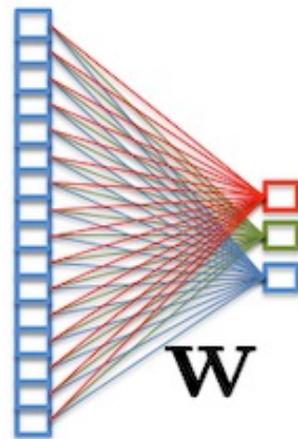


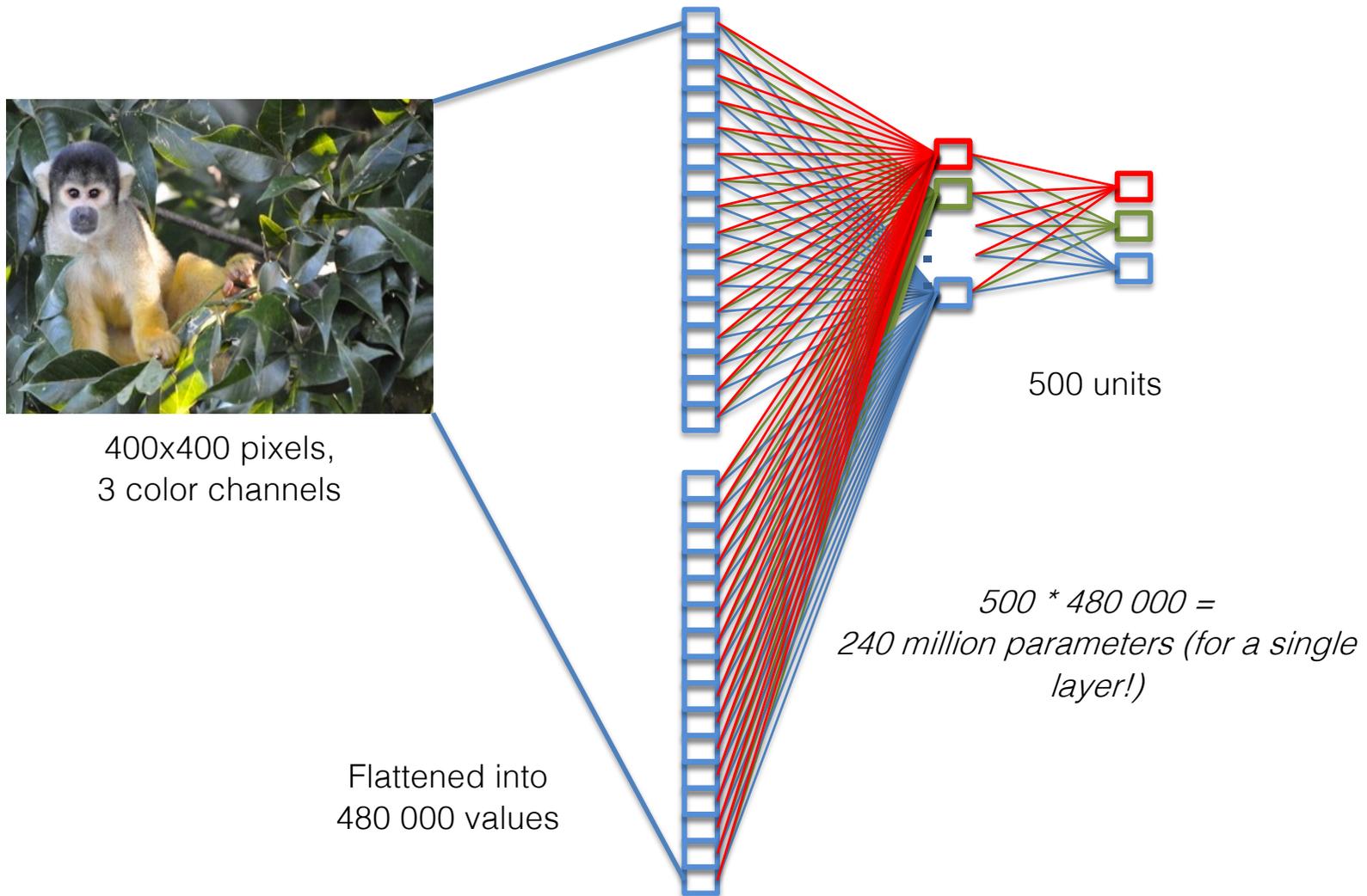
AI and Data Analytics

4.1 Shift equi/invariance and convolutions



Christian Wolf

Fully connected layers might be harmful



Fully connected layers might be harmful

Parameters learned for one part of the image do not generalize to other parts of the image.



\neq



Shift invariance

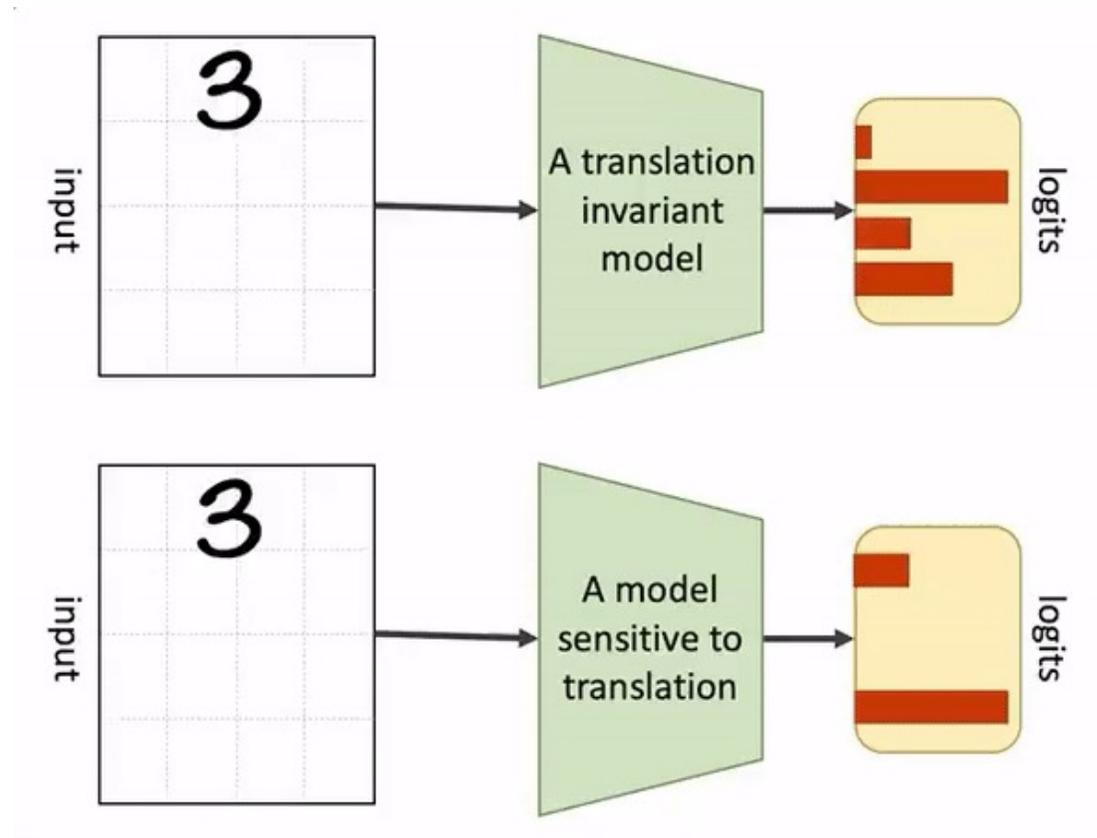
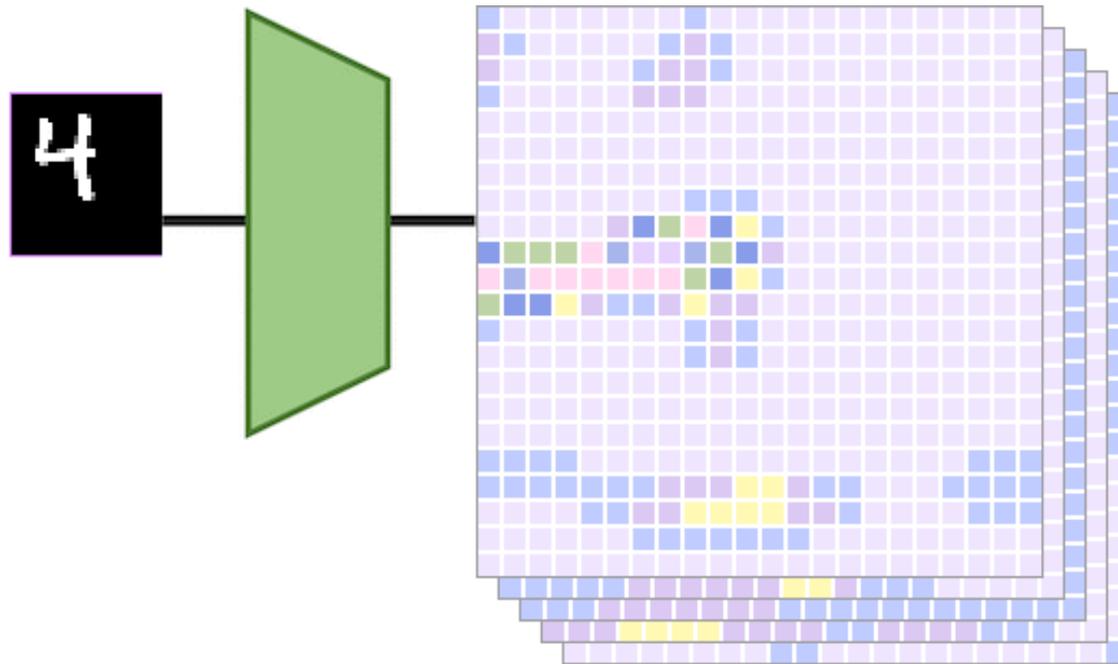


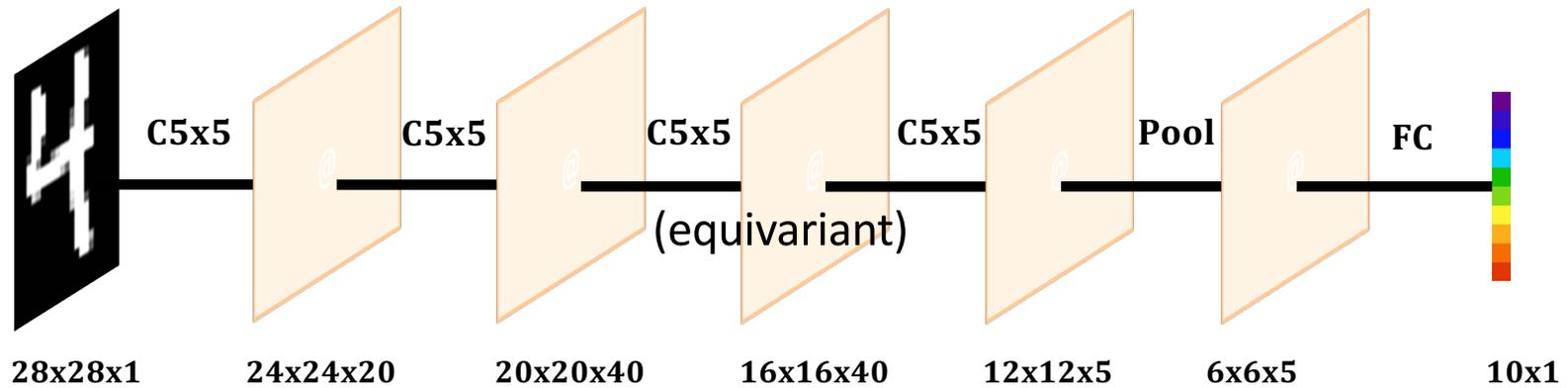
Figure by Samira Abnar (U. of Amsterdam)
<https://samiraabnar.github.io/>

Translation equivariance

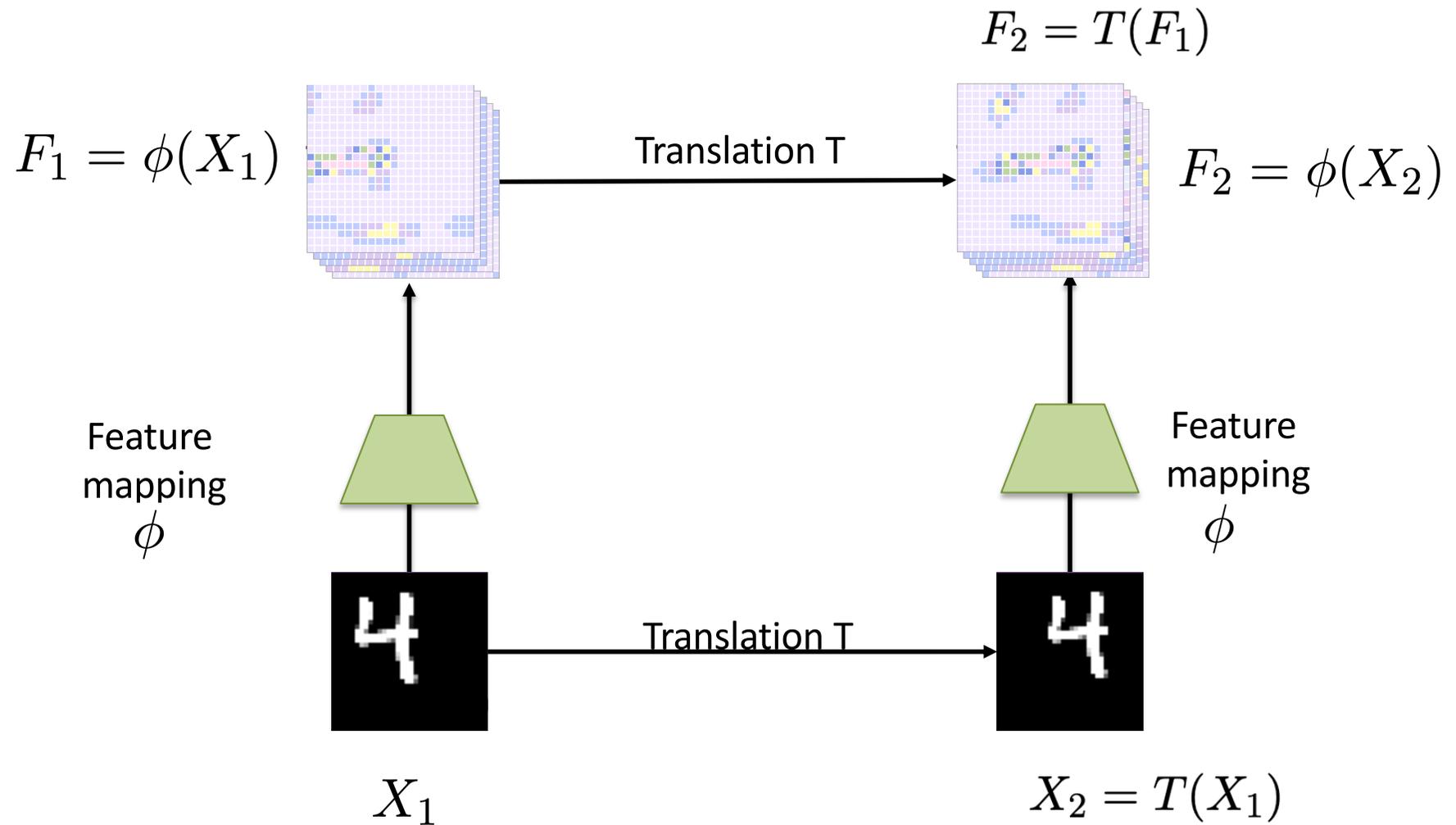
An equivariant mapping preserves the algebraic structure of the transformation.



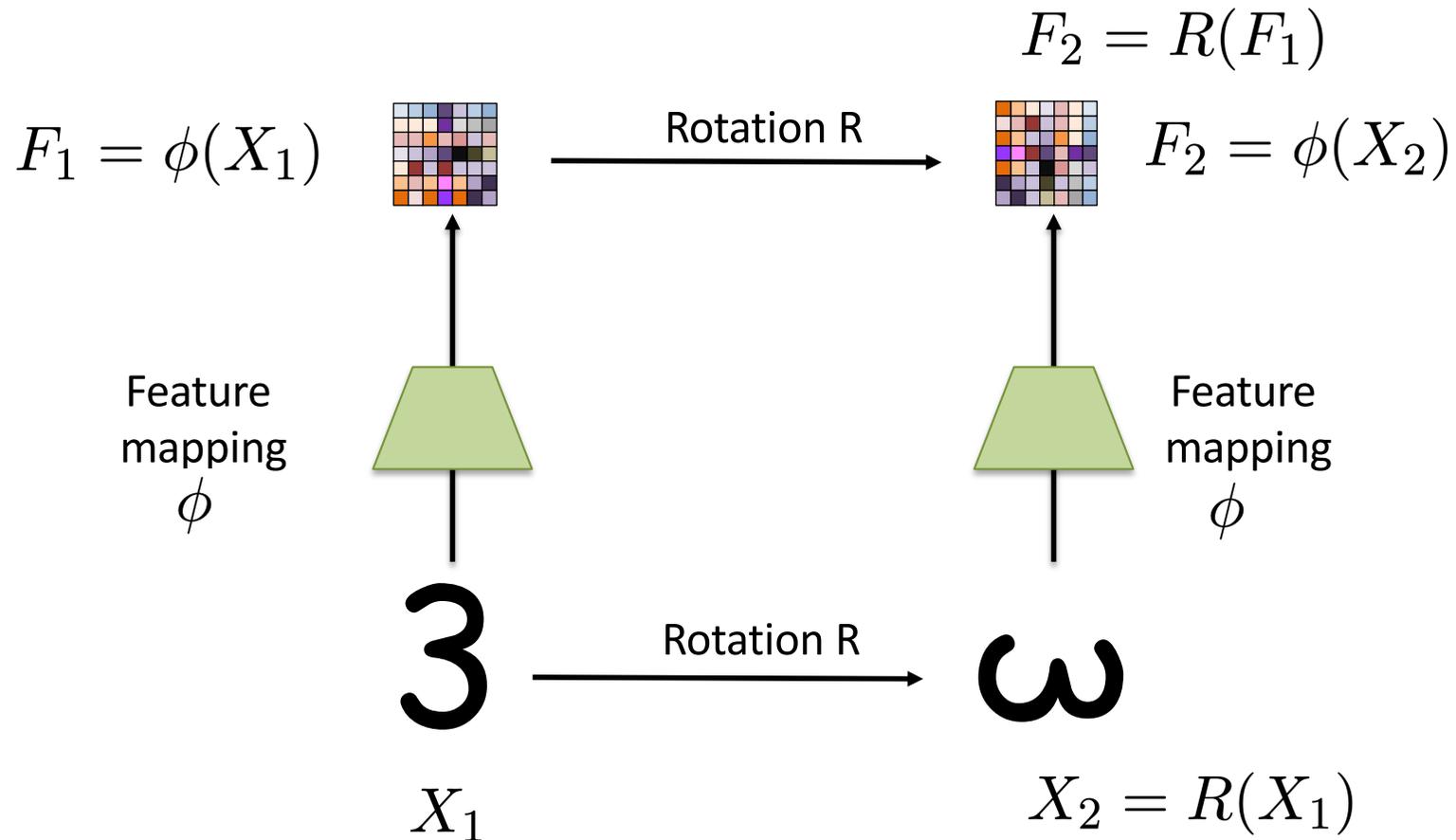
Translation Equivariance



Translation Equivariance



Rotation Equivariance



Shift-equivariant linear operations

Classical MLPs are sequences of linear layers followed by point-wise non-linearities. What kind of linear and shift-equivariant operators can there exist?

Let's consider an operator ϕ with following properties:

Linearity:

$$\phi(\alpha f + \beta f') = \alpha\phi(f) + \beta\phi(f')$$

Shift-equivariance:

$$\phi({}^{m,n}S(f)) = {}^{m,n}S(\phi(f))$$

where ${}^{m,n}S(f)$ shifts a signal by m,n .

Impulse response h :

$$h = \phi({}^{0,0}p)$$

where ${}^{0,0}p$ is a Dirac impulse centered at $0, 0$.

Shift-equivariant linear operations

We decompose the signal f into a series of Diracs ${}^{m,n}p$:

$$[\phi(f)](x, y) = \left[\phi \left(\sum_{m=-M/2}^{M/2} \sum_{n=-N/2}^{N/2} f(m, n) {}^{m,n}p \right) \right] (x, y)$$


 Dirac at position m, n

We use **linearity**:

$$= \left[\sum_{m=-M/2}^{M/2} \sum_{n=-N/2}^{N/2} f(m, n) \phi({}^{m,n}p) \right] (x, y)$$

We use **shift-equivariance**:

$$= \left[\sum_{m=-M/2}^{M/2} \sum_{n=-N/2}^{N/2} f(m, n) \phi({}^{m,n}S({}^{0,0}p)) \right] (x, y)$$

We use **linearity** again:

$$= \left[\sum_{m=-M/2}^{M/2} \sum_{n=-N/2}^{N/2} f(m, n) {}^{m,n}S(\phi({}^{0,0}p)) \right] (x, y)$$

Shift-equivariant linear operations

Change in notation: we replace $\phi^{(0,0)p}$ by h (through definition):

$$[\phi(f)](x, y) = \left[\sum_{m=-M/2}^{M/2} \sum_{n=-N/2}^{N/2} f(m, n)^{m,n} S(h) \right] (x, y)$$

Change in notation: make the shift-operator explicit:

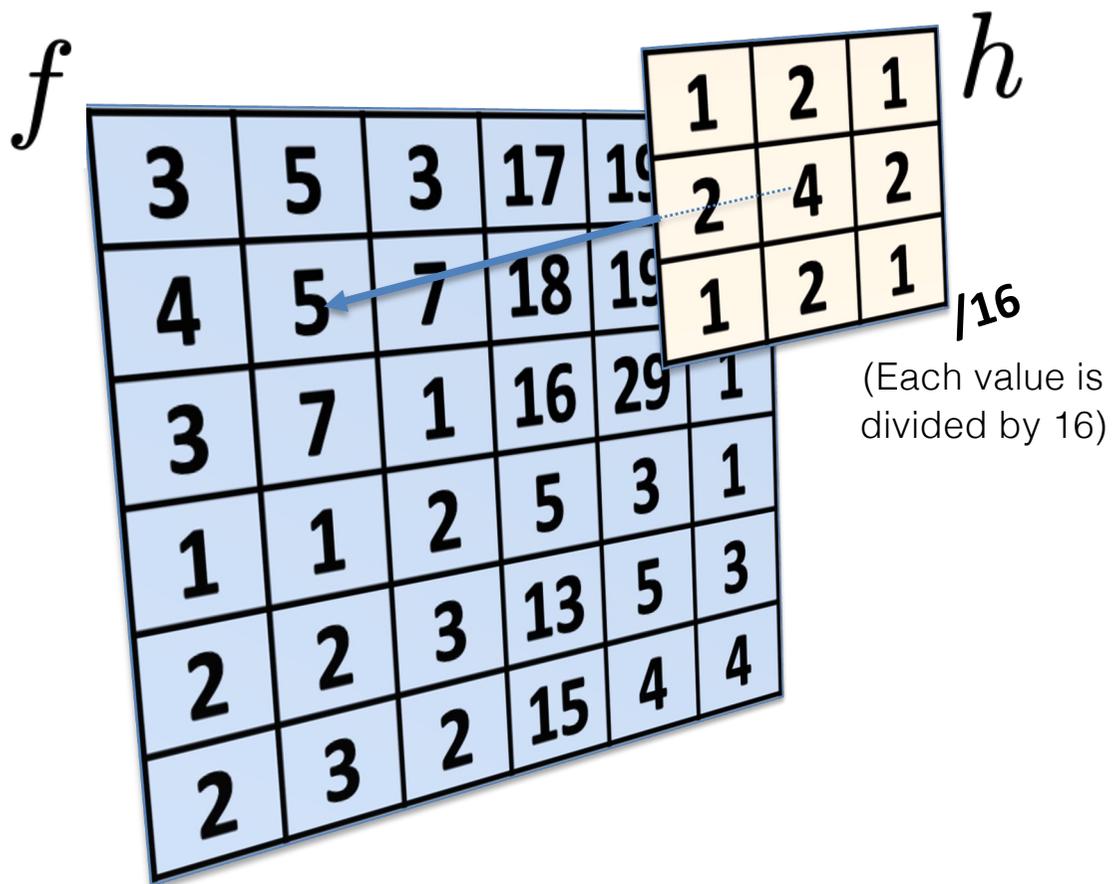
$$= \sum_{m=-M/2}^{M/2} \sum_{n=-N/2}^{N/2} f(m, n) h(x - m, y - n)$$

Change of variables:

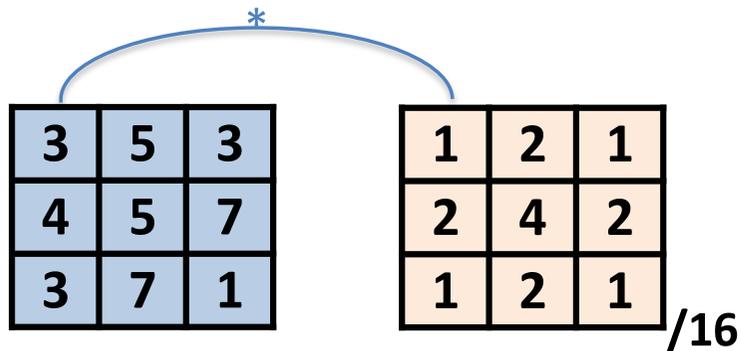
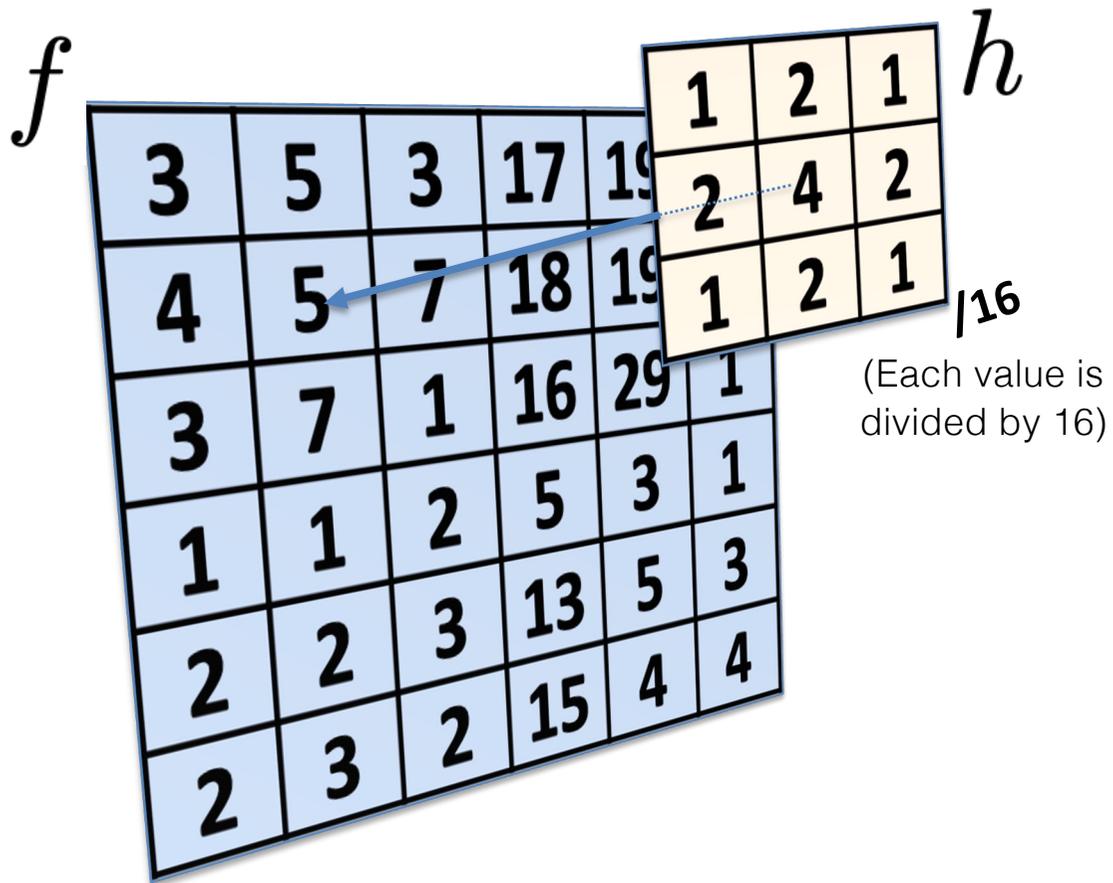
$$= \sum_{m'=-M/2}^{M/2} \sum_{n'=-N/2}^{N/2} f(x - m', y - n') h(m', n')$$

$(m' = x - m, n' = y - n) \Rightarrow$ **we get a convolution!**

Convolutions

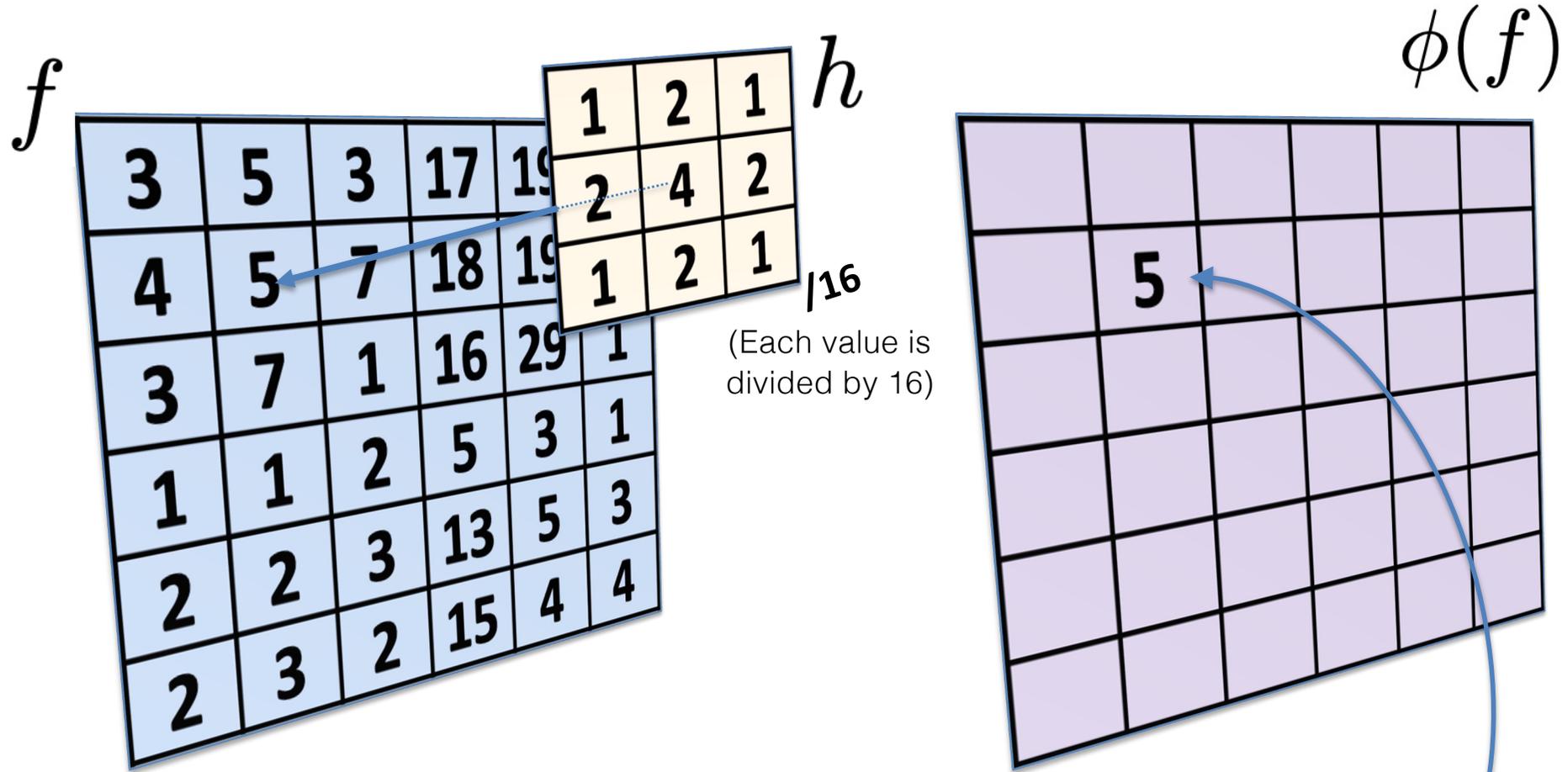


Convolutions

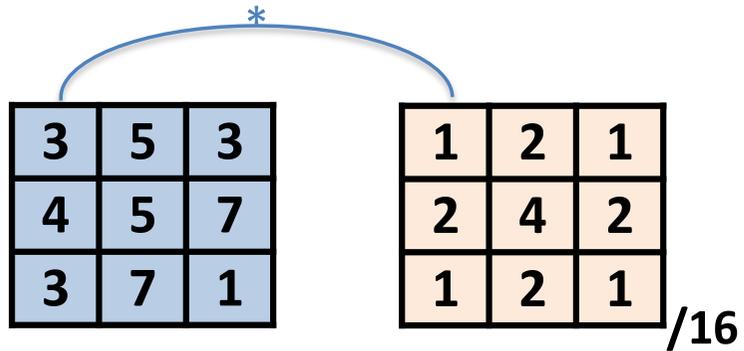


$$(3+2*5+3+2*4+5*5+2*7+3+2*7+1)/16 = 5$$

Convolutions

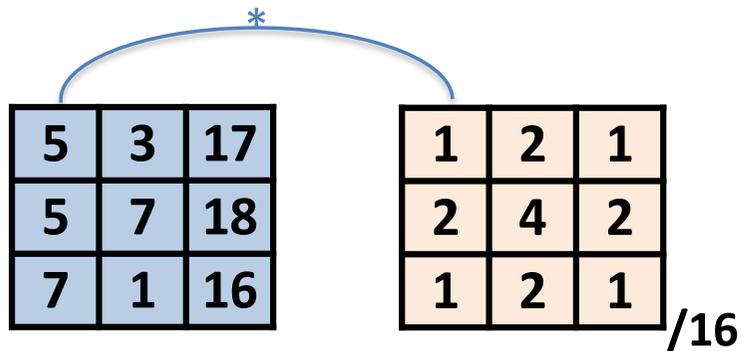
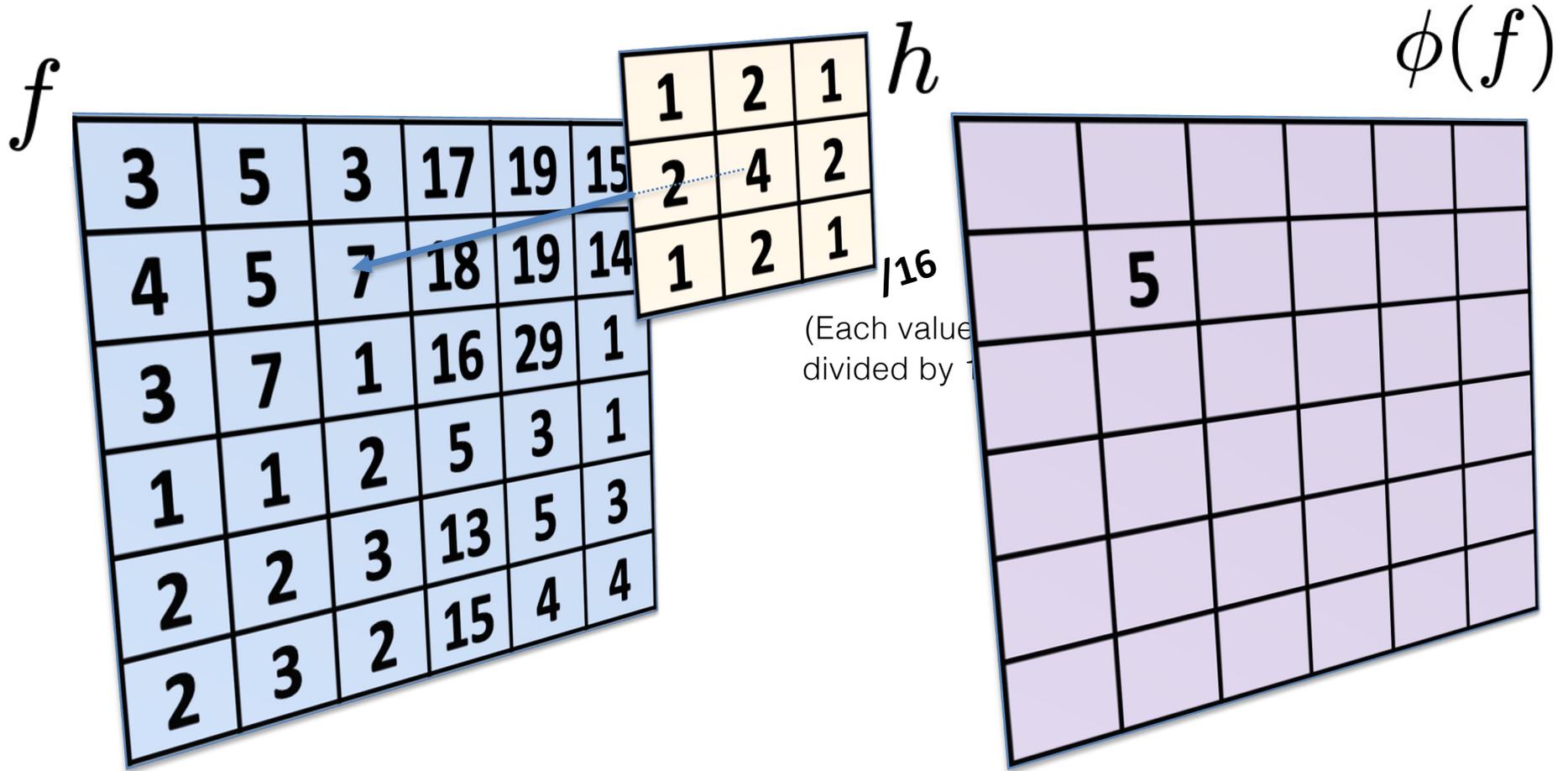


(Each value is divided by 16)

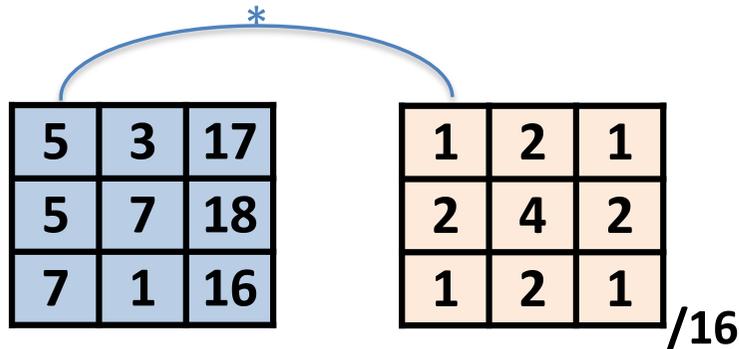
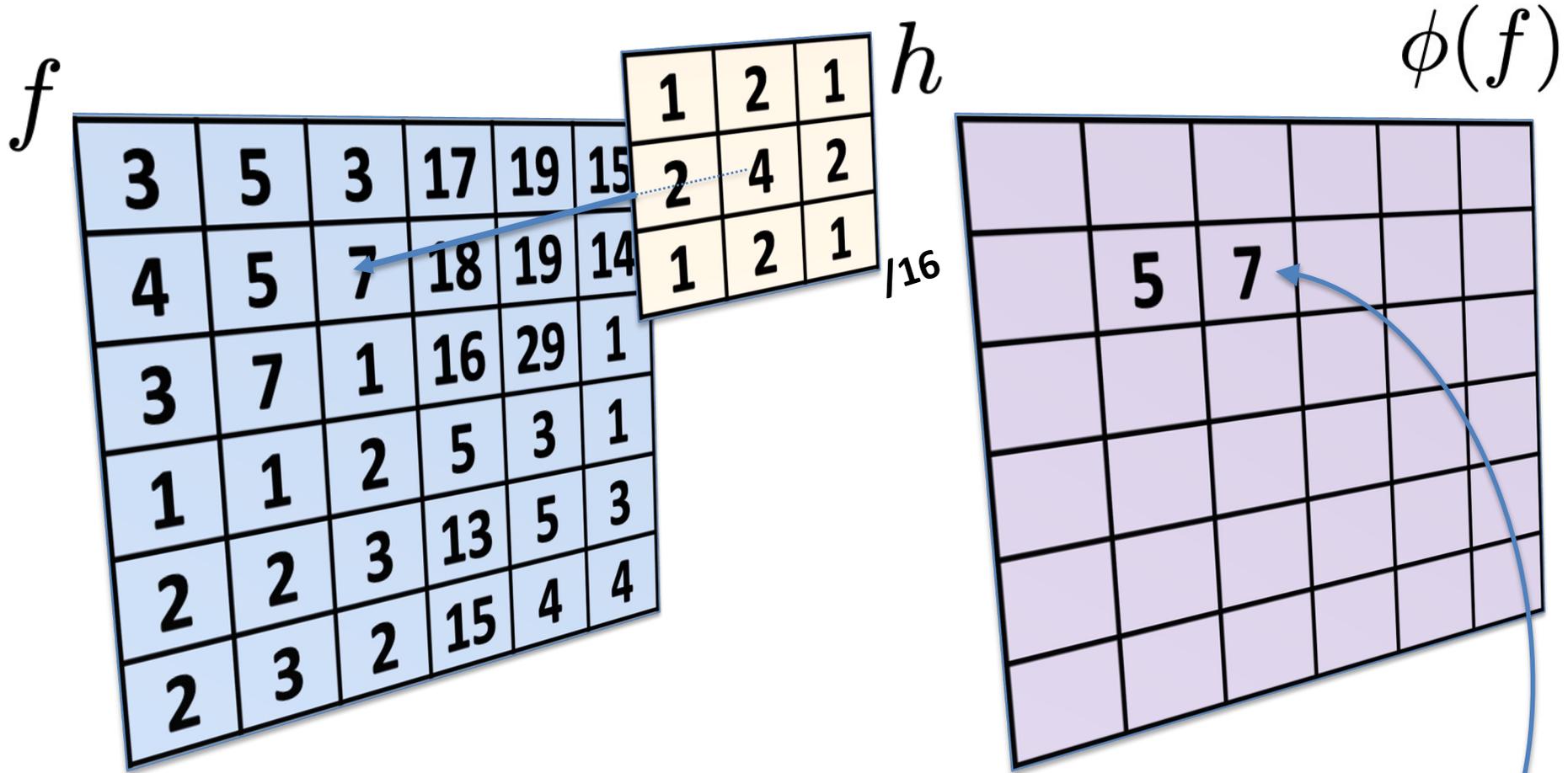


$$(3+2*5+3+2*4+5*5+2*7+3+2*7+1)/16 = 5$$

Convolutions

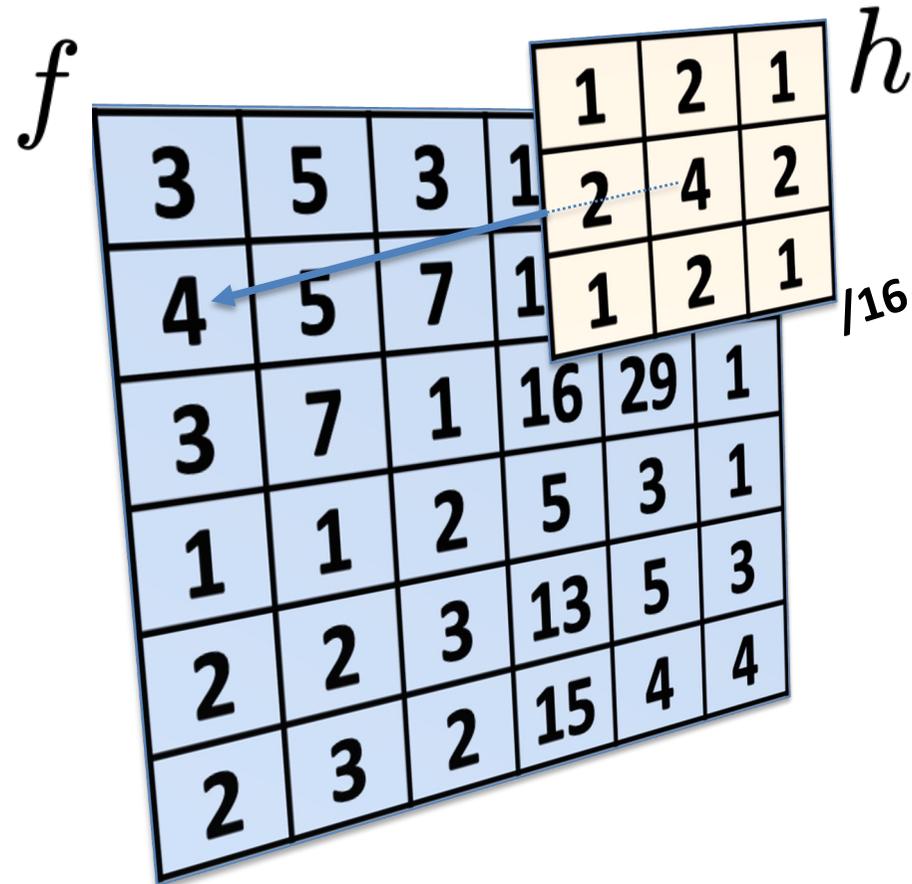


Convolutions

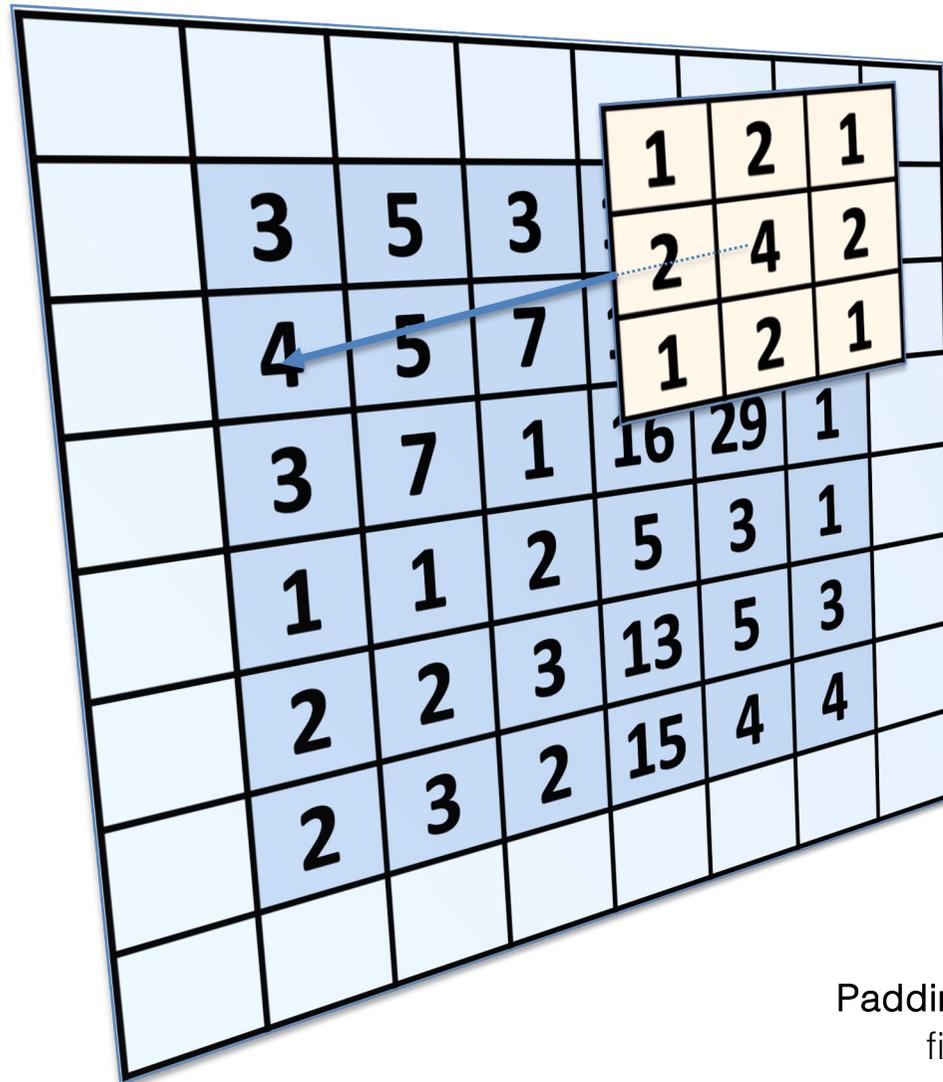


$$(5+2*3+17+2*5+4*7+2*18+7+2*1+16)/16 = 7$$

What happens on the borders?



Padding



Padding: create a virtual border filled with zero values

« Handcrafted « convolutions

Convolutions are standard operators in signal processing and image processing. In applications other than machine learning, their parameters (filter values) are not necessarily learned, but handcrafted.

Specifically chosen filters are frequently used for:

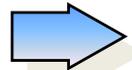
- Noise reduction
- Blurring / effect creation
- Deriving a signal
- Signal reconstruction (inverse filtering, Wiener filter)

The average filter

$$g = \frac{1}{4} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

$$g = \frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

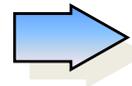
$$F(x,y) = f \otimes g \Leftrightarrow F(x,y) = \frac{1}{9} \sum_{u=-1}^{u=1} \sum_{v=-1}^{v=1} f(x+v, y+u)$$



Increasing the filter size

$$g = \frac{1}{25} \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 \end{bmatrix}$$

Bigger kernel = more blur



Computational complexity!

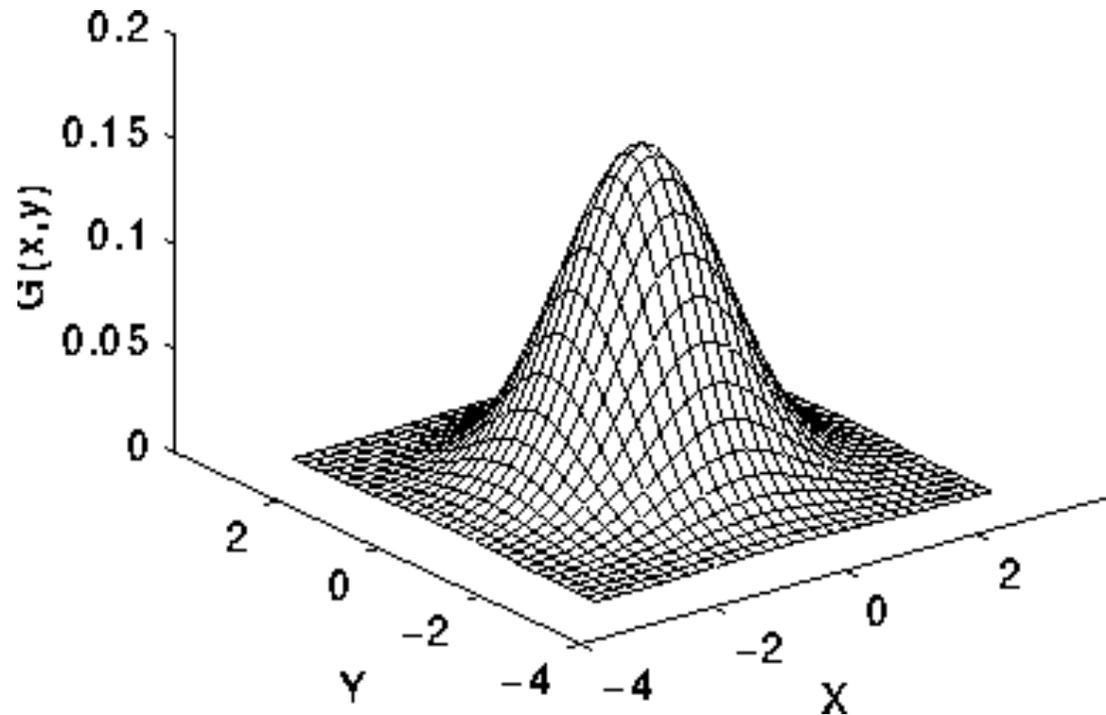
Gaussian filter

Discrete approximation of a 2D Gaussian filter:

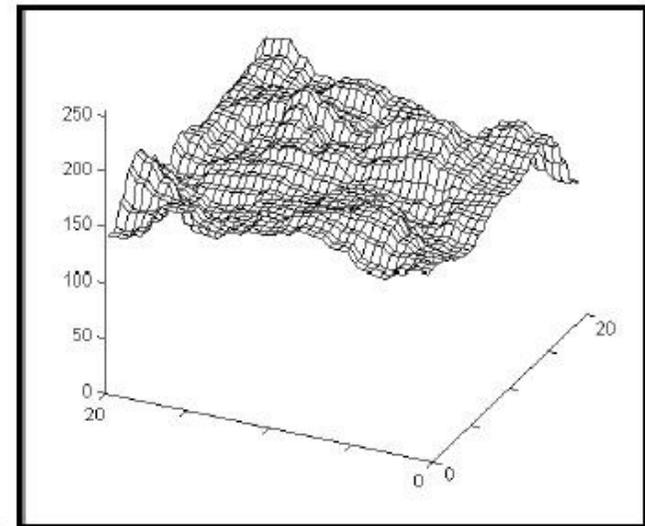
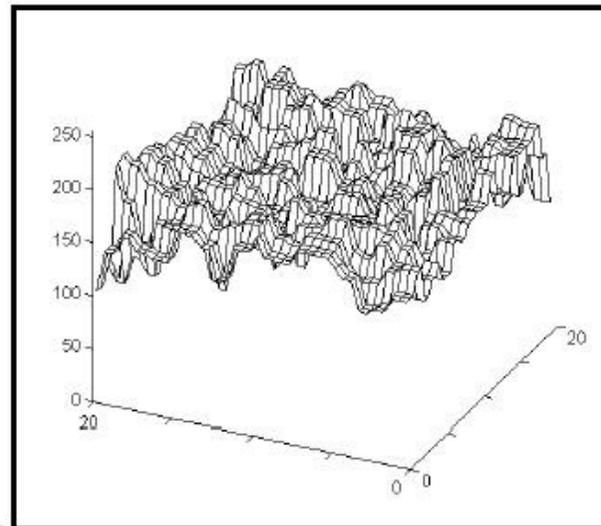
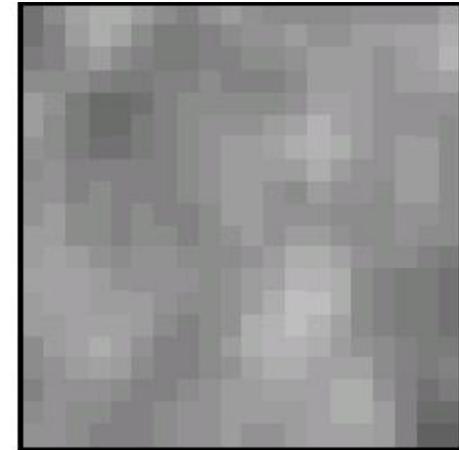
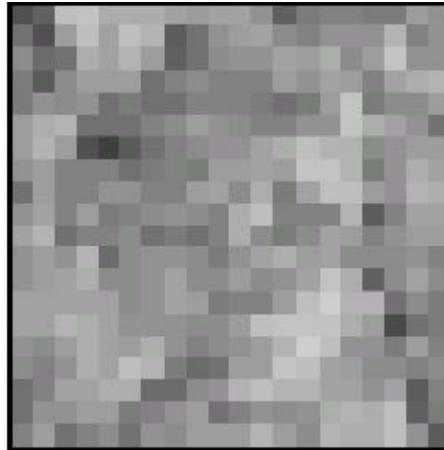
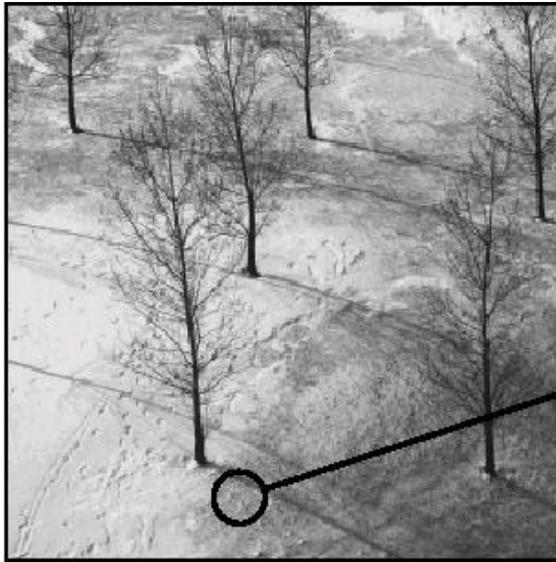
Weighted mean, the center pixels are more important than the border pixels

$$G(x, y) = \frac{1}{2\pi\sigma_x\sigma_y} e^{-\frac{1}{2}\left(\frac{x^2}{\sigma_x^2} + \frac{y^2}{\sigma_y^2}\right)}$$

$$g = \frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$$
$$g_5 = \frac{1}{246} \begin{bmatrix} 1 & 4 & 6 & 4 & 1 \\ 4 & 16 & 24 & 16 & 4 \\ 6 & 24 & 36 & 24 & 6 \\ 4 & 16 & 24 & 16 & 4 \\ 1 & 4 & 6 & 4 & 1 \end{bmatrix}$$



Gaussian low pass filter: details



Derive a signal: the Roberts filter

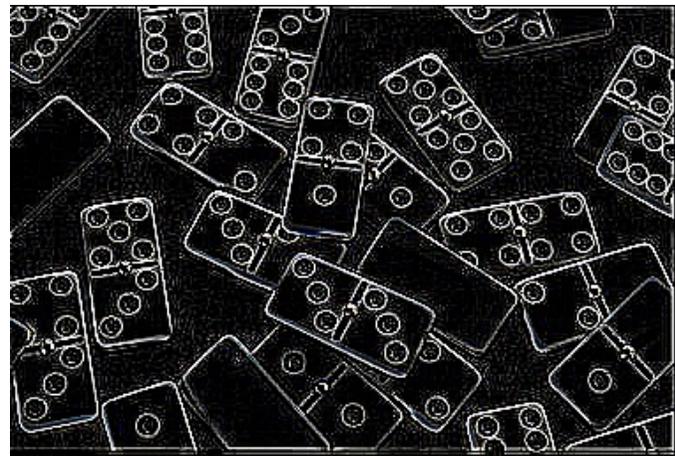
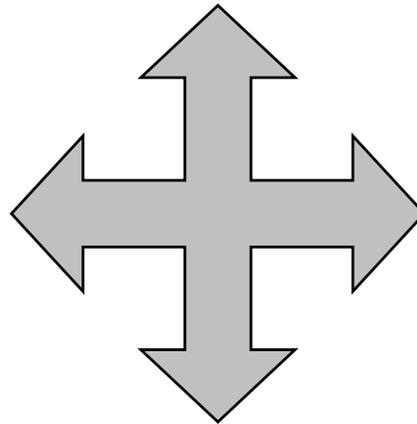
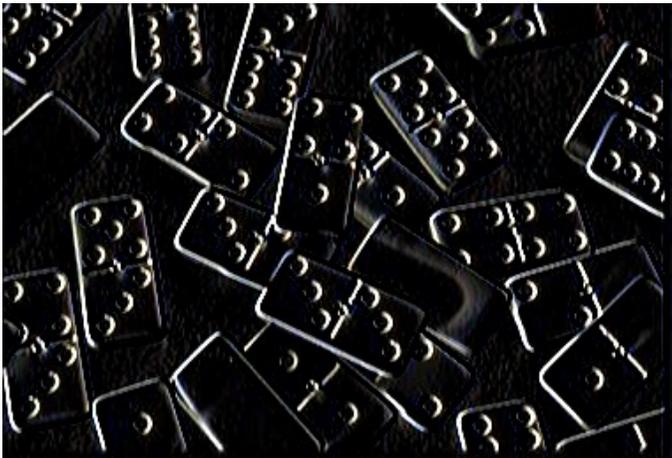
$$\frac{\partial f}{\partial x} \cong f \otimes [1 \quad -1] \quad \frac{\partial f}{\partial y} \cong f \otimes \begin{bmatrix} 1 \\ -1 \end{bmatrix} \quad \frac{\partial f}{\partial \vec{v}} \cong f \otimes \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}$$



$$\frac{\partial f}{\partial x} \cong f \otimes [1 \quad -1]$$



$$\frac{\partial f}{\partial y} \cong f \otimes \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$



$$\|\nabla f\| = \sqrt{\left(\frac{\partial f}{\partial x}\right)^2 + \left(\frac{\partial f}{\partial y}\right)^2}$$

Derive a signal: the Sobel filter

$$\frac{\partial f}{\partial x} \cong f \otimes \begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix} \quad \frac{\partial f}{\partial y} \cong f \otimes \begin{bmatrix} 1 & 2 & 1 \\ 0 & 0 & 0 \\ -1 & -2 & -1 \end{bmatrix}$$

$$\frac{\partial f}{\partial v1} \cong f \otimes \begin{bmatrix} 2 & 1 & 0 \\ 1 & 0 & -1 \\ 0 & -1 & -2 \end{bmatrix} \quad \frac{\partial f}{\partial v2} \cong f \otimes \begin{bmatrix} 0 & 1 & 2 \\ -1 & 0 & 1 \\ -2 & -1 & 0 \end{bmatrix}$$



Sharpen

Sharpened image = (A-1) original image + high pass filtered image

$$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 4 & -1 \\ 0 & -1 & 0 \end{bmatrix} \quad \begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix} \quad \text{Laplacian filter (high pass filter)}$$

A=1: standard Laplacian filter

A>1: a part of the original image is added to the high pass



A=2

$$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 9 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$

Noise!!!

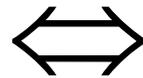
Iterative application

Iteratively convolving an image N times by a filter of size M corresponds to the application of a single convolution of a filter of size $M^* = 2 \times ((M-1) / 2) \times N + 1$

2 x a filter of size **3** = 1 x filter of size **5**
3 x a filter of size **3** = 1 x filter of size **7**

$$(f \otimes h) \otimes g = f \otimes (h \otimes g)$$

$$2 \times g_3 = \frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$$



$$g_5 = \frac{1}{256} \begin{bmatrix} 1 & 4 & 6 & 4 & 1 \\ 4 & 16 & 24 & 16 & 4 \\ 6 & 24 & 36 & 24 & 6 \\ 4 & 16 & 24 & 16 & 4 \\ 1 & 4 & 6 & 4 & 1 \end{bmatrix}$$

Convolutions in PyTorch

```
1 import torch
2 from torch.nn import functional as F
3
4 # Create a binary random tensor 5x5 tensor
5 # don't forget batch & channel dimensions
6 f = torch.rand((1,1,5,5))
7 f[f>0.5]=1
8 f[f<0.5]=0
9 print (f)
```

```
1 tensor([[[[1., 1., 1., 0., 0.],
2          [0., 0., 0., 0., 0.],
3          [0., 0., 1., 0., 0.],
4          [0., 0., 0., 1., 0.],
5          [0., 0., 1., 0., 0.]]]]])
```

```
1 # Create an average filter (all ones)
2 h = torch.ones((1,1,3,3))
3 print (h)
```

```
1 tensor([[[[1., 1., 1.],
2          [1., 1., 1.],
3          [1., 1., 1.]]]]])
```

Convolutions in PyTorch

```
1 print (f)
```

```
1 tensor([[[[1., 1., 1., 0., 0.],  
2         [0., 0., 0., 0., 0.],  
3         [0., 0., 1., 0., 0.],  
4         [0., 0., 0., 1., 0.],  
5         [0., 0., 1., 0., 0.]]]]])
```

```
1 # Convolve  
2 F.conv2d (f, h)
```

```
1 tensor([[[[4., 3., 2.],  
2         [1., 2., 2.],  
3         [2., 3., 3.]]]]])
```

The output is smaller: we lost 2 border rows and 2 border columns because we did not explicitly state that we want padding.

Convolutions in PyTorch

```
1 print (f)
```

```
1 tensor([[[[1., 1., 1., 0., 0.],  
2         [0., 0., 0., 0., 0.],  
3         [0., 0., 1., 0., 0.],  
4         [0., 0., 0., 1., 0.],  
5         [0., 0., 1., 0., 0.]]]]])
```

```
1 # Convolve with padding  
2 F.conv2d (f, h, padding=1)
```

```
1 tensor([[[[2., 3., 2., 1., 0.],  
2         [2., 4., 3., 2., 0.],  
3         [0., 1., 2., 2., 1.],  
4         [0., 2., 3., 3., 1.],  
5         [0., 1., 2., 2., 1.]]]]])
```

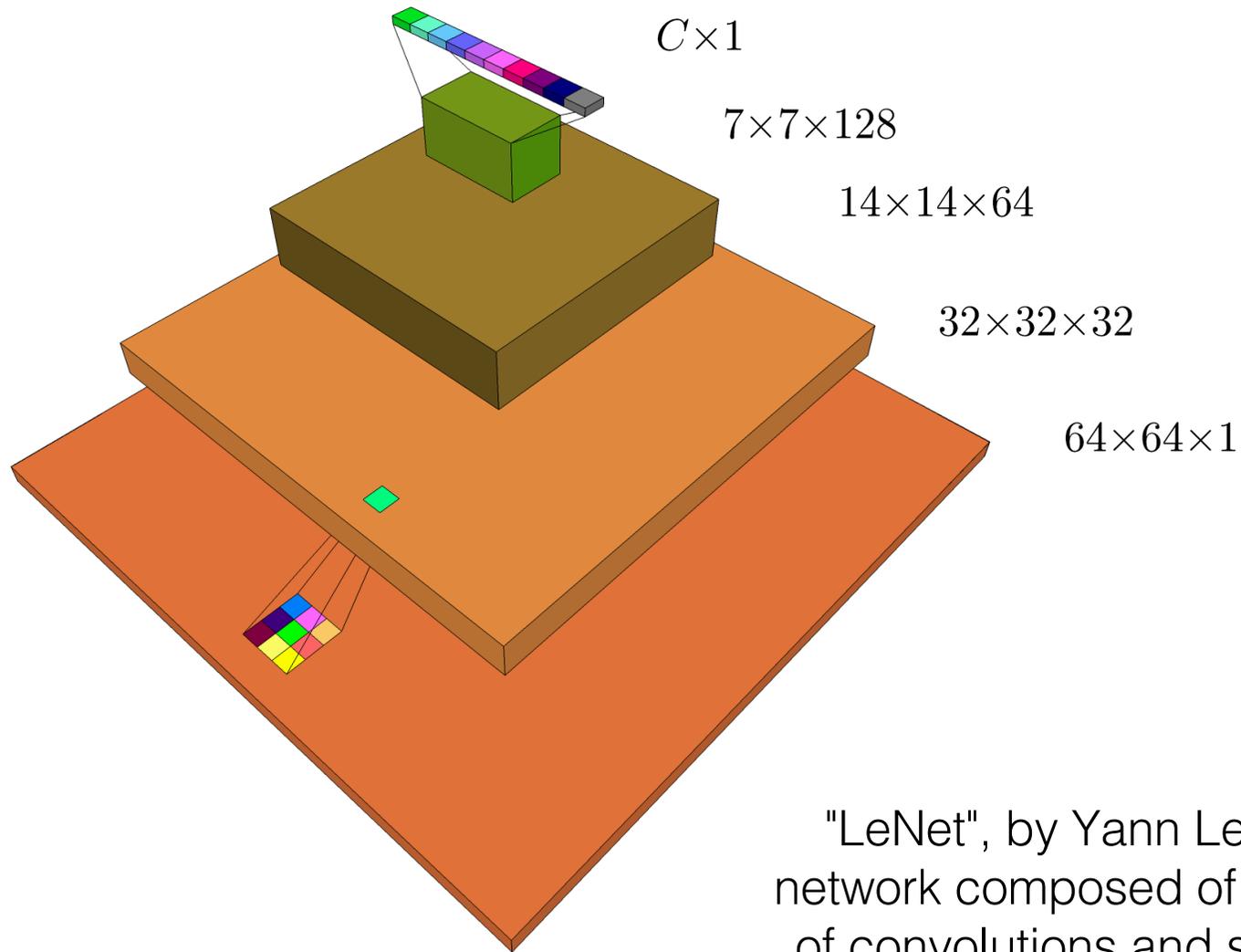
Learned convolutions

In Deep Learning, we use convolutions in the form of convolutional layers in neural networks.

They structure neural networks (“inductive bias”):

- Input layers can be multidimensional tensors (2D images, 3D videos) and are (usually) **not** flattened into vectors.
- Activations keep the structure: hidden layers can be spatially, temporally, or spatio-temporally structured.

« LeNet »



"LeNet", by Yann LeCun, is a network composed of sequences of convolutions and spatial size reductions (« pooling »).

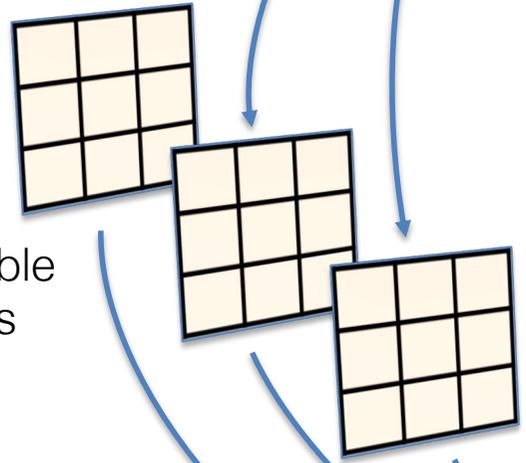
N (layer) input channels

3	5	3	17	19	15
4	5	7	18	19	14
3	7	1	16	29	1
1	1	2	5	3	1
2	2	3	13	5	3
2	3	2	15	4	4

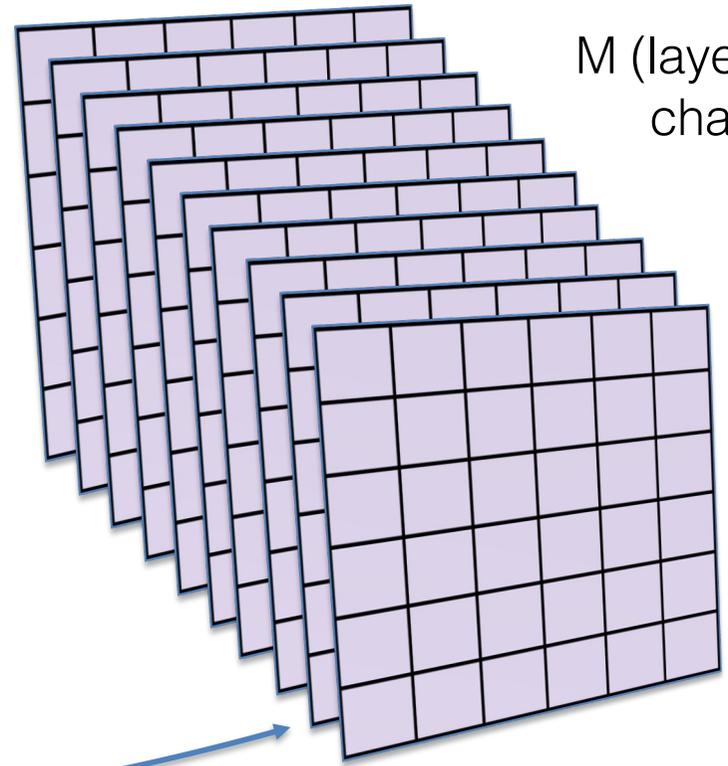
Each of the M output channels is a linear combination of the N filtered input images.

We need N*M filters!

Trainable filters



M (layer) output channels



Σ

+bias

Or, in equations ...

In matrix/vector notation:

$$y_i = \sum_j f_j * h_{ij} + b_i$$

- * is the convolution operation
- y_i is output channel i (a 2D tensor for each sample of the batch)
- h_{ij} is filter j for output channel i (a 2D tensor for each sample of the batch)
- b is a bias vector (one value for each output channel)

Summary/recall: tensor dimensions

Recall: PyTorch functions operate on multi-dimensional tensors and follow conventions on the order of dimensions.

- The first dimension is the batch dimension
 - ▶ Use 1 if you don't use batches (= batches of size 1).
 - ▶ Losses are reduced (sum or mean) over samples in a batch
 - ▶ We don't convolve over batches!
- the second dimension is the channel dimension
 - ▶ Use 1 if you don't use channels (= single channels).
 - ▶ We don't convolve over channels!
 - ▶ Each output channel receives filtered responses of each input channel.
- the following dimensions are application dependant:
 - ▶ e.g. rows, columns in images.
 - ▶ Convolutions can be done over these dimensions.
 - ▶ We can convolve over sub-dimensions (eg 1D convolutions of image rows).

Max pooling

Pooling reduces the spatial resolution of a tensor, combining multiple values into a single one.

3	5	3	17	19	15
4	5	7	18	19	14
3	7	1	16	29	1
1	1	2	5	3	1
2	2	3	13	5	3
2	3	2	15	4	4

7	29
3	15

Max pooling: take the maximum of a HxW neighborhood.

4.3	16.4
1.88	5.88

Average pooling

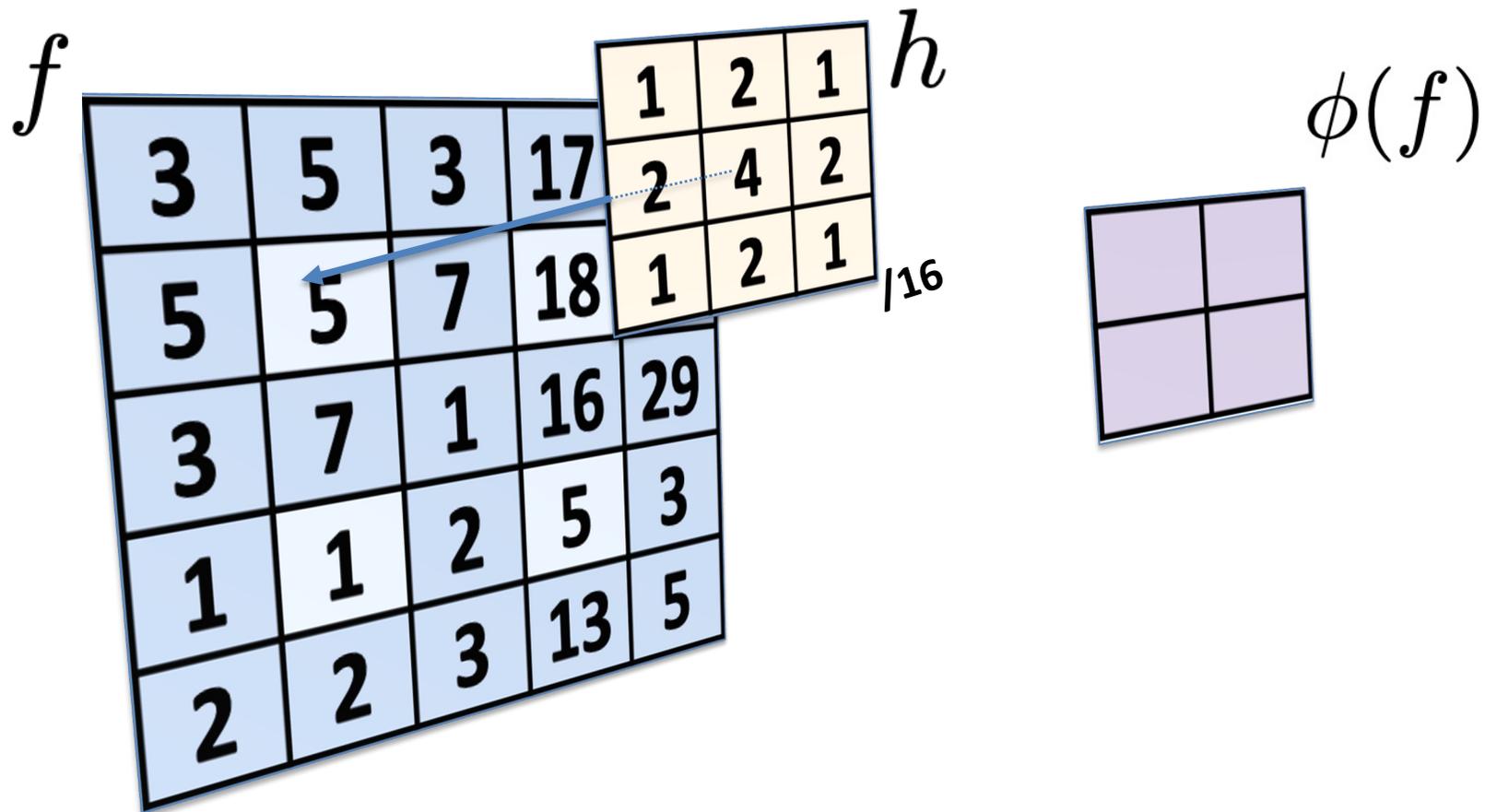
Example: kernel=3x3

Convolutions + pooling ...

- Convolutions introduce shift-equivariance
- Convolutions + pooling create shift-invariant models

Convolutions with stride

The stride parameter allows to move the kernel in multiples. Shown below: stride = 2 (and no padding).



LeNet in PyTorch

Lecture code:
lenet.py

```
1 class LeNet(torch.nn.Module):
2     def __init__(self):
3         super(LeNet, self).__init__()
4         # 1 input channel, 20 output channels,
5         # 5x5 filters, stride=1, no padding
6         self.conv1 = torch.nn.Conv2d(1, 20, 5, 1, 0)
7         self.conv2 = torch.nn.Conv2d(20, 50, 5, 1, 0)
8         self.fc1 = torch.nn.Linear(4*4*50, 500)
9         self.fc2 = torch.nn.Linear(500, 10)
10
11     def forward(self, x):
12         x = F.relu(self.conv1(x))
13         # Max pooling with a filter size of 2x2
14         # and a stride of 2
15         x = F.max_pool2d(x, 2, 2)
16         x = F.relu(self.conv2(x))
17         x = F.max_pool2d(x, 2, 2)
18         x = x.view(-1, 4*4*50)
19         x = F.relu(self.fc1(x))
20         return self.fc2(x)
```

The Conv2D function

Description of the `torch.nn.Conv2d` method:

```
1 torch.nn.Conv2d(  
2     in_channels,           # number of input channels  
3     out_channels,        # number of input channels  
4     kernel_size,        # scalar=square, 5 means 5x5  
5                           # or tuple eg (5,3)  
6     stride=1,           # stride (scalar or tuple)  
7     padding=0,         # padding (scalar or tuple)  
8     dilation=1,  
9     groups=1,  
10    bias=True,          # add bias, True or False  
11    padding_mode='zeros') # zeroes, replicate, reflect
```

Conv2D: functional vs. NN layer

There are two different functions:

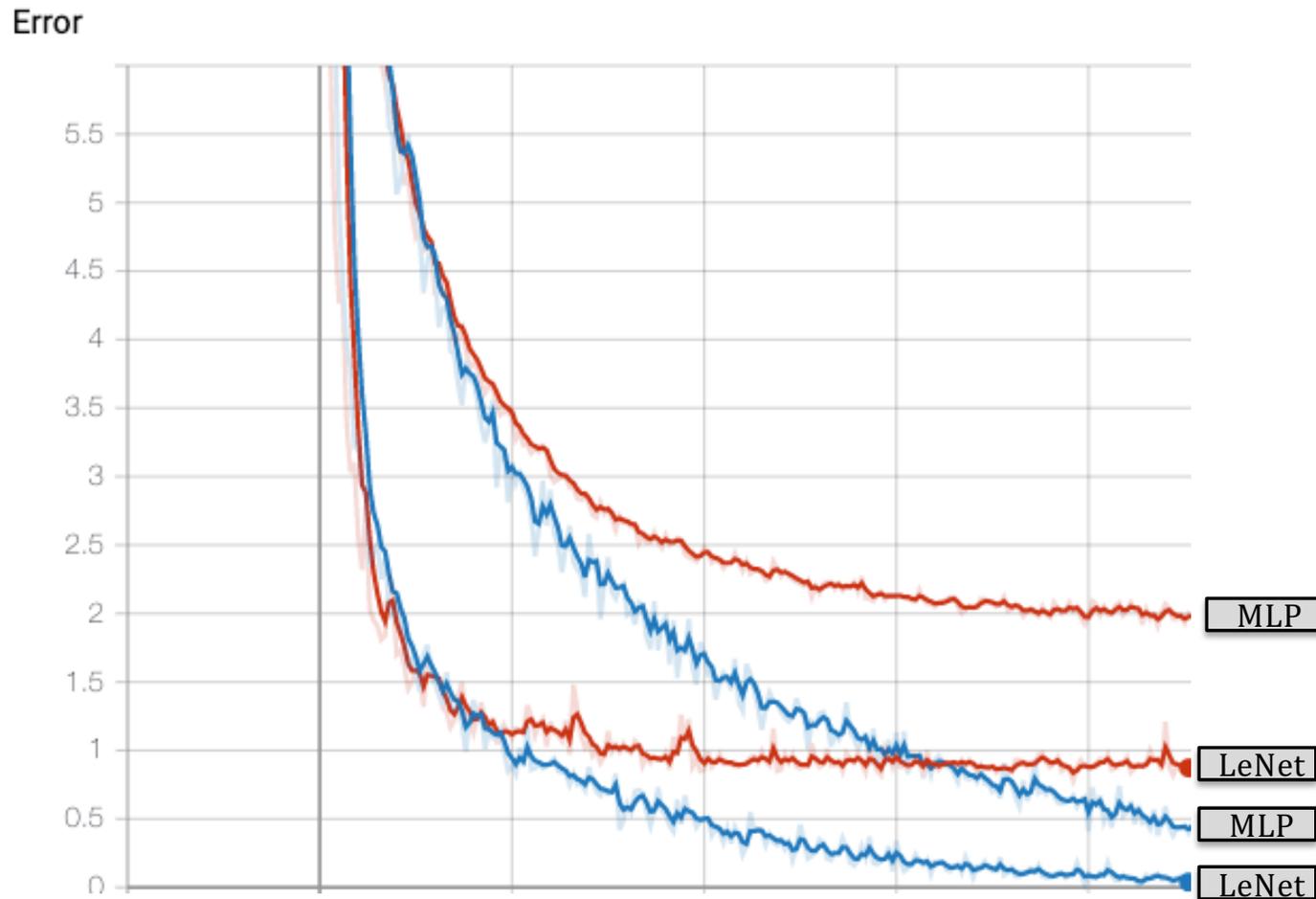
- `torch.nn.functional.conv2d` is a function which takes two tensors and returns the convolved result tensor:

```
1 from torch.nn import functional as F
2 (...)
3 out = F.conv2d (f, h)
```

- `torch.nn.Conv2d` is a method which creates a 2D convolution layer with trainable filters and eventually a bias. The result is the configured convolution function:

```
1 import torch
2 (...)
3 c = torch.nn.Conv2d(1, 20, 5, 1, 0)
4 out = c(f)
```

MLP vs LeNet on MNIST



Blue: training, Red: validation

Template matching

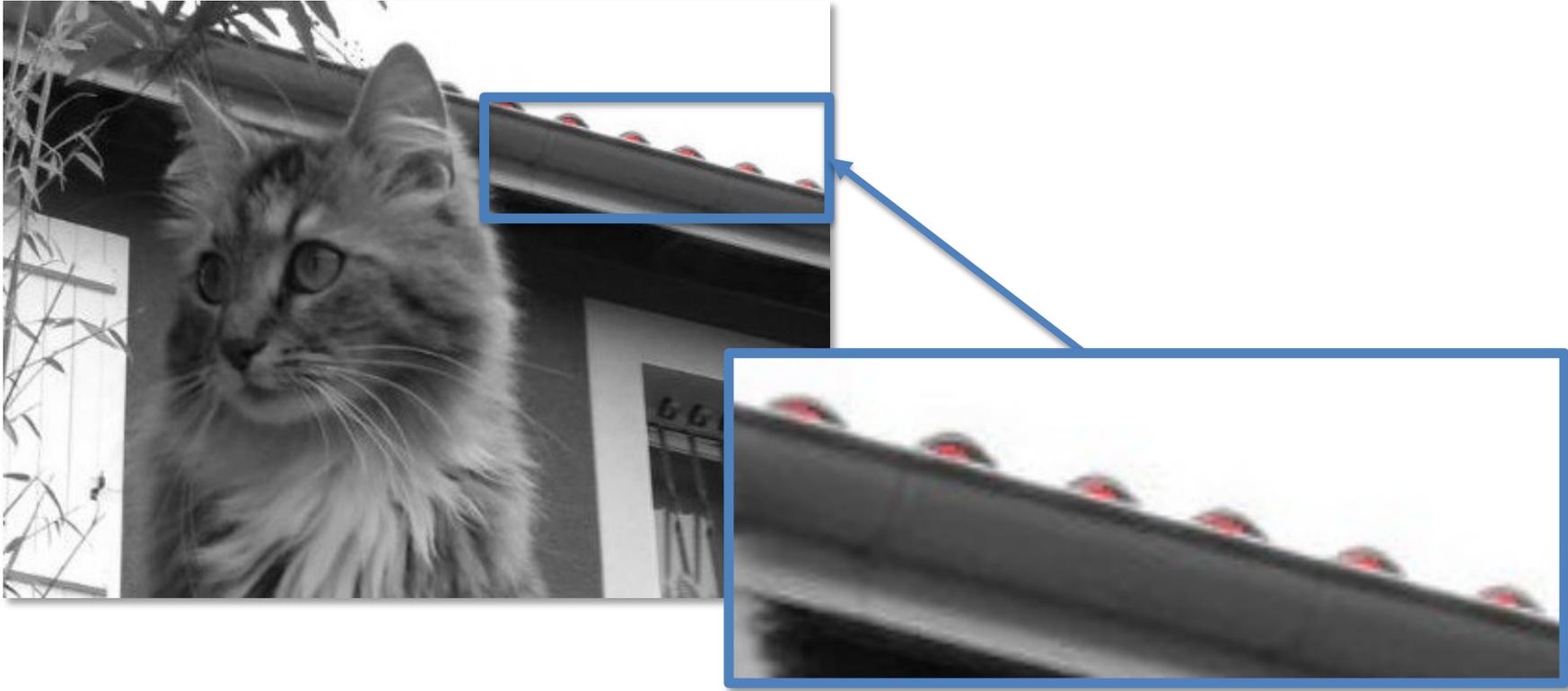
Input image:



Filter kernel

- We normalize input and kernel (subtract mean, divide by standard deviation).
- We convolve the input with the kernel
- We threshold response and superimpose it on the red channel.

Template matching



Convolutions can match patterns:

- either in input space
- or in intermediate network layers (which still have a spatial organization).