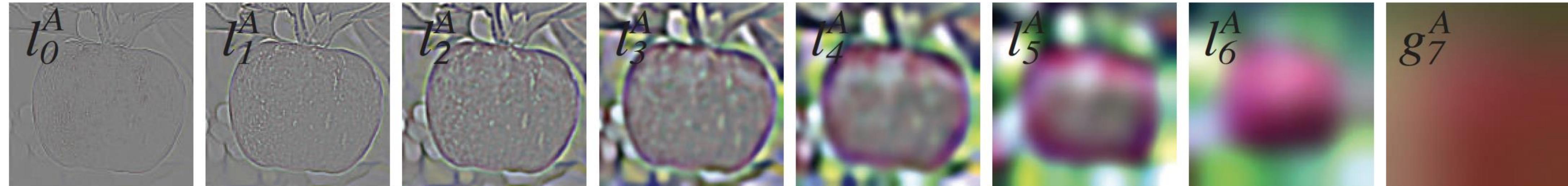


Image Blending

 I^A I^B m I

$$I = m * I^A + (1 - m) * I^B$$

Image Blending with the Laplacian Pyramid



$$l_k = l_k^A * m_k + l_i^B * (1 - m_k)$$

Image Blending with the Laplacian Pyramid



Image Blending with the Laplacian Pyramid

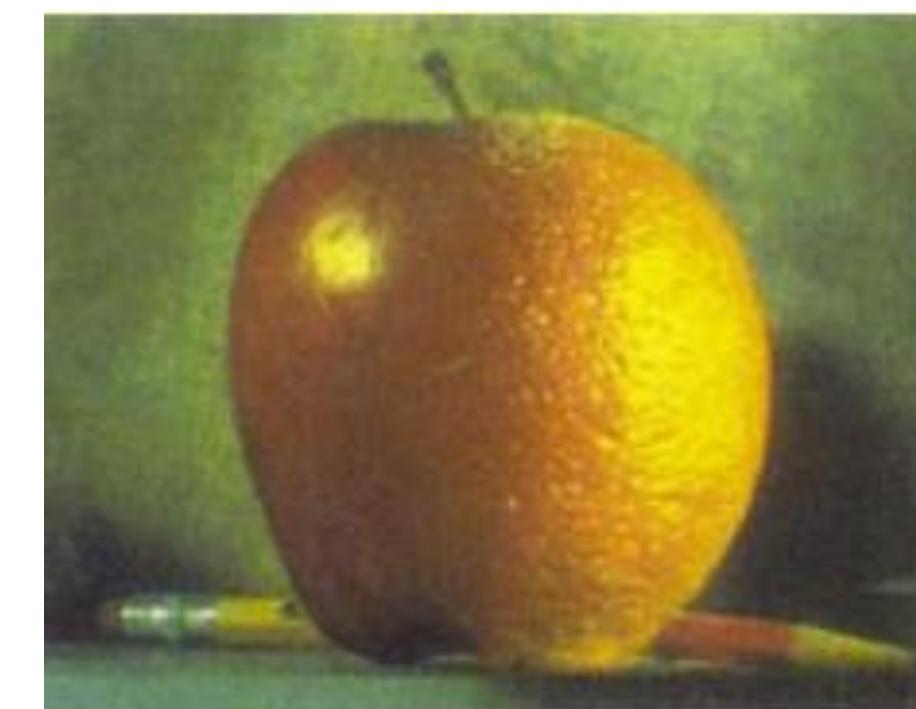
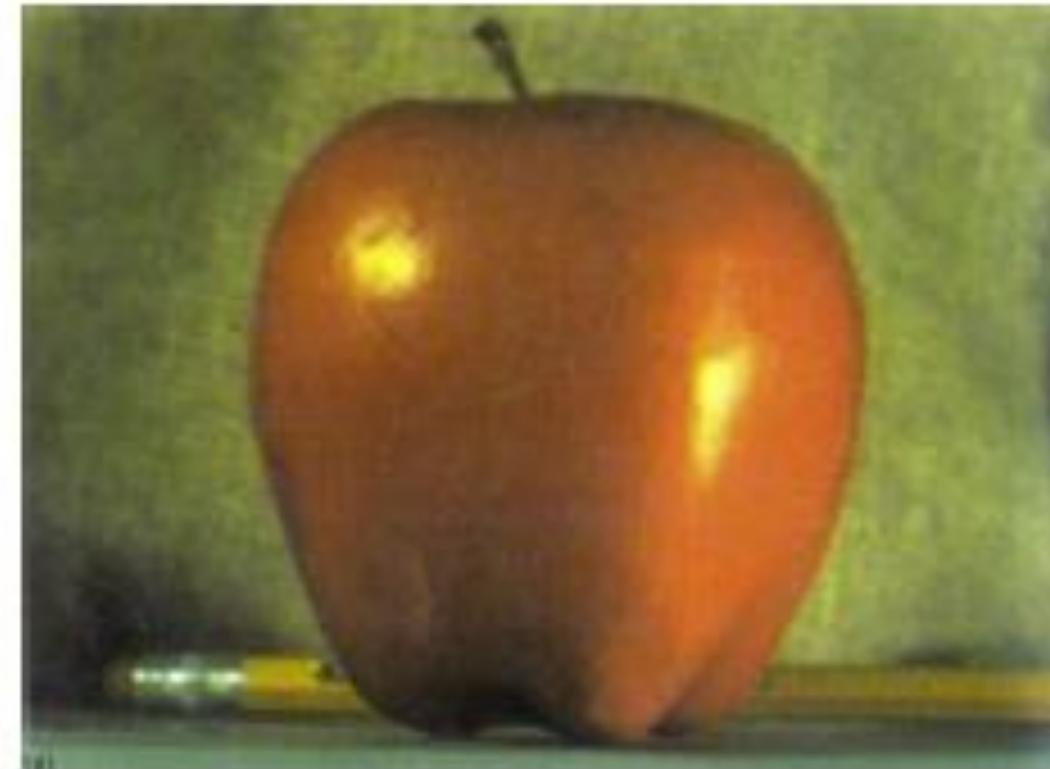
532

IEEE TRANSACTIONS ON COMMUNICATIONS, VOL. COM-31, NO. 4, APRIL 1983

The Laplacian Pyramid as a Compact Image Code

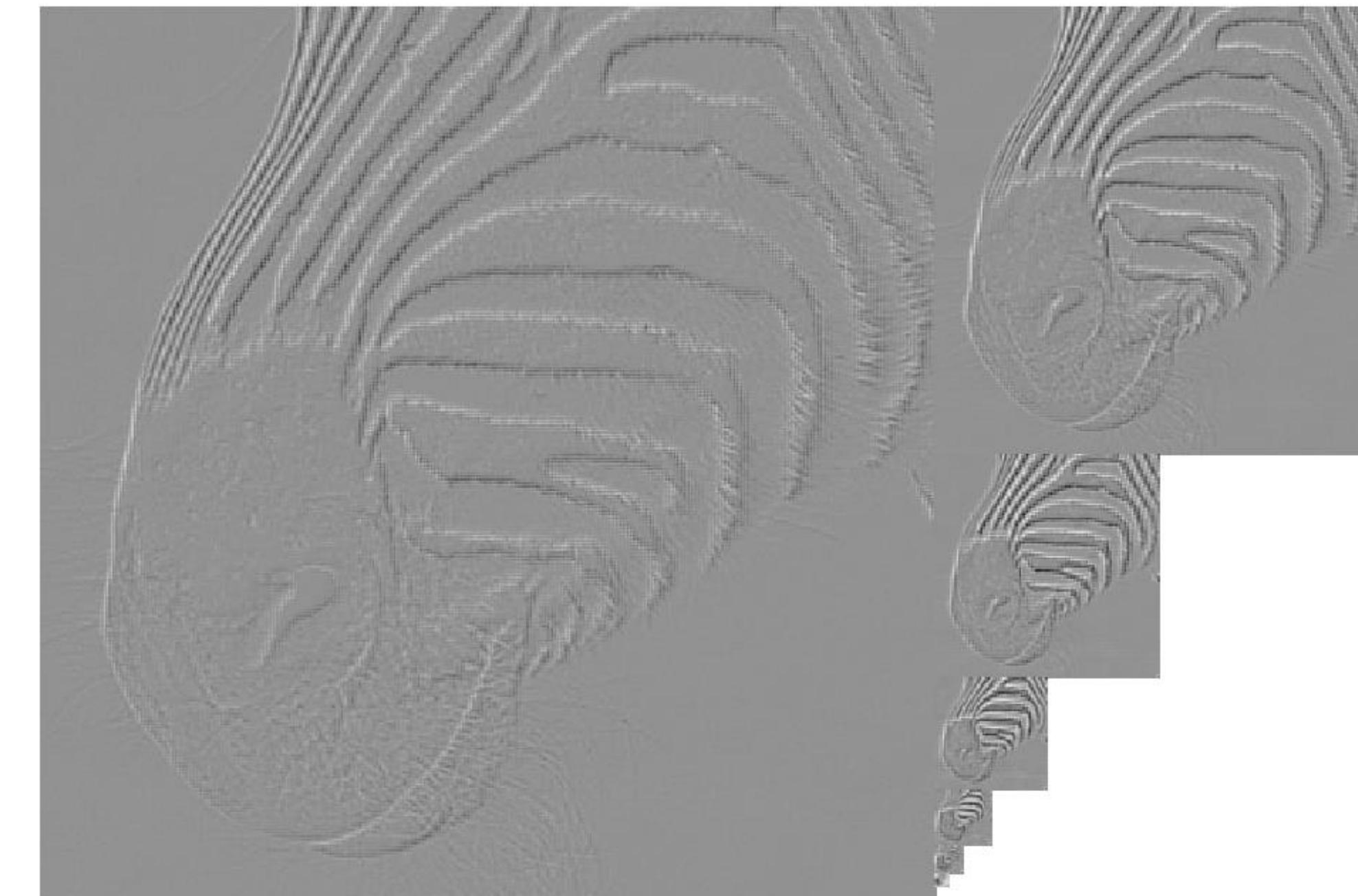
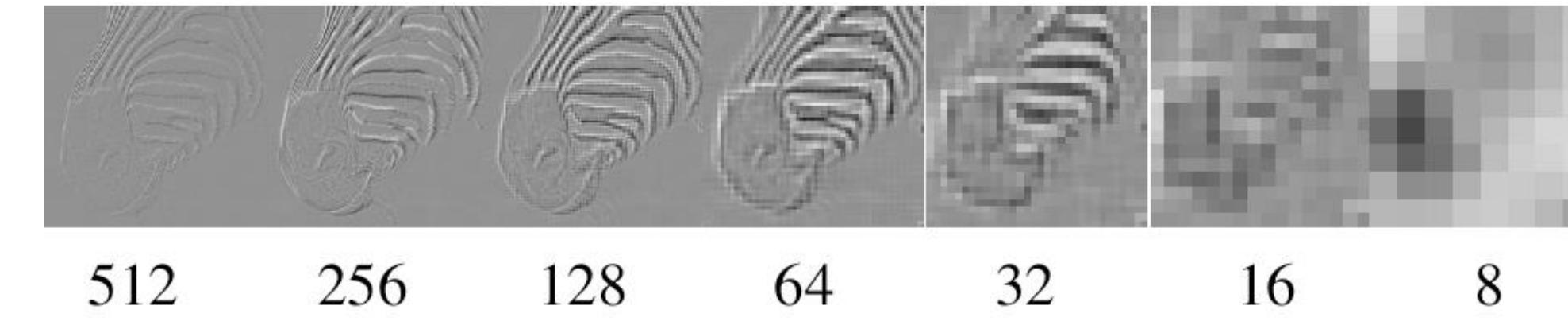
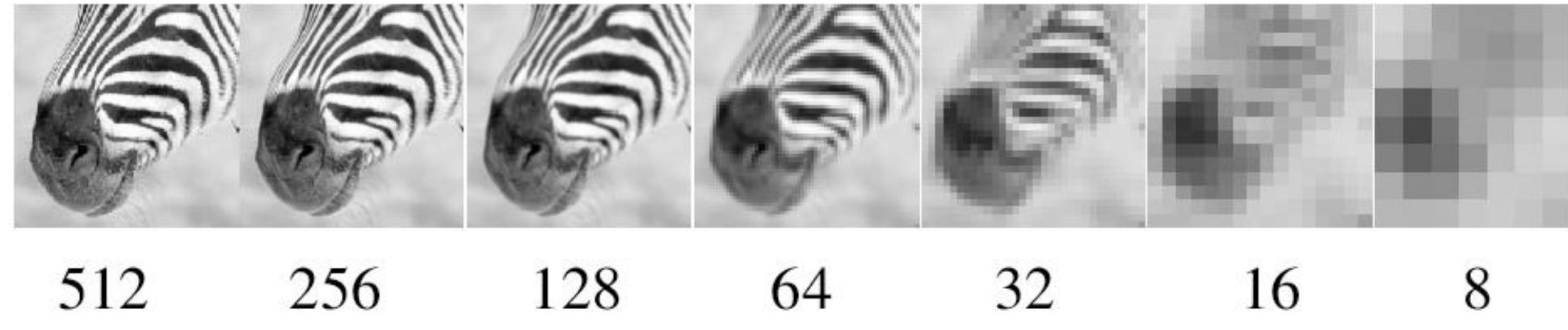
PETER J. BURT, MEMBER, IEEE, AND EDWARD H. ADELSON

- Build Laplacian pyramid for both images: LA, LB
- Build Gaussian pyramid for mask: G
- Build a combined Laplacian pyramid:
- Collapse L to obtain the blended image



http://persci.mit.edu/pub_pdfs/pyramid83.pdf

Image pyramids

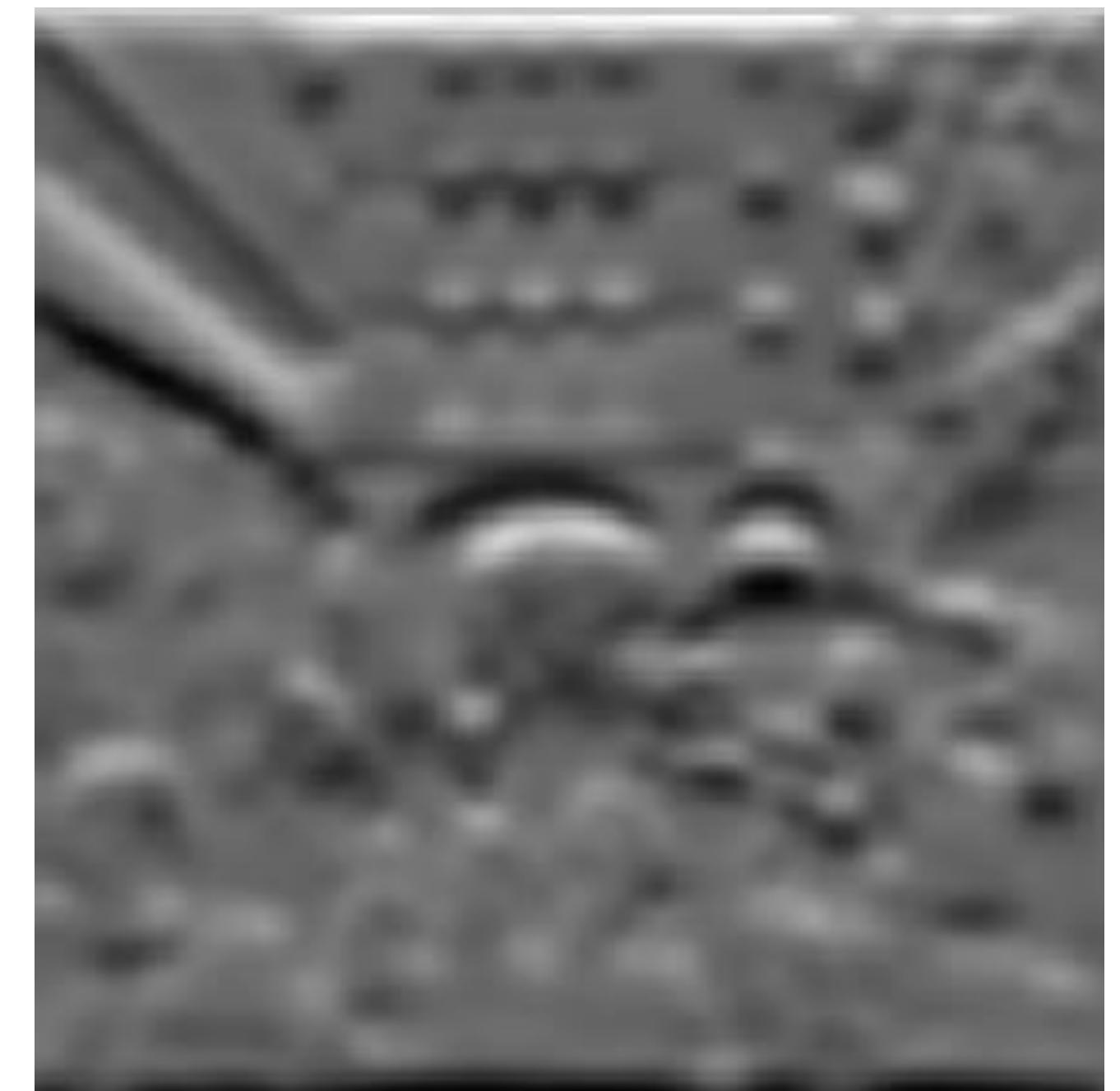
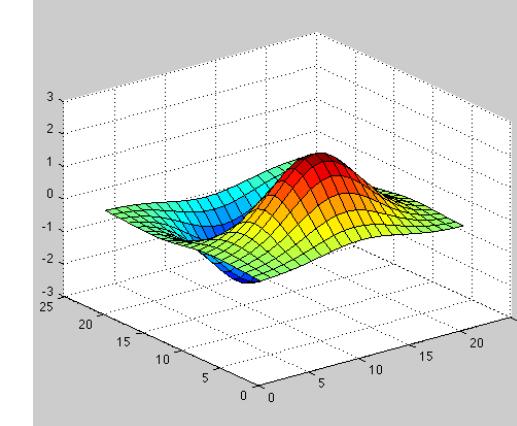
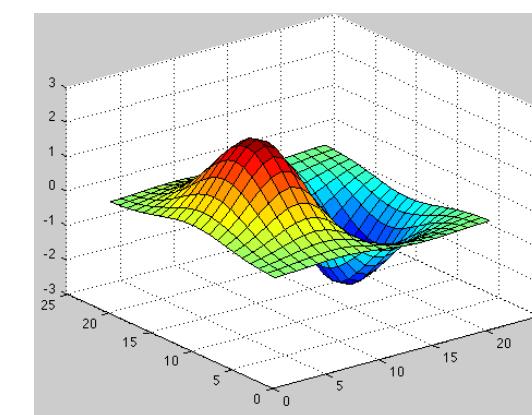


Gaussian Pyr

Laplacian Pyr

And many more: QMF, steerable, ...Convnets!

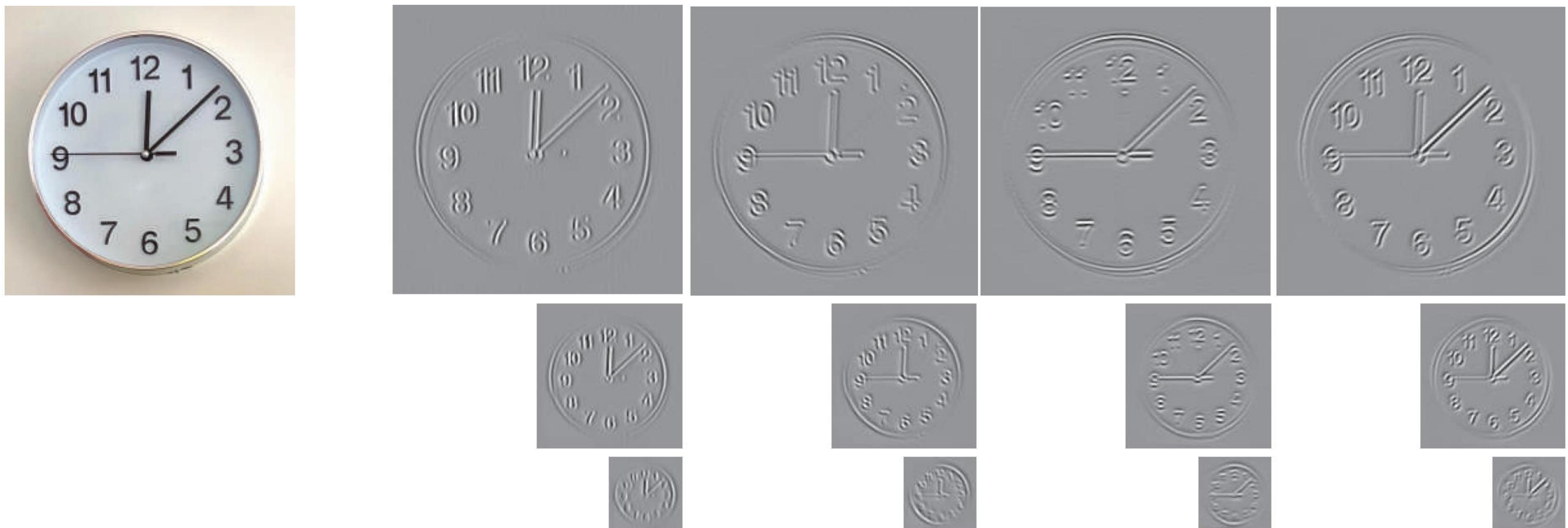
Orientations



Orientations

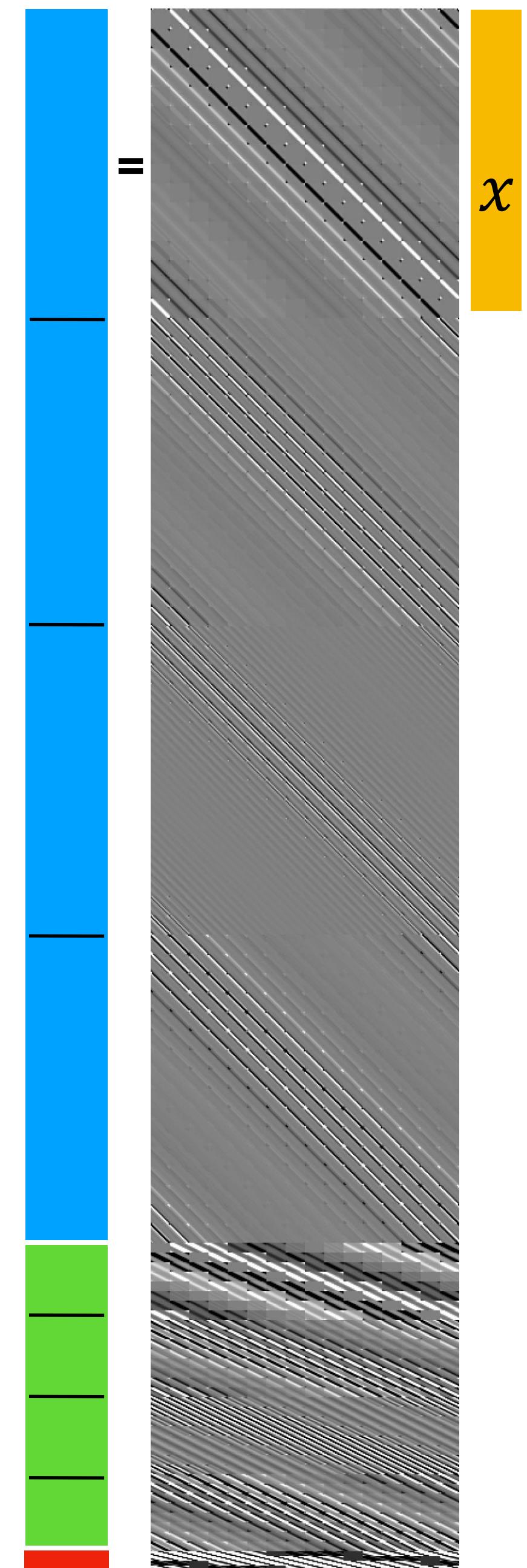
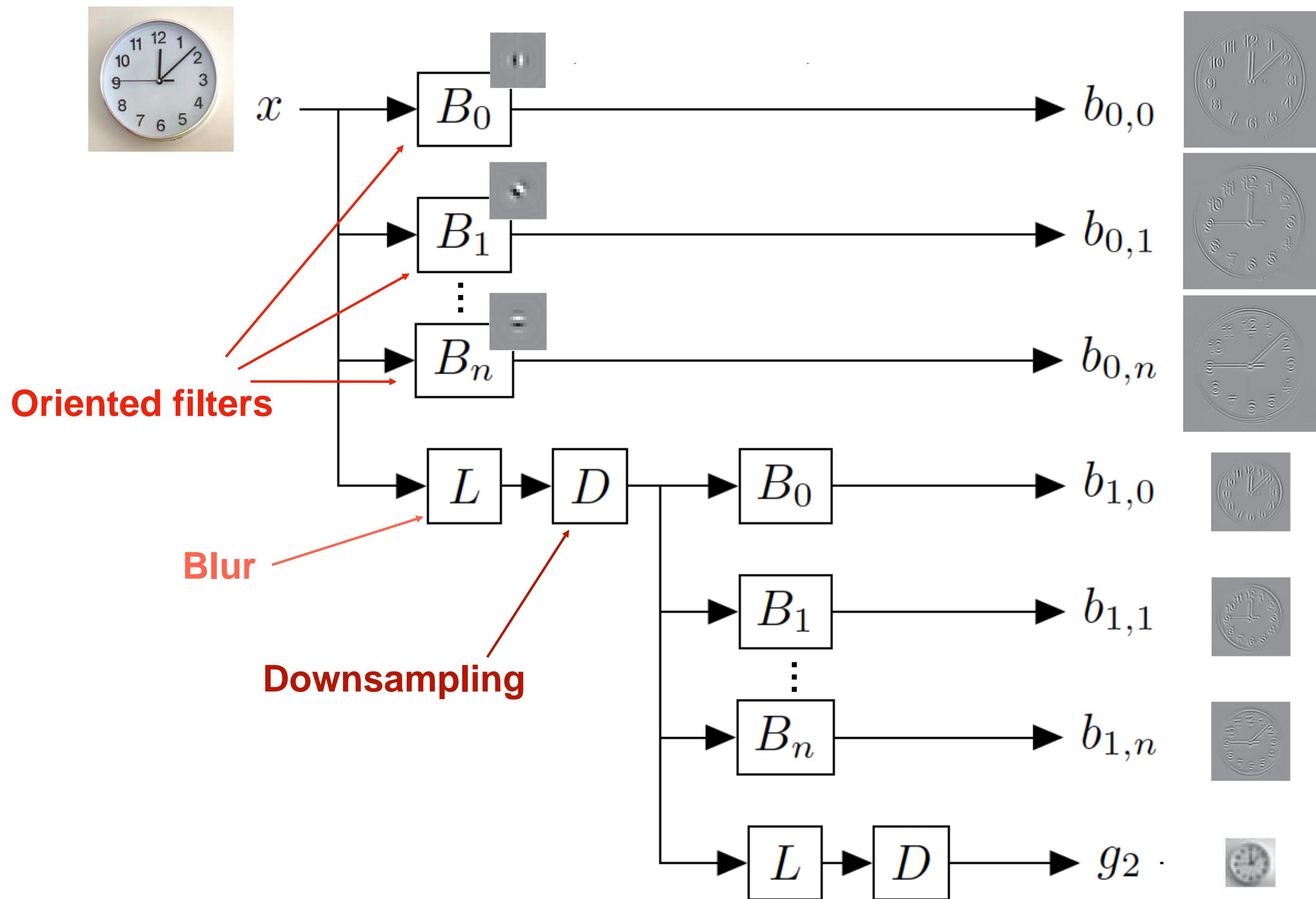


Steerable Pyramid

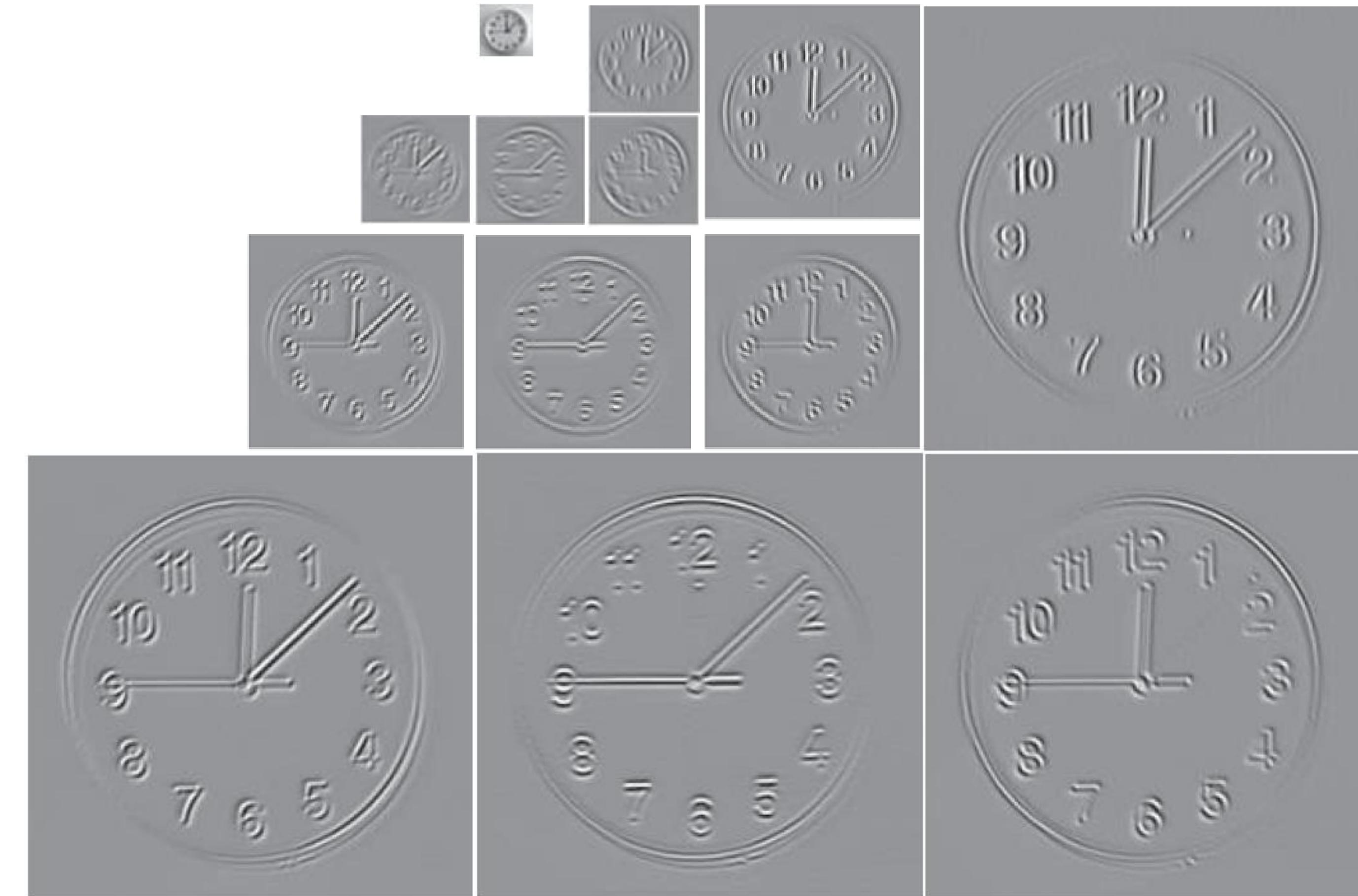
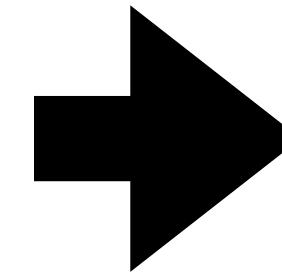


Low pass residual

Steerable Pyramid

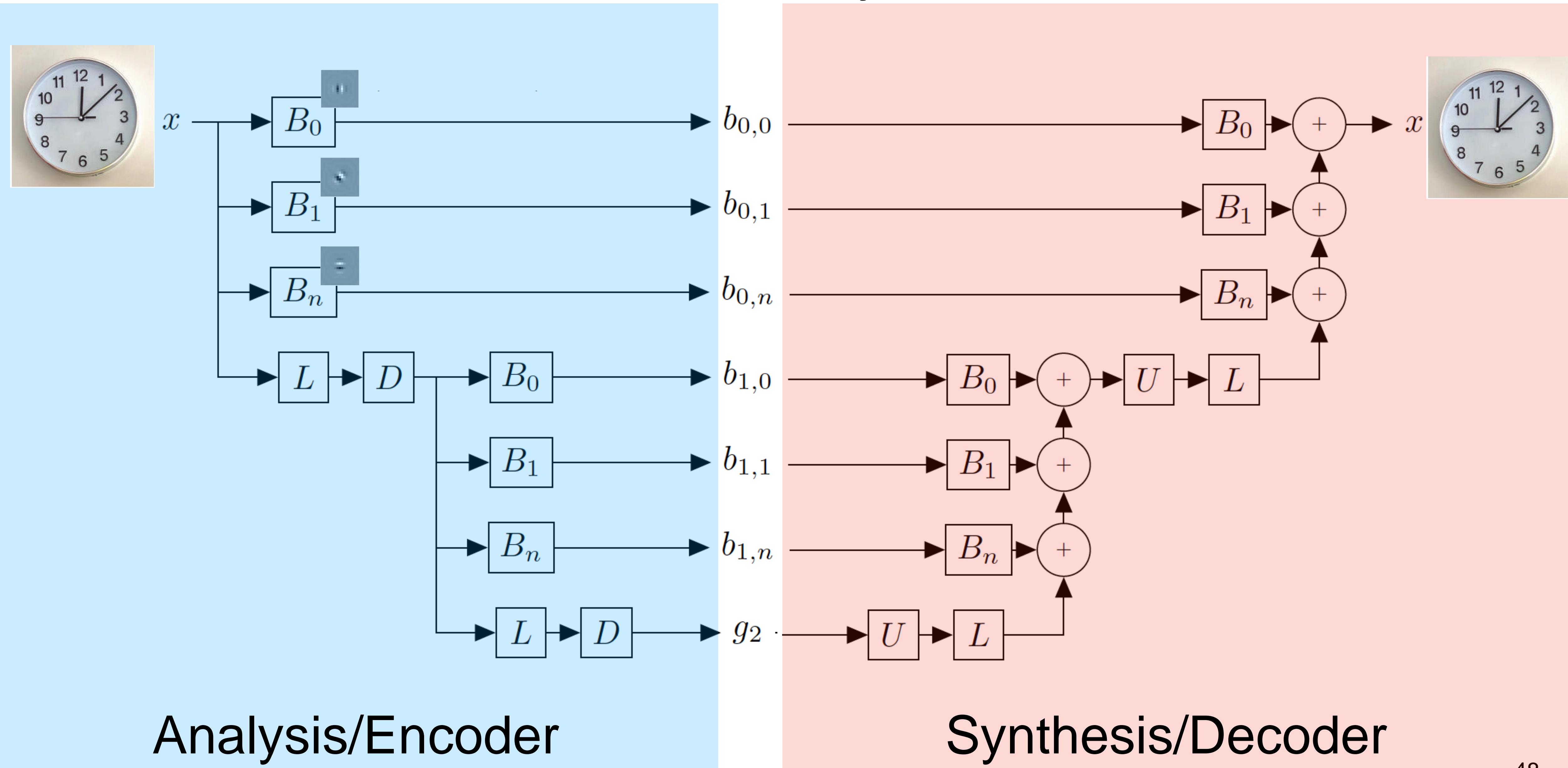


Steerable Pyramid



3 scales, 4 orientations

Steerable Pyramid



Analysis/Encoder

Synthesis/Decoder

What about learning the filters?

Wavelet-like receptive fields emerge from a network that learns sparse codes for natural images.

Bruno A. Olshausen¹ and David J. Field

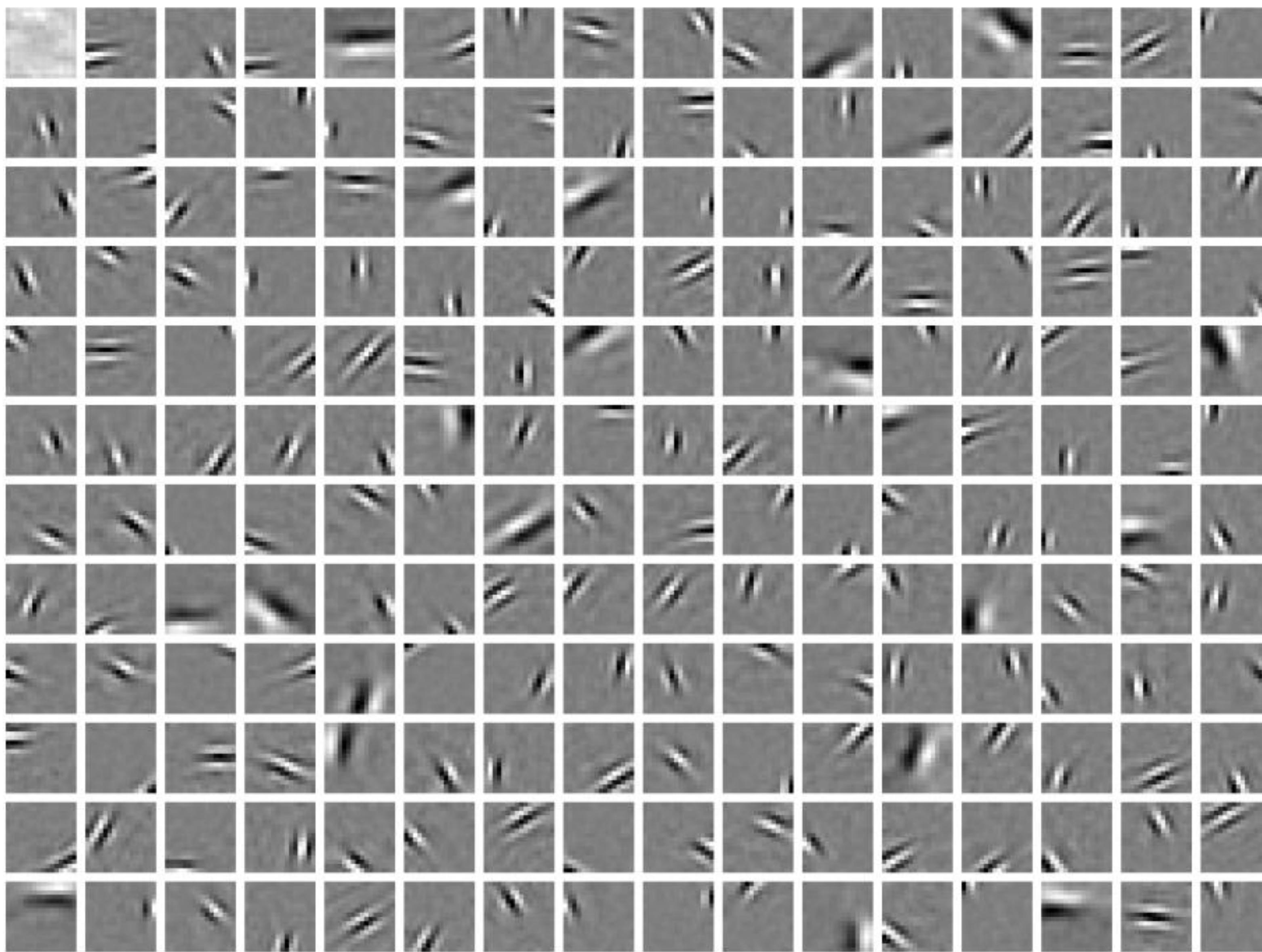
$$E = -[\text{preserve information}] - \lambda[\text{sparseness of } a_i], \quad (2)$$

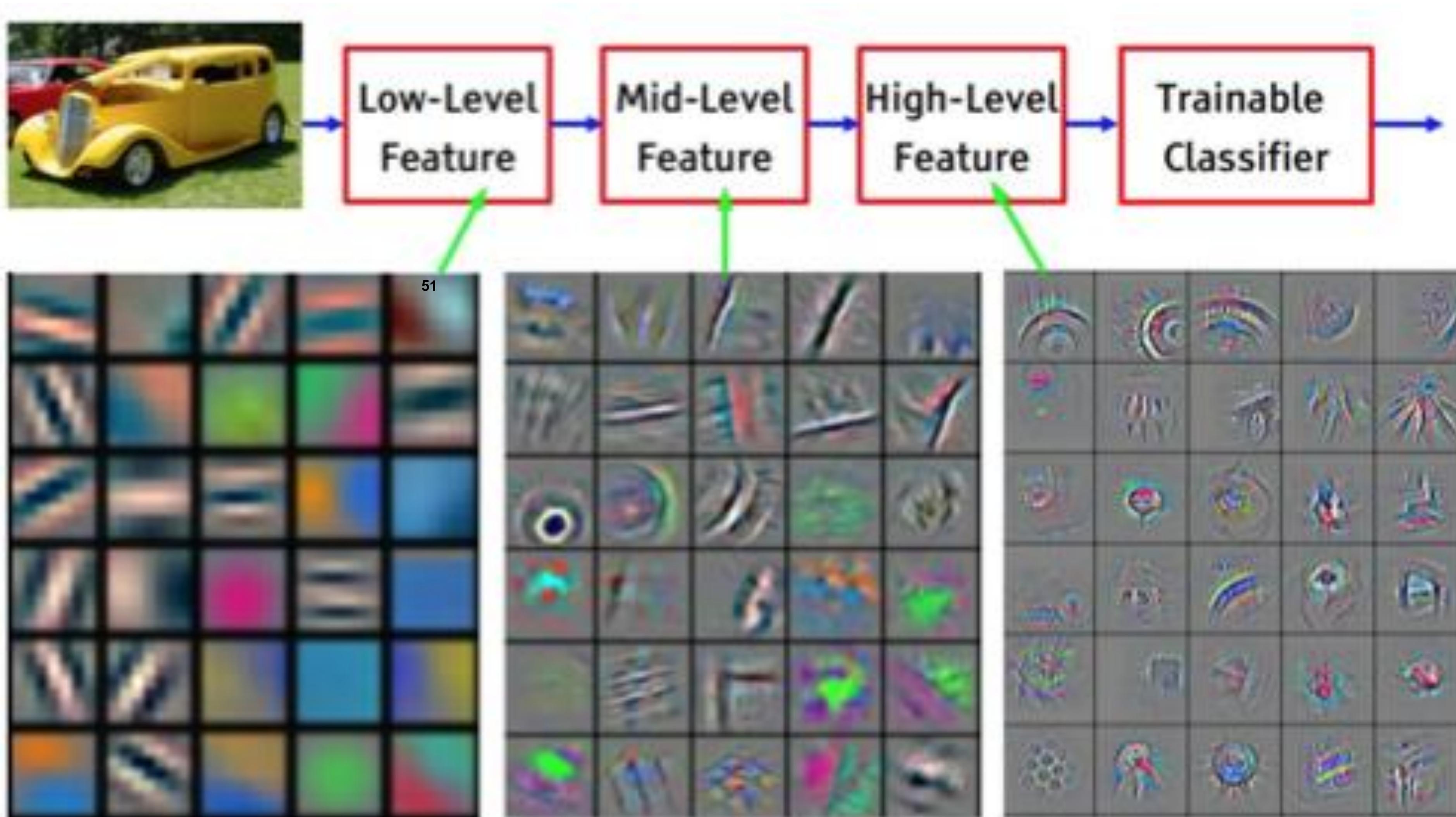
where λ is a positive constant that determines the importance of the second term relative to the first. The first term measures how well the code describes the image, and we choose this to be the mean square of the error between the actual image and the reconstructed image:

$$[\text{preserve information}] = -\sum_{x,y} \left[I(x,y) - \sum_i a_i \phi_i(x,y) \right]^2. \quad (3)$$

Learned filters

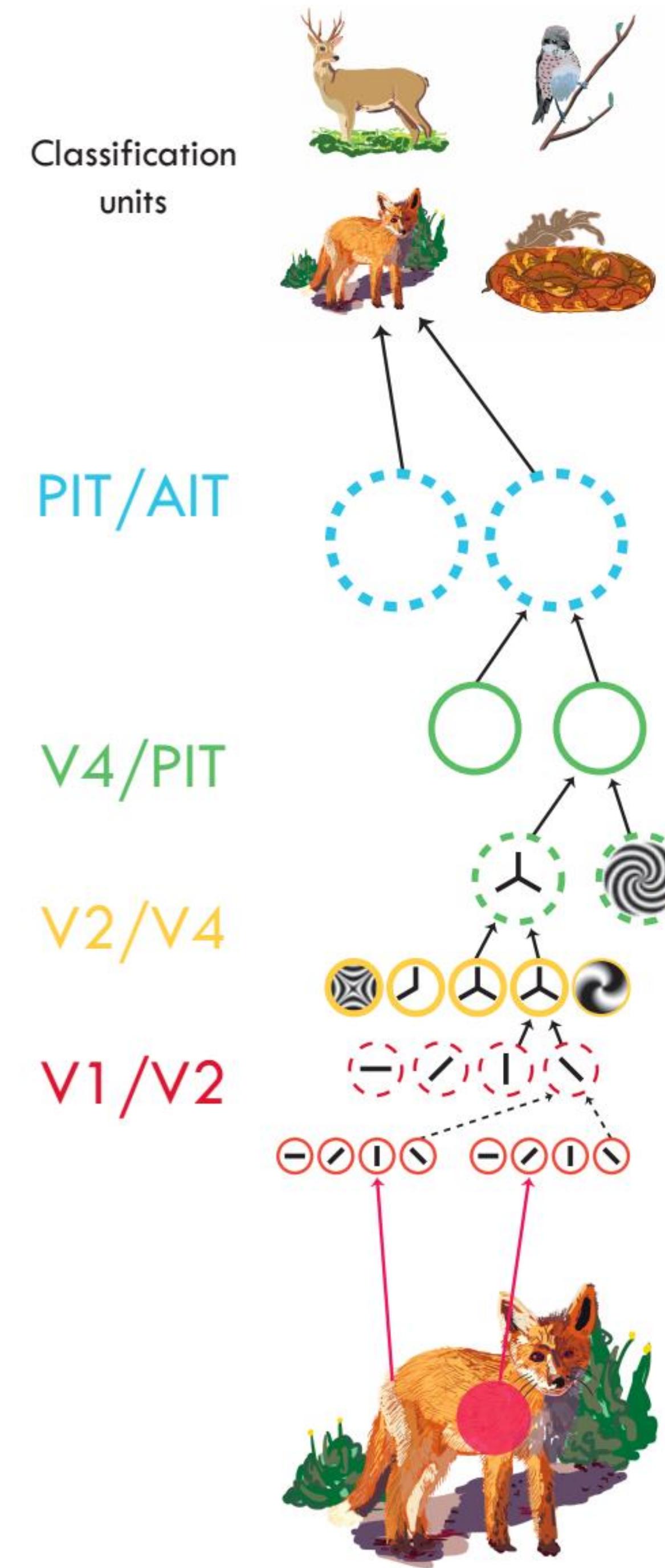
a.





Feature visualization of convolutional net trained on ImageNet from [Zeiler & Fergus 2013]

[slide courtesy Yann LeCun]



Serre, 2014

Image classification

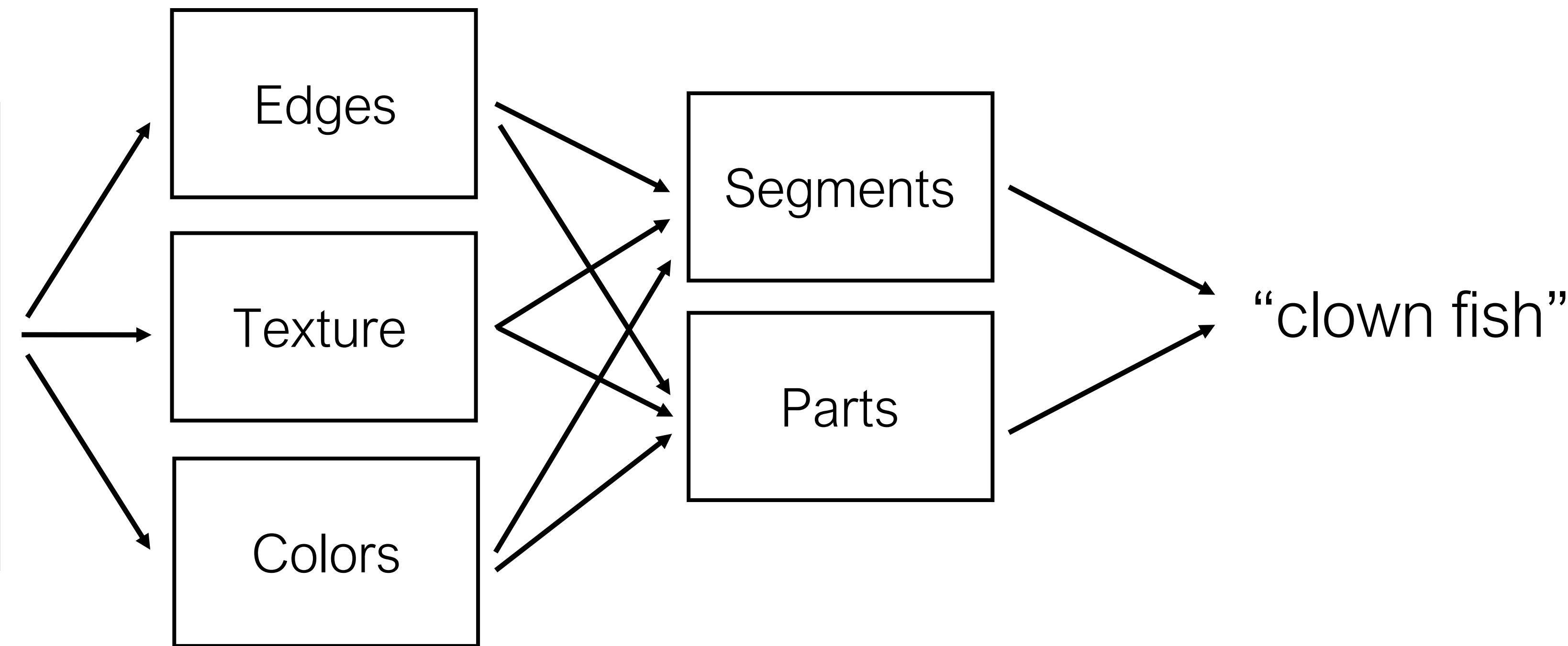


Image classification

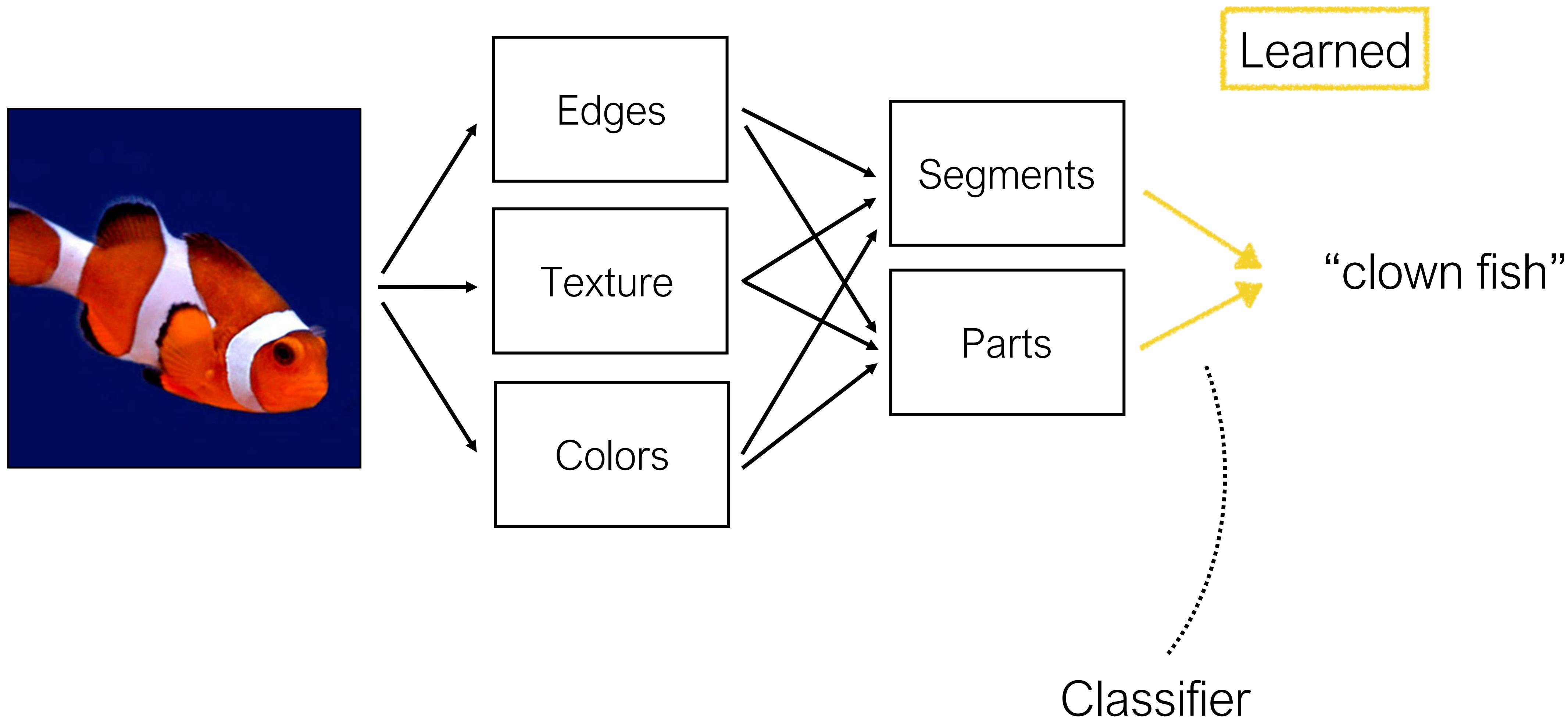


Image classification

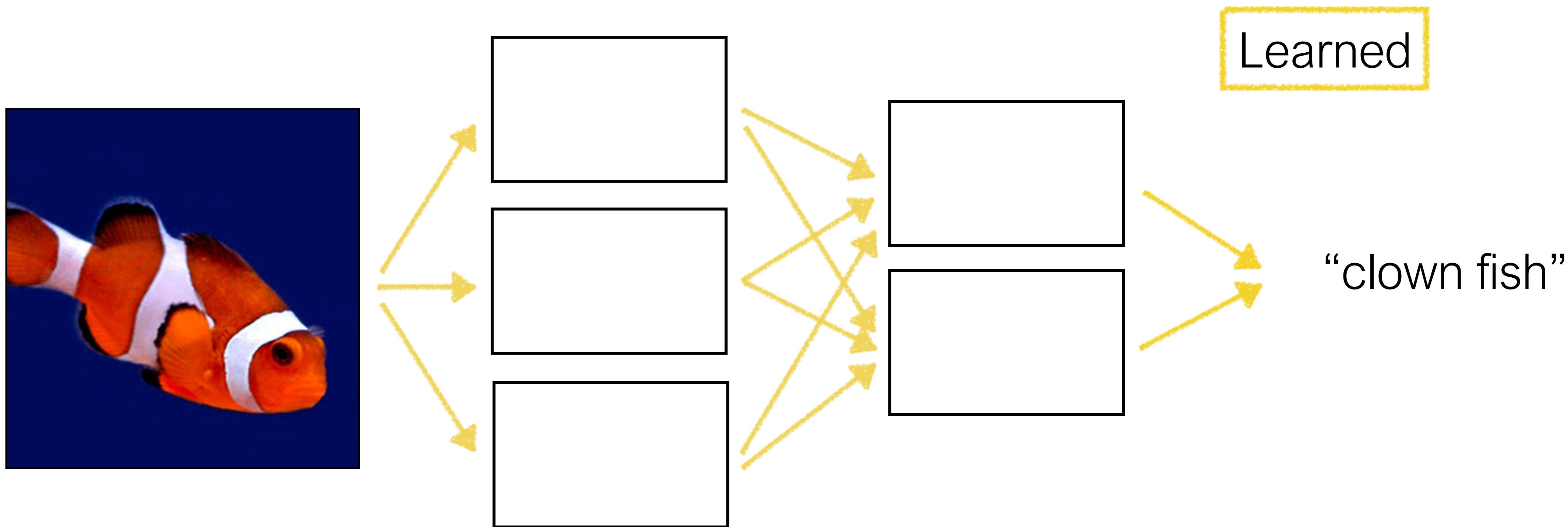
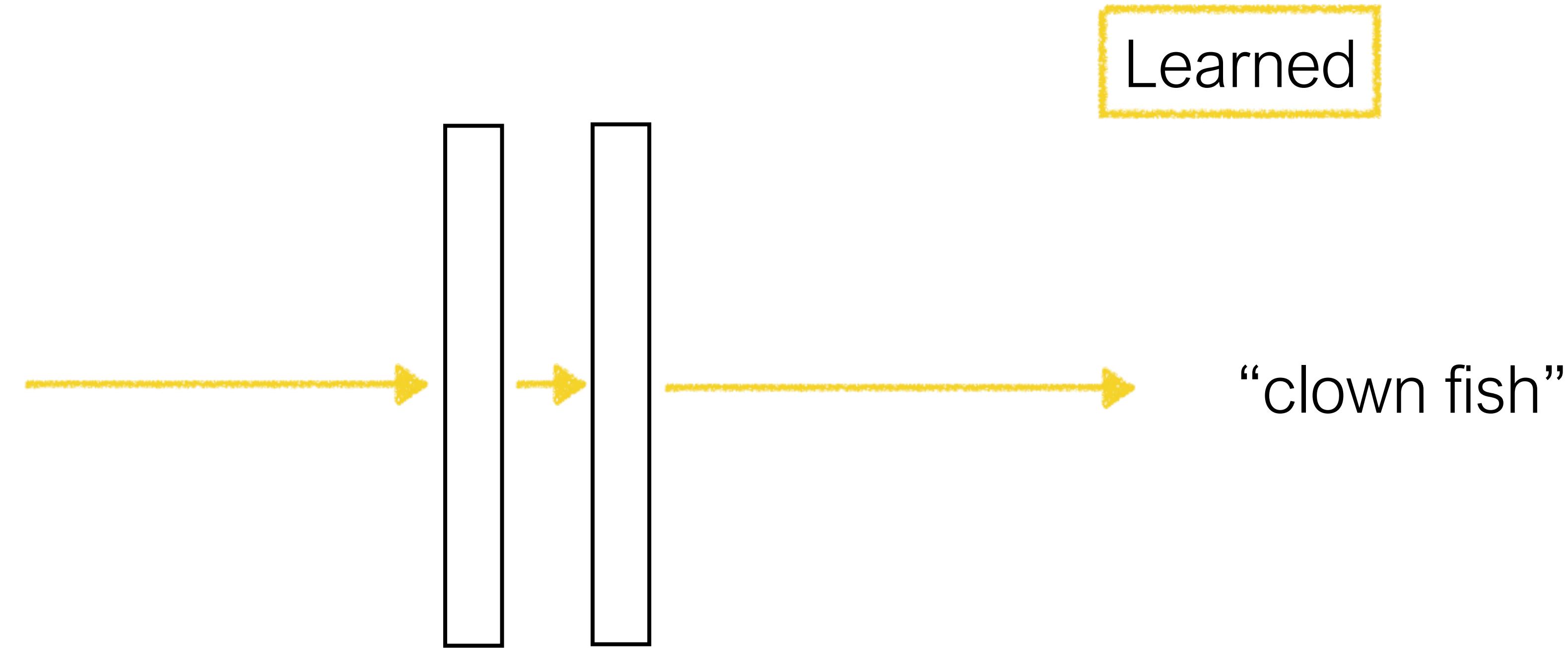
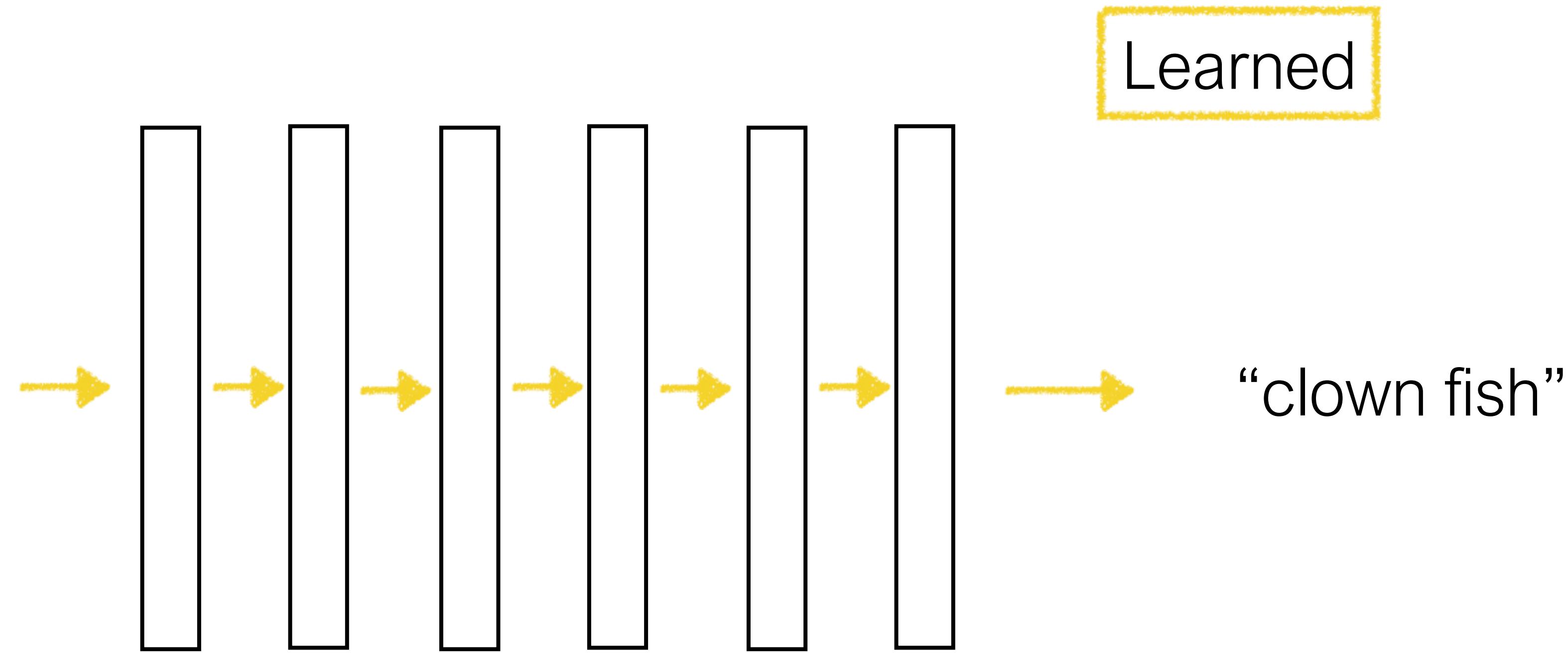


Image classification



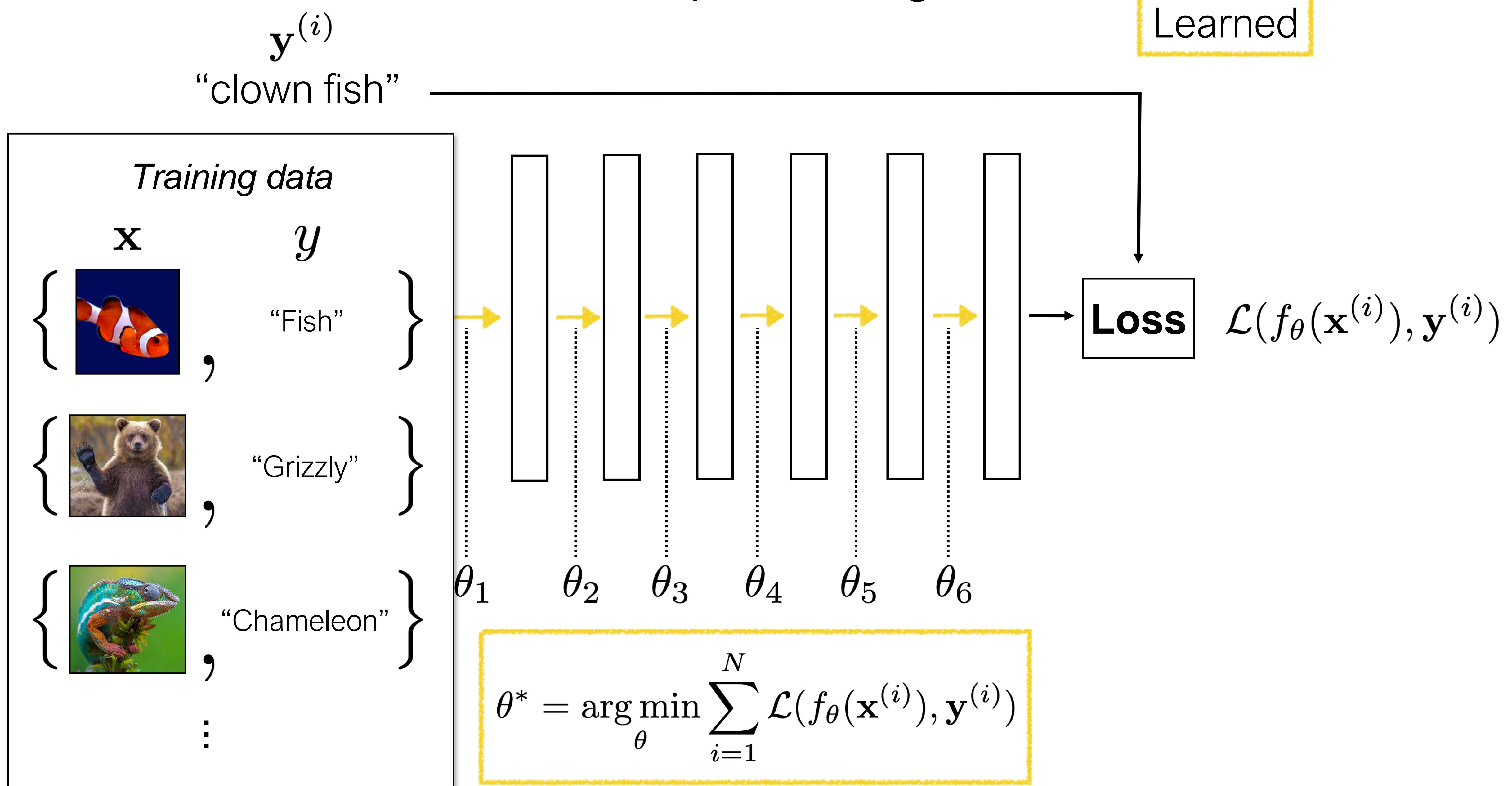
Neural net

Image classification

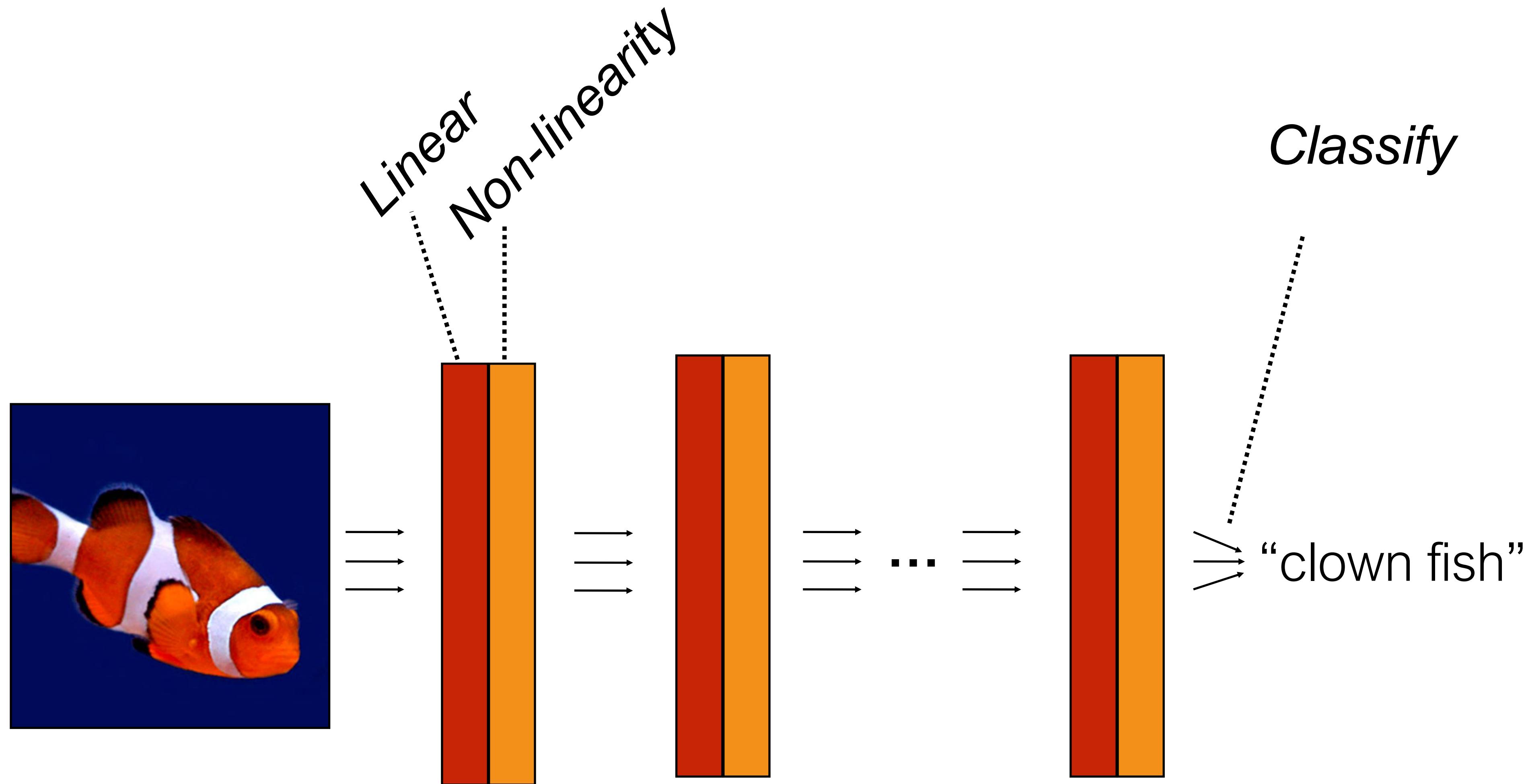


Deep neural net

Deep learning

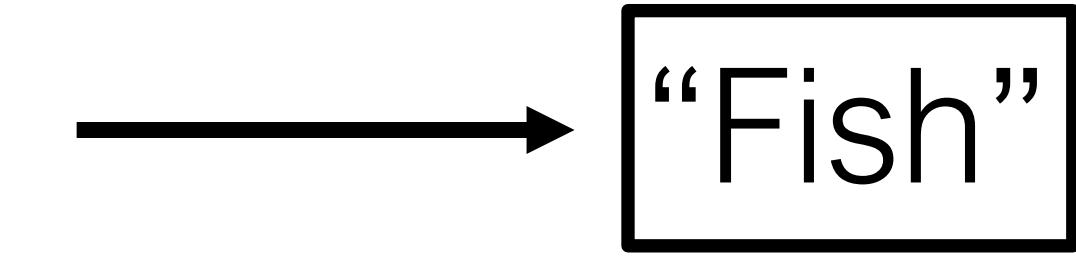
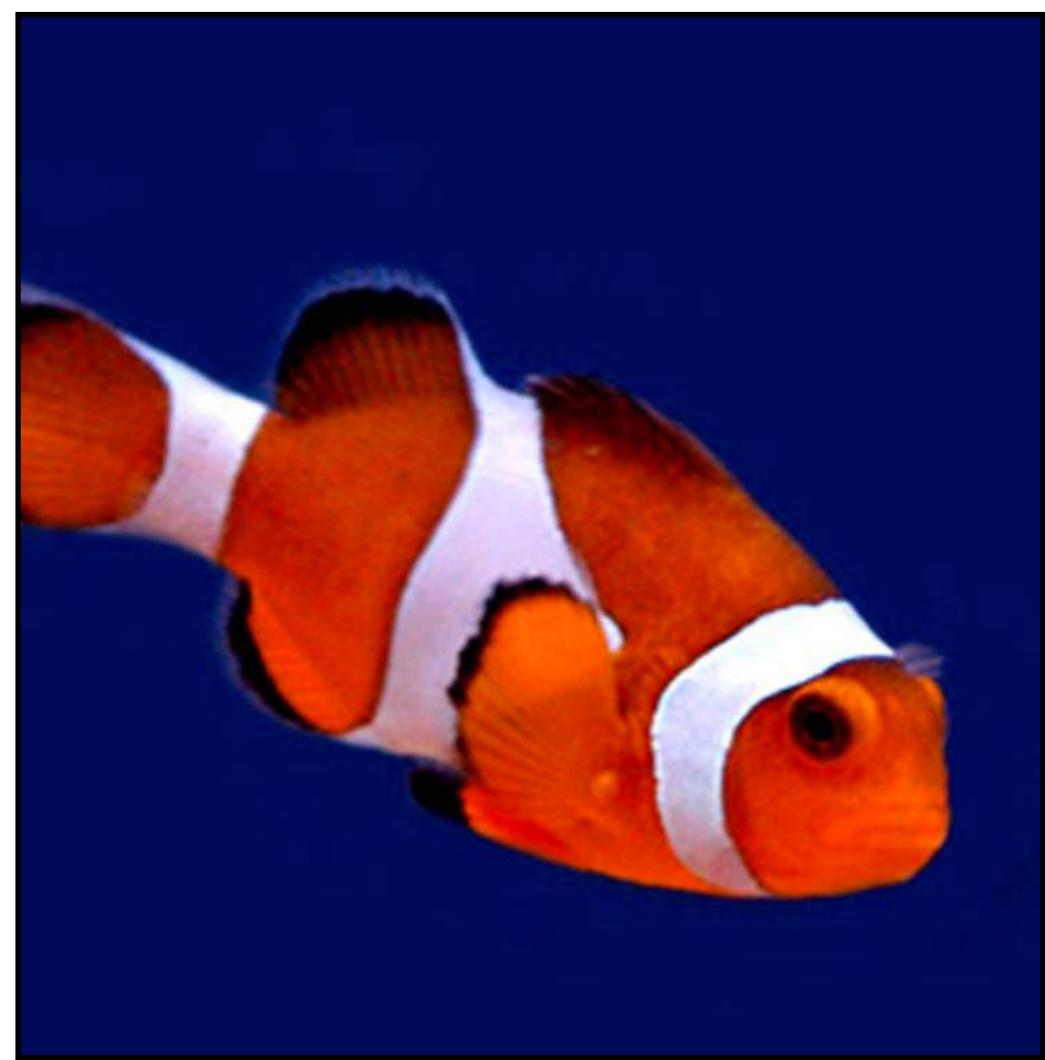


Deep nets



$$f(\mathbf{x}) = f_L(f_{L-1}(\dots f_2(f_1(\mathbf{x}))))$$

Image classification is boring!



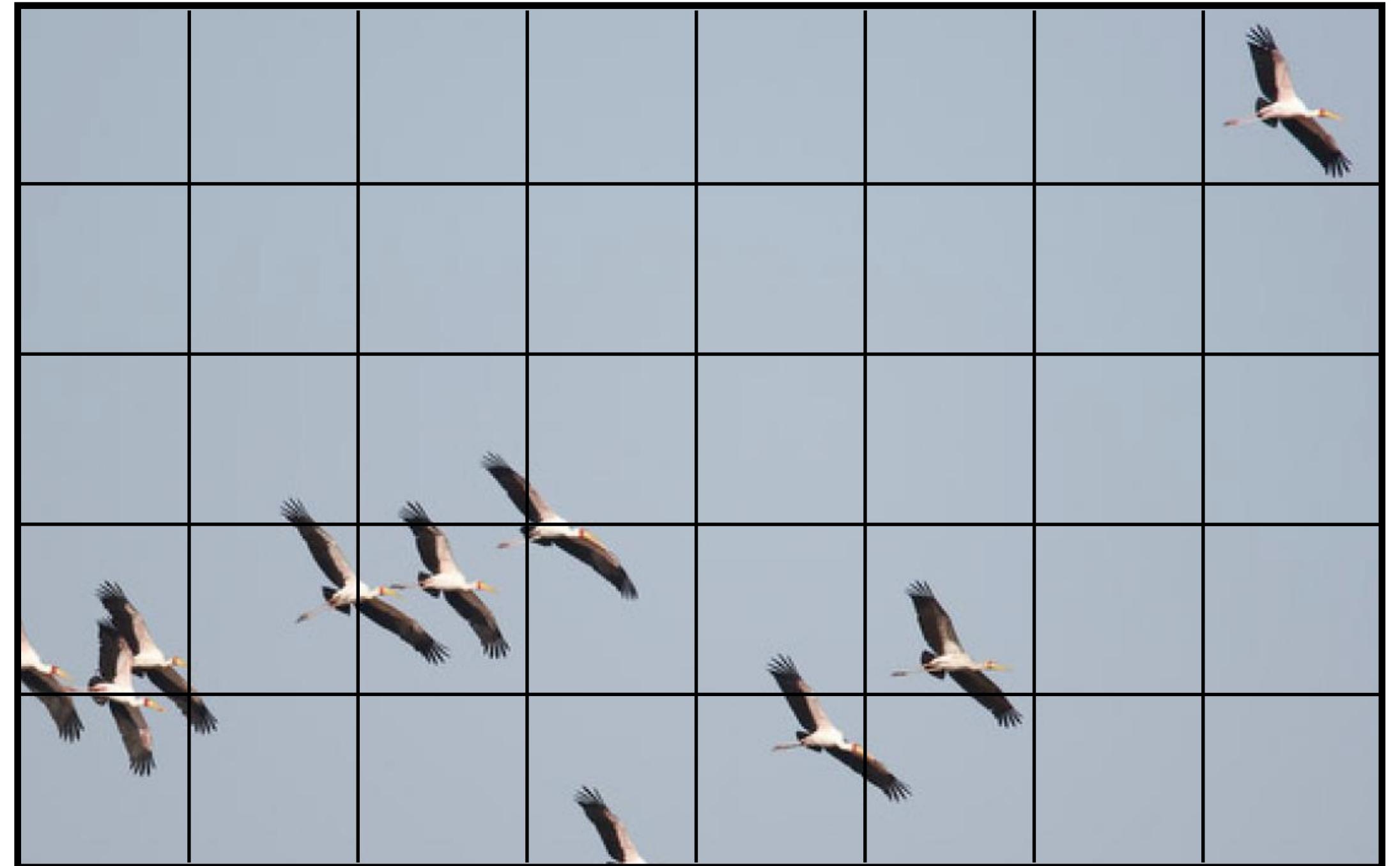
“Fish”

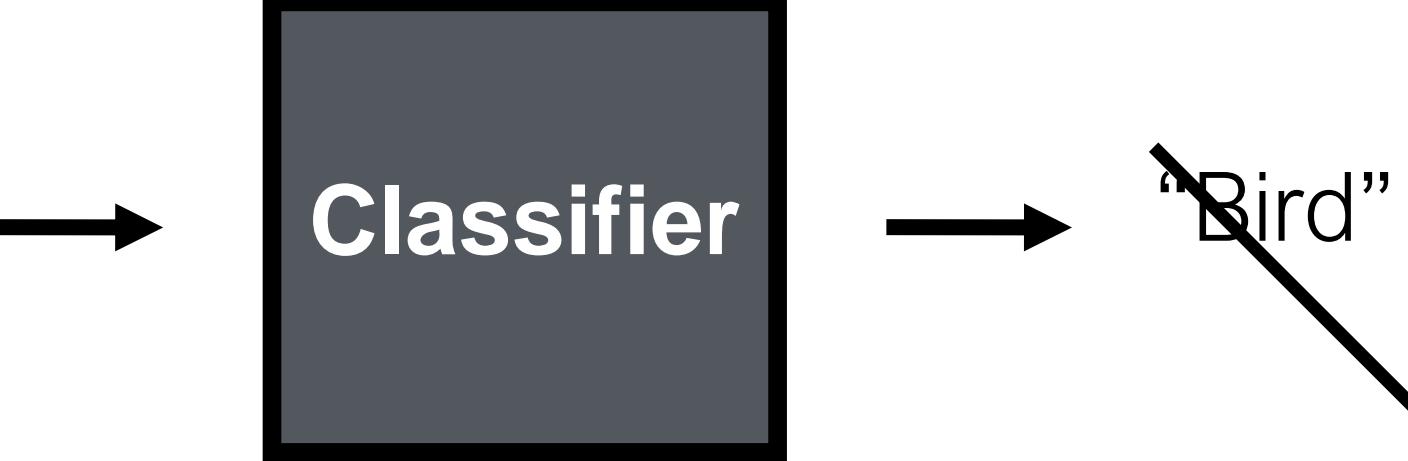
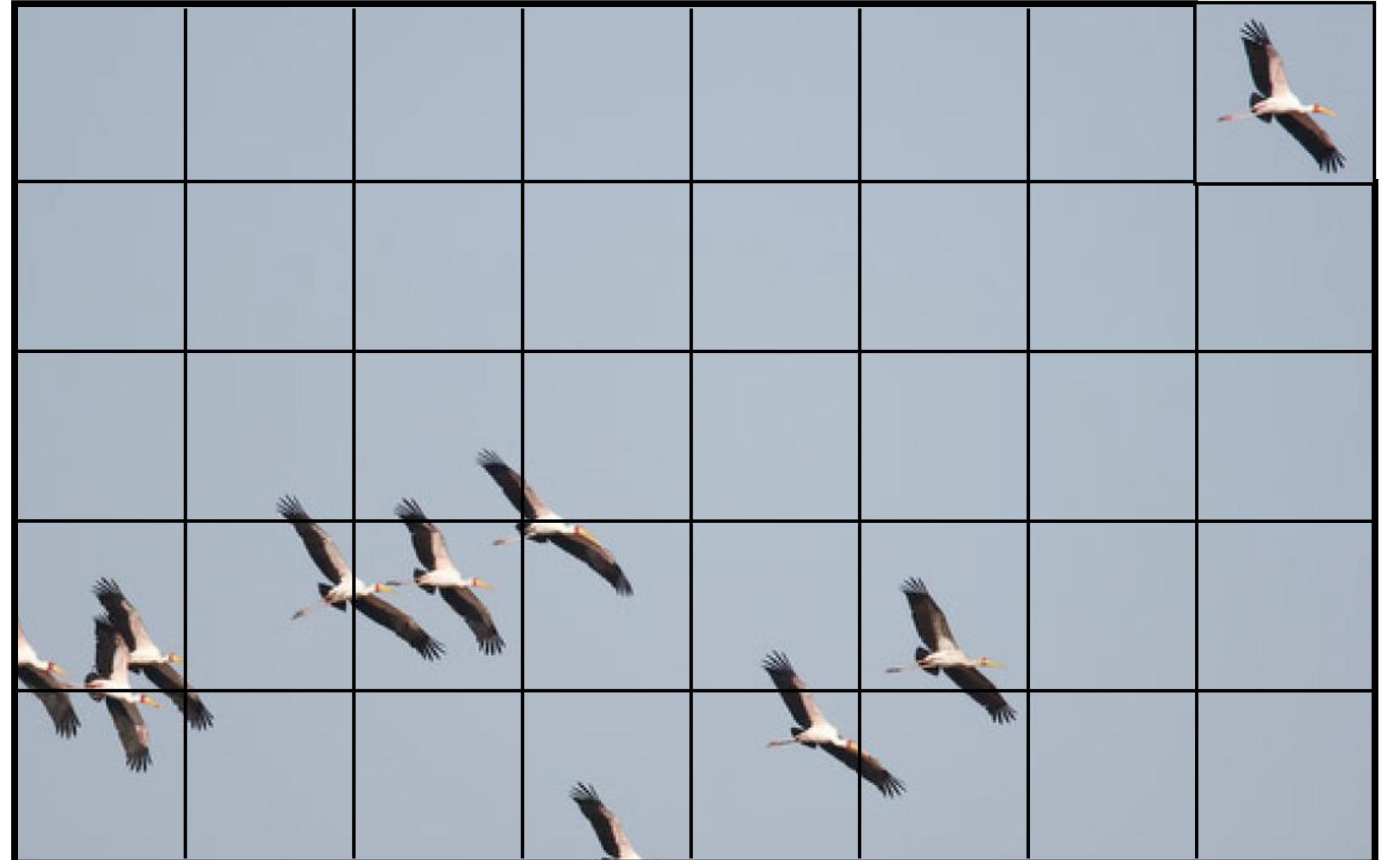
image **x**

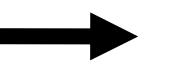
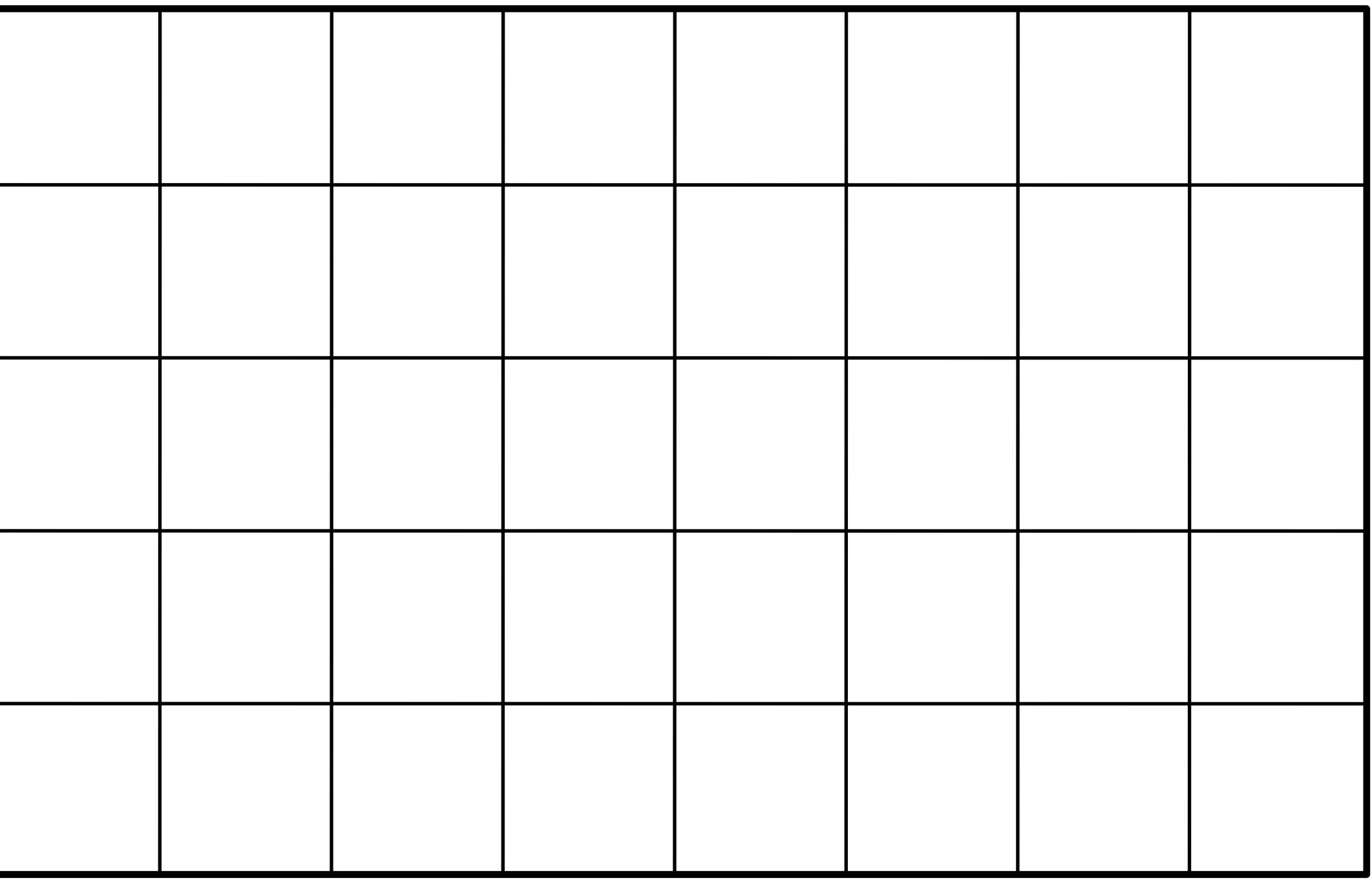
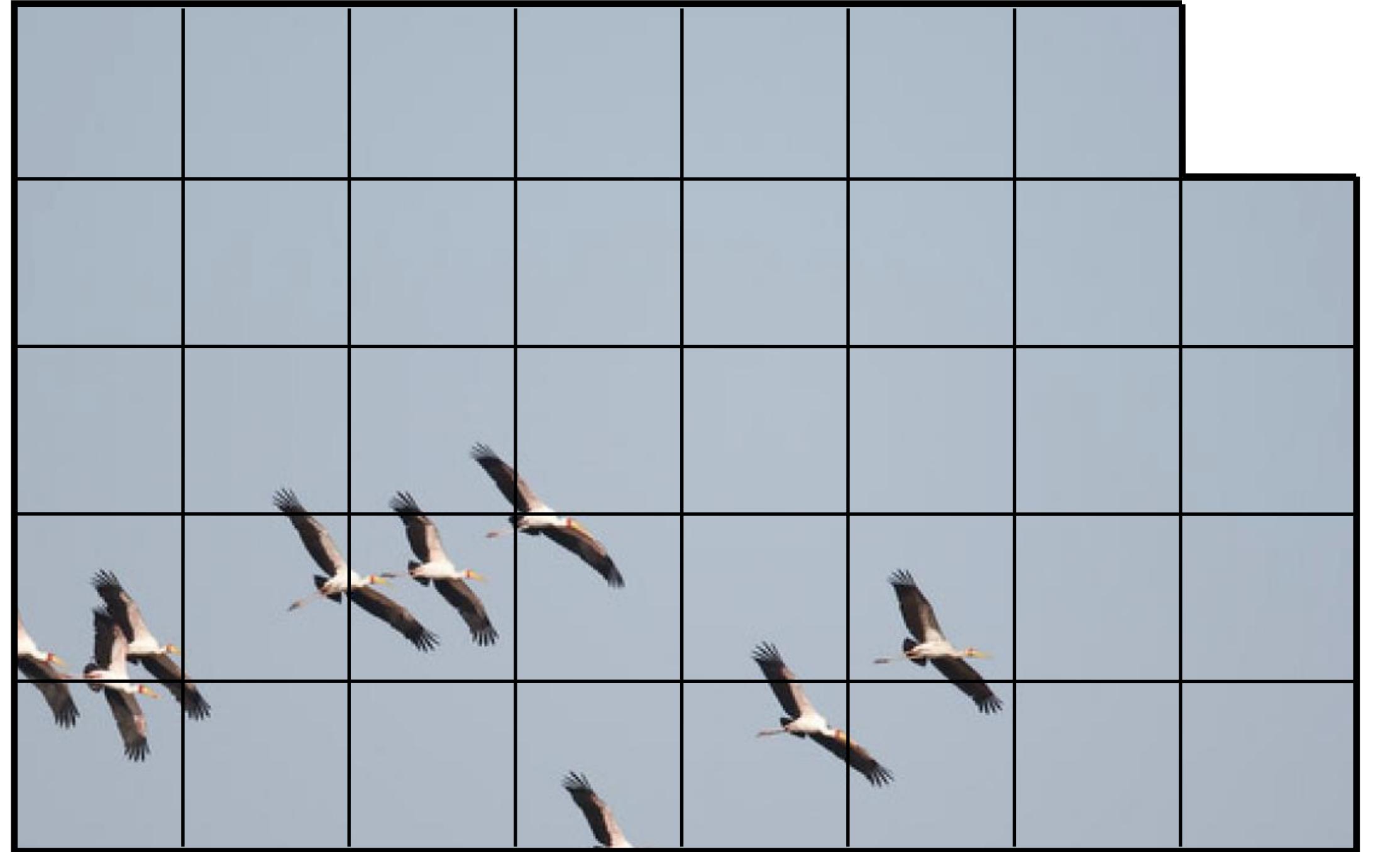
label **y**



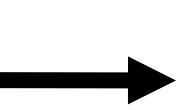
Photo credit: Fredo Durand



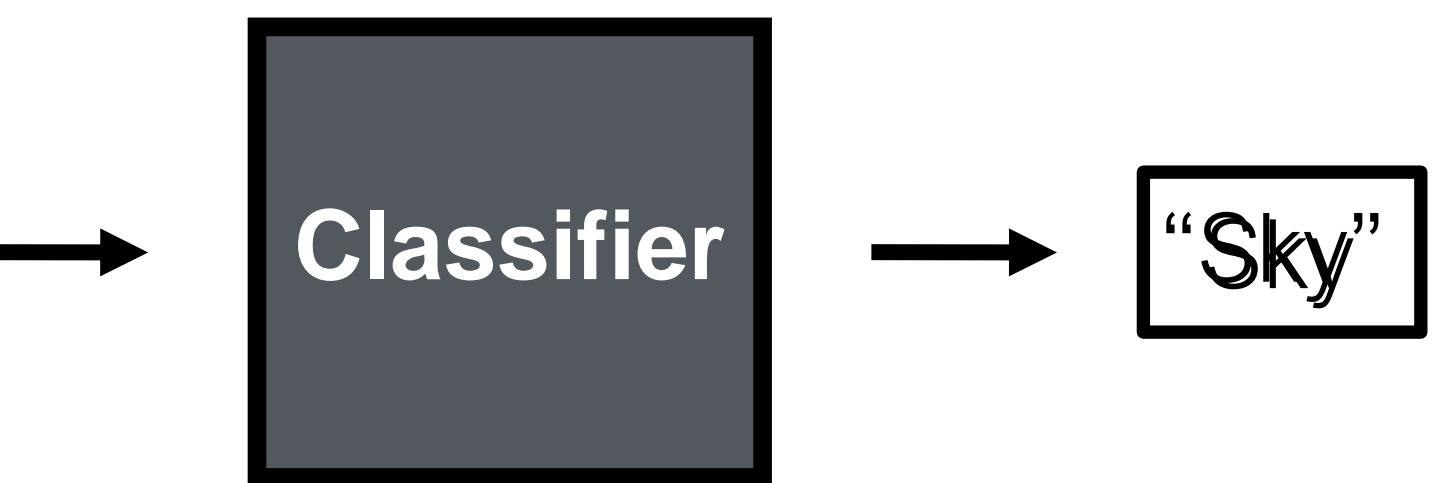
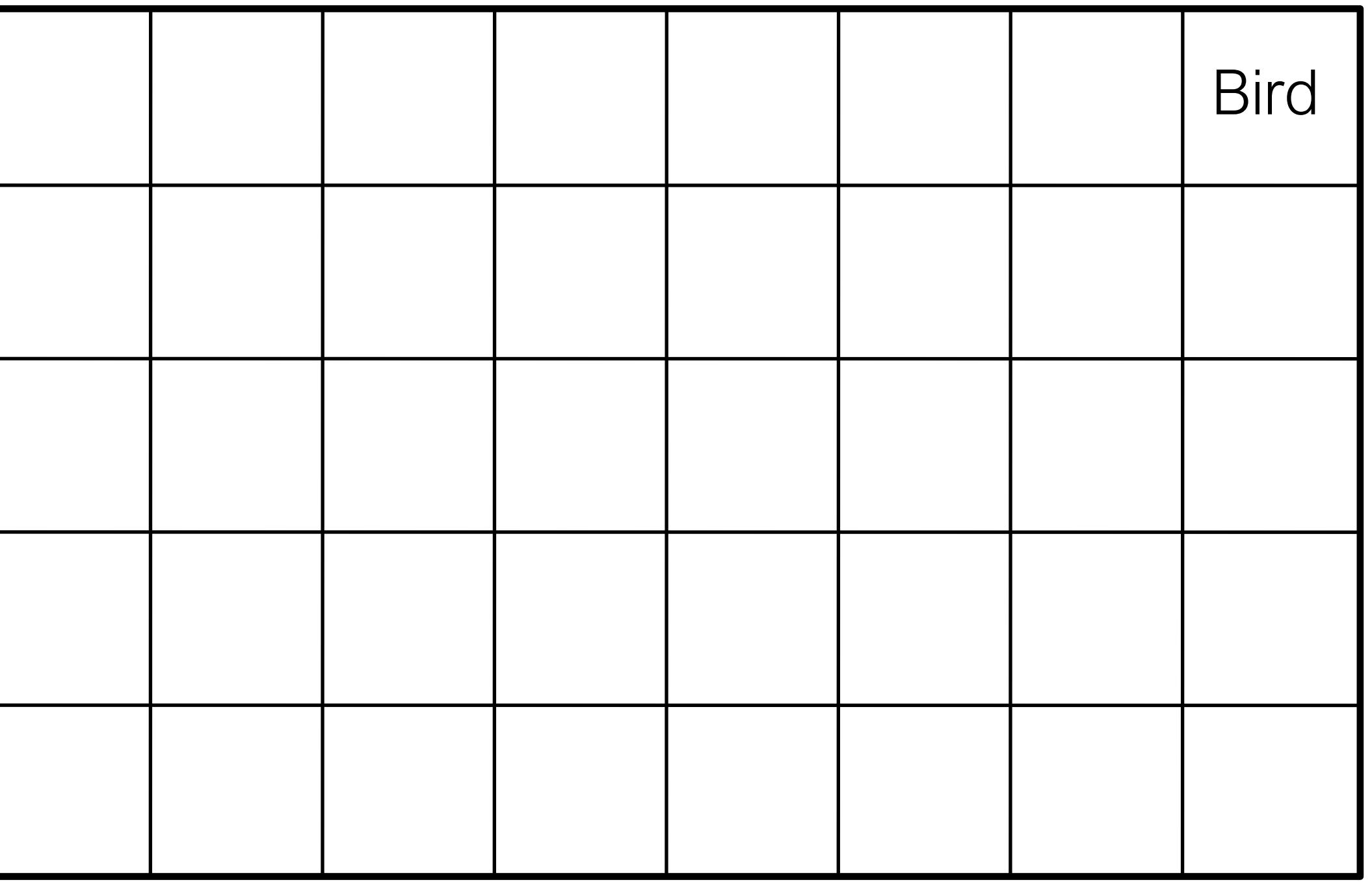
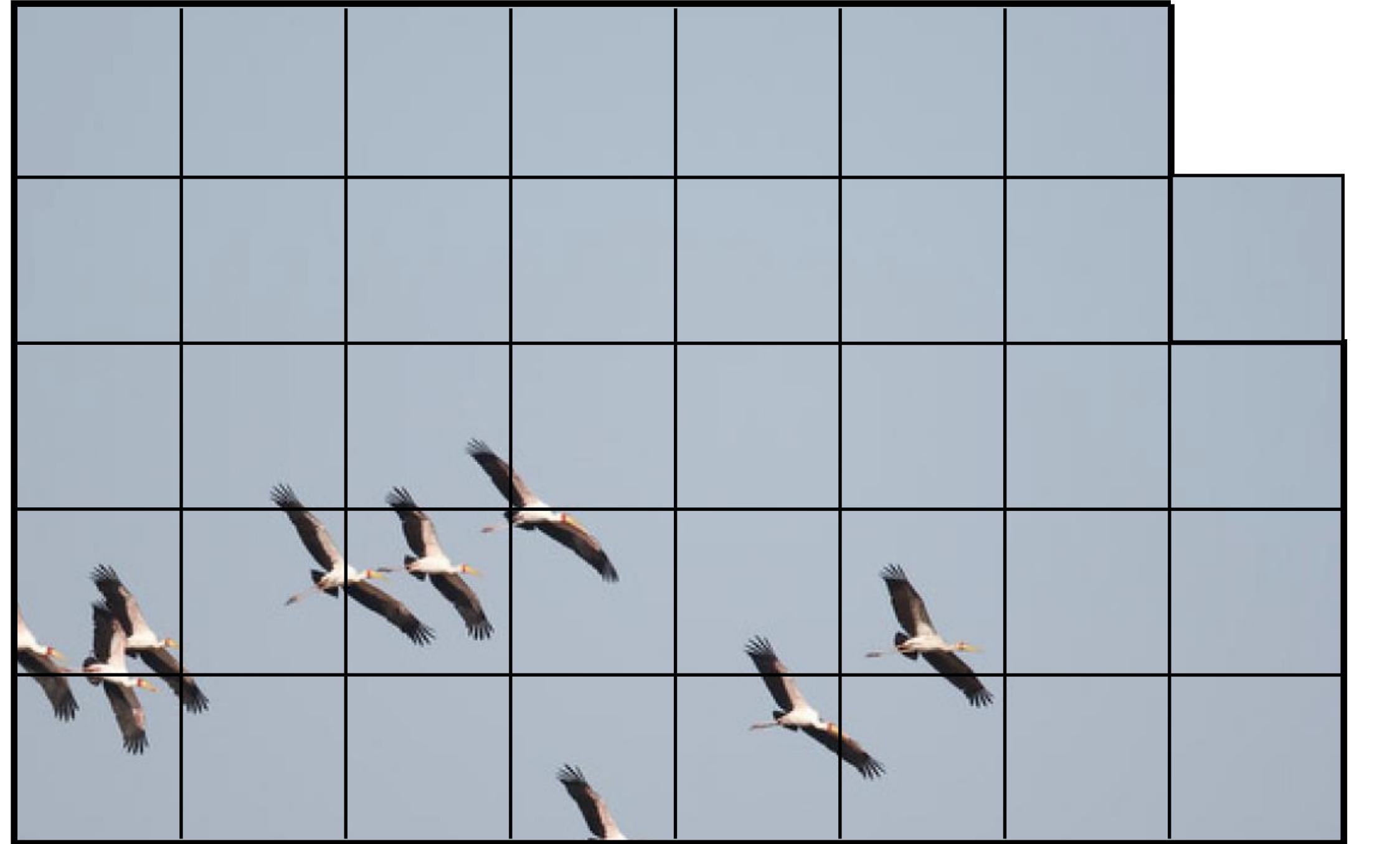


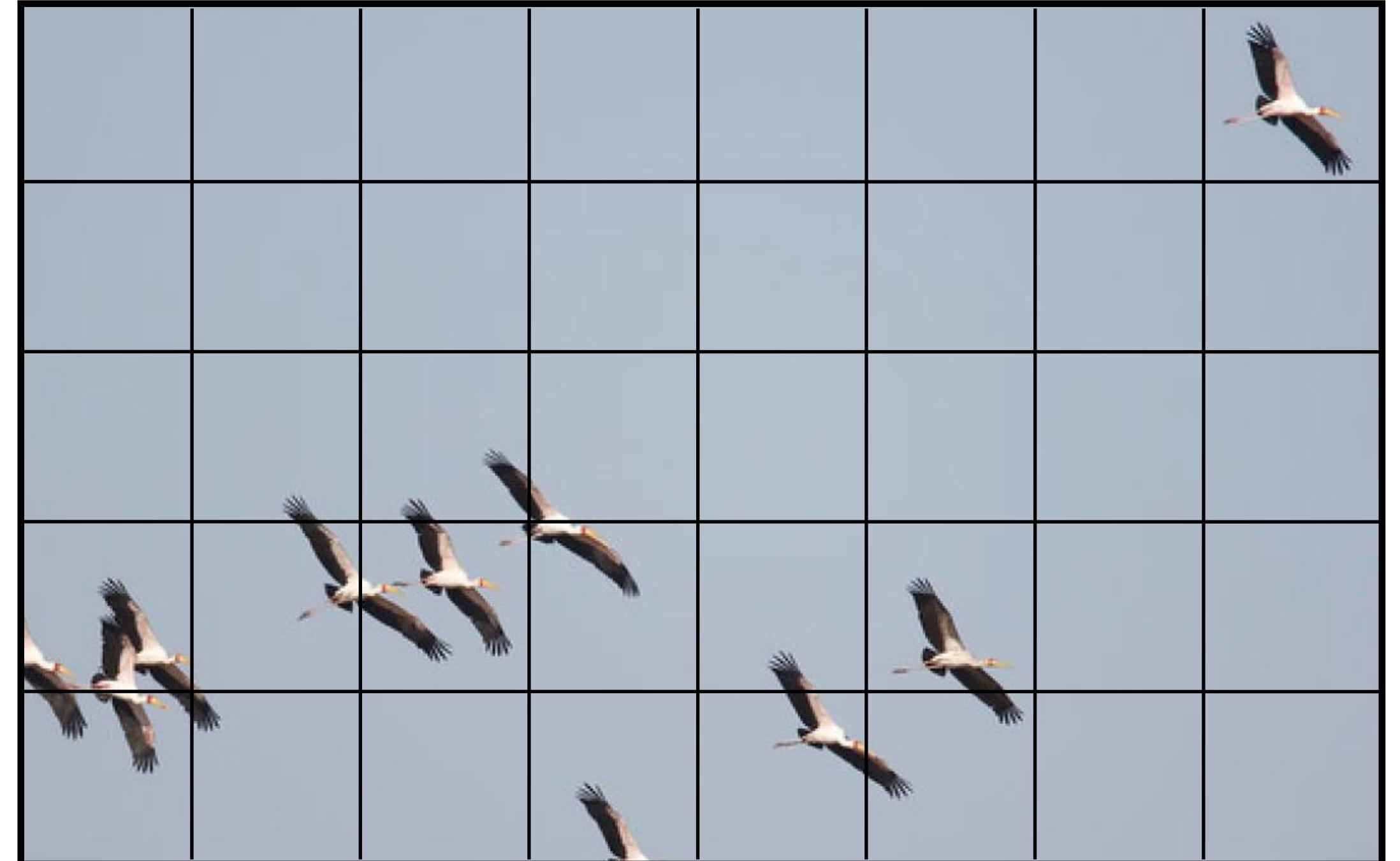


Classifier

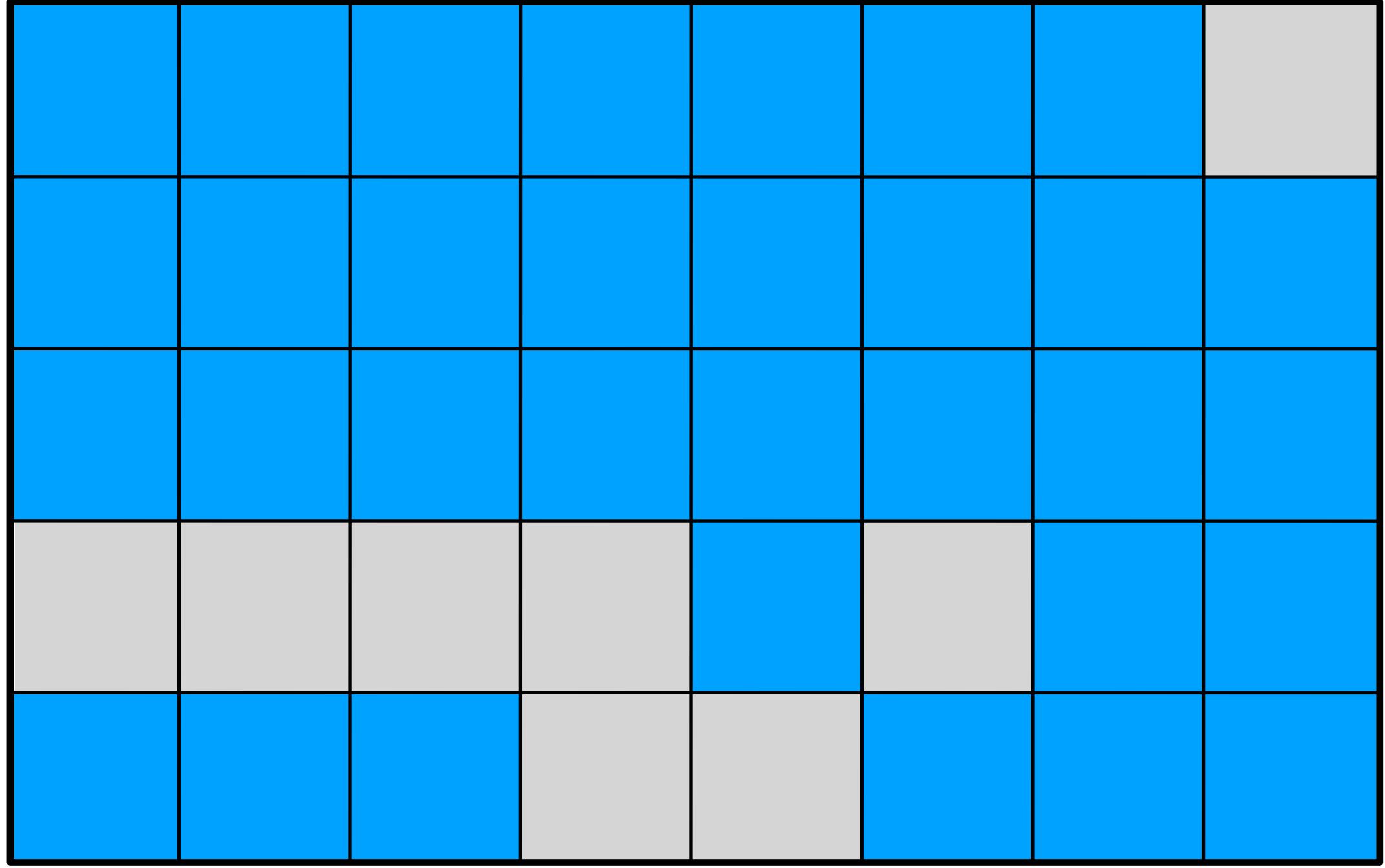
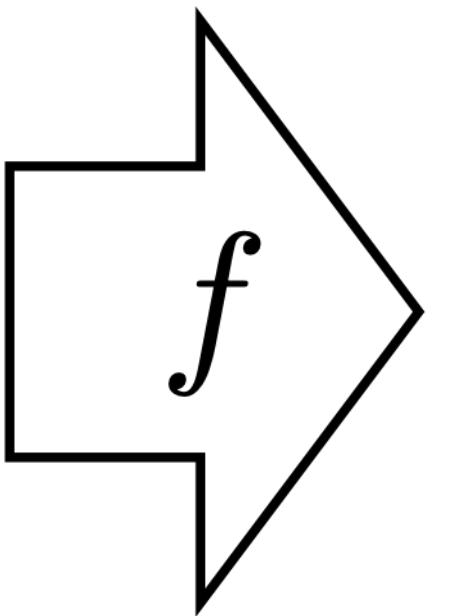
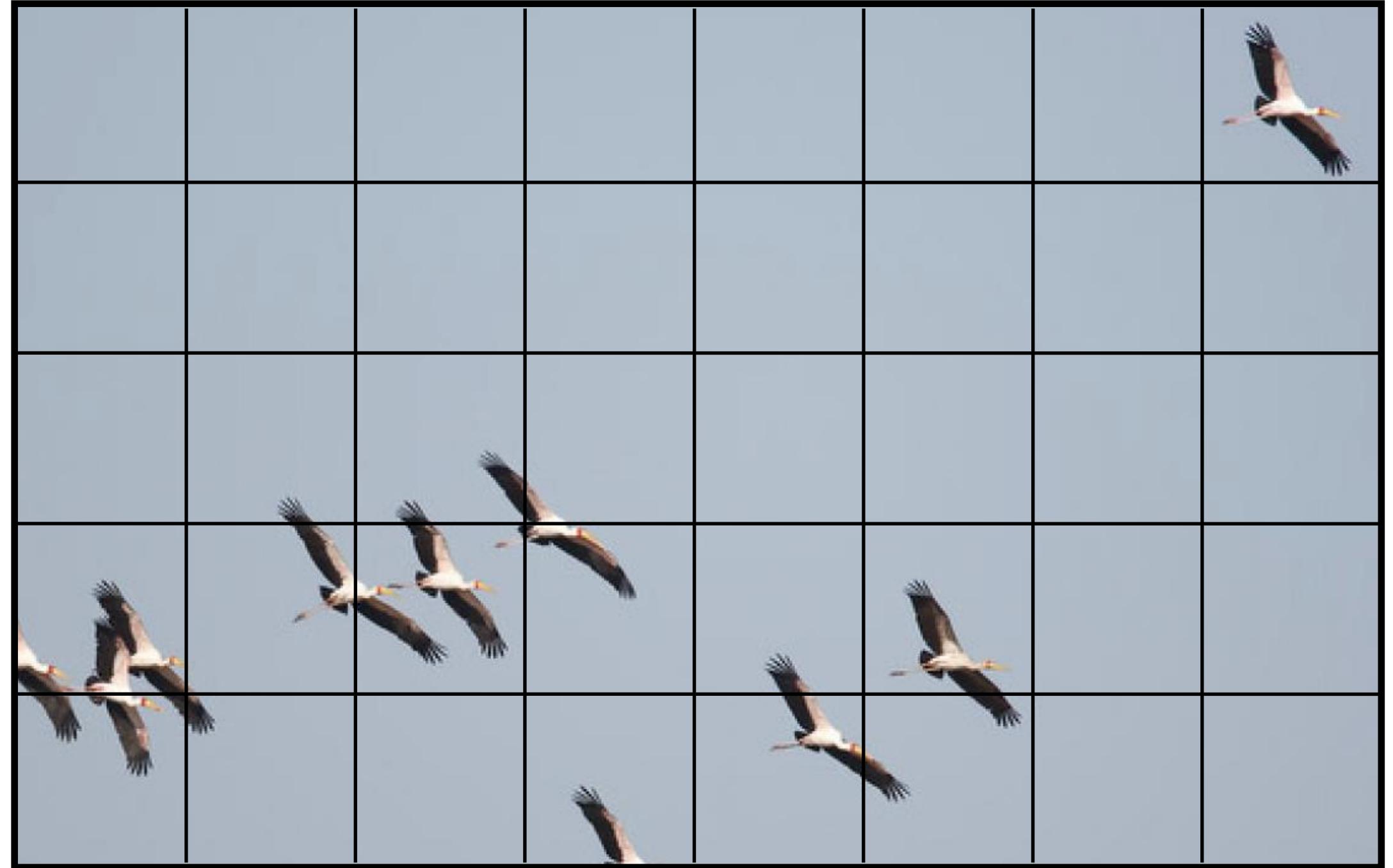


“Bird”





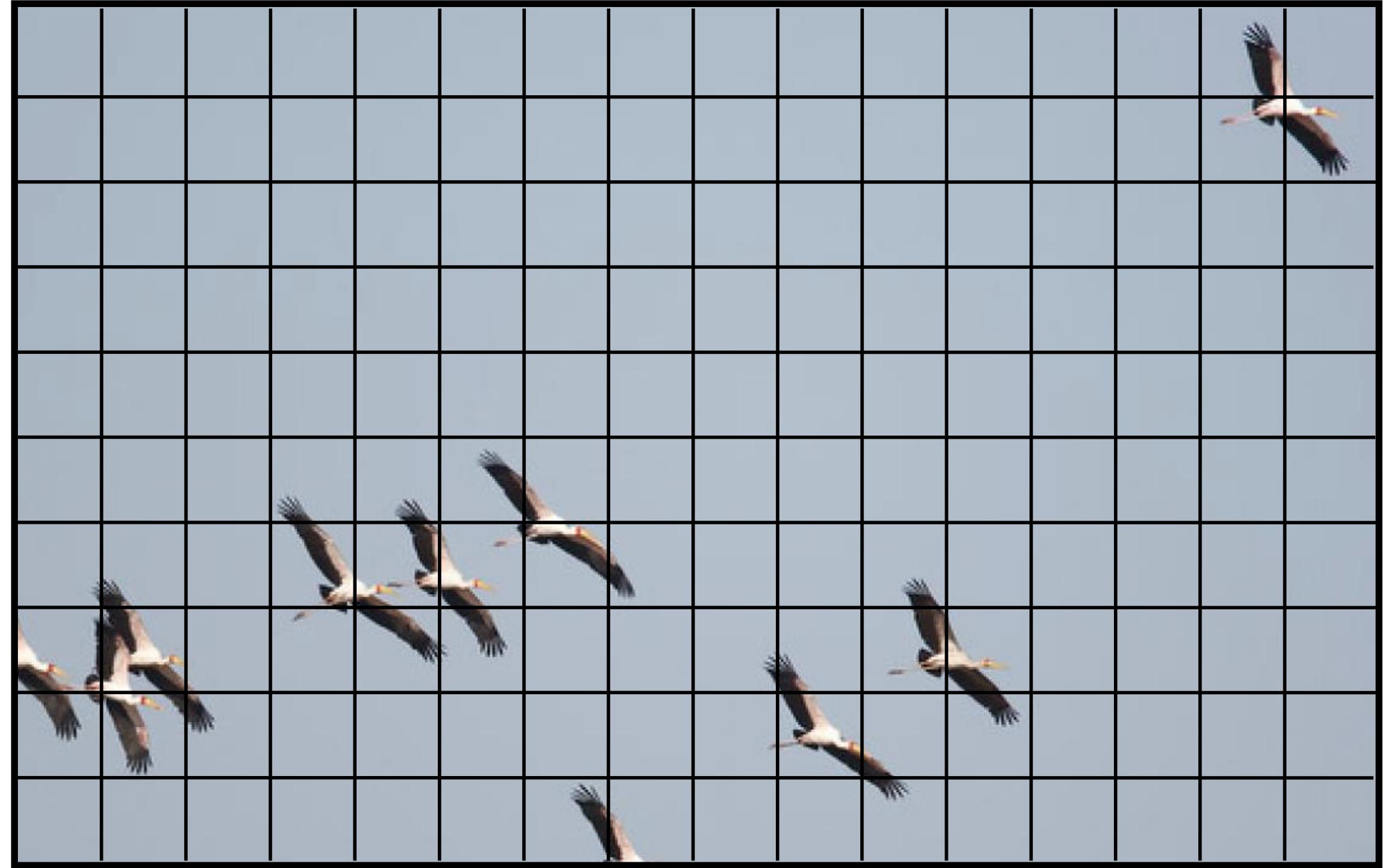
Sky	Sky	Sky	Sky	Sky	Sky	Sky	Bird
Sky	Sky	Sky	Sky	Sky	Sky	Sky	Sky
Sky	Sky	Sky	Sky	Sky	Sky	Sky	Sky
Bird	Bird	Bird	Sky	Bird	Sky	Sky	Sky
Sky	Sky	Sky	Bird	Sky	Sky	Sky	Sky



Problem:

What happens to objects that are bigger?

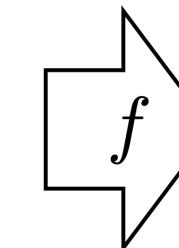
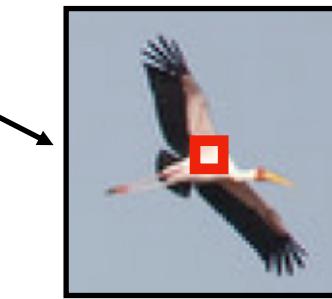
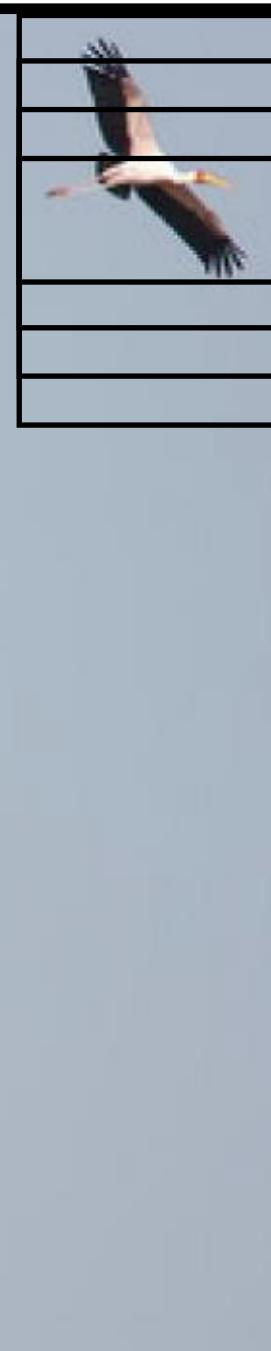
What if an object crosses multiple cells?



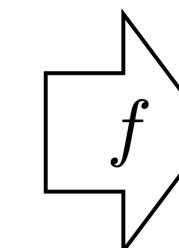
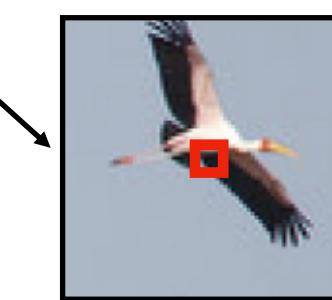
“Cell”-based approach is limited.

What can we do instead?

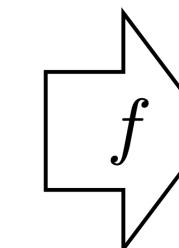
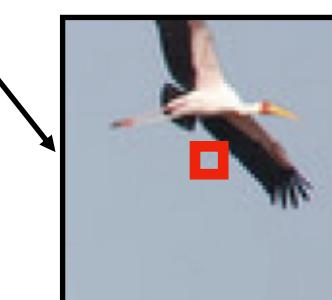
What's the object class of the center pixel?



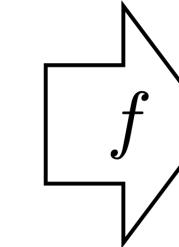
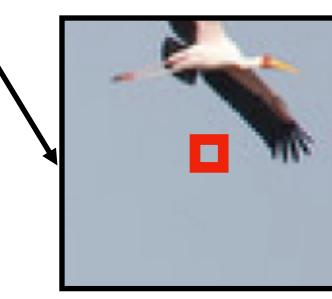
"Bird"



"Bird"



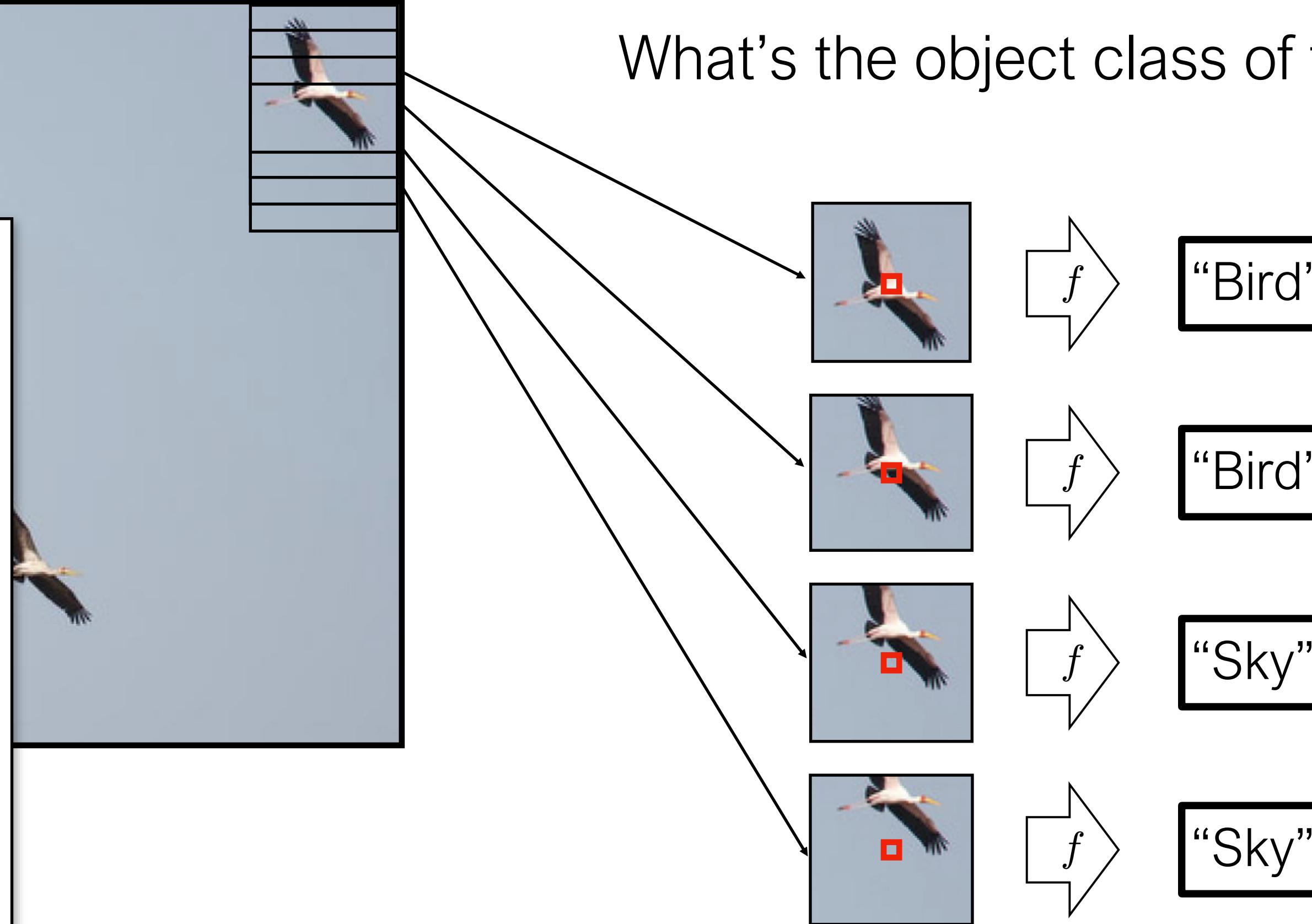
"Sky"

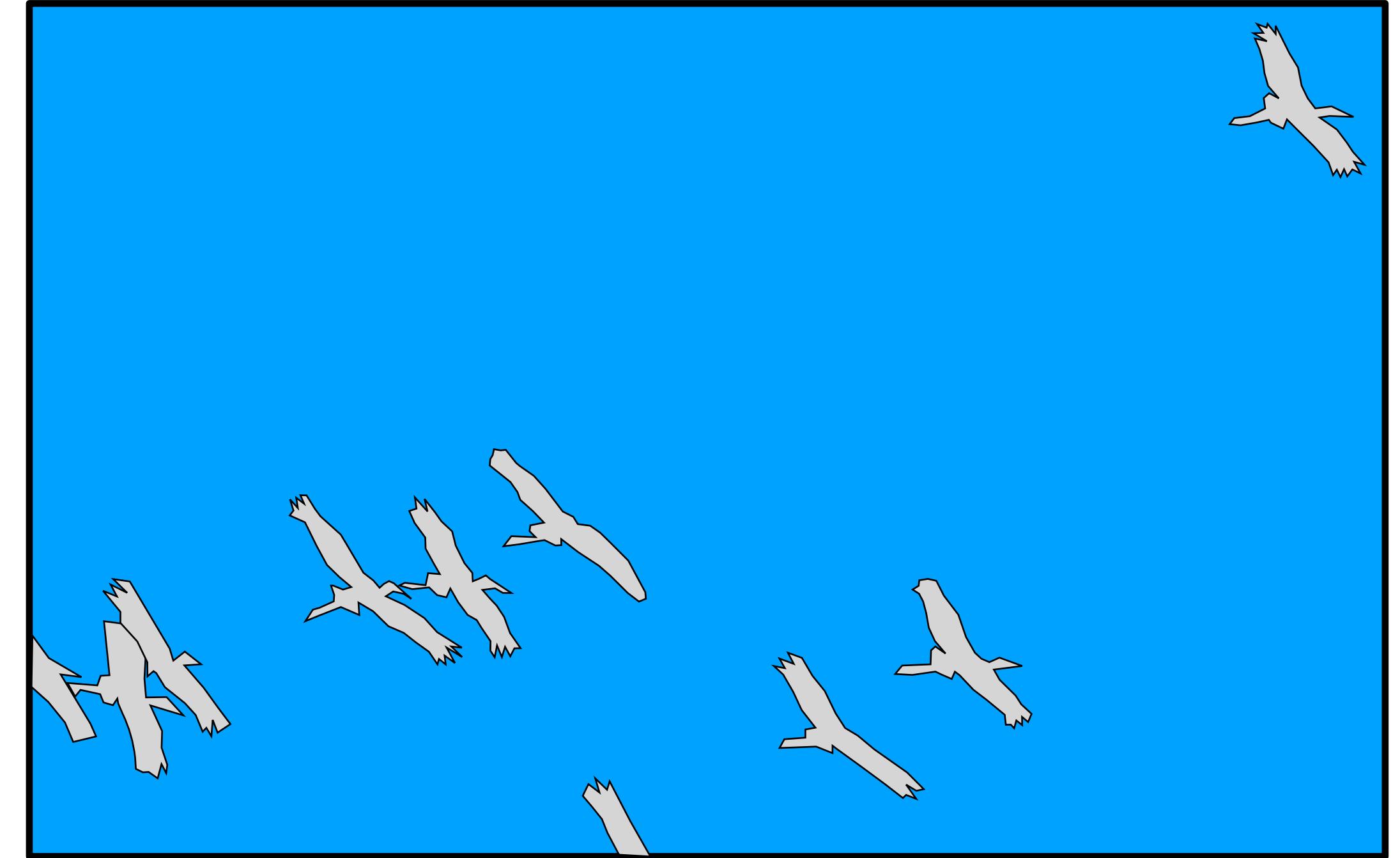
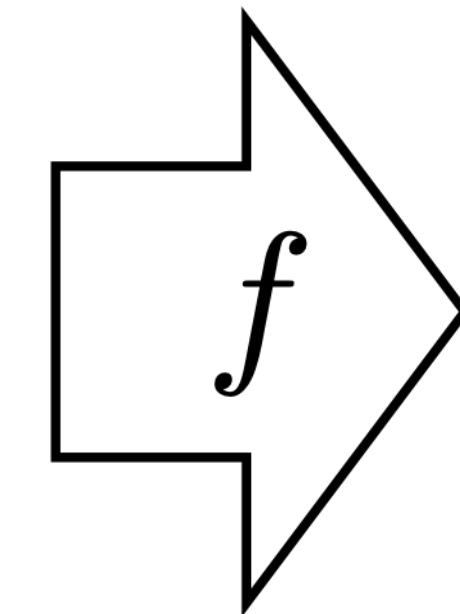


"Sky"

What's the object class of the center pixel?

<i>Training data</i>	
<i>x</i>	<i>y</i>
{  , “Bird” } , , “Bird” } , , “Sky” } , : 	



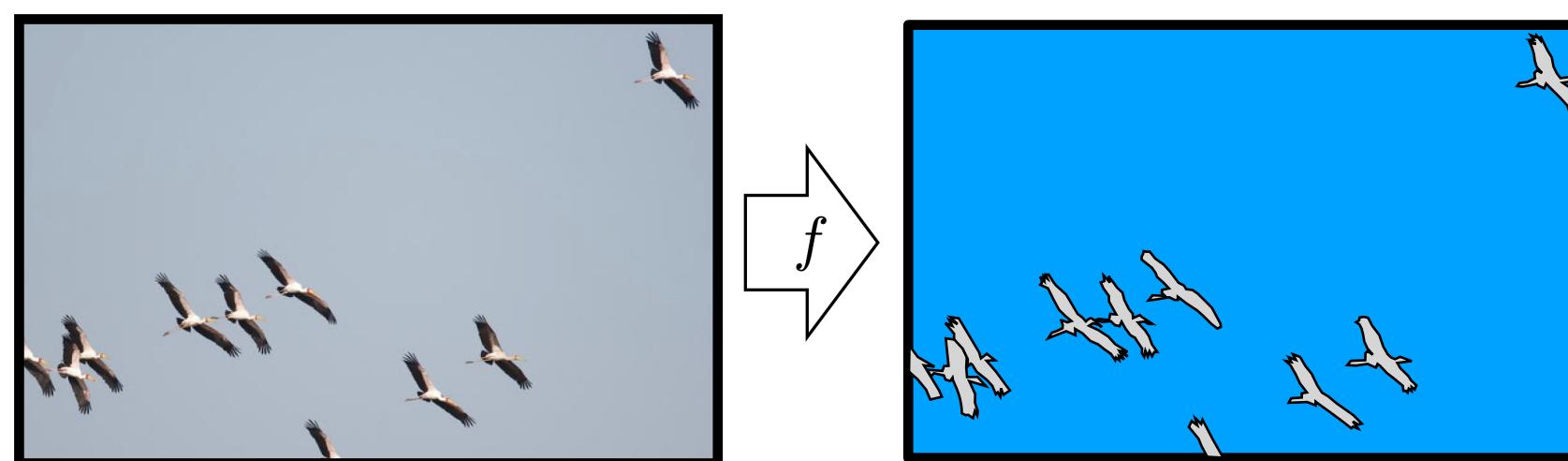
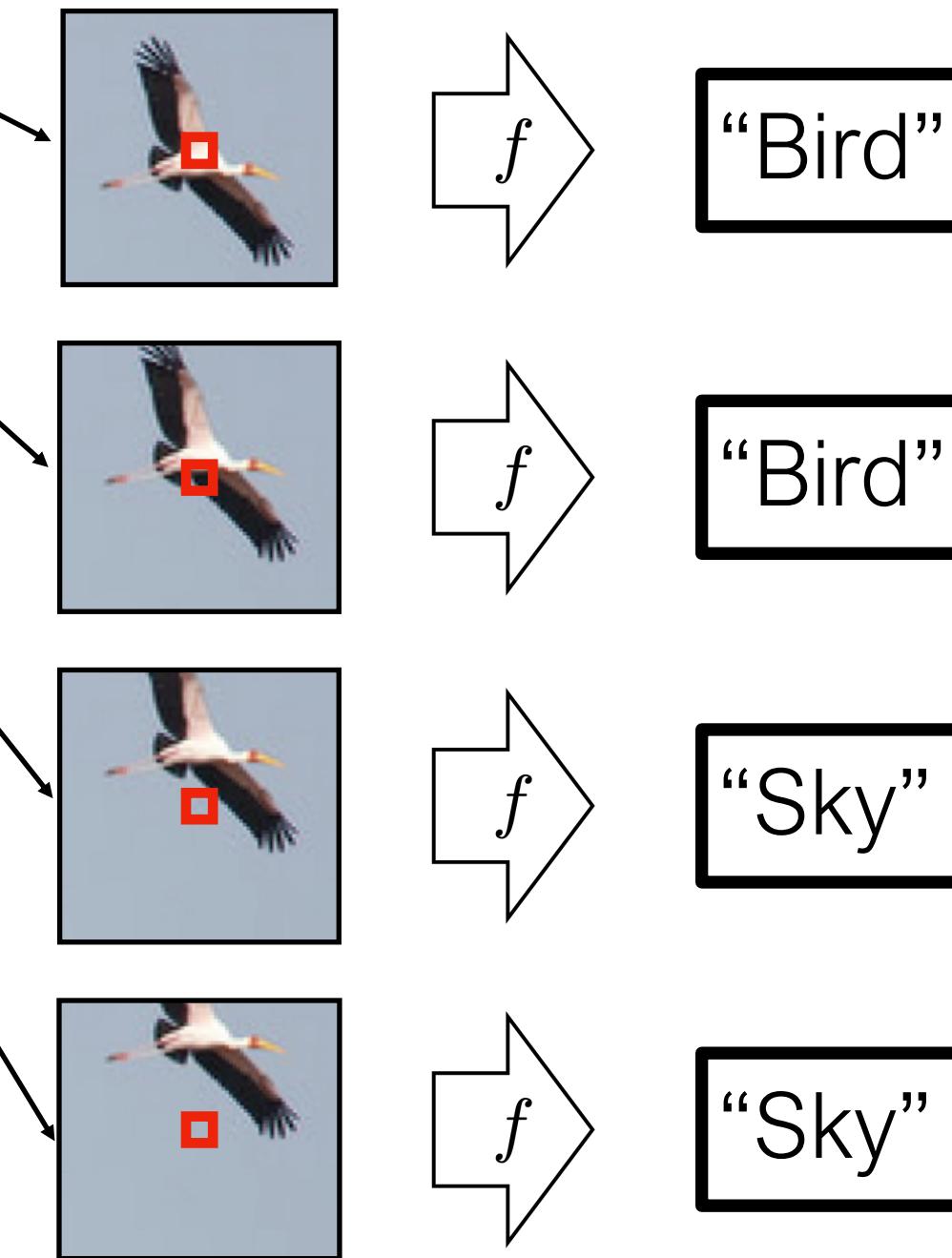


(Colors represent one-hot codes)

This problem is called **semantic segmentation**



What's the object class of the center pixel?

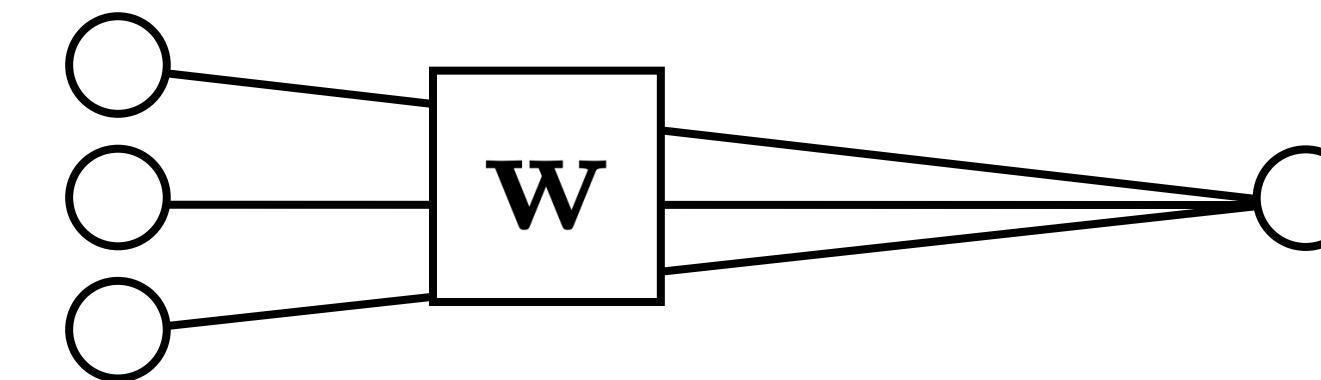


An *equivariant* mapping:

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

Translation invariance: process each patch in the same way.

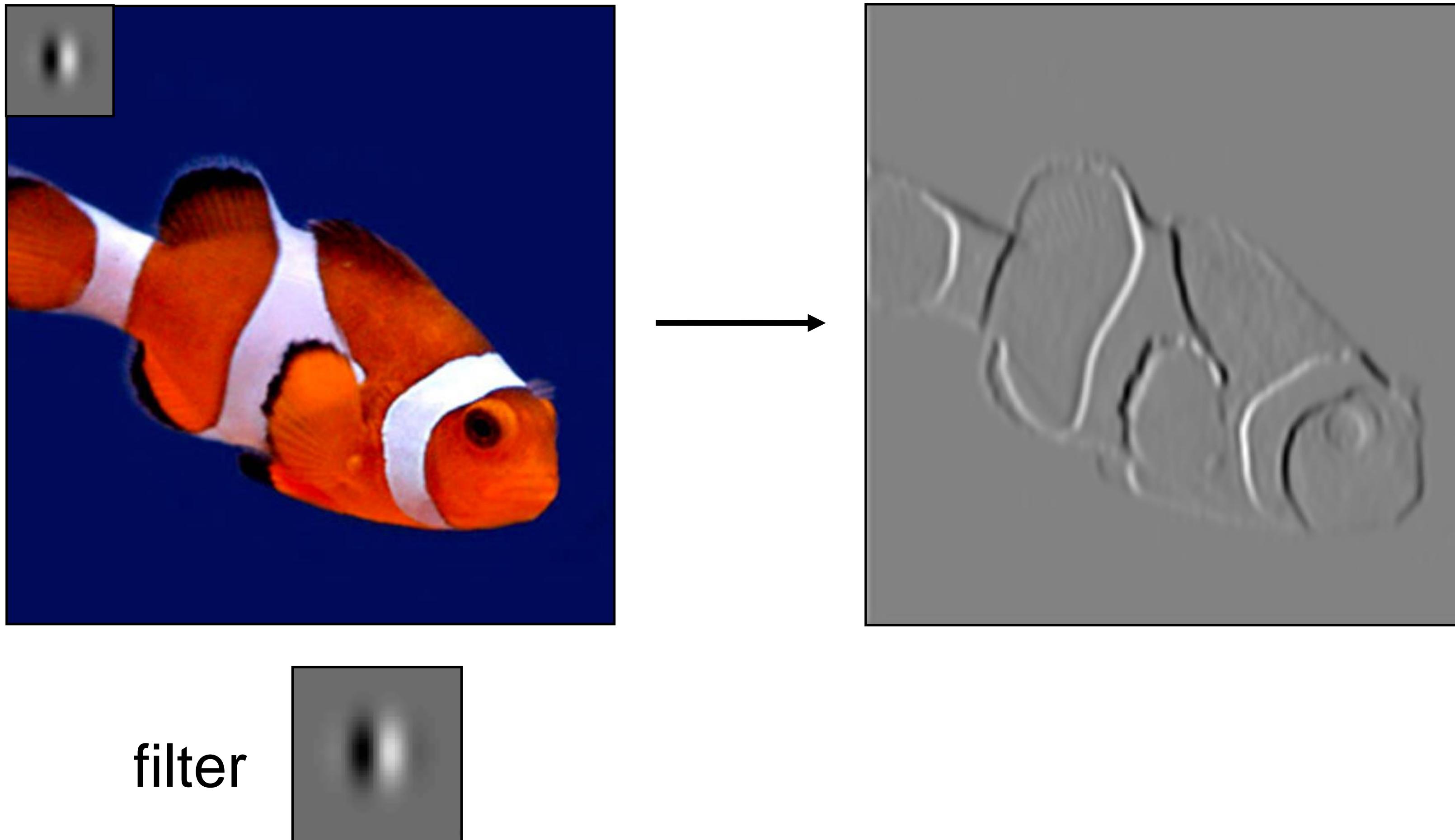
W computes a weighted sum of all pixels in the patch



W is a convolutional kernel applied to the full image!

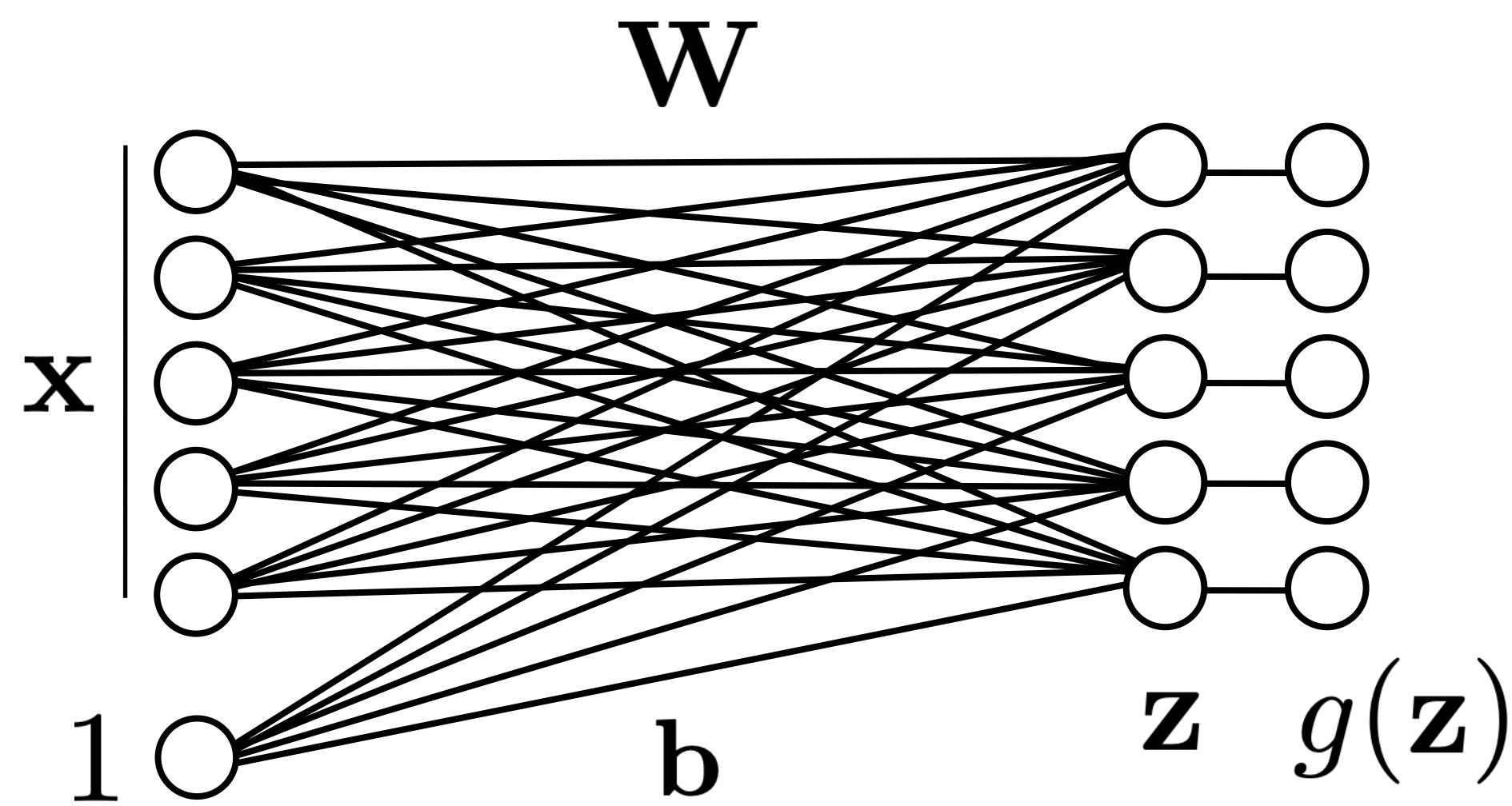


Convolution

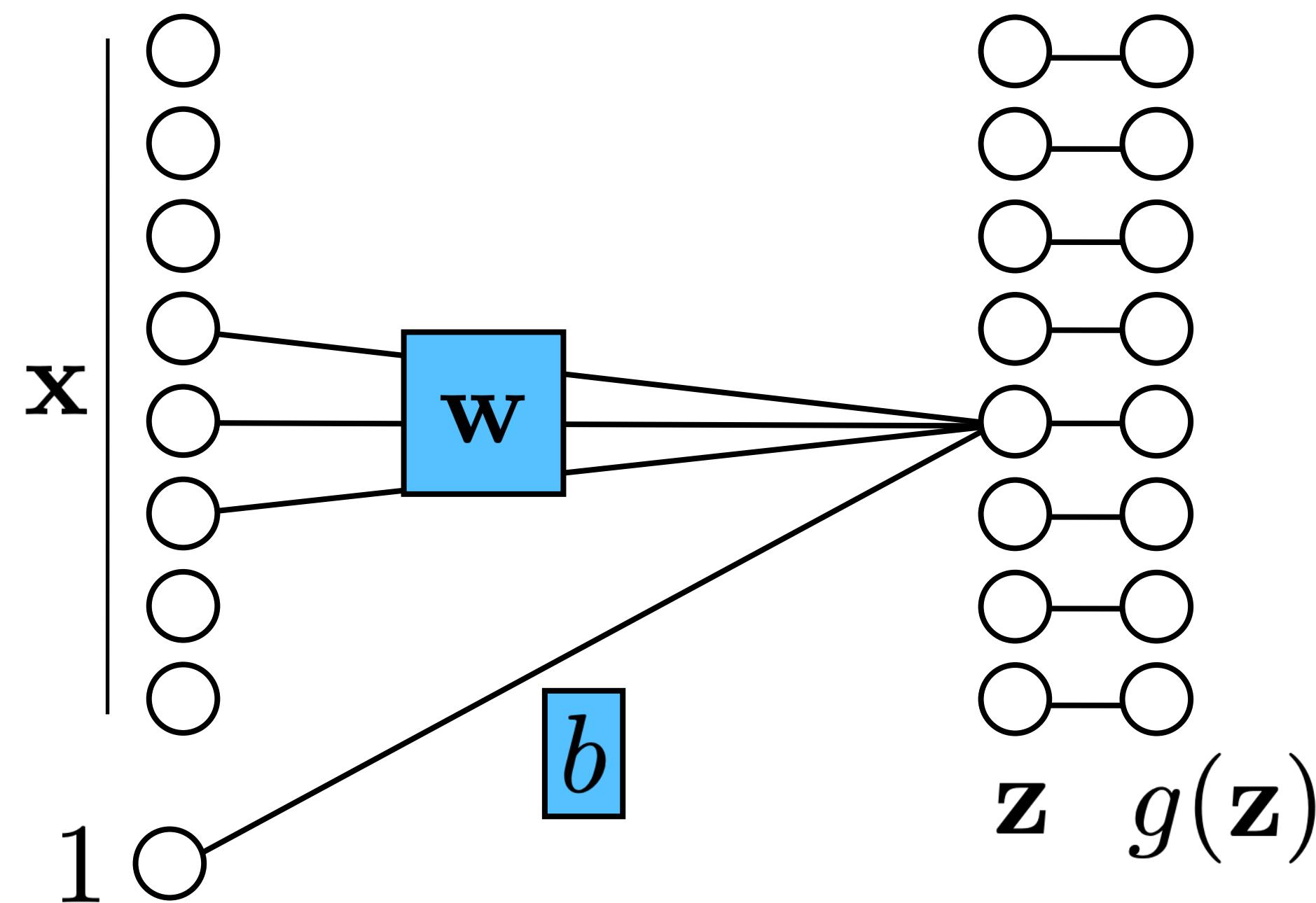


Fully-connected network

Fully-connected (fc) layer



Locally connected network

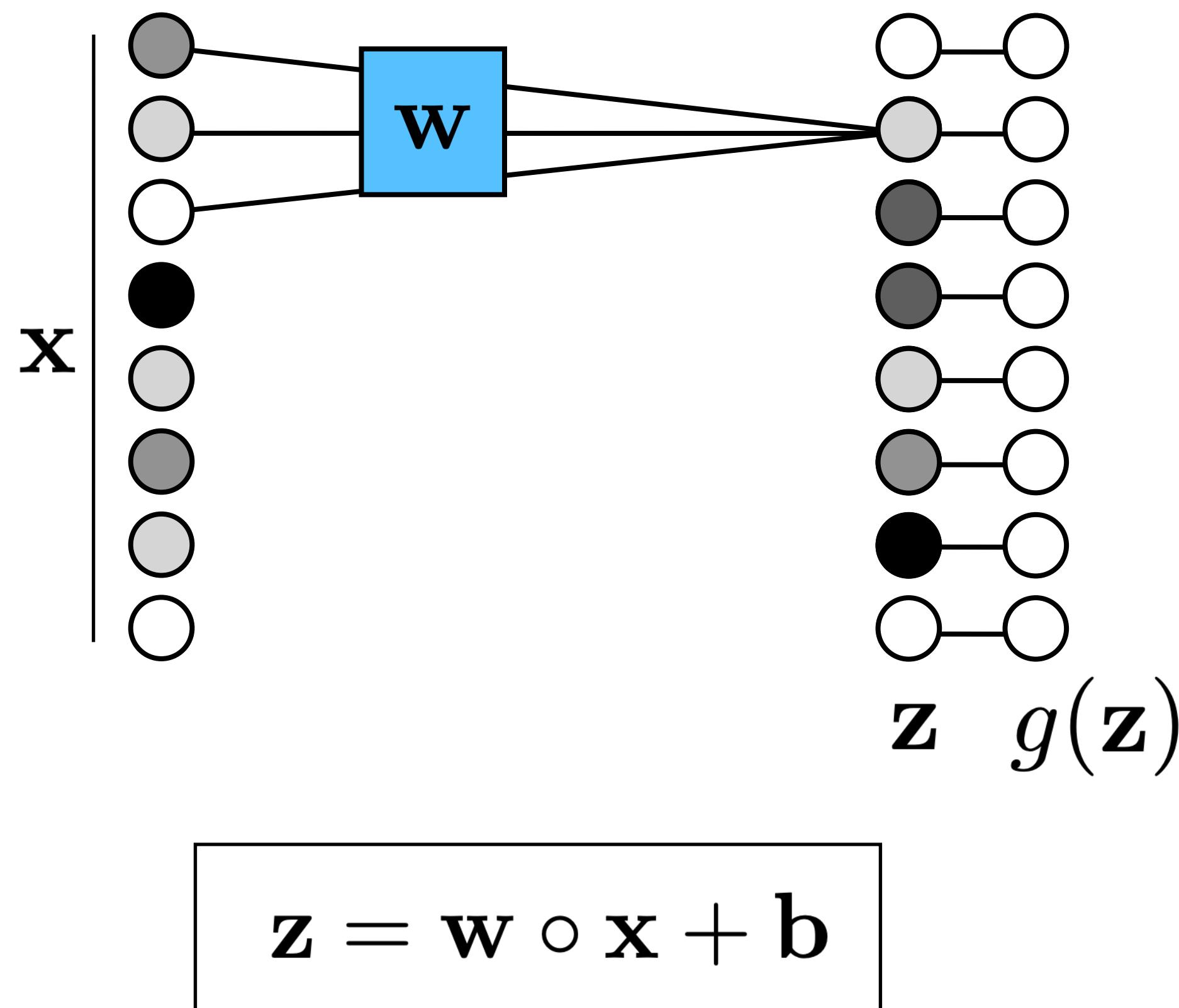


Often, we assume output is a **local** function of input.

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

Convolutional neural network

Conv layer

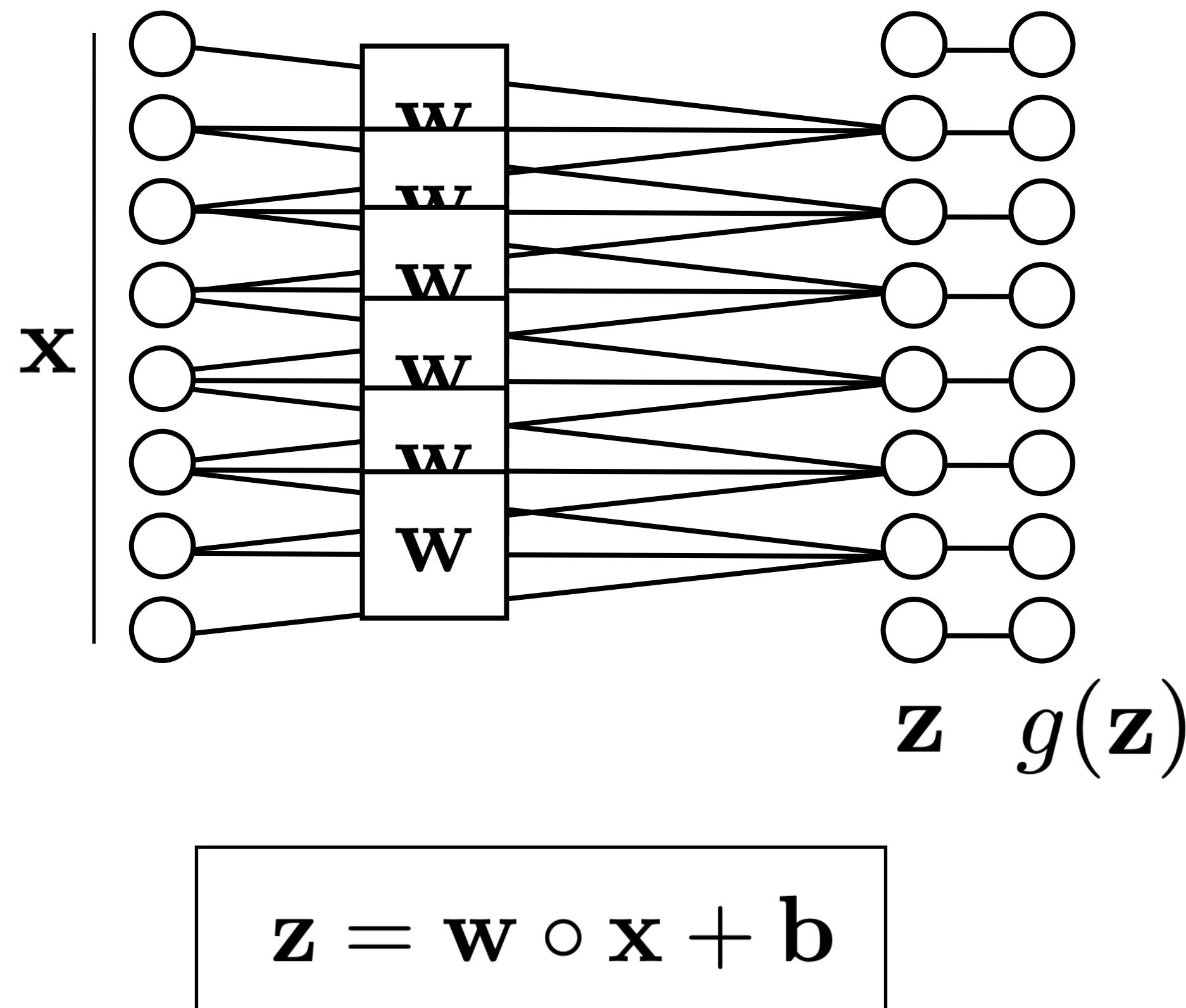


Often, we assume output is a **local** function of input.

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.

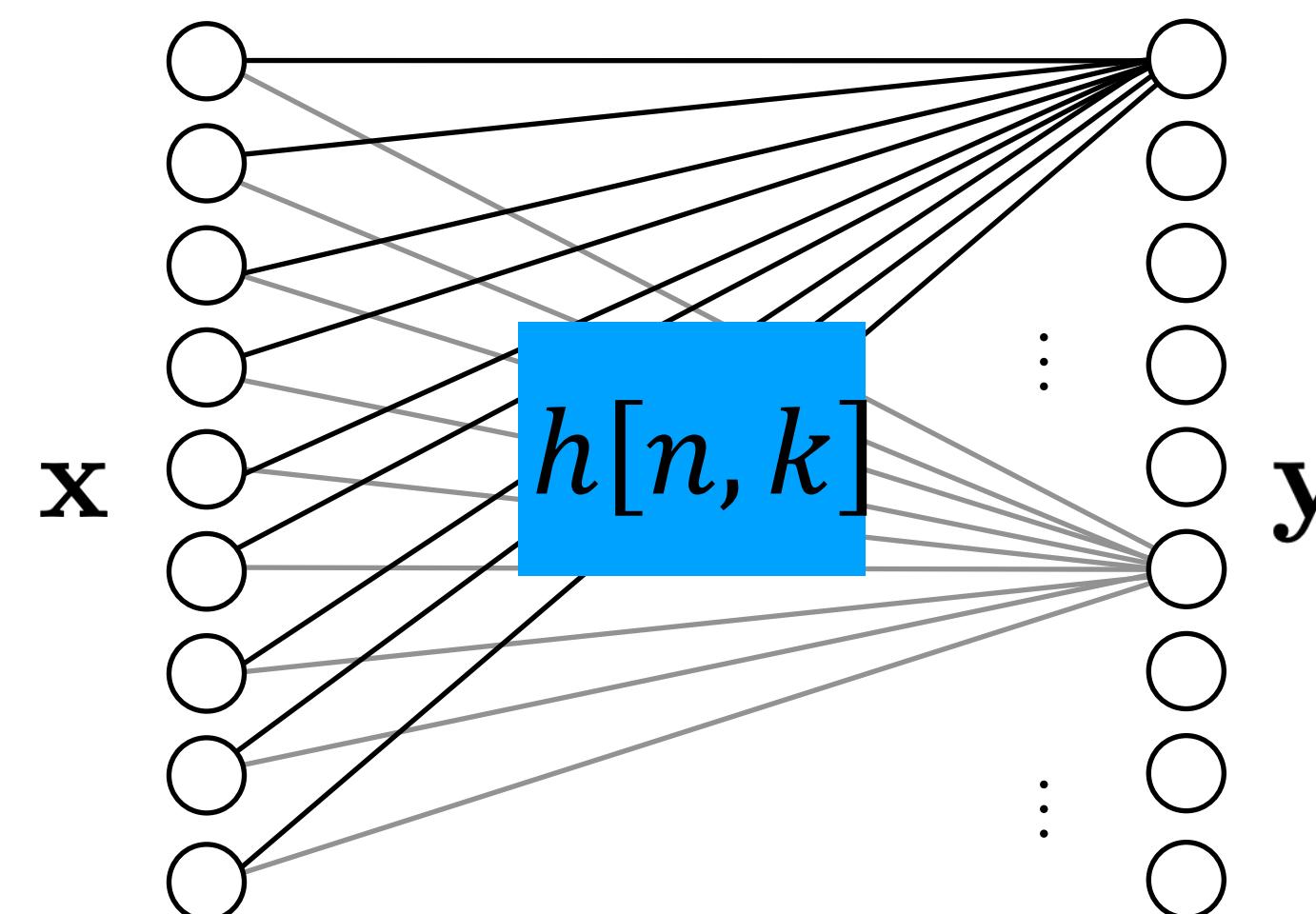
Linear system: $y = f(x)$

A linear function f can be written as a matrix multiplication:

y = $X \cdot h[n, k]$

n indexes rows,
k indexes columns

It can also be represented as a fully connected linear neural network



$$y[n] = \sum_{k=0}^{N-1} h[n, k]x[k]$$

$h[n, k]$ Is the strength of the connection between $x[k]$ and $y[n]$

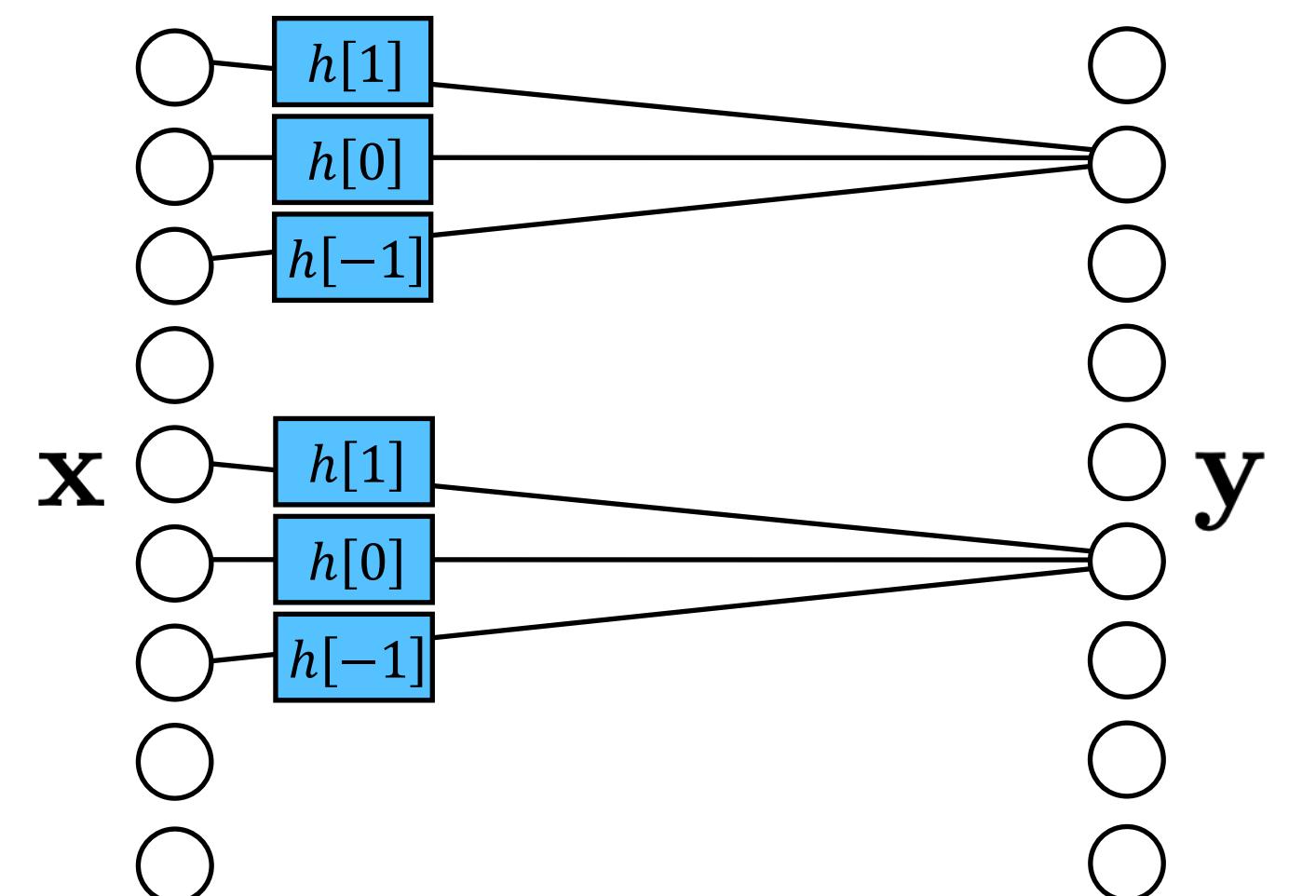
Convolution

A linear shift invariant (LSI) function f can be written as a matrix multiplication:

$y =$ X

$h[n - k]$ n indexes rows,
 k indexes columns

It can also be represented as a convolutional layer of neural net:



$$y \rightarrow y[n] = \sum_{k=-1}^1 h[k]x[n-k]$$

$h[n - k]$ is the strength of the connection between $x[k]$ and $y[n]$

$$\mathbf{y} = \mathbf{W} * \mathbf{x}$$

e.g., pixel image

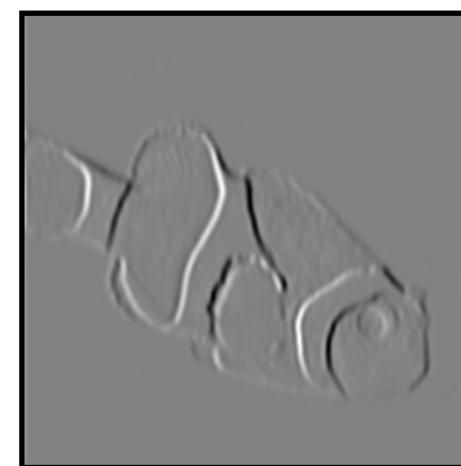
- Constrained linear layer
- Fewer parameters —> easier to learn, less overfitting

Five views on convolutional layers

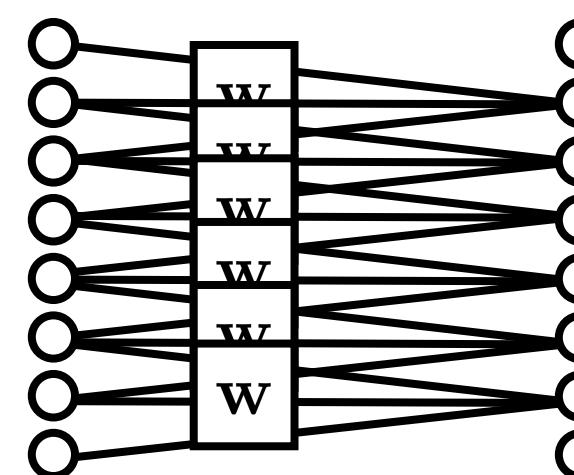
1. Equivariant with translation $f(\text{translate}(x)) = \text{translate}(f(x))$

2. Patch processing

3. Image filter



4. Parameter sharing



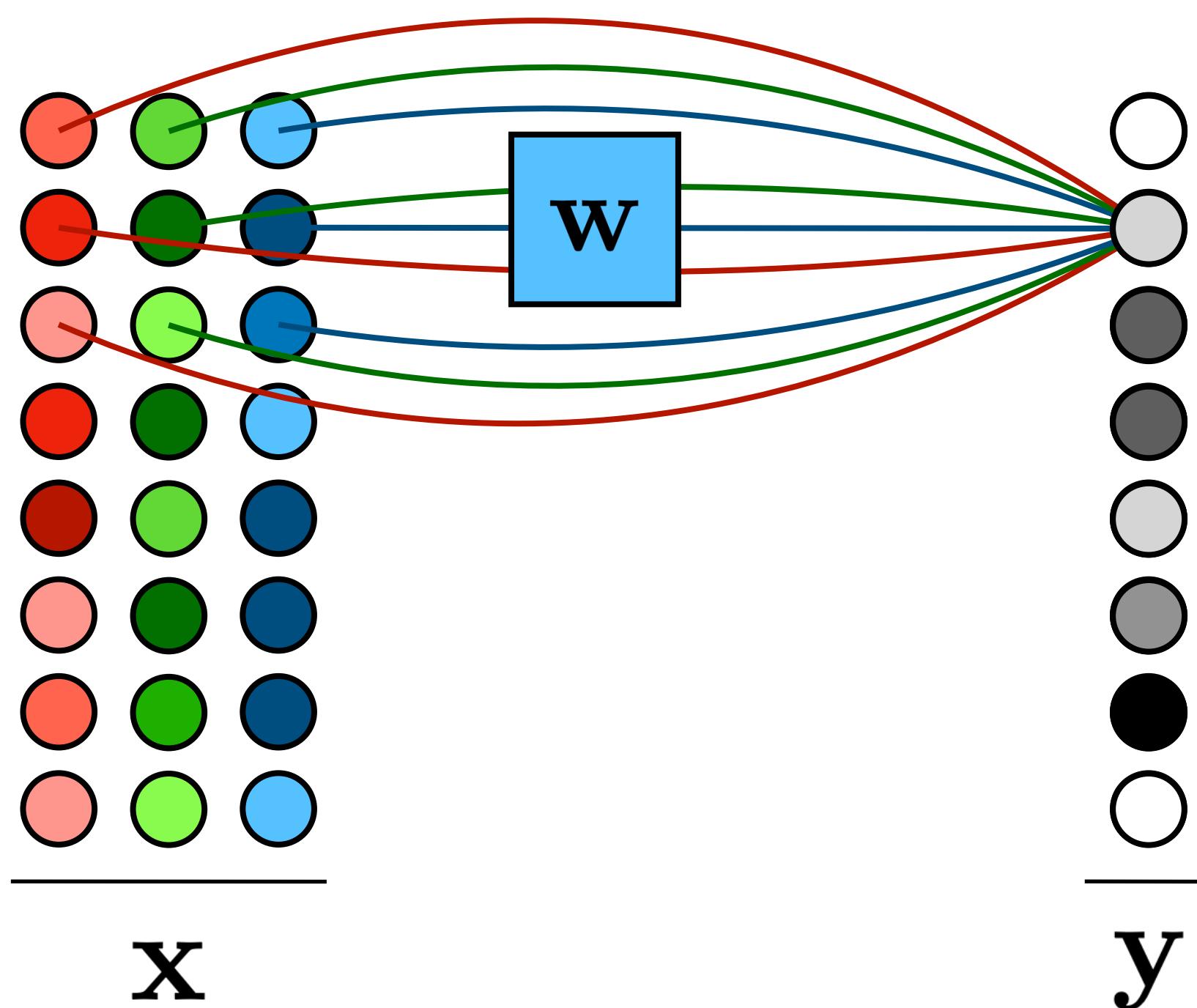
5. A way to process variable-sized tensors

What if we have color?

(aka multiple input channels?)

Multiple channel inputs

Conv layer

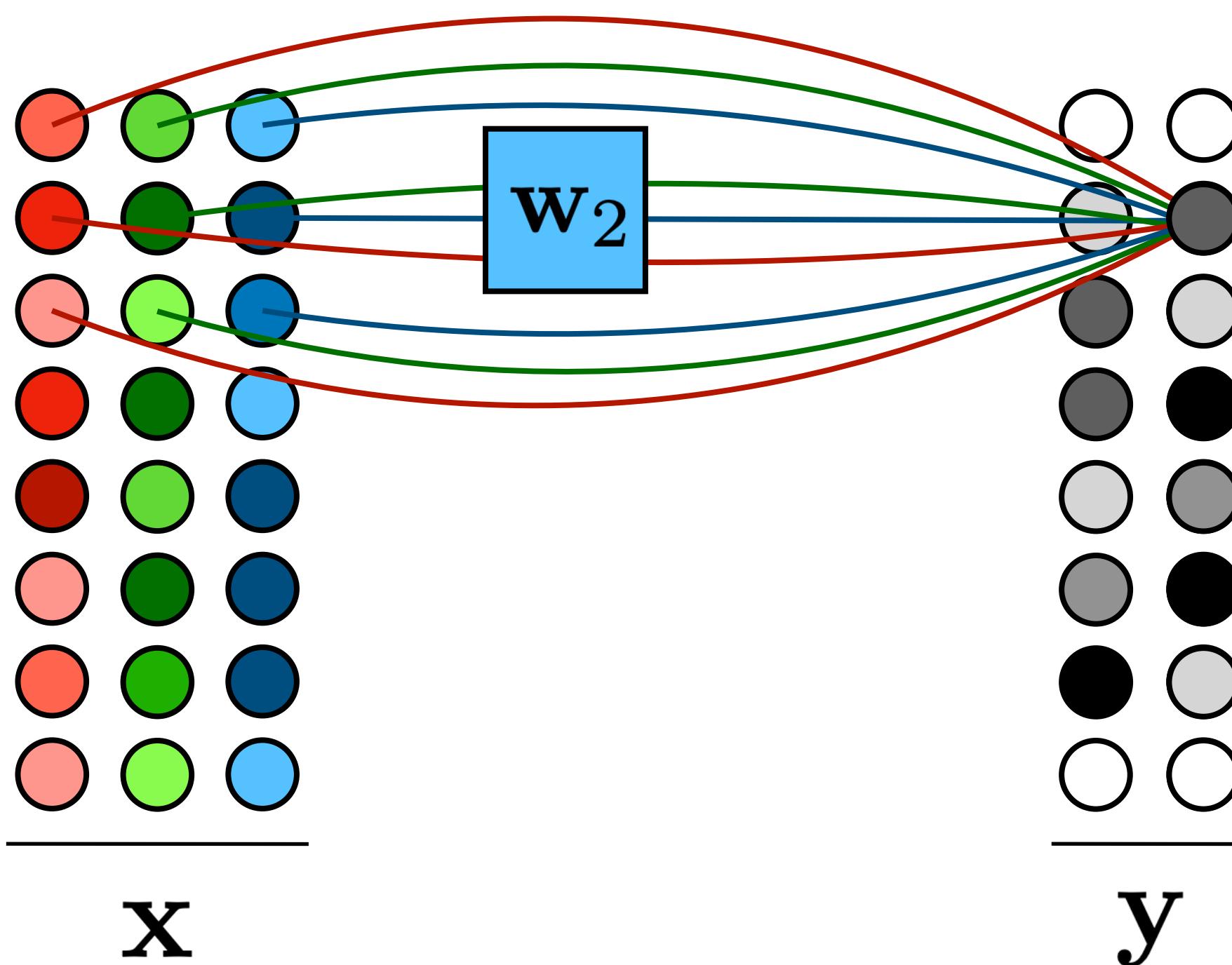


$$\mathbf{y} = \sum_c \mathbf{w}_c \circ \mathbf{x}_c$$

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

Multiple channel *outputs*

Conv layer

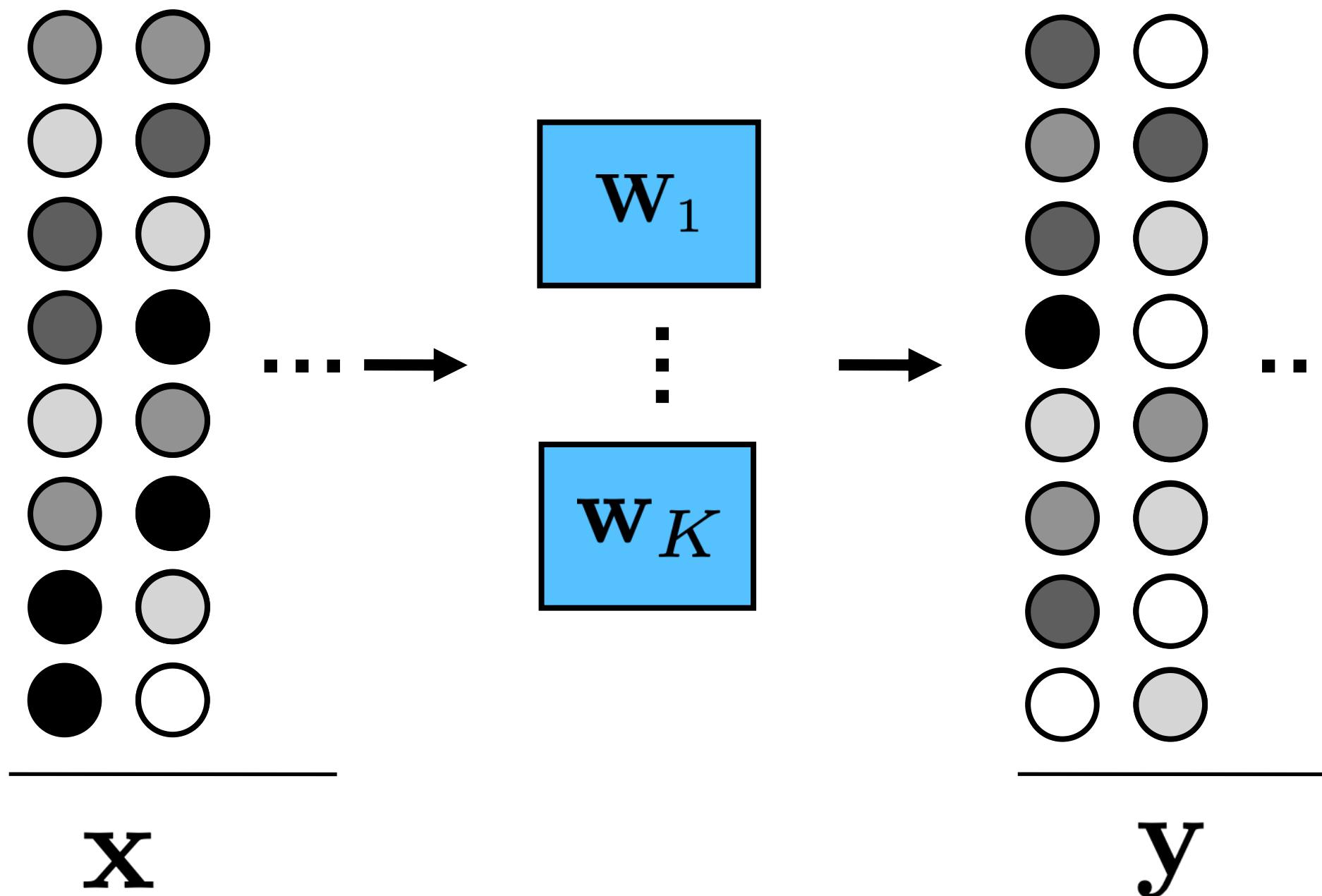


$$\mathbf{y}_k = \sum_c \mathbf{w}_{k_c} \circ \mathbf{x}_c$$

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

Multiple channels

Conv layer

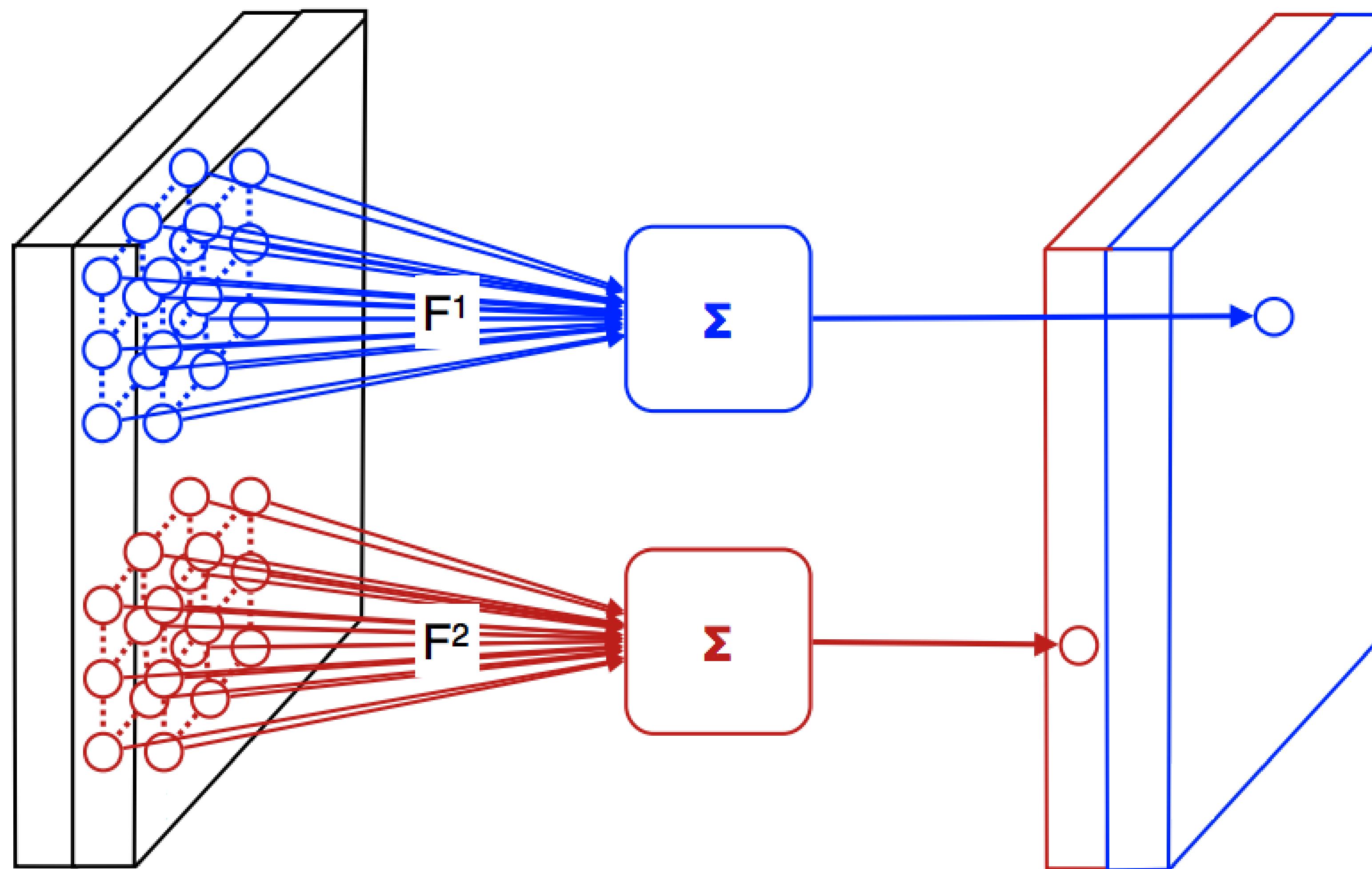


$$\mathbf{y}_k = \sum_c \mathbf{w}_{k_c} \circ \mathbf{x}_c$$
$$\mathbb{R}^{N \times C} \rightarrow \mathbb{R}^{N \times K}$$

Input features

A bank of 2 filters

2-dimensional output
feature maps



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

[Figure modified from Andrea Vedaldi]

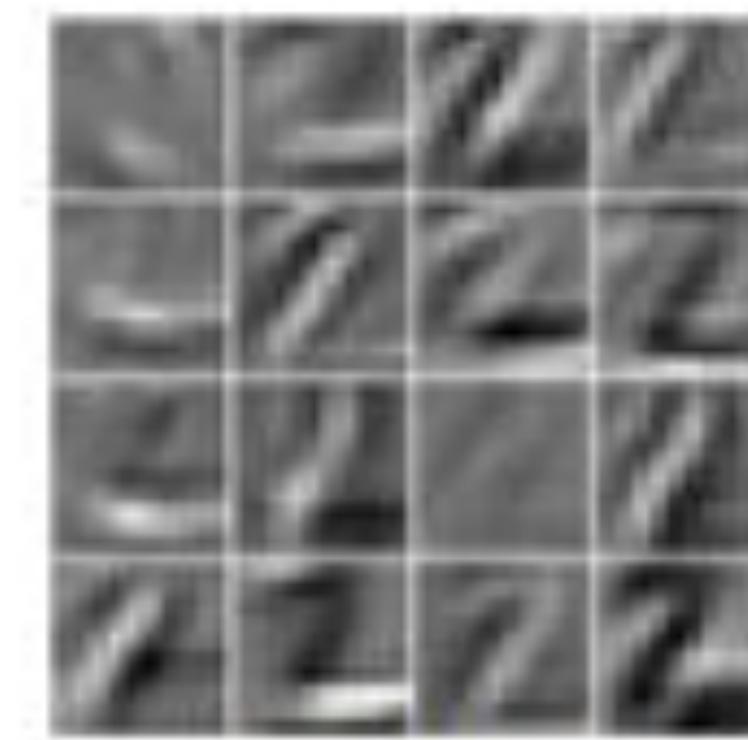
Feature maps



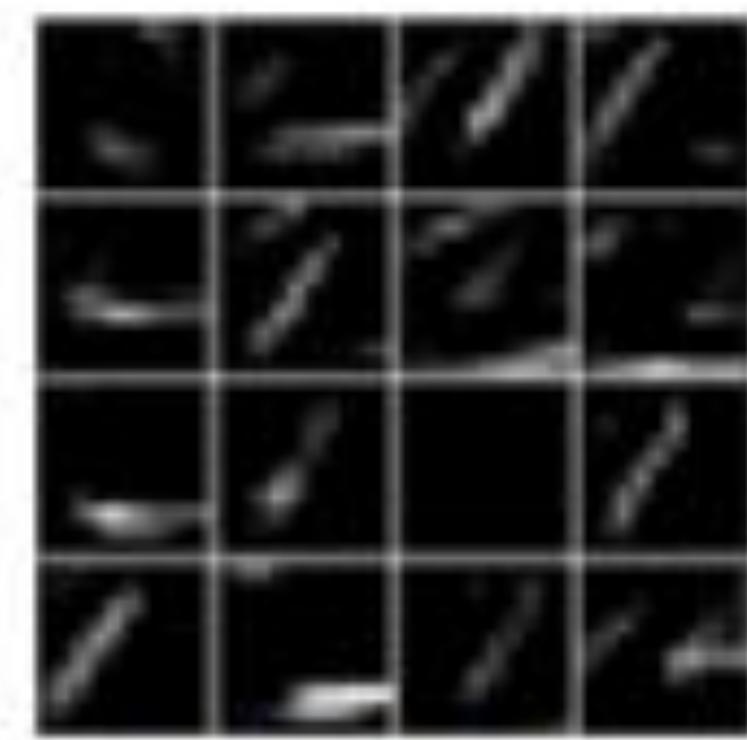
conv1



relu1



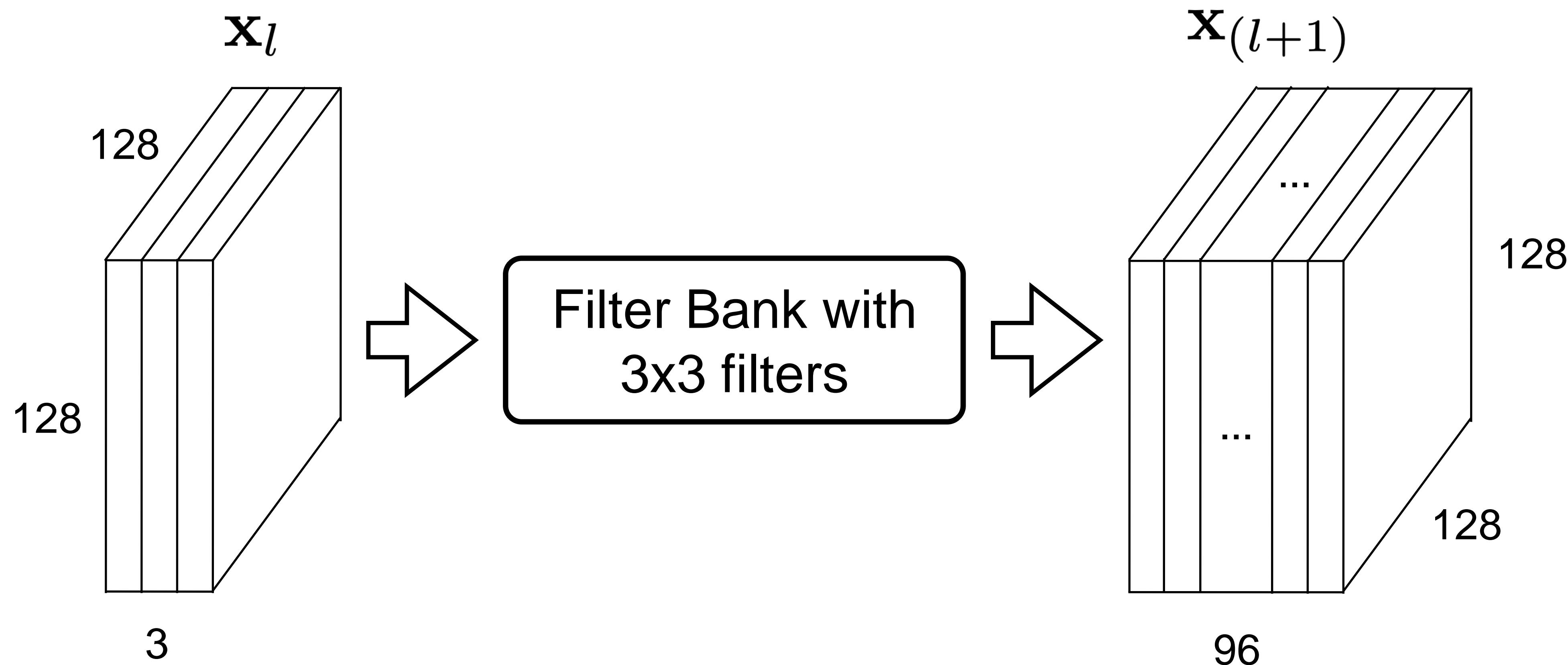
conv2



relu2

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

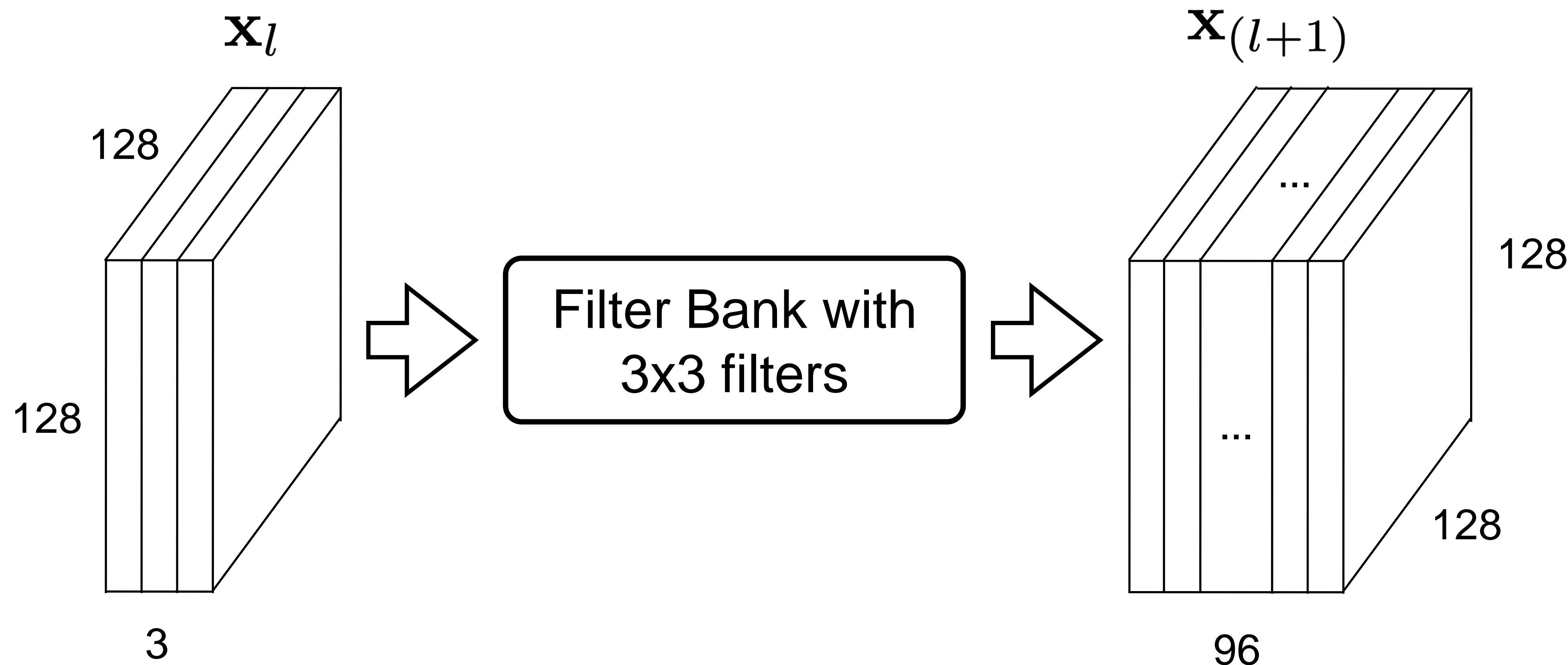
Multiple channels: Example



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