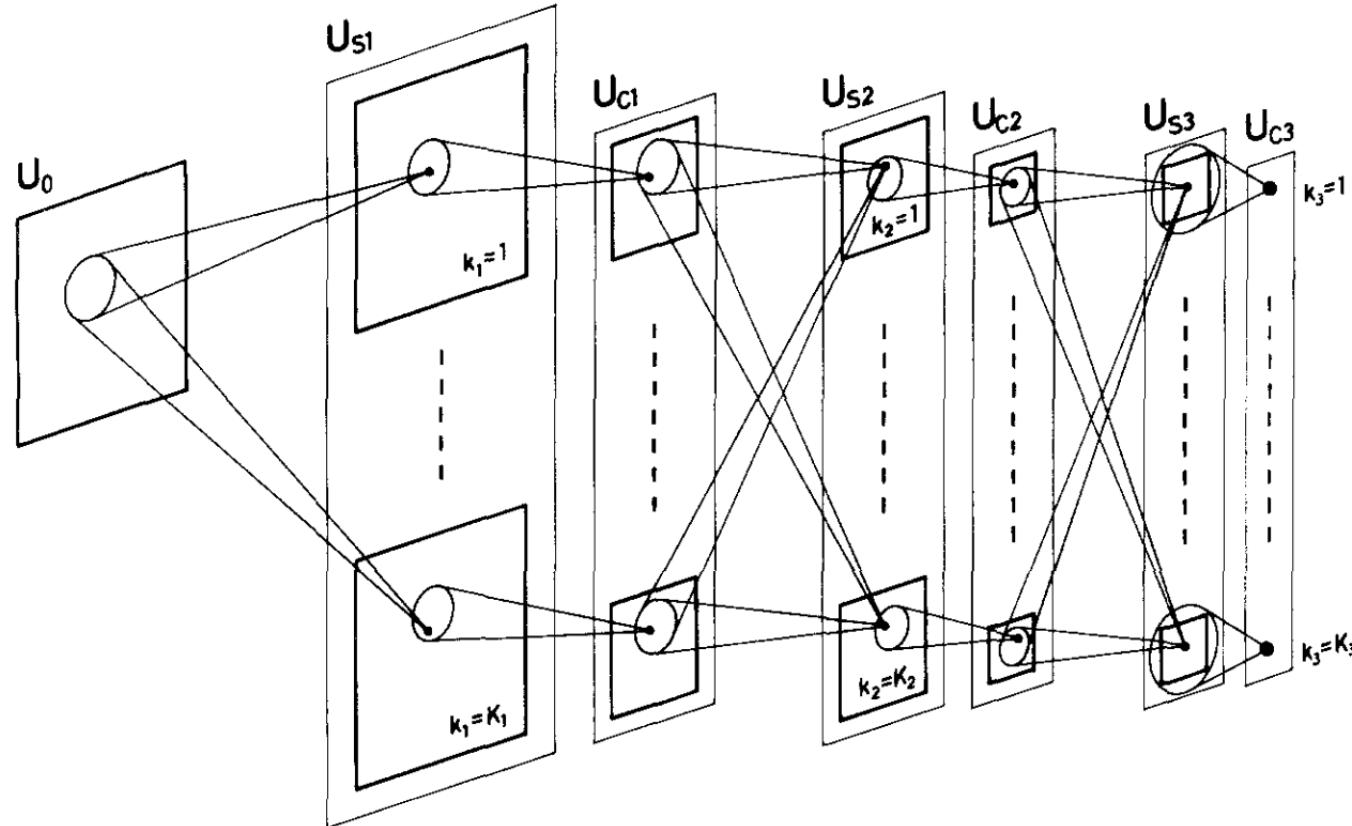


Lecture 9

Convolutional Neural Networks



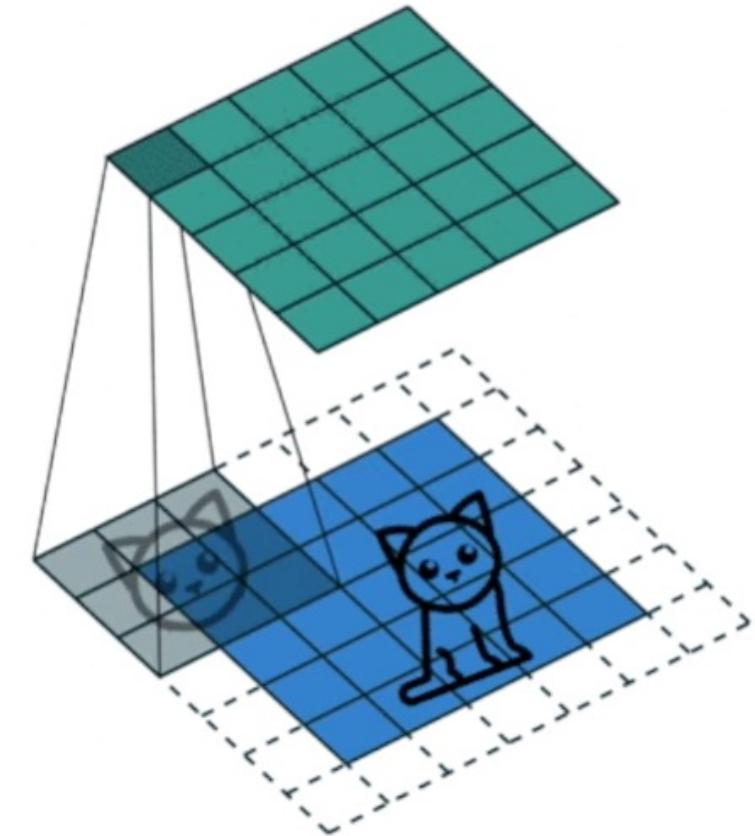
Overview

- Convolution
- Convolutional Neural Networks
- Case Study: LeNet, AlexNet, VGG
- Visualization

Convolution

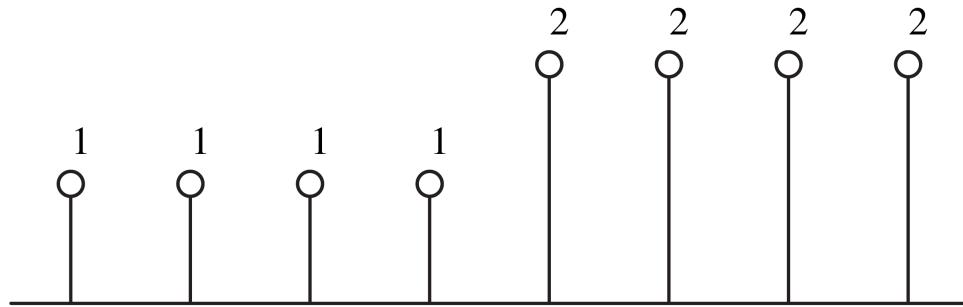
Convolution is template matching ...

- with a sliding window
- abstract templates
- similarity measured by dot product
- stronger activation, better matching



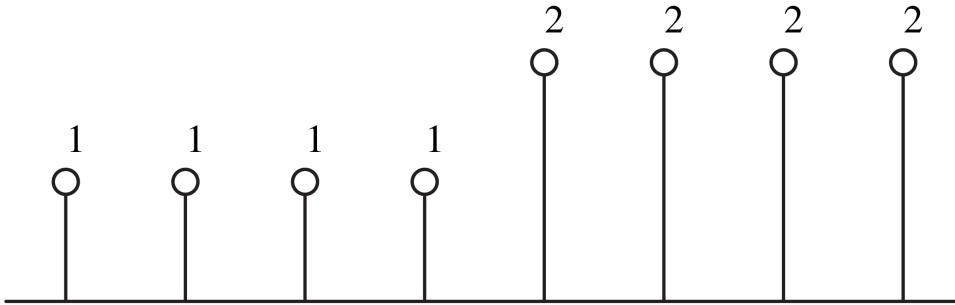
Convolution: a 1-D example

input

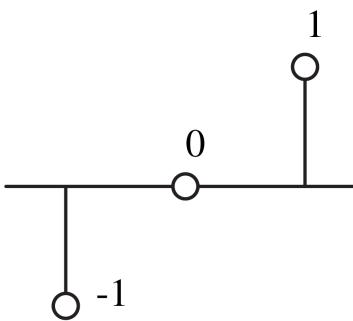


Convolution: a 1-D example

input



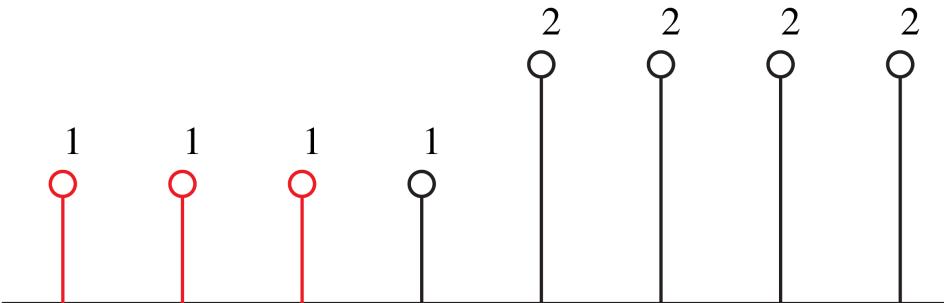
filter



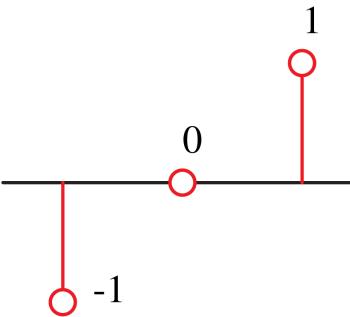
Convolution: a 1-D example

- sliding window
- dot product

input



filter



$$-1 \times 1 + 0 \times 1 + 1 \times 1$$

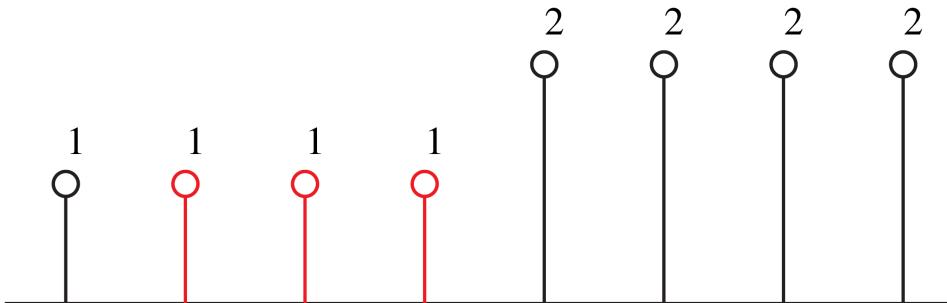
output



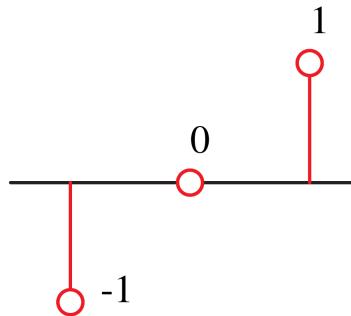
Convolution: a 1-D example

- sliding window
- dot product

input



filter



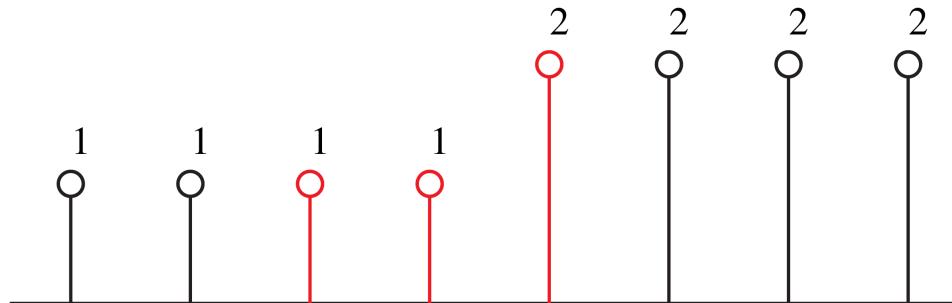
output



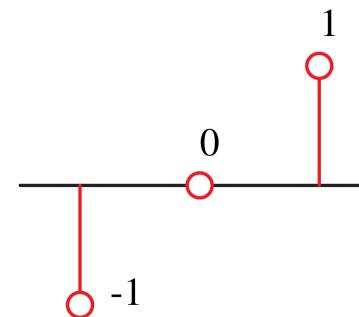
Convolution: a 1-D example

- sliding window
- dot product

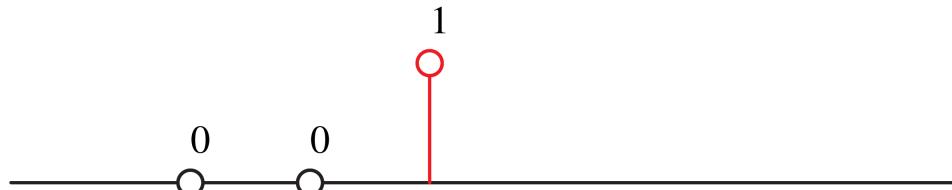
input



filter



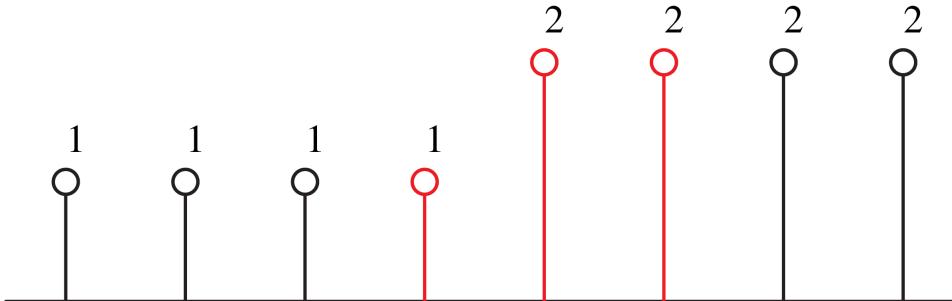
output



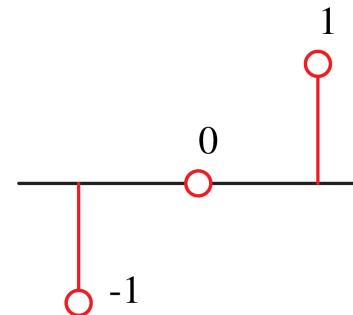
Convolution: a 1-D example

- sliding window
- dot product

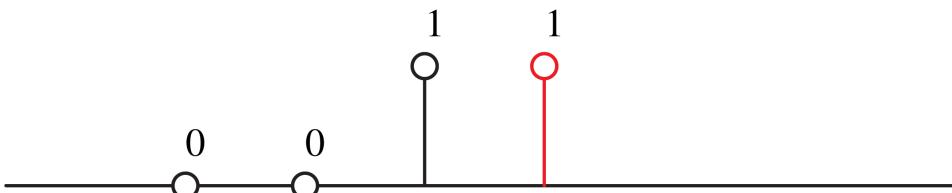
input



filter



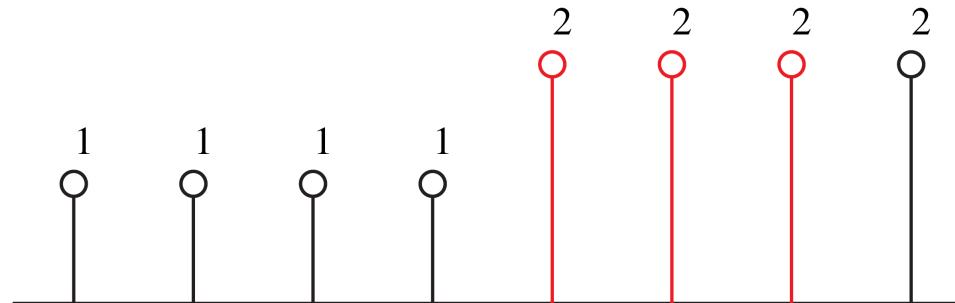
output



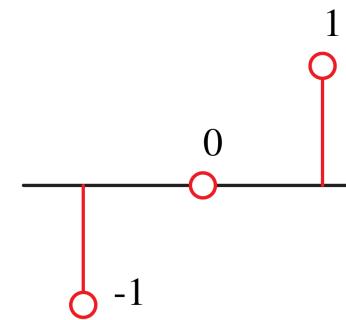
Convolution: a 1-D example

- sliding window
- dot product

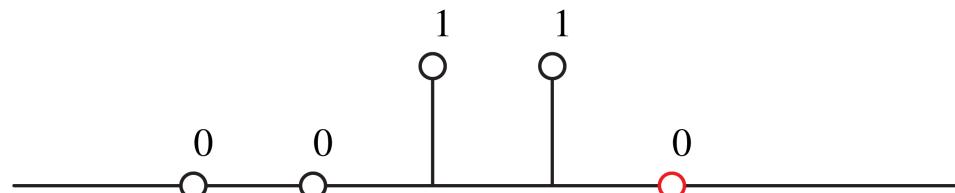
input



filter



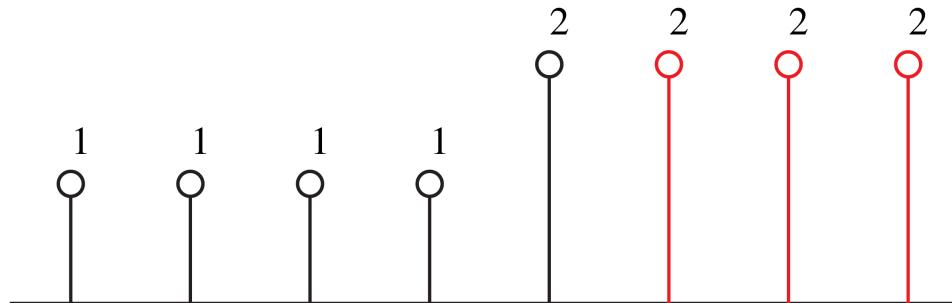
output



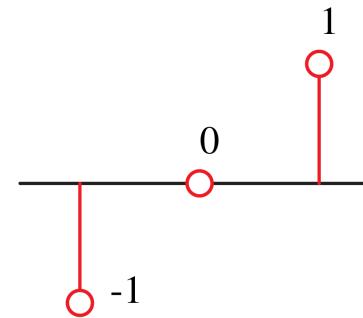
Convolution: a 1-D example

- sliding window
- dot product

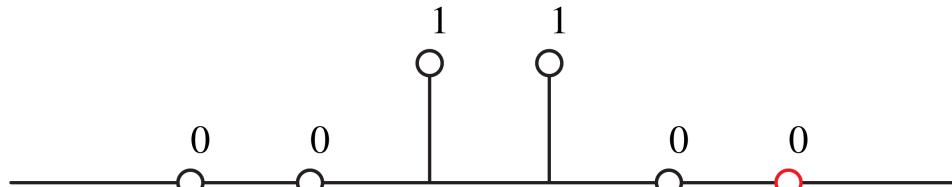
input



filter



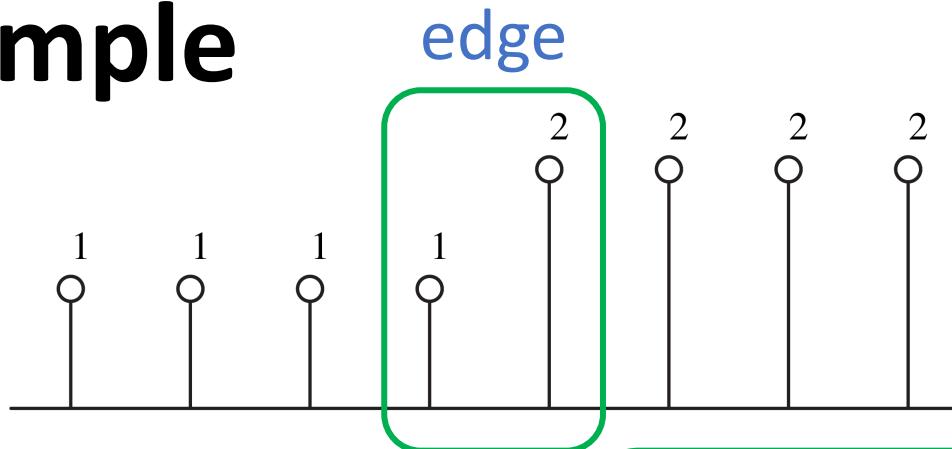
output



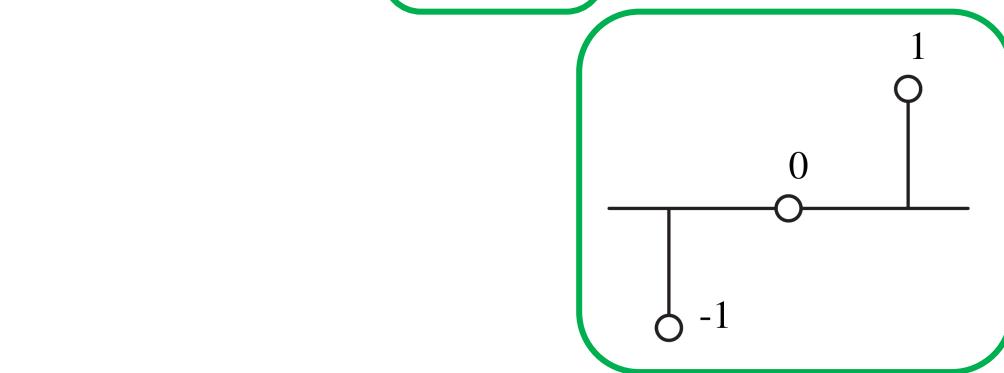
Convolution: a 1-D example

- sliding window
- dot product

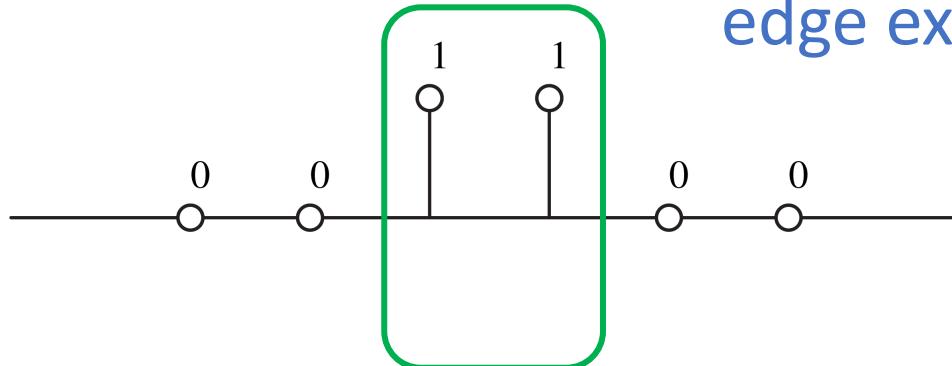
input



filter



output



indicator of edge ("feature")

Convolution: a 2-D example

input

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

output

filter

1	2	1
0	0	0
-1	-2	-1

Convolution: a 2-D example

input

0	1	0	2	0	1	0	0	0	0
0	0	0	0	0	0	1	1	1	0
0	-1	1	-2	1	-1	1	1	1	0
0	1	1	1	1	1	1	1	1	0
0	1	1	1	1	1	1	1	1	0
0	0	1	1	1	1	0	0	0	0
0	0	1	1	1	1	0	0	0	0
0	0	0	0	0	0	0	0	0	0

output

-3						

filter

1	2	1
0	0	0
-1	-2	-1

- sliding window
- dot product

Convolution: a 2-D example

input

0	0	1	0	2	0	1	0	0	0	0
0	0	0	0	0	0	0	1	1	1	0
0	1	-1	1	-2	1	-1	1	1	1	0
0	1	1	1	1	1	1	1	1	1	0
0	1	1	1	1	1	1	1	1	1	0
0	0	1	1	1	1	0	0	0	0	0
0	0	1	1	1	1	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0

output

-3	-4						

filter

1	2	1
0	0	0
-1	-2	-1

- sliding window
- dot product

Convolution: a 2-D example

input

0	0	0	1	0	2	0	1	0	0	0	0
0	0	0	0	0	0	0	0	1	1	0	0
0	1	1	-1	1	-2	1	-1	1	1	0	0
0	1	1	1	1	1	1	1	1	1	0	0
0	1	1	1	1	1	1	1	1	1	0	0
0	0	1	1	1	0	0	0	0	0	0	0
0	0	1	1	1	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0

output

-3	-4	-4				

filter

1	2	1
0	0	0
-1	-2	-1

- sliding window
- dot product

Convolution: a 2-D example

input

0	0	0	0	1	0	1	0	0	0
0	0	0	0	0	0	1	0	1	0
0	1	1	1	-1	1	-2	1	-1	1
0	1	1	1	1	1	1	1	1	0
0	1	1	1	1	1	1	1	1	0
0	0	1	1	1	0	0	0	0	0
0	0	1	1	1	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

output

-3	-4	-4	-4			

filter

1	2	1
0	0	0
-1	-2	-1

- sliding window
- dot product

Convolution: a 2-D example

input

0	0	0	0	0	0	0	0
0	0	0	0	0	1	1	0
0	1	1	1	1	1	1	0
0	1	1	1	1	1	1	0
0	1	1	1	1	1	1	0
0	0	1	1	1	01	02	01
0	0	1	1	1	00	00	00
0	0	0	0	0	-1	-2	-1

filter

1	2	1
0	0	0
-1	-2	-1

- sliding window
- dot product

output

-3	-4	-4	-4	-4	-3
-3	-4	-4	-3	-1	0
0	0	0	0	0	0
2	1	0	1	3	3
2	1	0	1	3	3
1	3	4	3	1	0

Convolution: a 2-D example

$$y[n, m] = \sum_{i=-r}^r \sum_{j=-r}^r w[i, j] x[n + i, m + j]$$

output map

filter weights

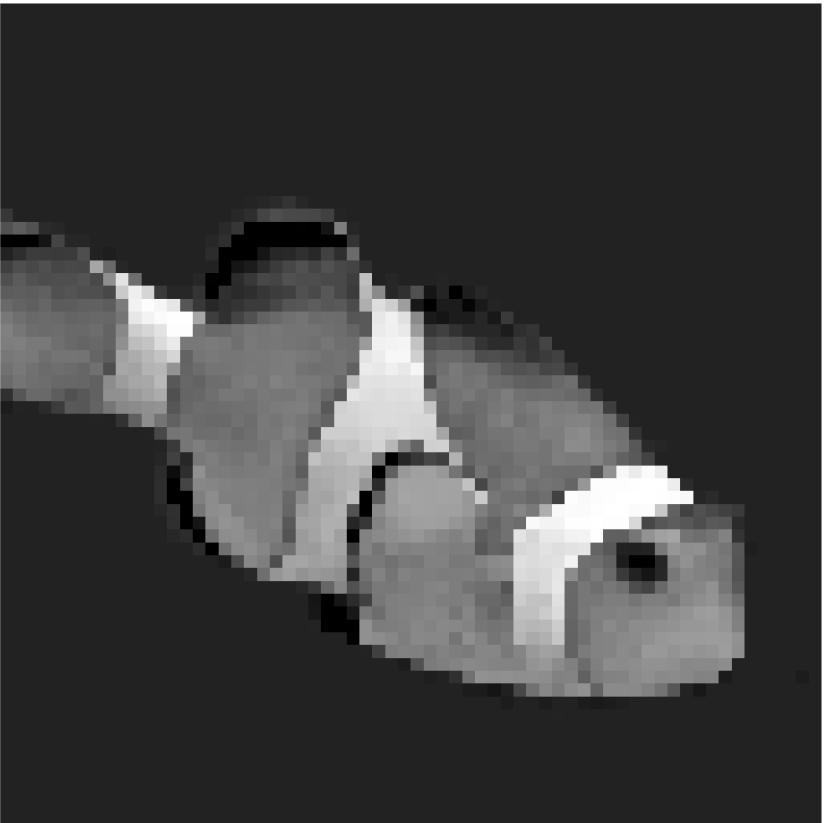
input map

coordinates in a local window

r : kernel radius
kernel size = $2r + 1$

* In ConvNets, convolution is often implemented as **cross-correlation** (no flipping)

Convolution: 2-D

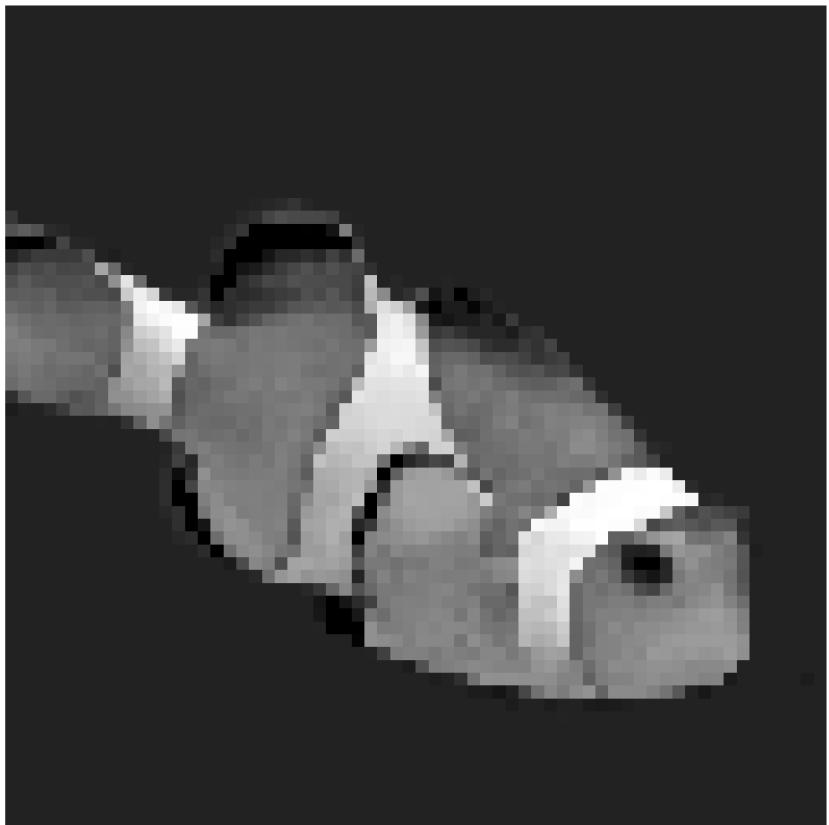


$$\text{filter} \quad \begin{array}{|c|c|c|} \hline 1 & 2 & 1 \\ \hline 0 & 0 & 0 \\ \hline -1 & -2 & -1 \\ \hline \end{array} =$$



This example is a Sobel edge extractor of horizontal edges.

Convolution: 2-D



$$\text{input} * \text{filter} = \text{output}$$

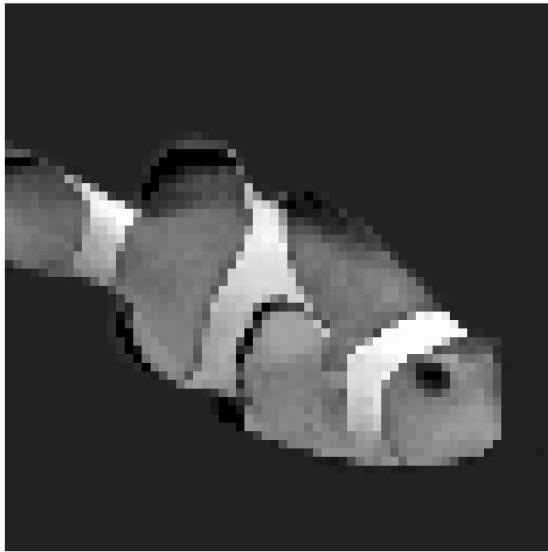
The diagram illustrates the convolution process. On the left is the input image. In the center is a 3x3 filter kernel labeled "filter". The asterisk (*) indicates the convolution operation, and the equals sign (=) indicates the resulting output image on the right.

Light Gray	White	Light Gray
Dark Gray	Black	Dark Gray
Dark Gray	Black	Dark Gray



This example is a Sobel edge extractor of horizontal edges.

Convolution: Multi-channel outputs



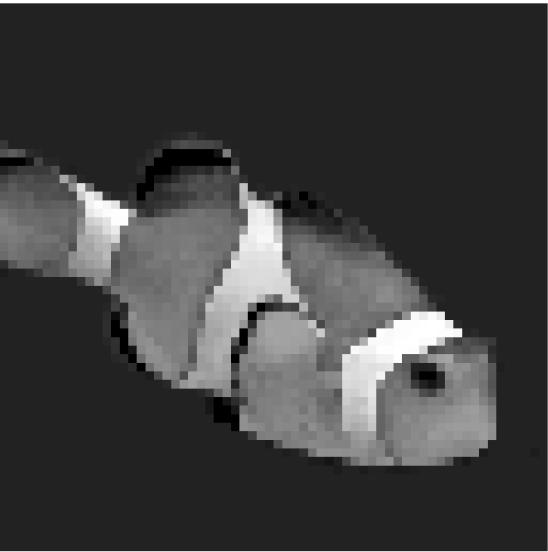
$$\begin{matrix} * & \begin{array}{|c|c|c|} \hline 1 & 2 & 1 \\ \hline 0 & 0 & 0 \\ \hline -1 & -2 & -1 \\ \hline \end{array} & = \end{matrix}$$

one filter, one feature

$$\begin{matrix} * & \begin{array}{|c|c|c|} \hline 1 & 0 & -1 \\ \hline 2 & 0 & -2 \\ \hline 1 & 0 & -1 \\ \hline \end{array} & = \end{matrix}$$



Convolution: Multi-channel outputs



$$\text{input} * \begin{matrix} \text{filter} \\ \text{matrix} \end{matrix} = \text{output}$$

The input image is multiplied by a 3x3 convolutional filter matrix. The filter consists of nine gray squares arranged in a 3x3 grid, with the central square being black.

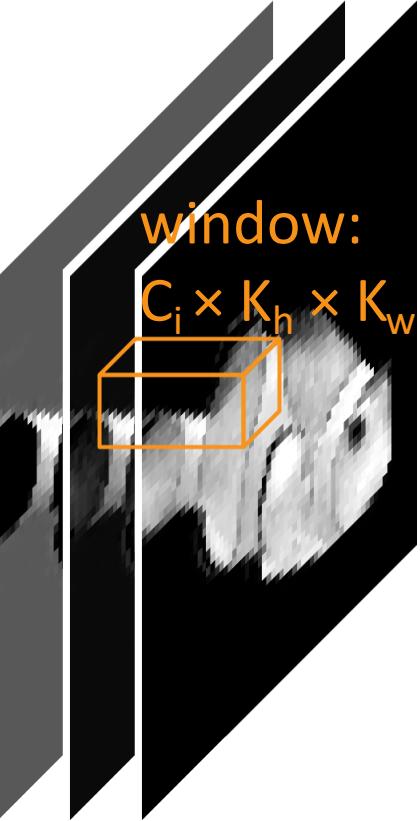
one filter, one feature

$$\text{input} * \begin{matrix} \text{filter} \\ \text{matrix} \end{matrix} = \text{output}$$

The input image is multiplied by a 3x3 convolutional filter matrix. The filter consists of nine gray squares arranged in a 3x3 grid, with the bottom-right square being black.



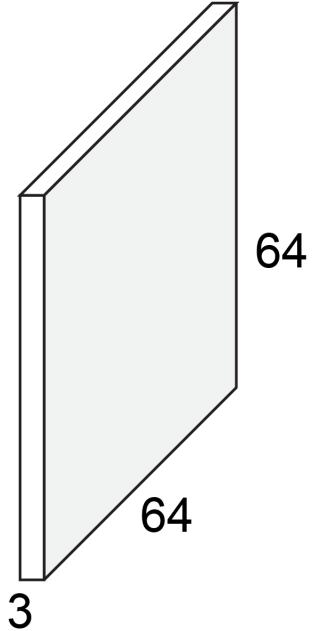
Convolution: Multi-channel inputs



$$\text{window: } C_i \times K_h \times K_w \quad * \quad \text{filter: } C_i \times K_h \times K_w =$$
A diagram illustrating a convolution operation. It shows a 3D filter kernel (represented by a 3x3 grid of squares) being multiplied (*) with a 3D input window (represented by a 3x3 grid of squares). The result of the multiplication is indicated by an equals sign (=).

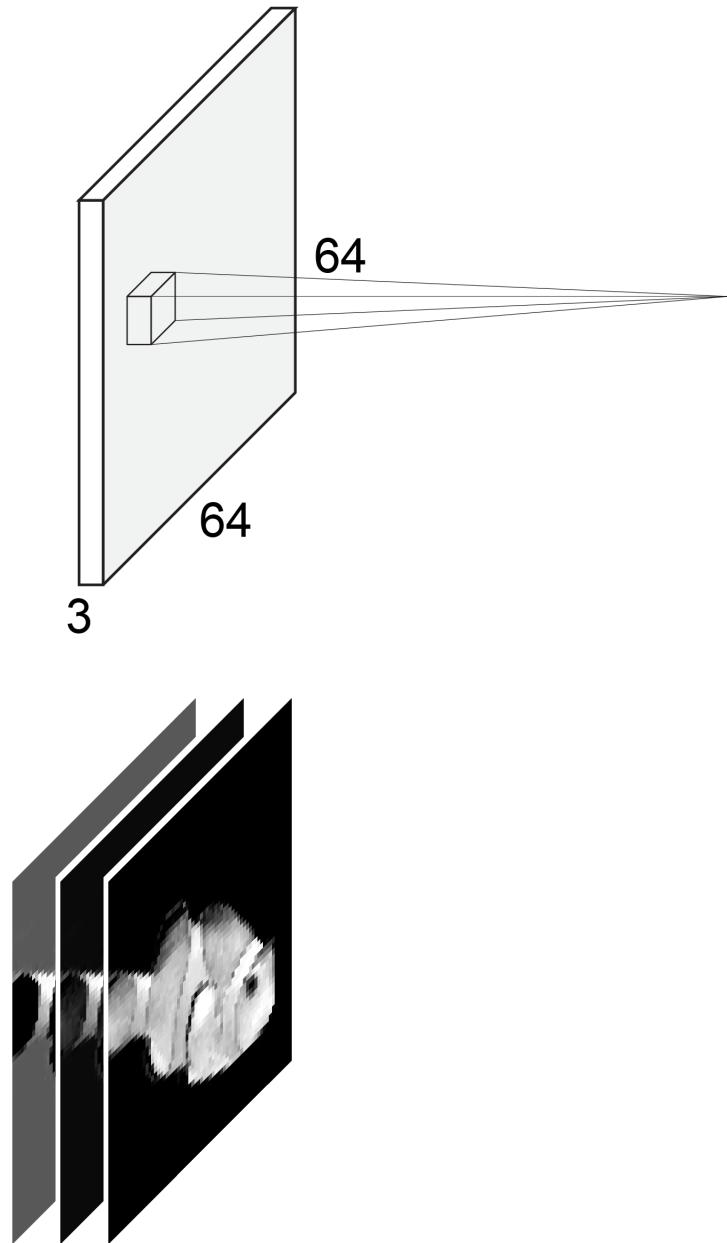


Convolution: tensor views



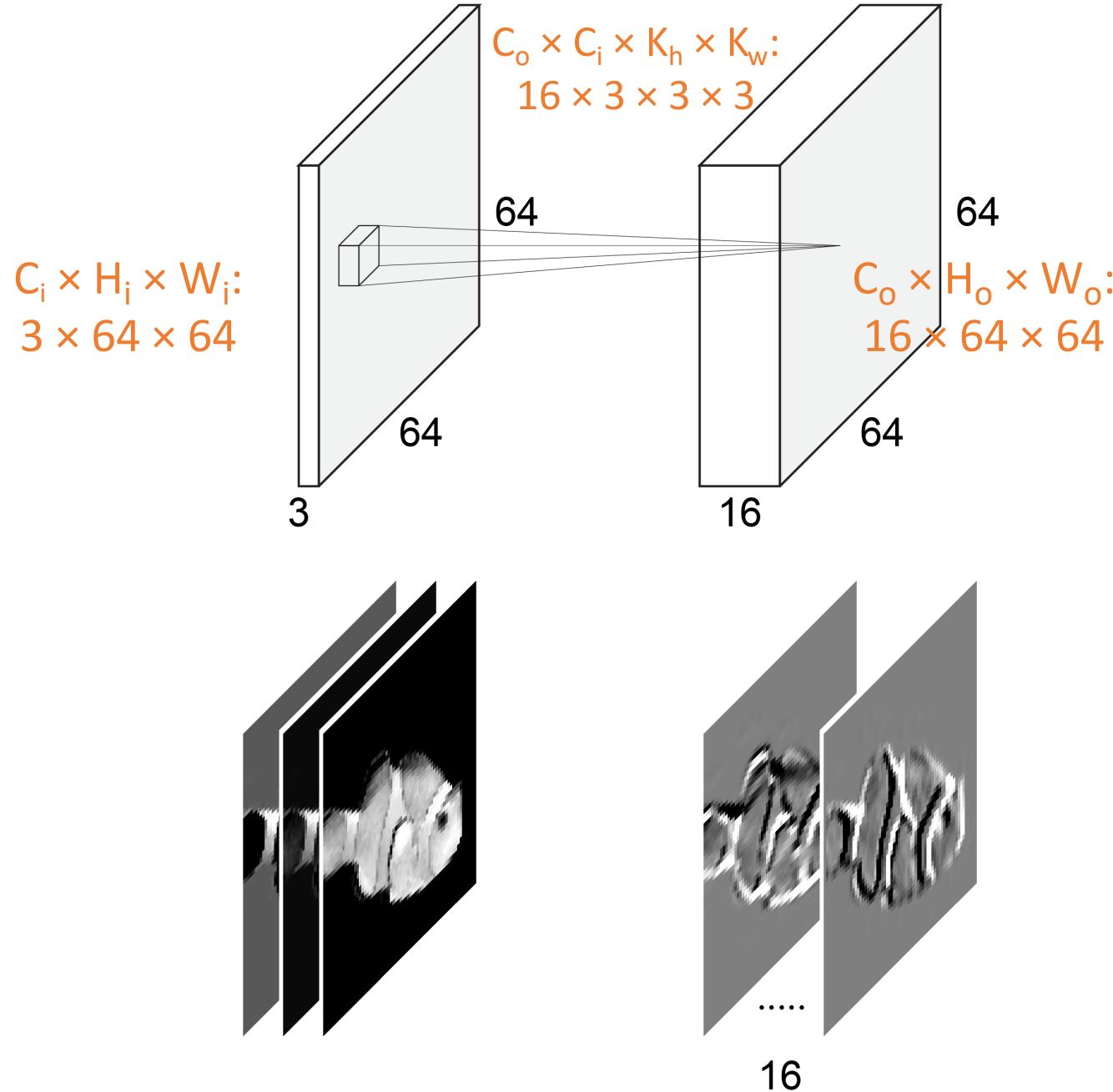
- Tensor: high-dimension array
- feature maps
 - 3-D tensor: $C \times H \times W$
 - C: channels
 - H: height
 - W: width

Convolution: tensor views



- Tensor: high-dimension array
- feature maps
 - 3-D tensor: $C \times H \times W$
 - C : channels
 - H : height
 - W : width
- filters
 - 4-D tensor: $C_o \times C_i \times K_h \times K_w$
 - C_o : output channels
 - C_i : input channels
 - K_h, K_w : filter height, width

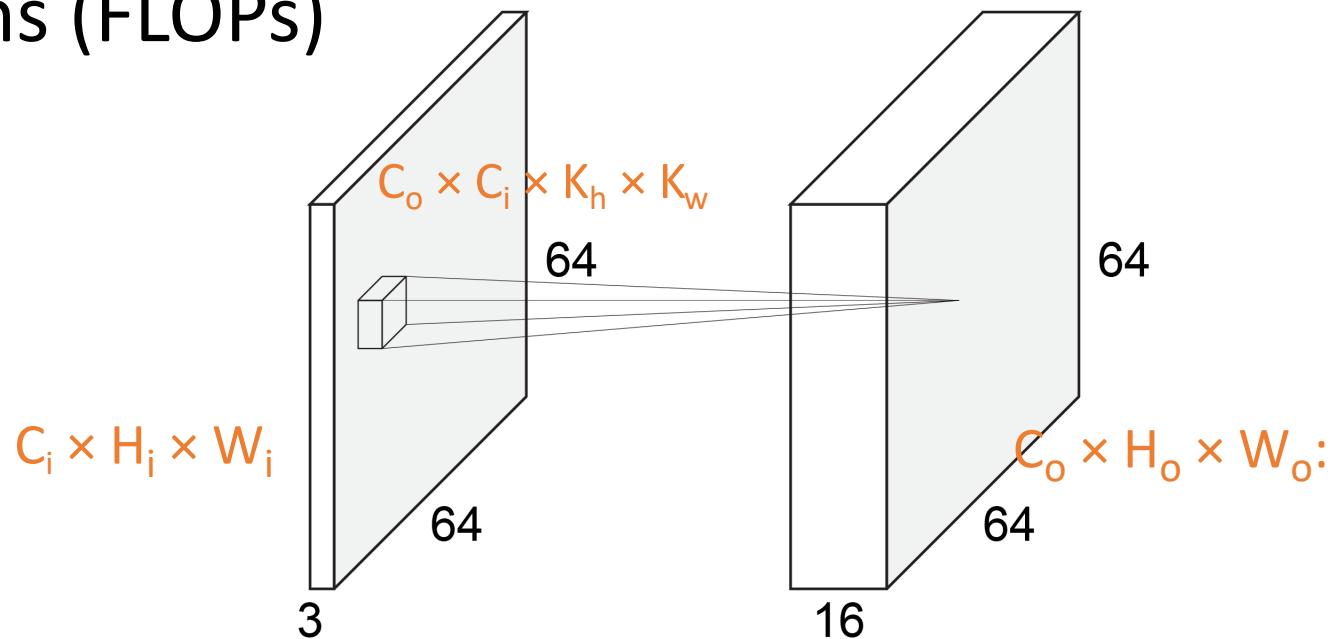
Convolution: tensor views



- Tensor: high-dimension array
- feature maps
 - 3-D tensor: $C \times H \times W$
 - C : channels
 - H : height
 - W : width
- filters
 - 4-D tensor: $C_o \times C_i \times K_h \times K_w$
 - C_o : output channels
 - C_i : input channels
 - K_h, K_w : filter height, width

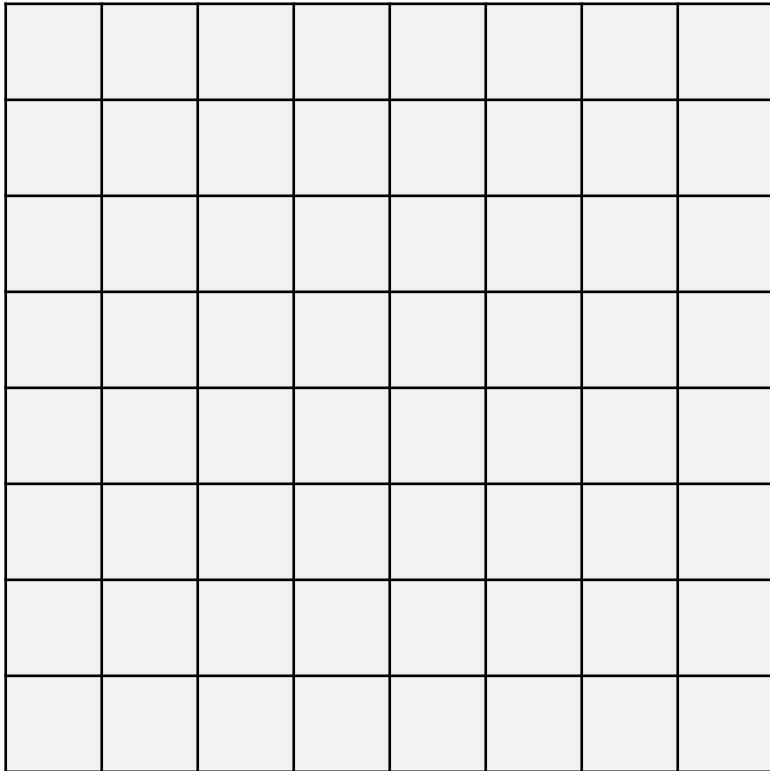
Convolution: # parameters and # operations

- # parameters
 - weights: $C_o \times C_i \times K_h \times K_w$
 - bias: C_o
- # floating-point operations (FLOPs)
 - # params $\times H_o \times W_o$
- # params per FLOP
 - $1 / \text{spatial_size}$
 - parameter-efficient

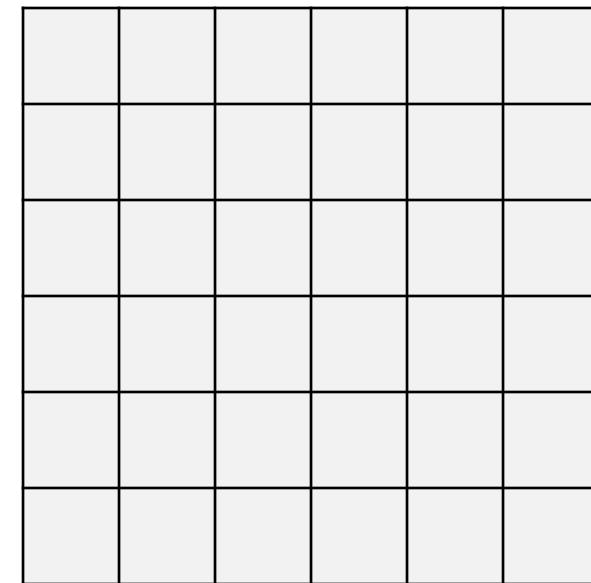
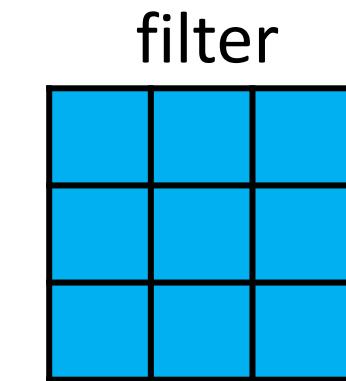


Convolution: padding

input: $H \times W = 8 \times 8$

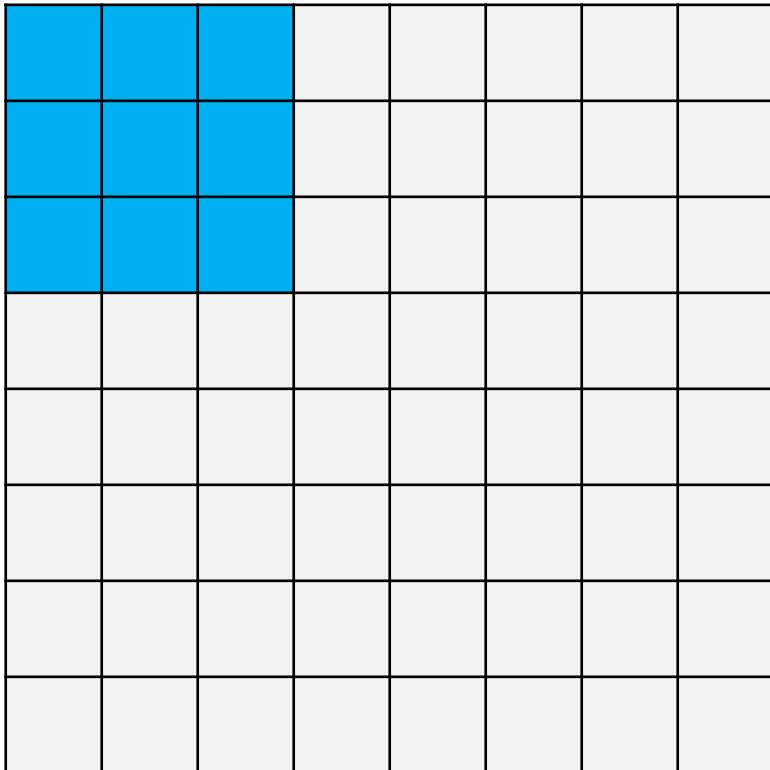


output: $H \times W = 6 \times 6$

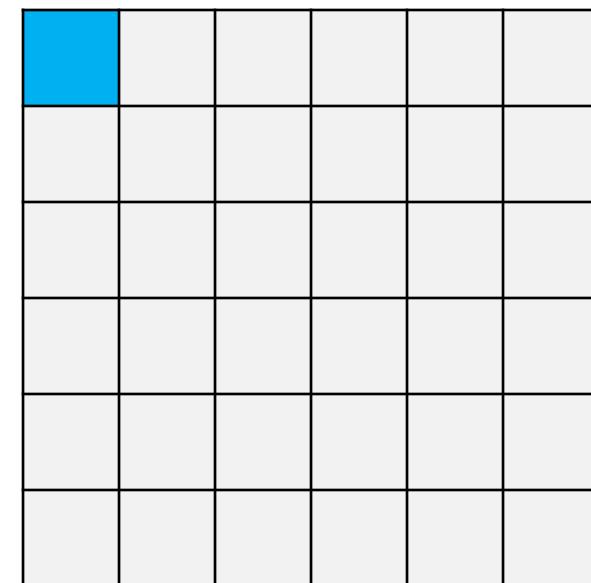
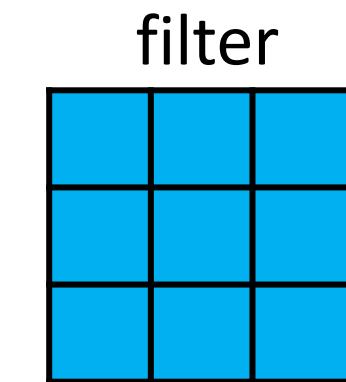


Convolution: padding

input: $H \times W = 8 \times 8$

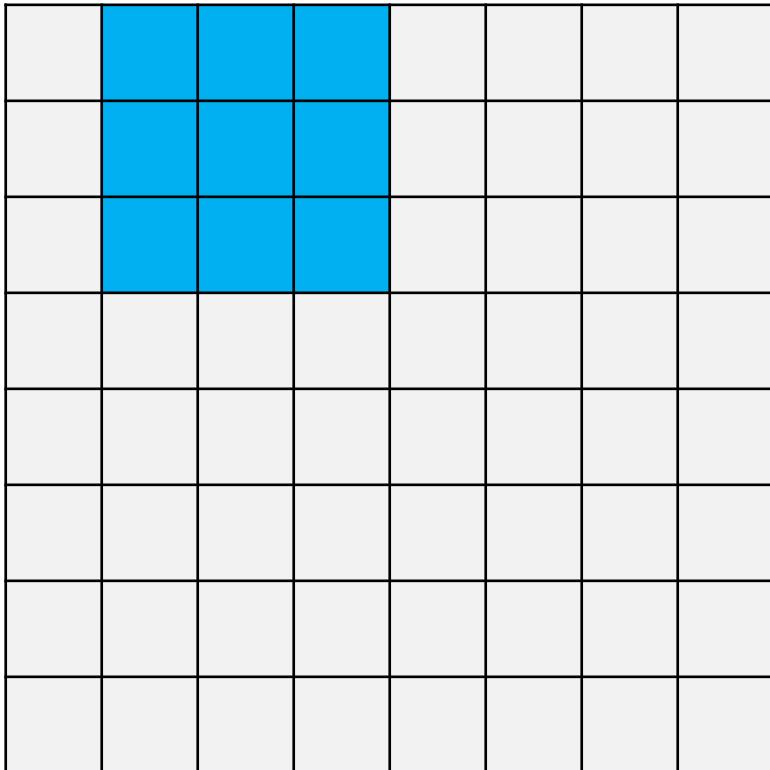


output: $H \times W = 6 \times 6$

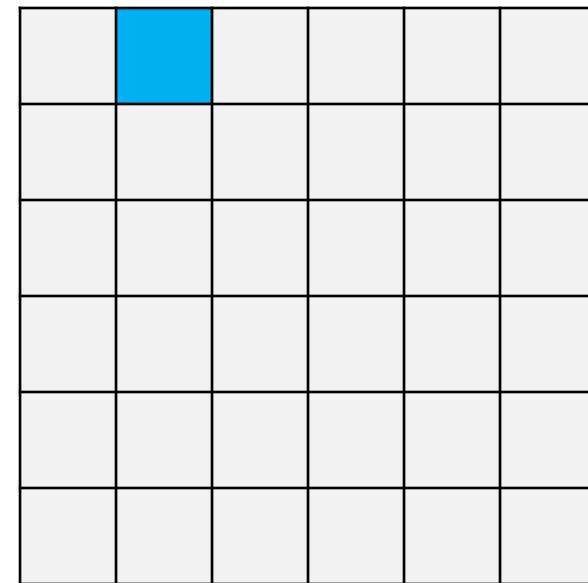
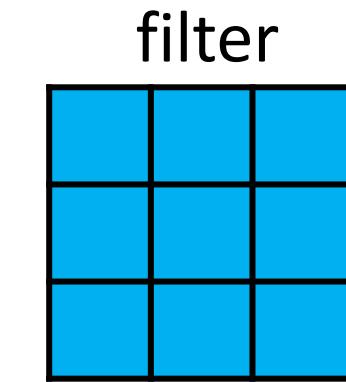


Convolution: padding

input: $H \times W = 8 \times 8$

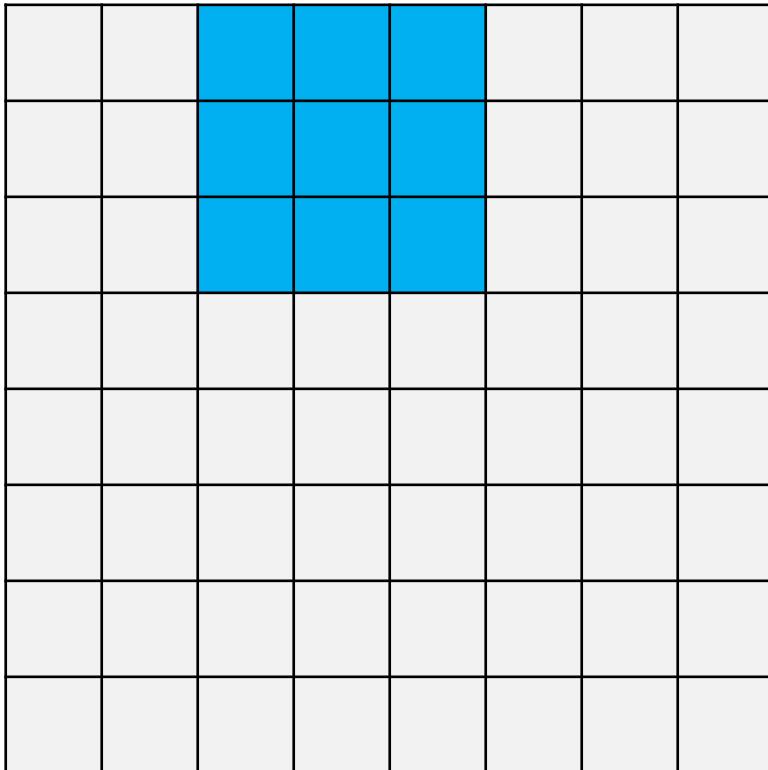


output: $H \times W = 6 \times 6$

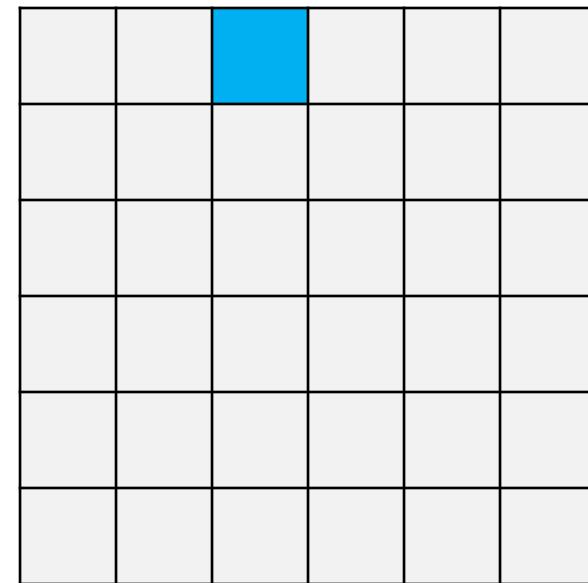


Convolution: padding

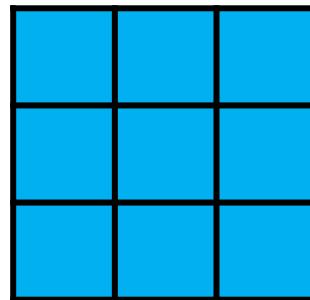
input: $H \times W = 8 \times 8$



output: $H \times W = 6 \times 6$

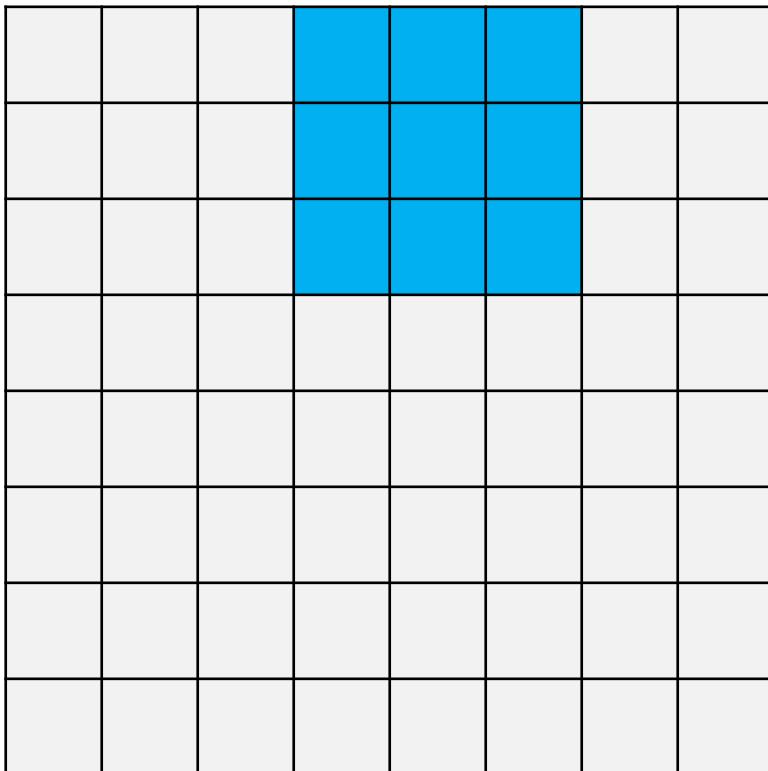


filter

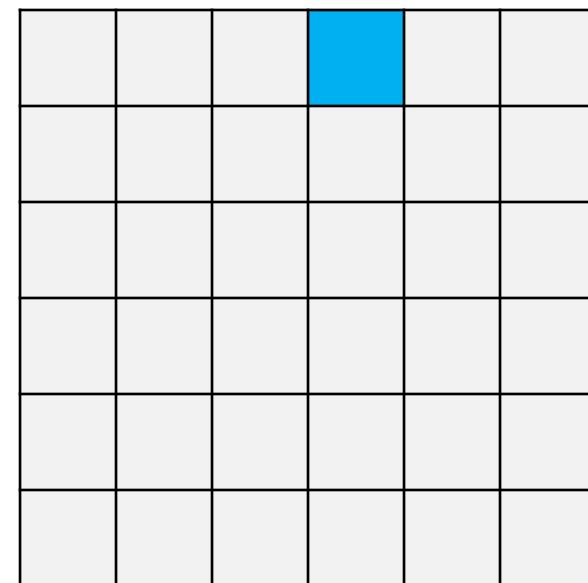


Convolution: padding

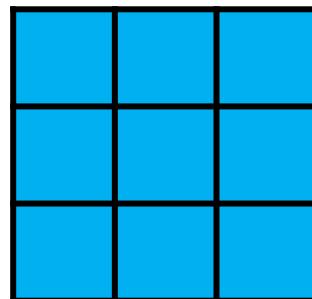
input: $H \times W = 8 \times 8$



output: $H \times W = 6 \times 6$

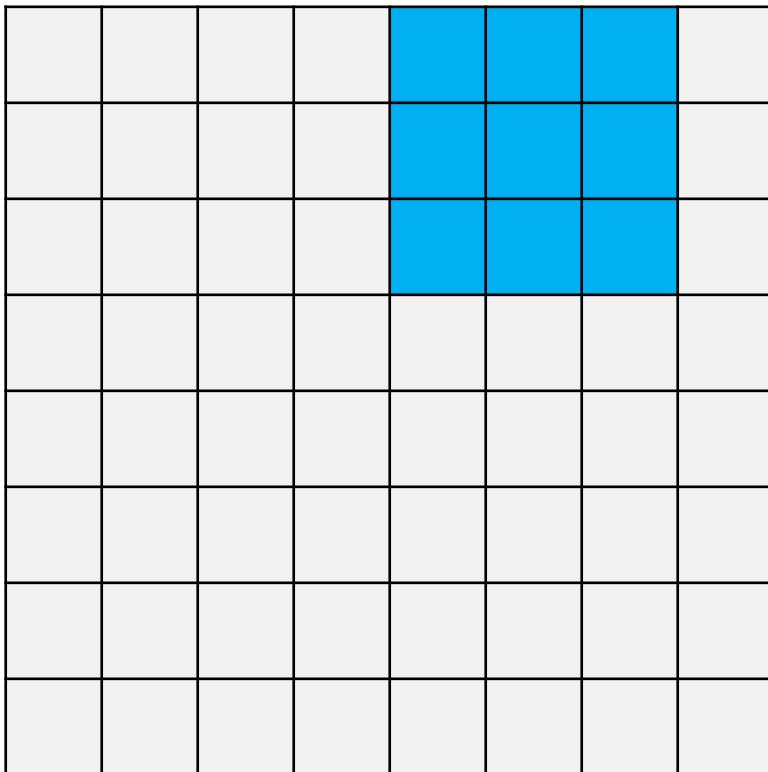


filter

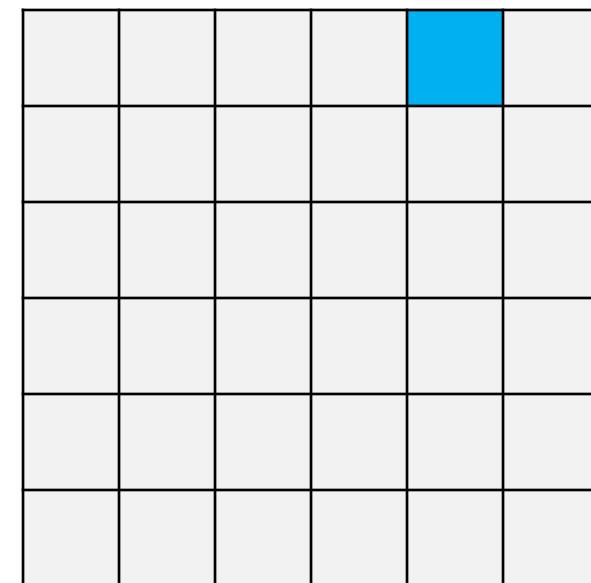
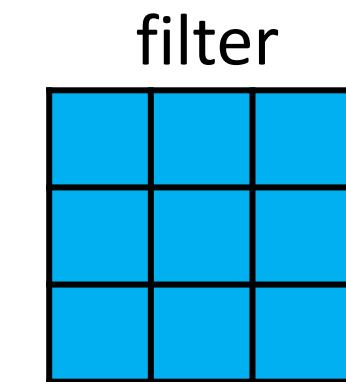


Convolution: padding

input: $H \times W = 8 \times 8$

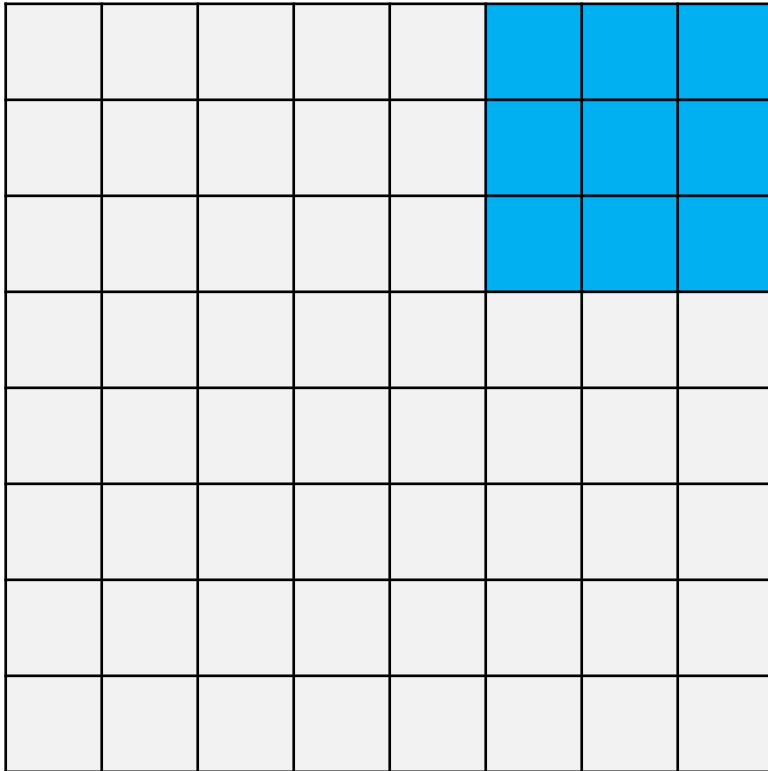


output: $H \times W = 6 \times 6$

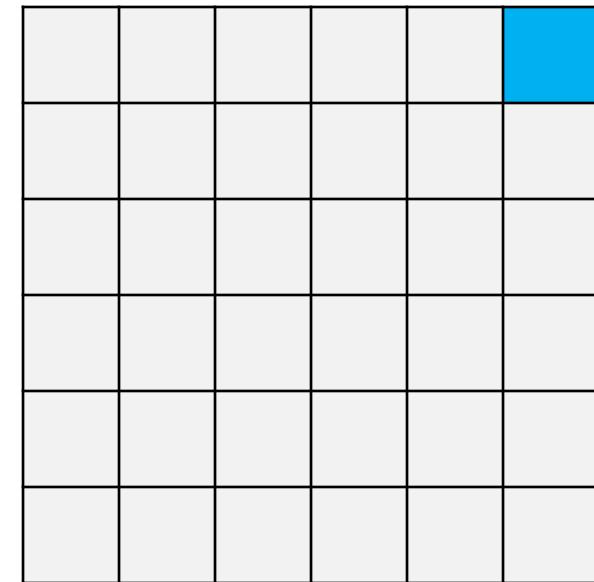
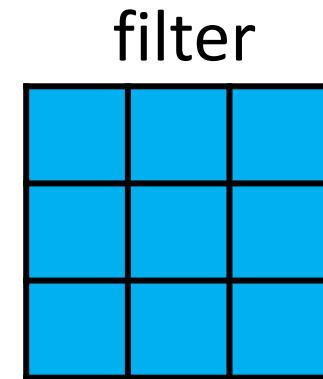


Convolution: padding

input: $H \times W = 8 \times 8$



output: $H \times W = 6 \times 6$

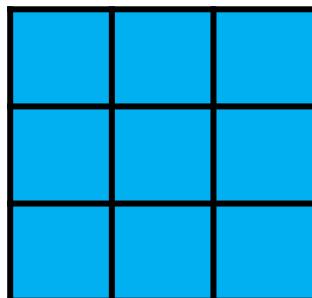


$$H_{\text{out}} = H_{\text{in}} - K_h + 1$$

Convolution: padding

input: 8×8 , + pad

filter

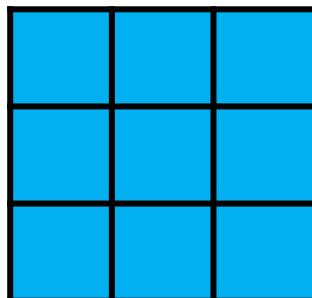


output: $H \times W = 8 \times 8$

Convolution: padding

input: 8×8 , + pad

filter

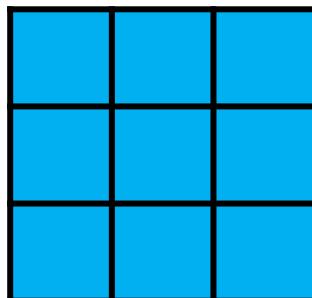


output: $H \times W = 8 \times 8$

Convolution: padding

input: 8×8 , + pad

filter

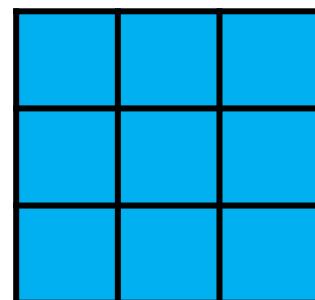


output: $H \times W = 8 \times 8$

Convolution: padding

input: 8×8 , + pad

filter

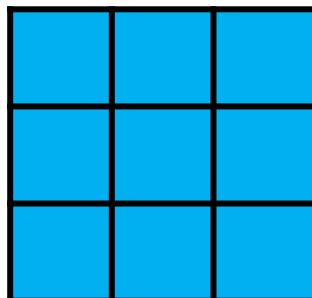


output: $H \times W = 8 \times 8$

Convolution: padding

input: 8×8 , + pad

filter

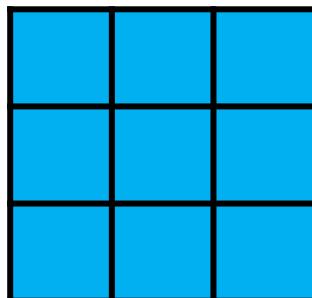


output: $H \times W = 8 \times 8$

Convolution: padding

input: 8×8 , + pad

filter

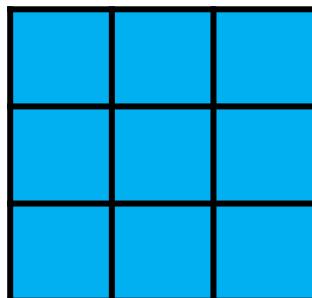


output: $H \times W = 8 \times 8$

Convolution: padding

input: 8×8 , + pad

filter

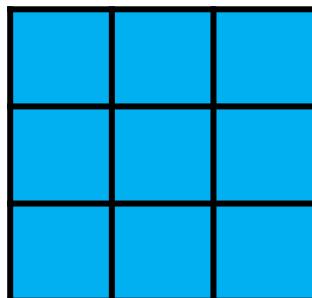


output: $H \times W = 8 \times 8$

Convolution: padding

input: 8×8 , + pad

filter

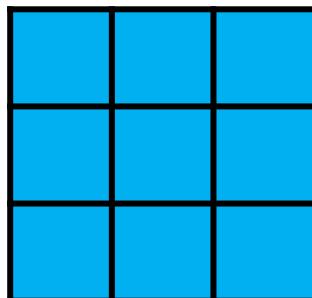


output: $H \times W = 8 \times 8$

Convolution: padding

input: 8×8 , + pad

filter



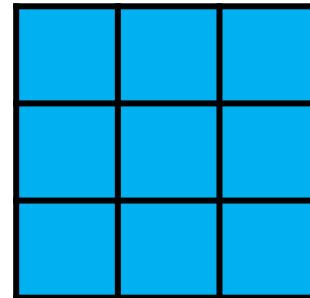
output: $H \times W = 8 \times 8$

Convolution: padding

input: 8×8 , + pad

0	0	0	0	0	0	0	0	0	0
0									0
0									0
0									0
0									0
0									0
0									0
0									0
0	0	0	0	0	0	0	0	0	0

filter



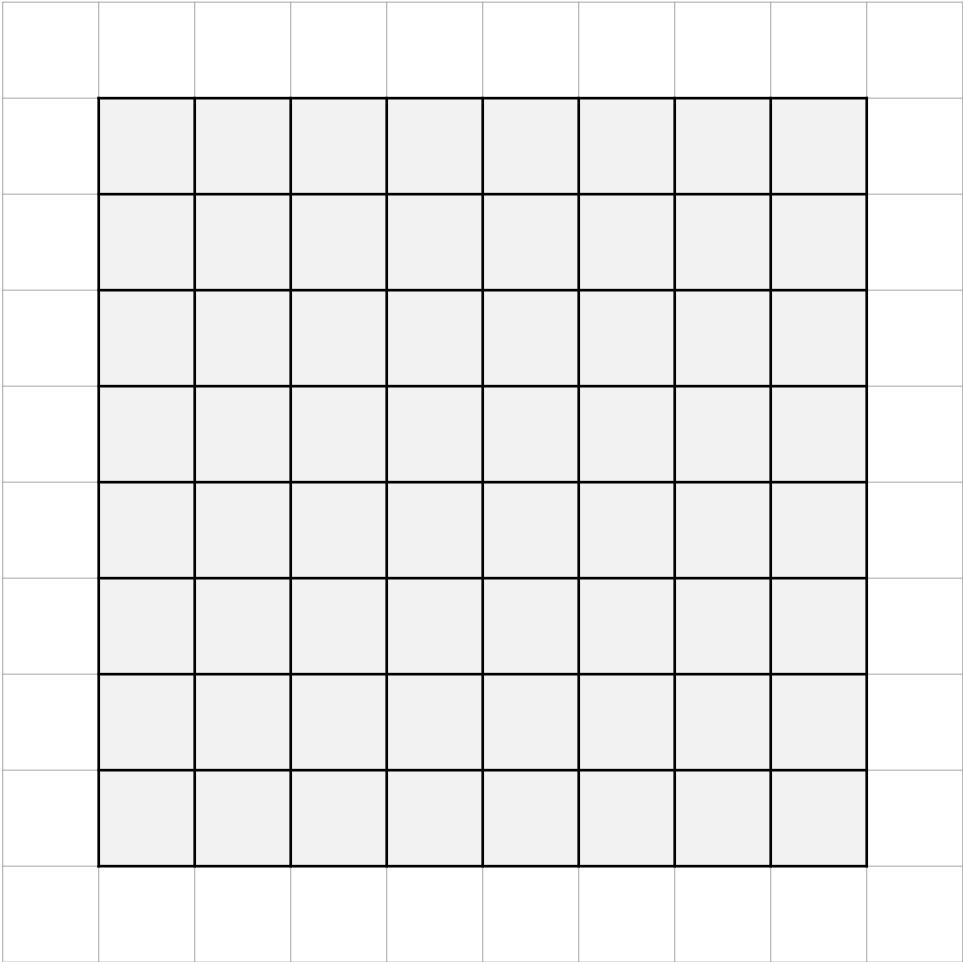
- $\text{pad} = [\text{kernel_size} / 2]$
- maintains feature map size

output: $H \times W = 8 \times 8$

$$H_{\text{out}} = H_{\text{in}} + 2\text{pad}_h - K_h + 1$$

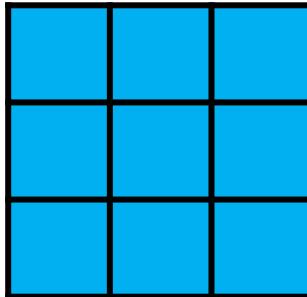
Convolution: stride

input

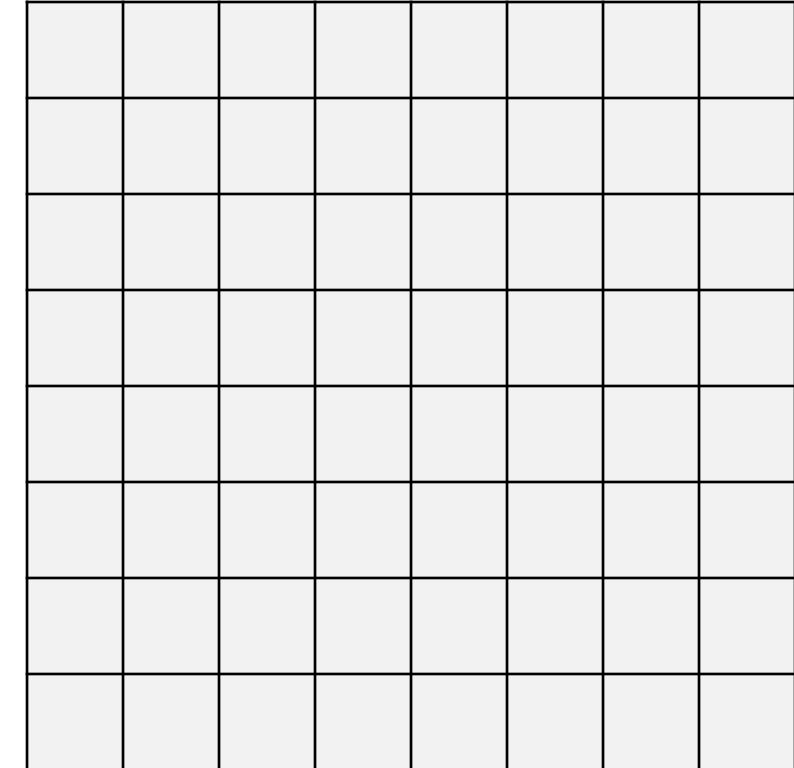


stride = 2

filter

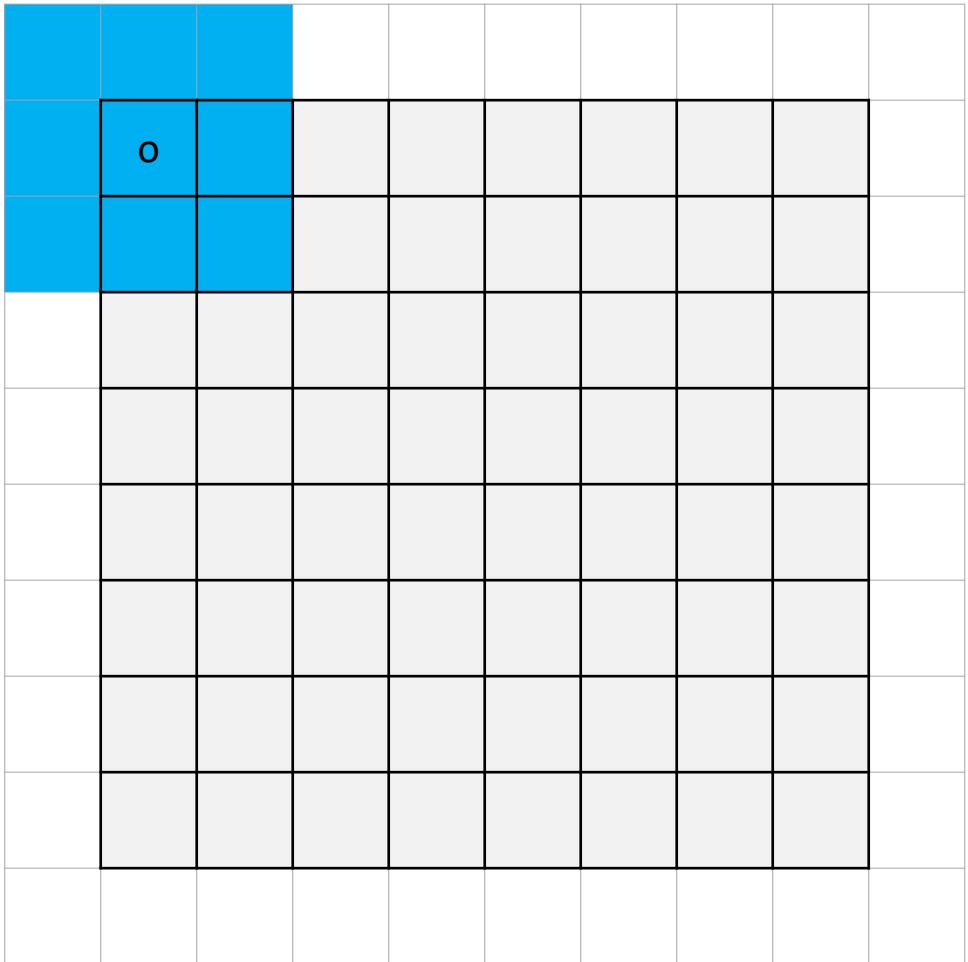


output



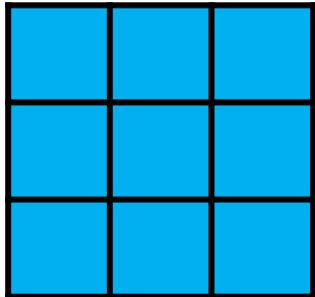
Convolution: stride

input

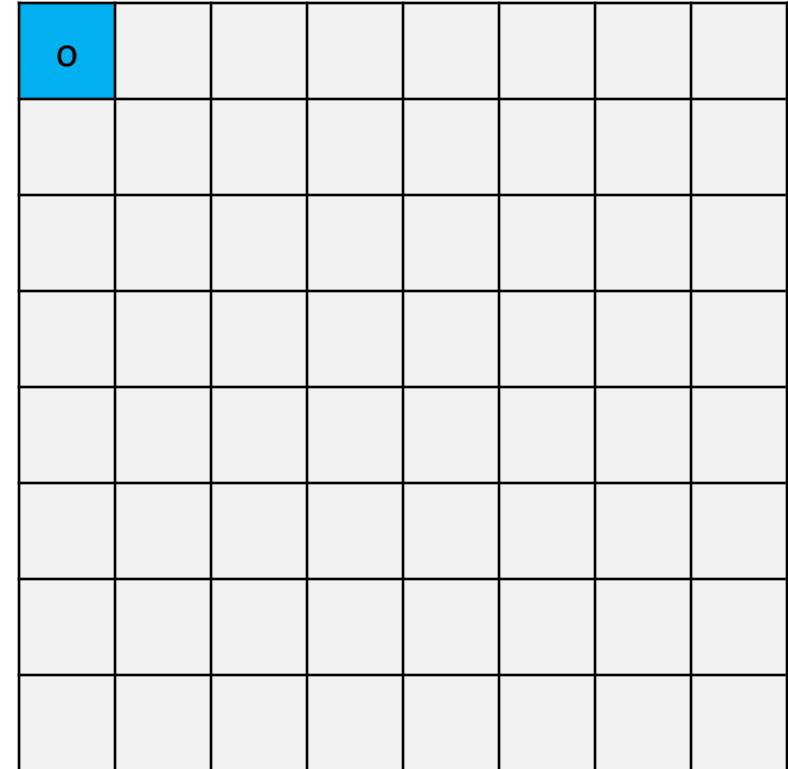


stride = 2

filter

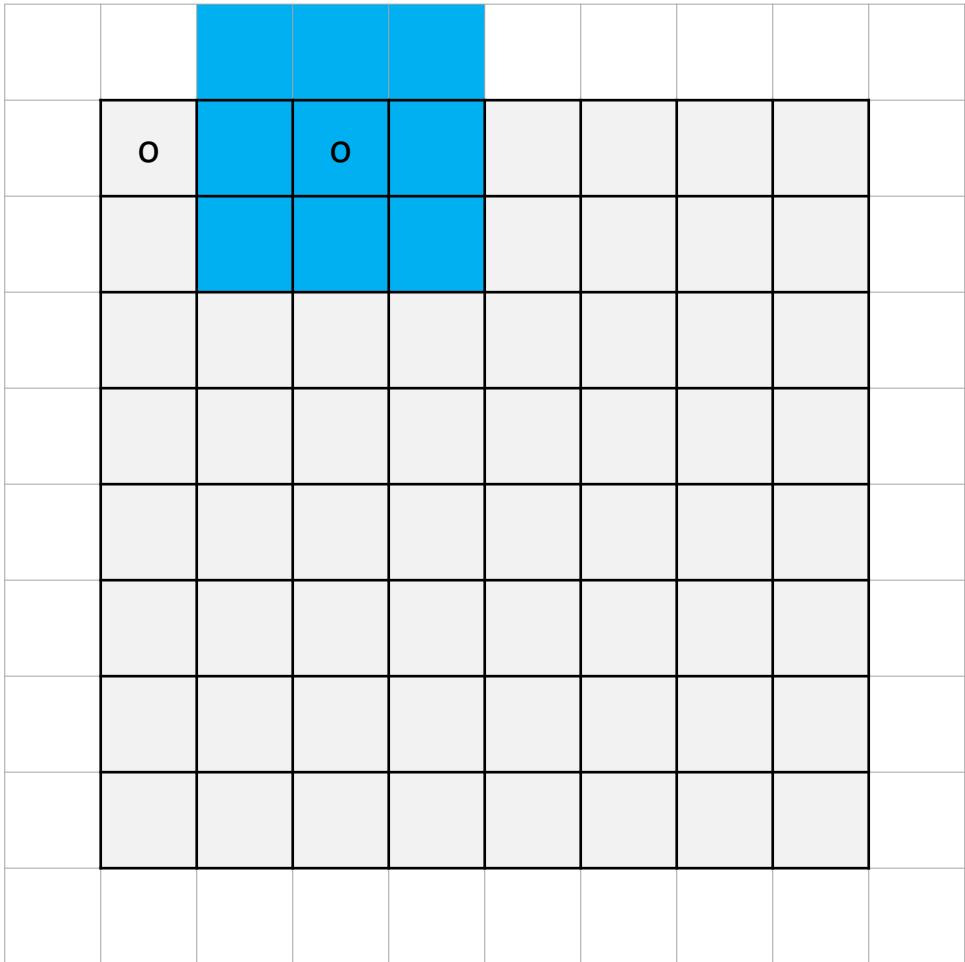


output



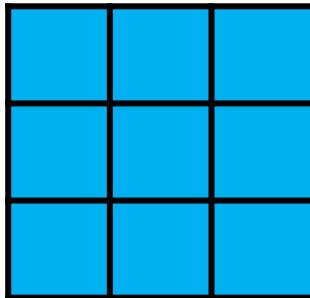
Convolution: stride

input

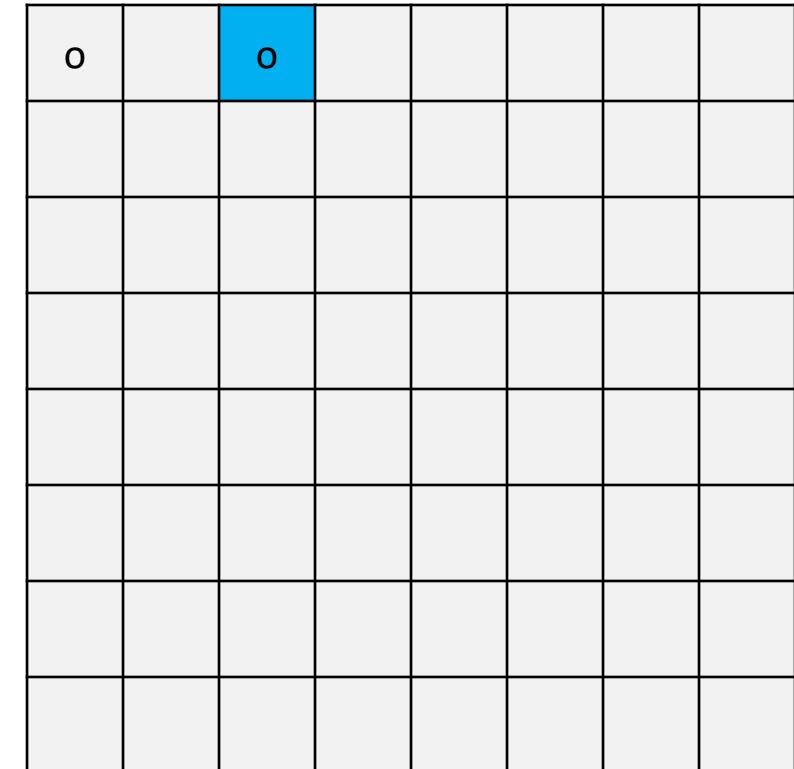


stride = 2

filter



output

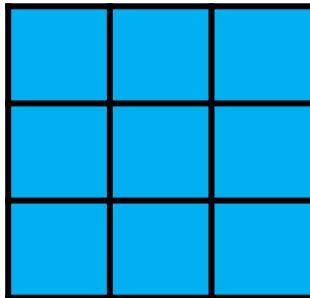


Convolution: stride

input

stride = 2

filter



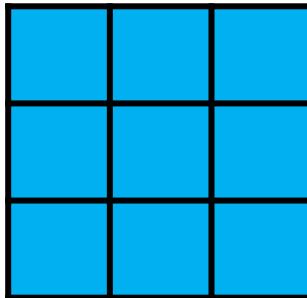
output

Convolution: stride

input

stride = 2

filter



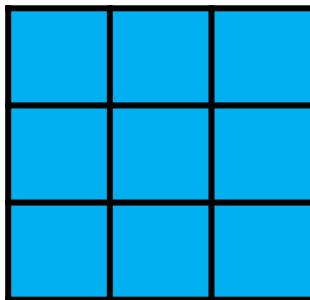
output

Convolution: stride

input

stride = 2

filter



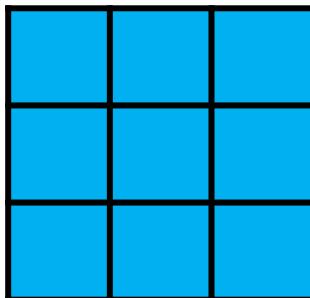
output

Convolution: stride

input

o		o		o		o		
o		o		o		o		
o		o		o		o		
o		o		o	o	o	o	
o		o		o		o		

filter



output

o		o		o		o	
o		o		o		o	
o		o		o		o	
o		o		o	o	o	o

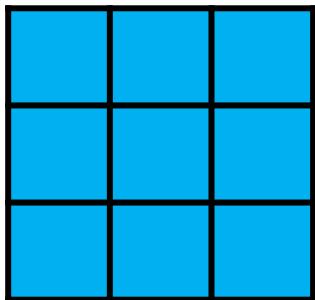
Convolution: stride

input: $H \times W = 8 \times 8$

o		o		o		o		
o		o		o		o		
o		o		o		o		
o		o		o		o		

stride = 2

filter



output: $H \times W = 4 \times 4$

o		o		o		o	
o		o		o		o	
o		o		o		o	
o		o		o		o	

Convolution: stride

input: $H \times W = 8 \times 8$

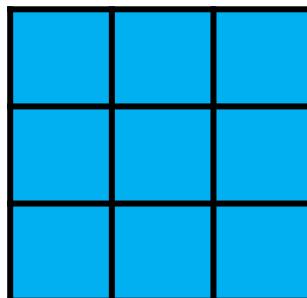
o		o		o		o		
o		o		o		o		
o		o		o		o		
o		o		o		o		

stride = 2

- reduces feature map size
- compress and abstract

output: $H \times W = 4 \times 4$

filter



o	o	o	o
o	o	o	o
o	o	o	o
o	o	o	o

$$H_{out} = \lfloor (H_{in} + 2\text{pad}_h - K_h) / \text{str} \rfloor + 1$$

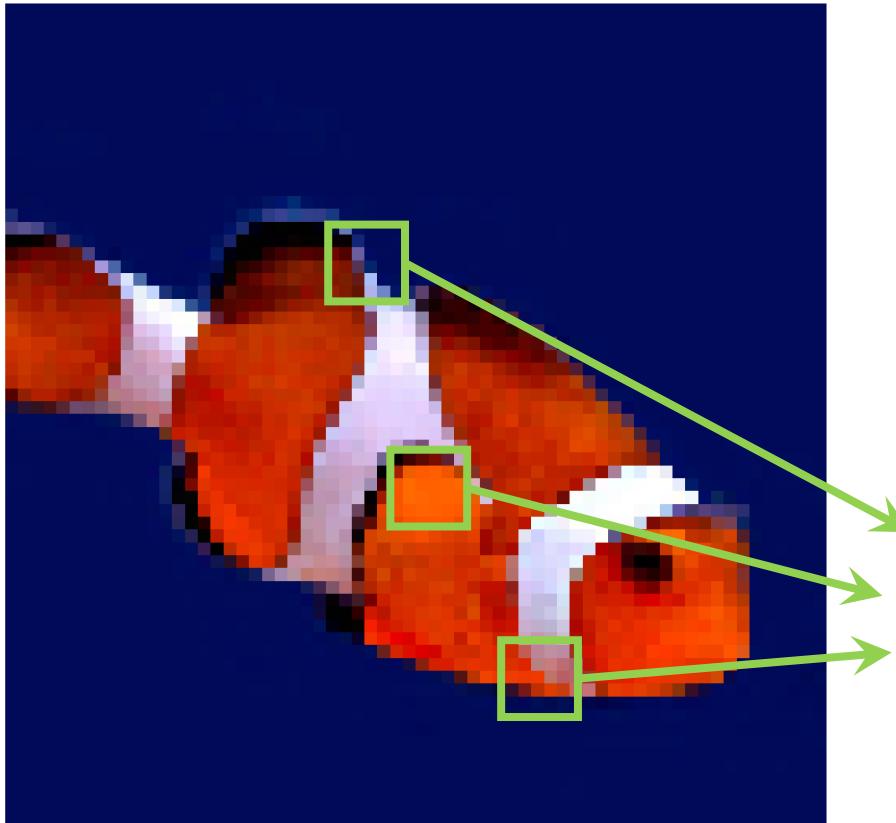
*rounding operation depends on libraries

Convolution: arguments

- kernel size: $K_h \times K_w$ (e.g., 3×3)
- output channels: C_o
 - input channels: C_i determined by input
- padding (e.g., 1)
- stride (e.g., 1 or 2)

Convolution: translation-invariance

- Process each window in the same way



apply the same weights regardless of
the window location

Convolution: translation-invariance

- Process each window in the same way

$$y[n, m] = \sum_{i=-r}^r \sum_{j=-r}^r w[i, j] x[n + i, m + j]$$

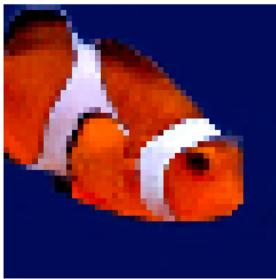
apply the same weights regardless of
the window location

Convolution: translation-equivariance

- Equivariance: $f(\mathcal{T}(x)) = \mathcal{T}(f(x))$

Convolution: translation-equivariance

- Equivariance: $f(\mathcal{T}(x)) = \mathcal{T}(f(x))$

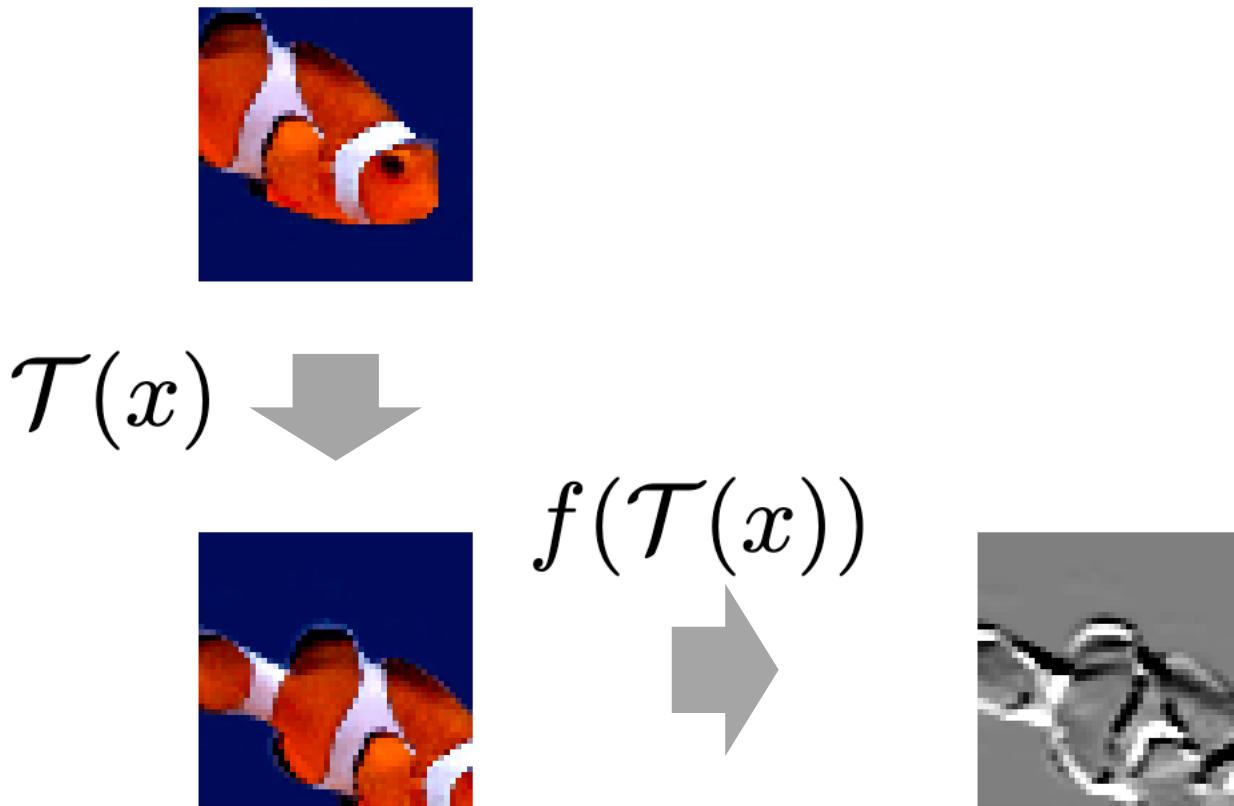


$$\mathcal{T}(x) \downarrow$$



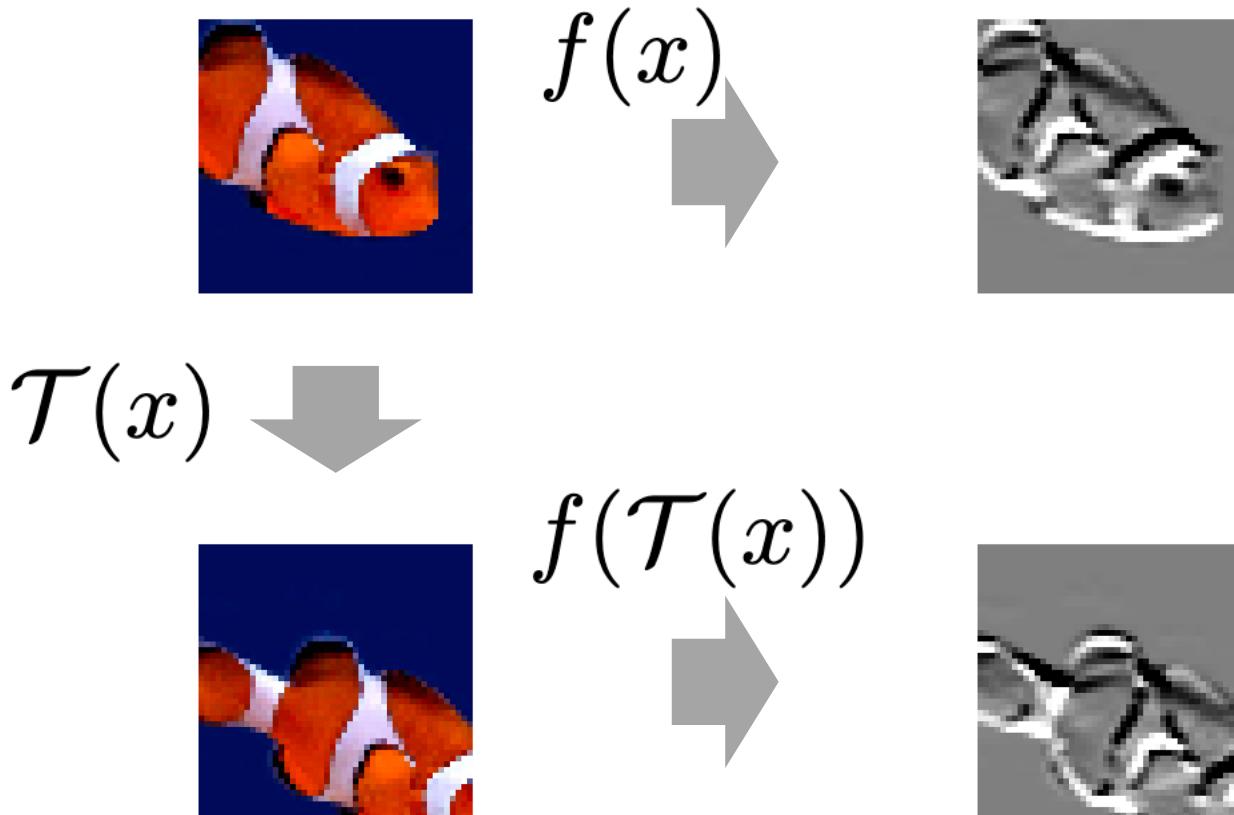
Convolution: translation-equivariance

- Equivariance: $f(\mathcal{T}(x)) = \mathcal{T}(f(x))$



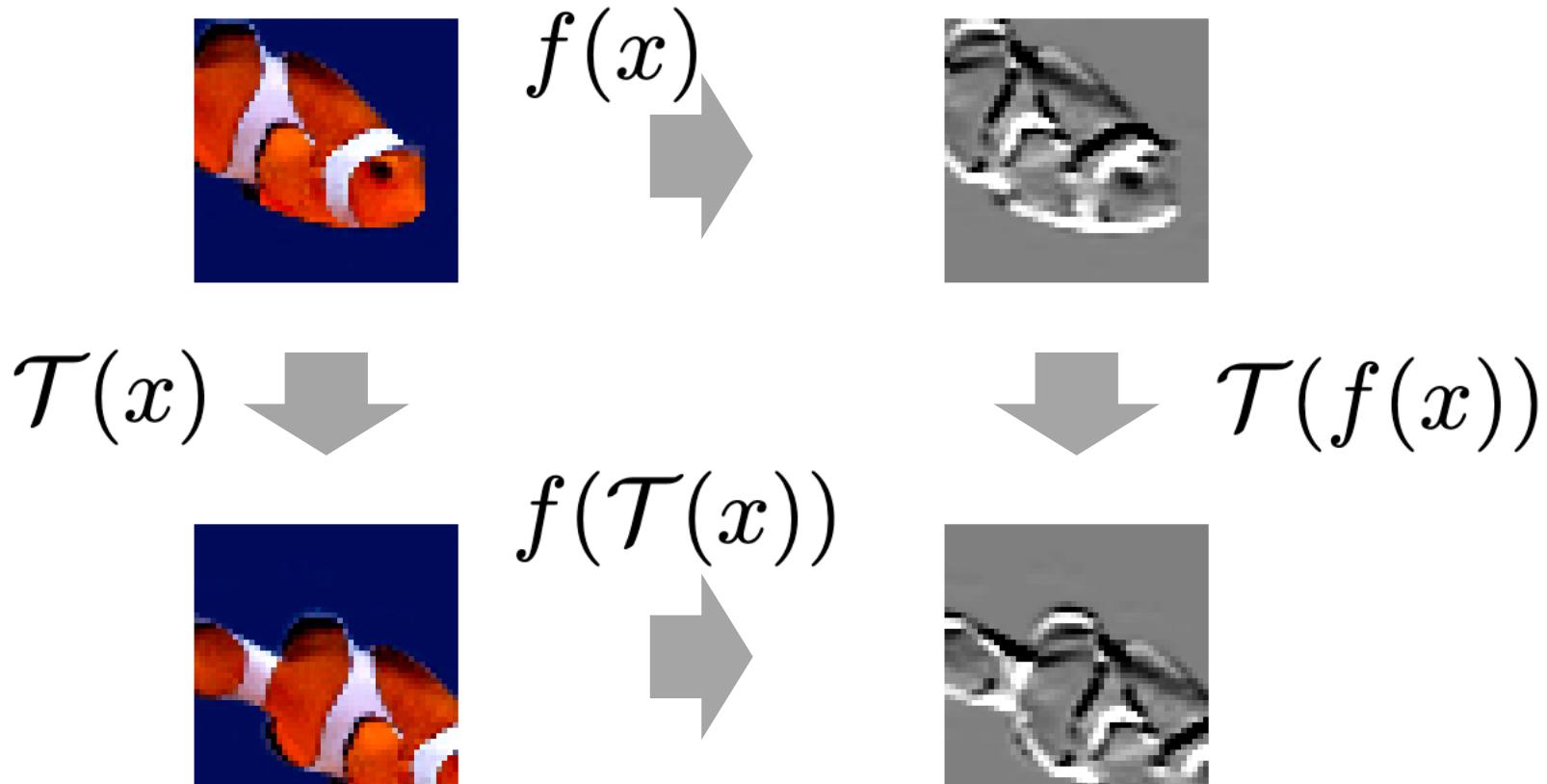
Convolution: translation-equivariance

- Equivariance: $f(\mathcal{T}(x)) = \mathcal{T}(f(x))$



Convolution: translation-equivariance

- Equivariance: $f(\mathcal{T}(x)) = \mathcal{T}(f(x))$

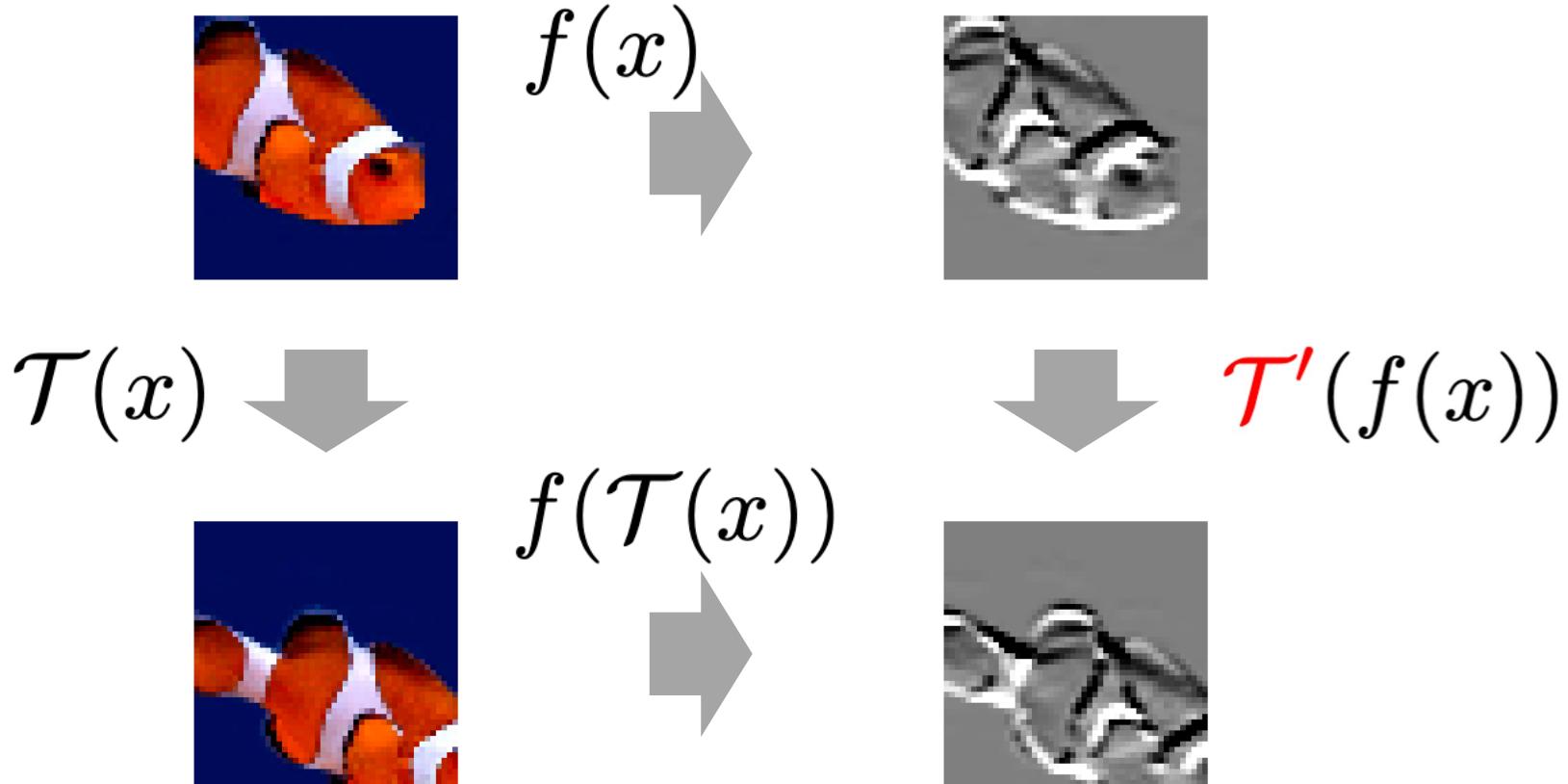


* assuming infinite canvas

Convolution: translation-equivariance

- Equivariance: $f(\mathcal{T}(x)) = \mathcal{T}'(f(x))$

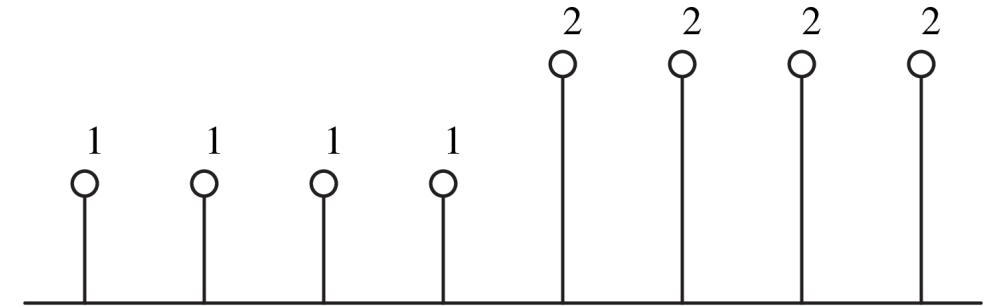
\mathcal{T}' does not depend on x



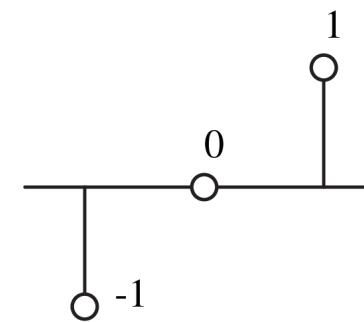
* assuming infinite canvas

Convolution = local connection + weight-sharing

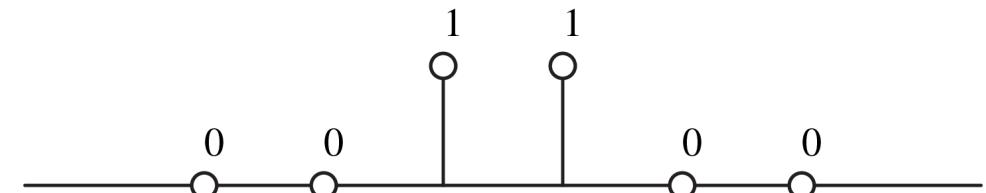
input



filter

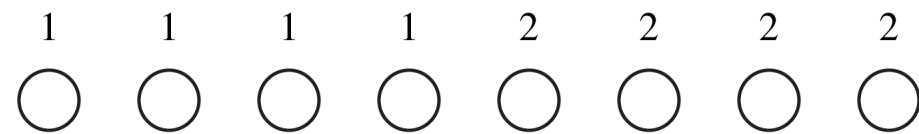


output



Convolution = local connection + weight-sharing

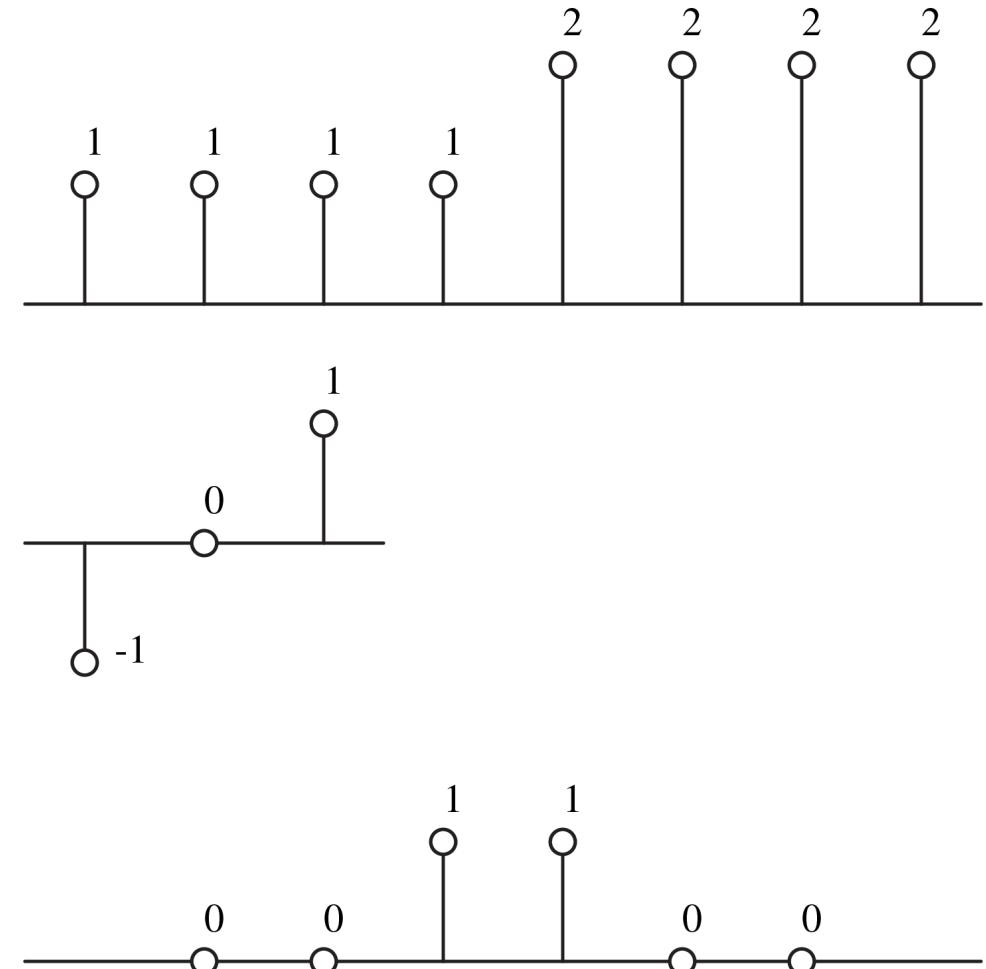
Let's plot it in a different way!



input

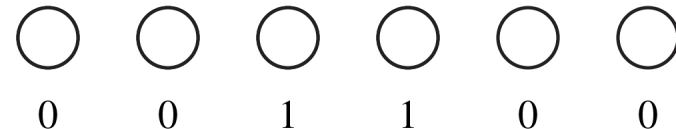
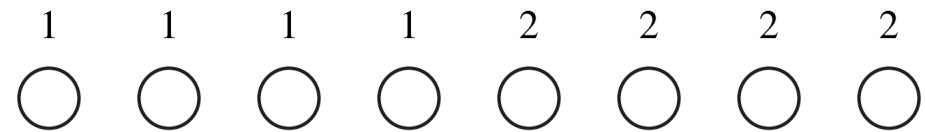
filter

output

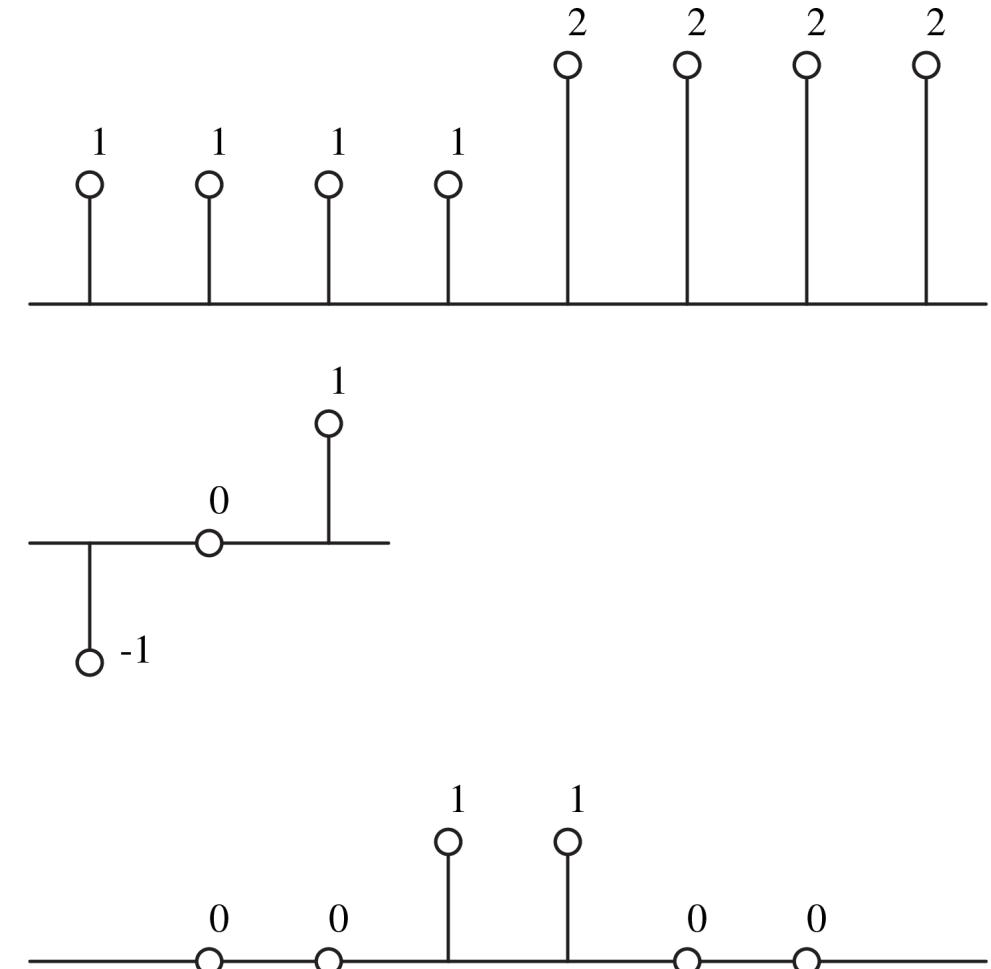


Convolution = local connection + weight-sharing

Let's plot it in a different way!

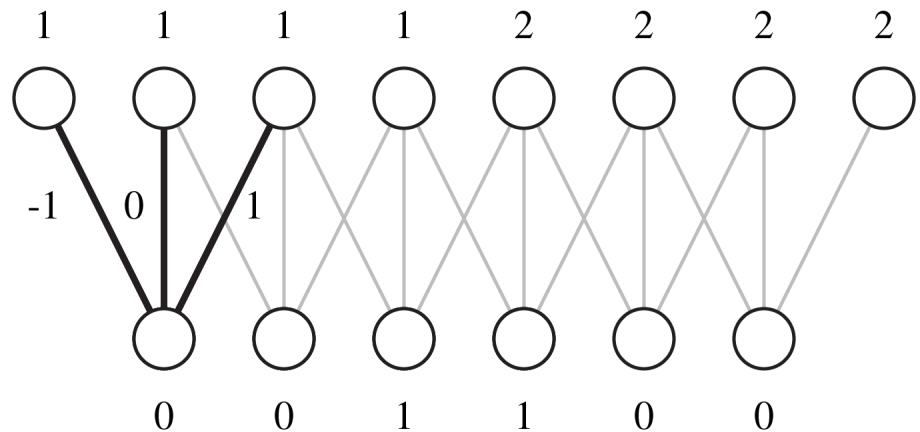


output



Convolution = local connection + weight-sharing

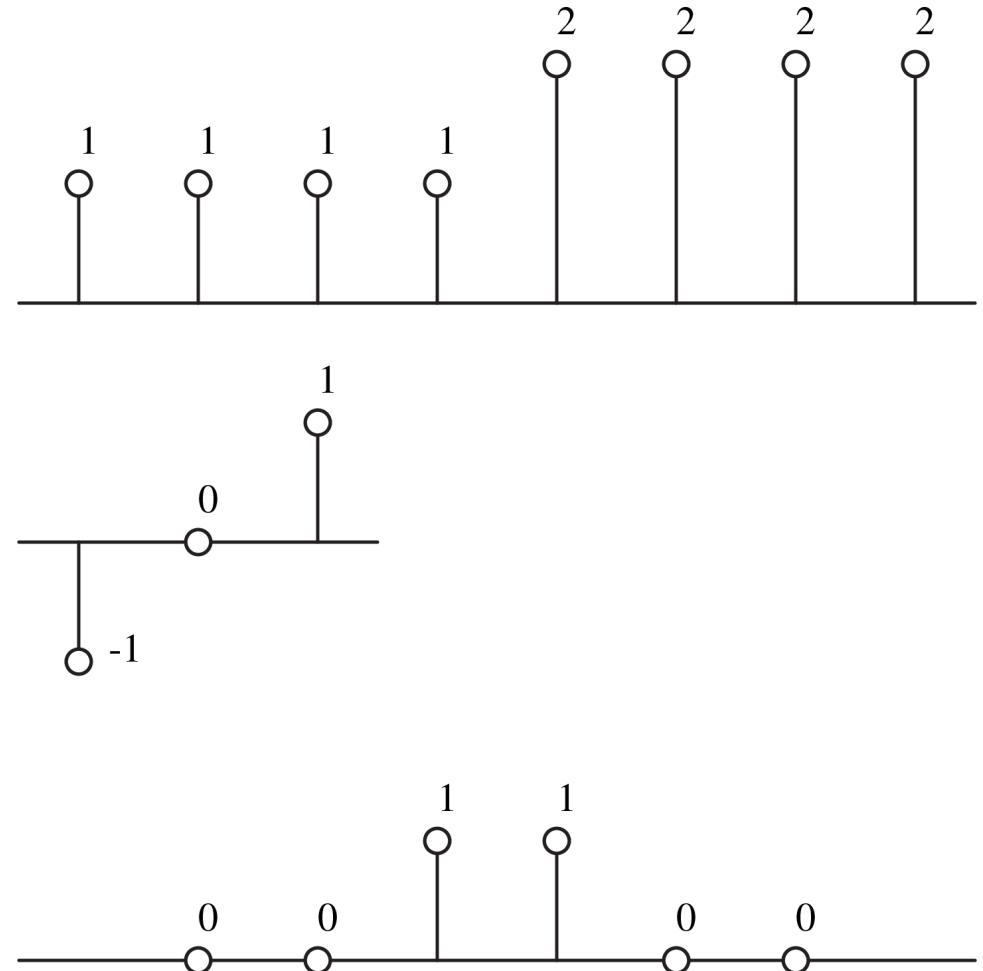
Let's plot it in a different way!



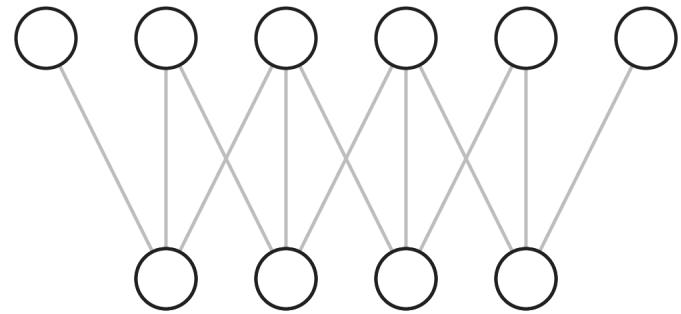
input

filter

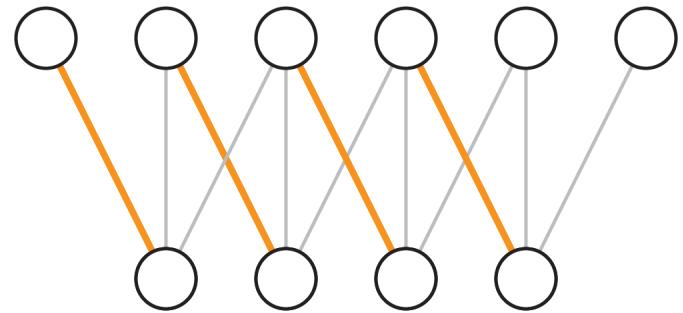
output



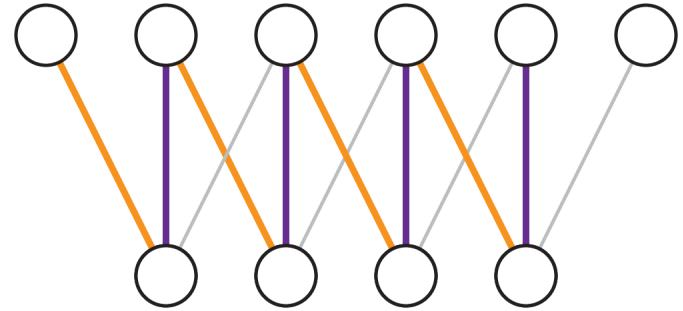
Convolution = local connection + weight-sharing



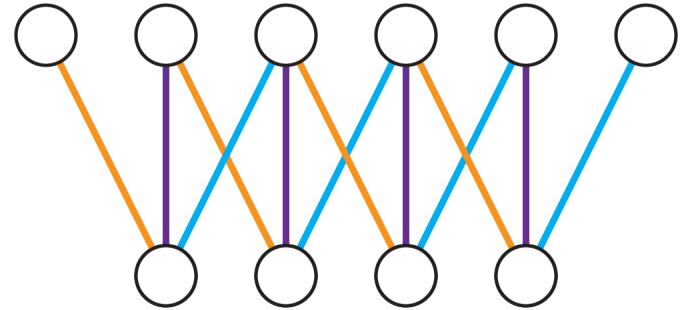
Convolution = local connection + weight-sharing



Convolution = local connection + weight-sharing

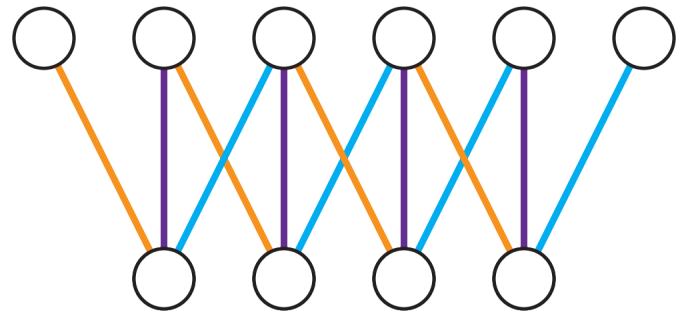


Convolution = local connection + weight-sharing



Convolution = local connection + weight-sharing

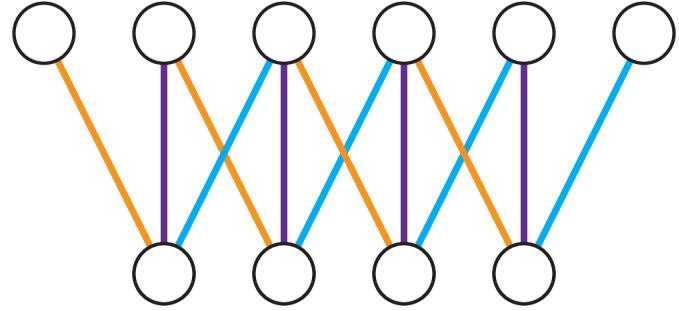
convolution



params = 3 here

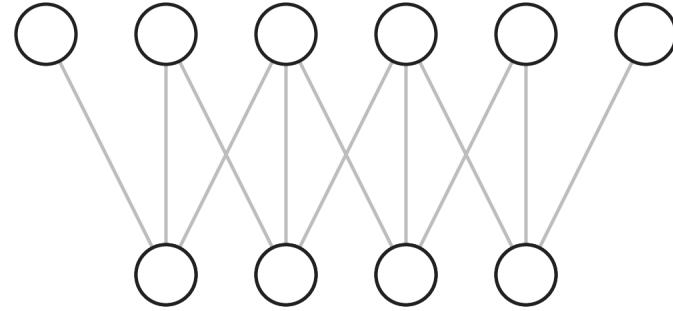
Convolution = local connection + weight-sharing

convolution



params = 3 here

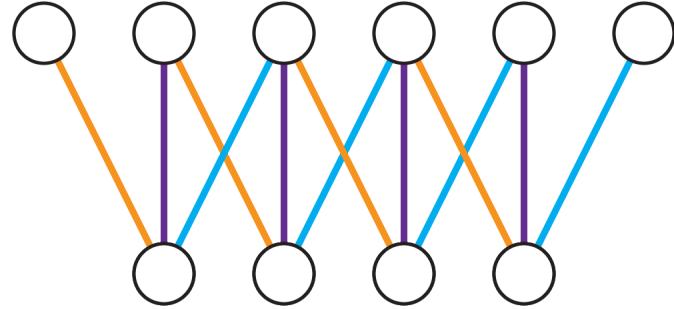
locally connected



params = 12 here
no weight-sharing

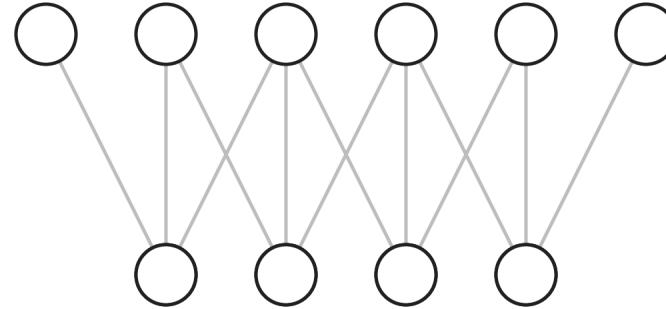
Convolution = local connection + weight-sharing

convolution



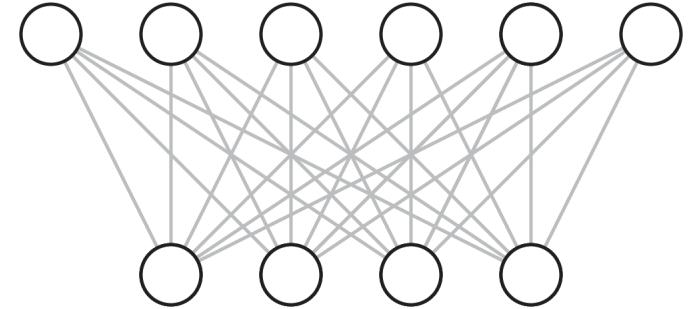
params = 3 here

locally connected



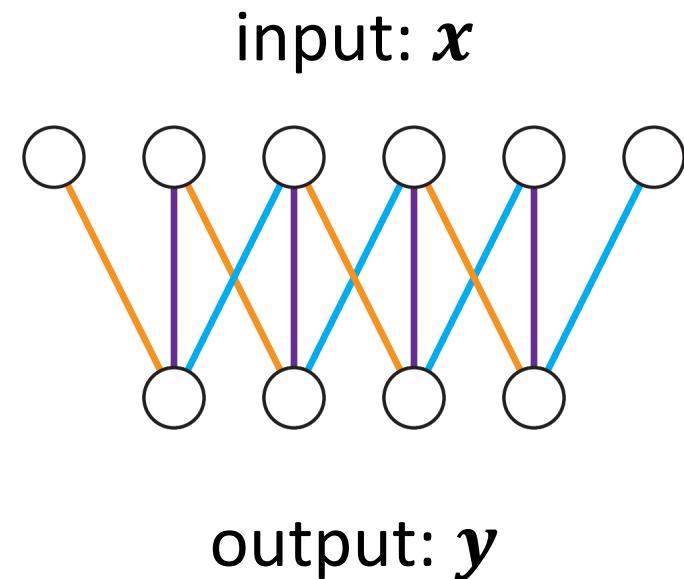
params = 12 here
no weight-sharing

fully connected

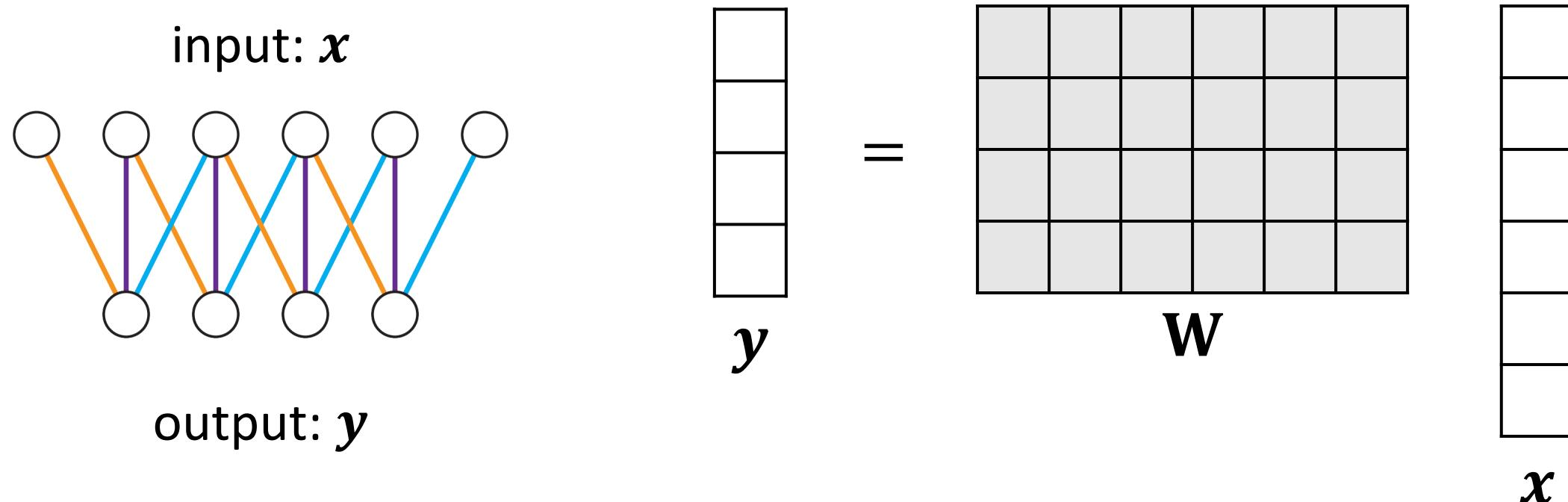


params = 24 here
no weight-sharing

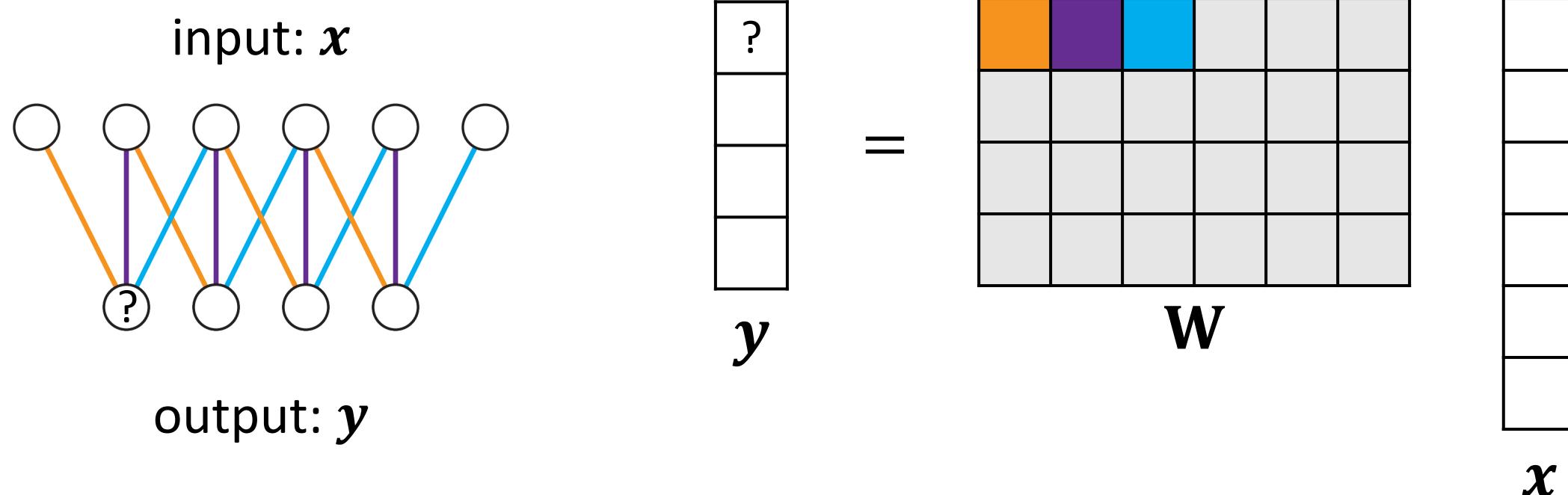
Convolution: linear transform



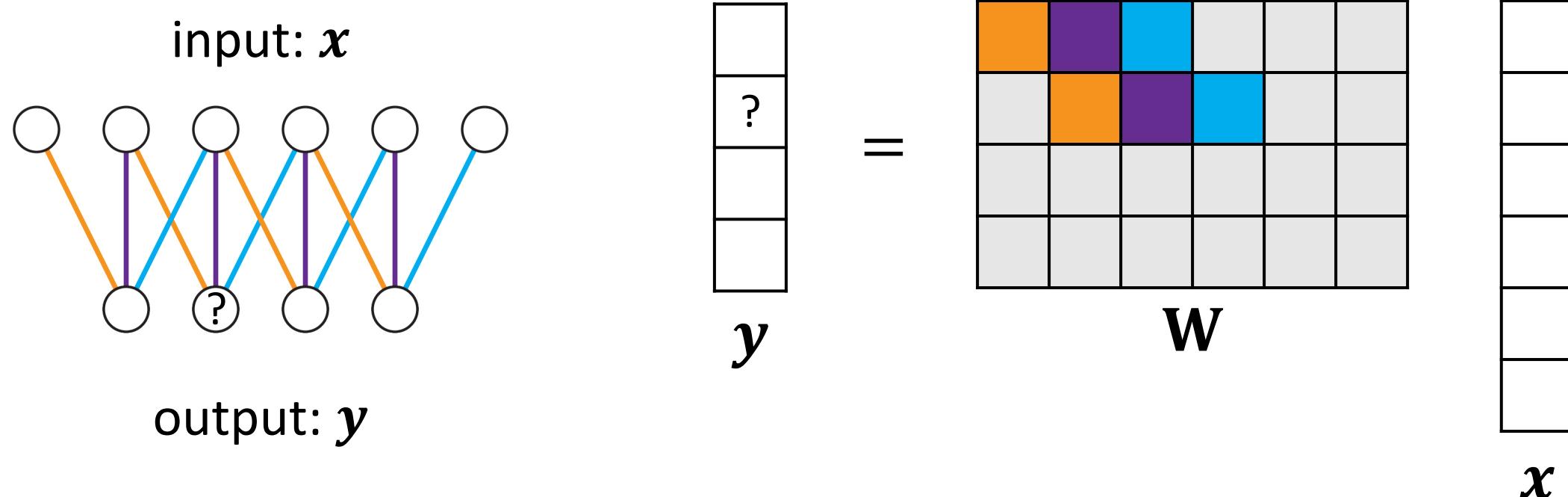
Convolution: linear transform



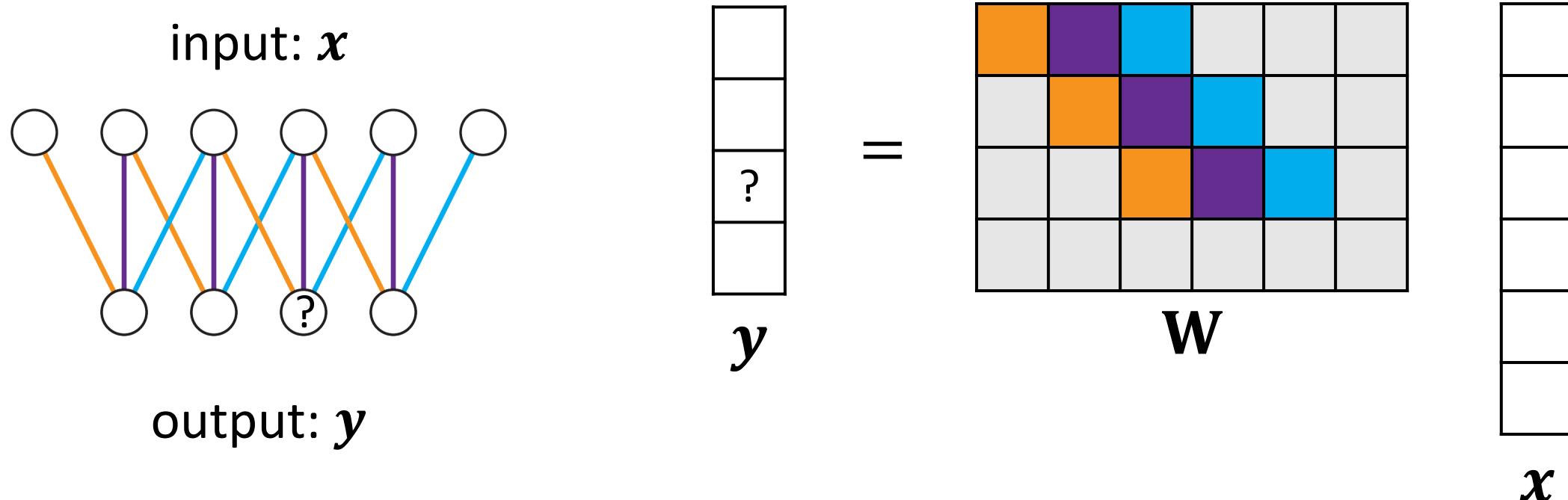
Convolution: linear transform



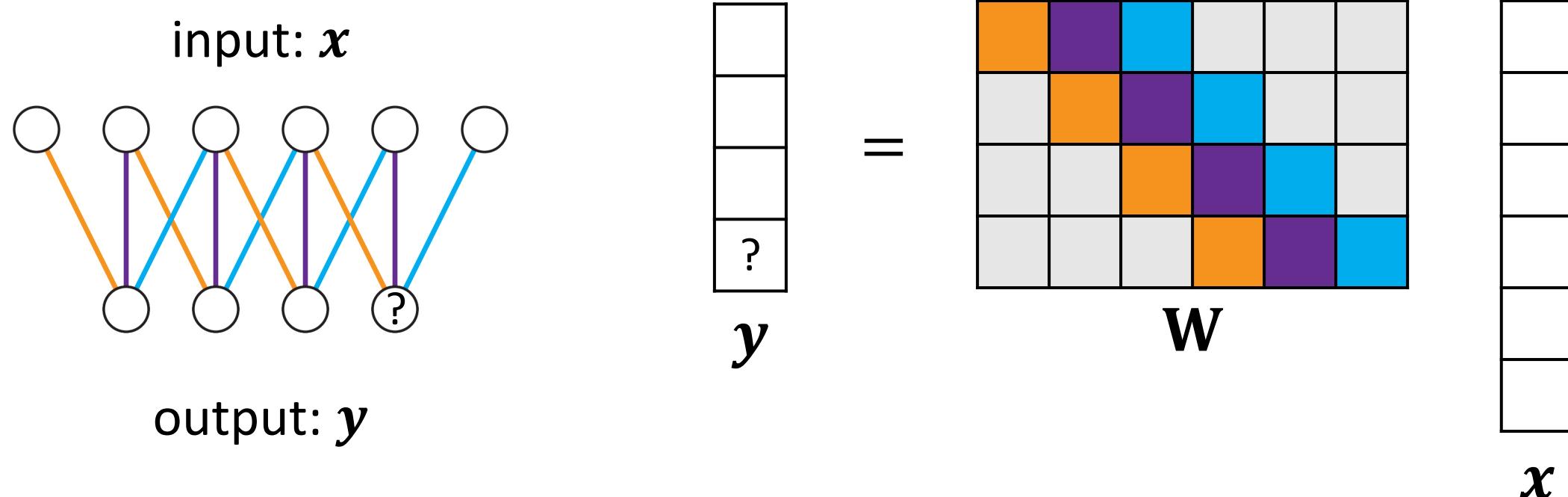
Convolution: linear transform



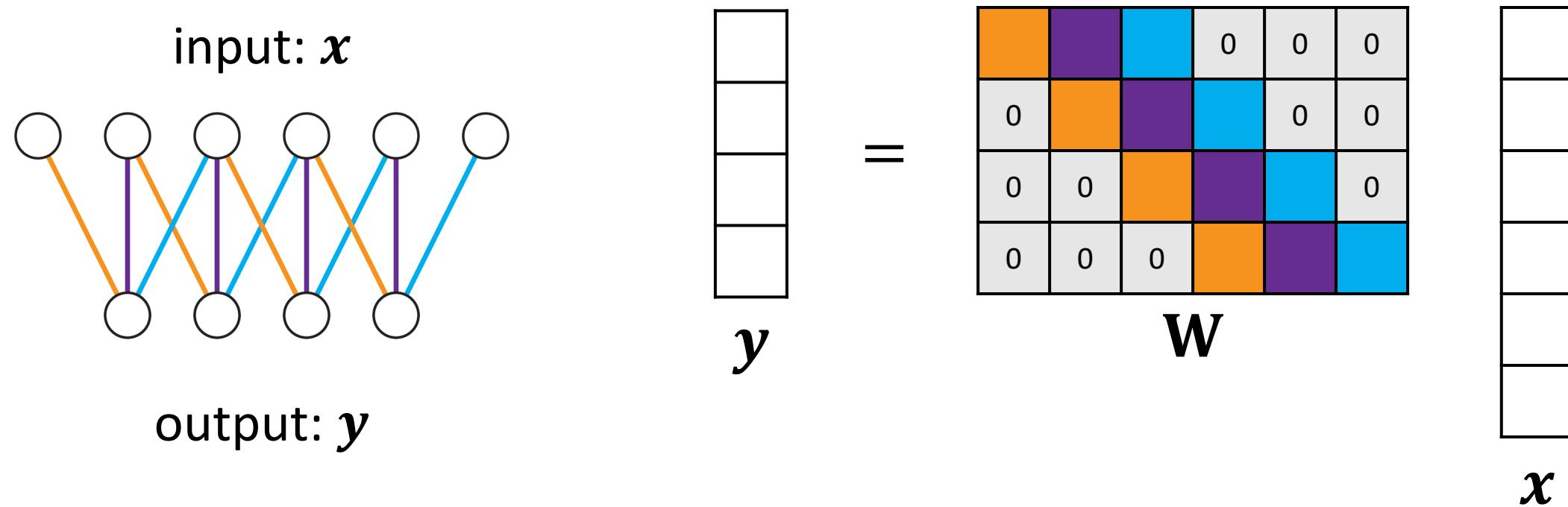
Convolution: linear transform



Convolution: linear transform

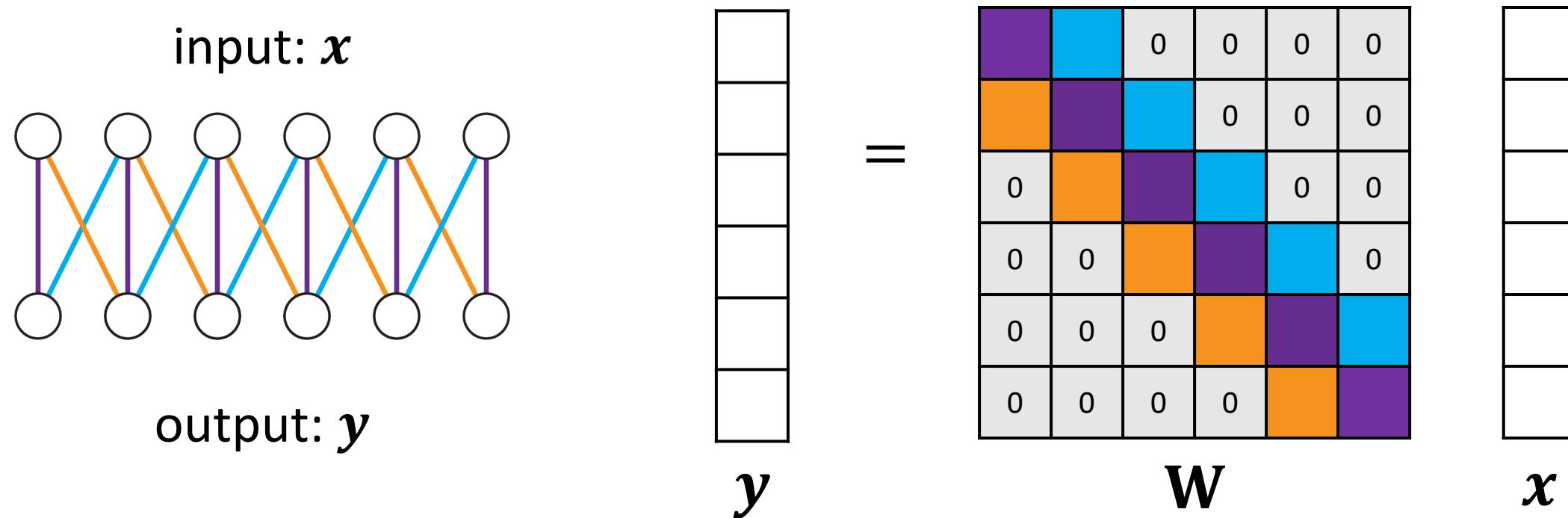


Convolution: linear transform



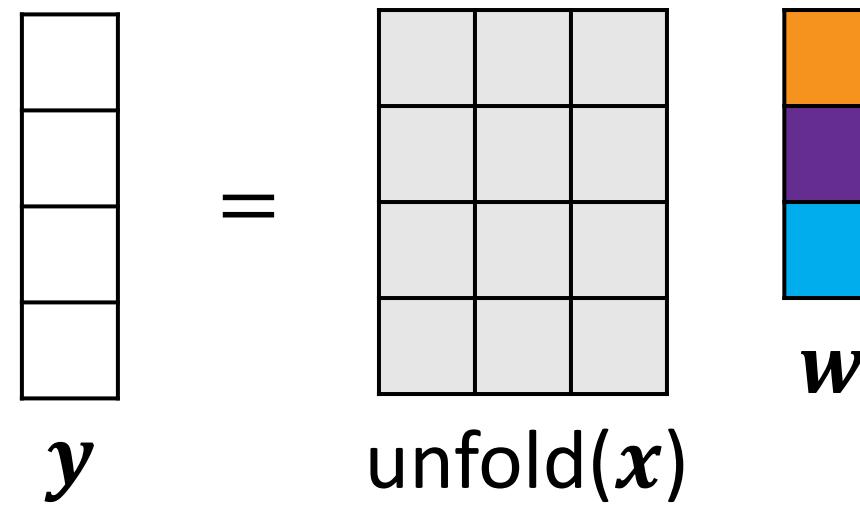
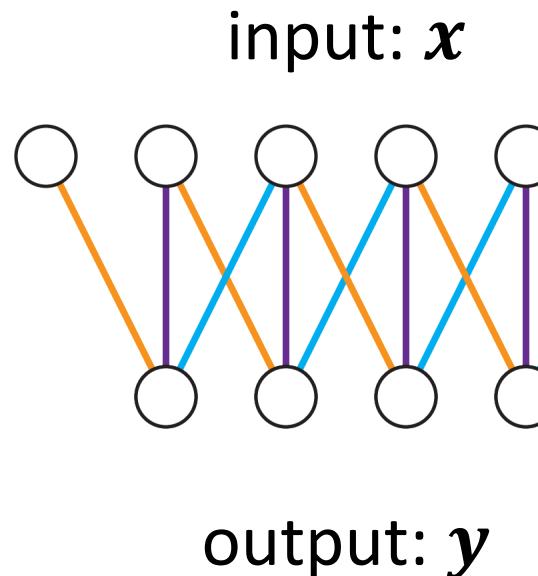
- **local connection:** W is sparse
- **weight-sharing:** # params = 3 here in W

Convolution: linear transform



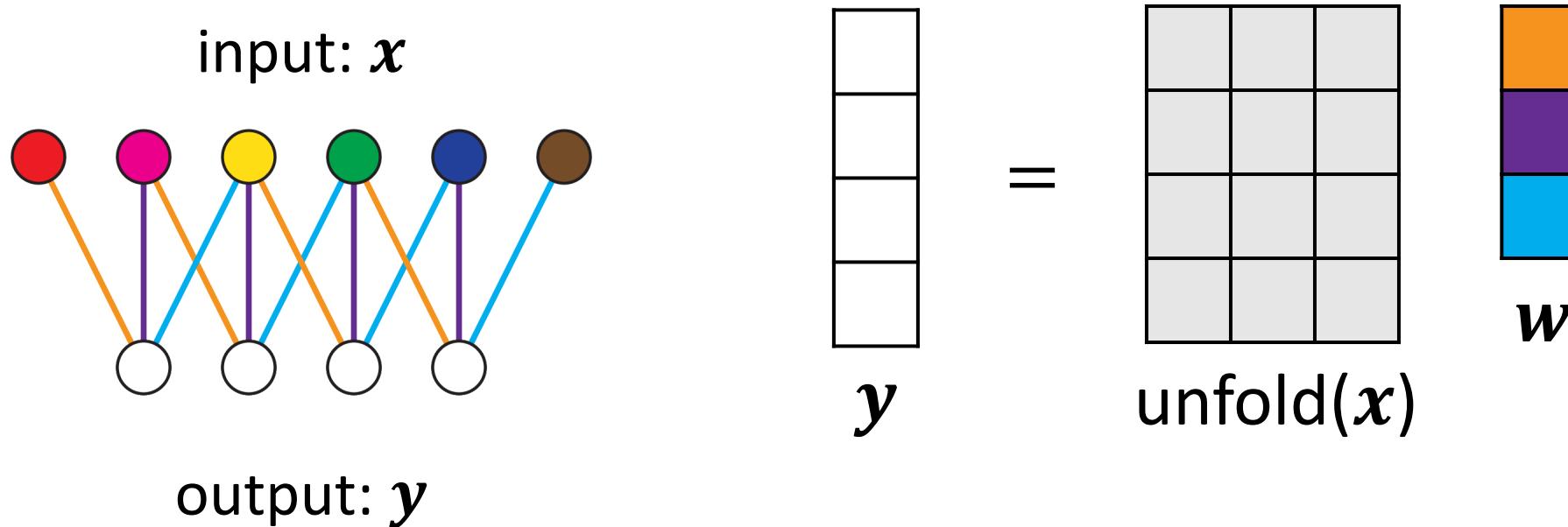
- This is the case with padding
- *Exercise:* what about stride = 2?
- *Exercise:* what about 2-D convolution?

Convolution: unfold / im2col

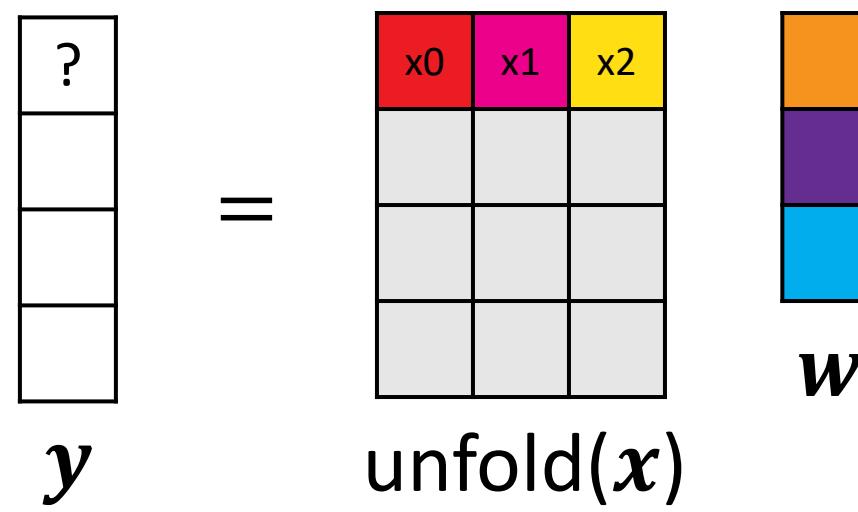
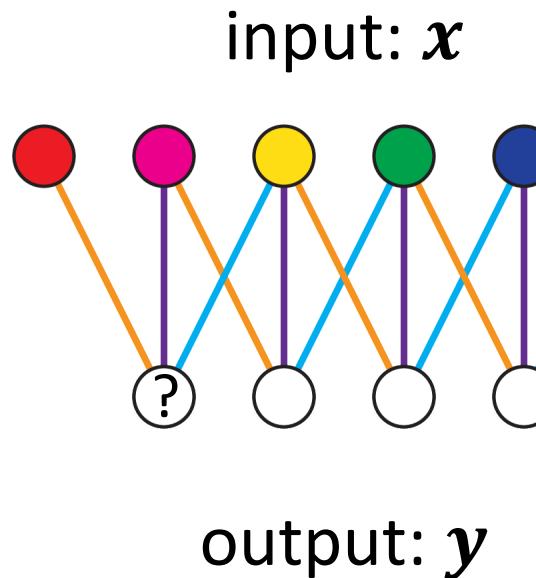


Can we find another linear transform,
such that we don't need sparse matmul?

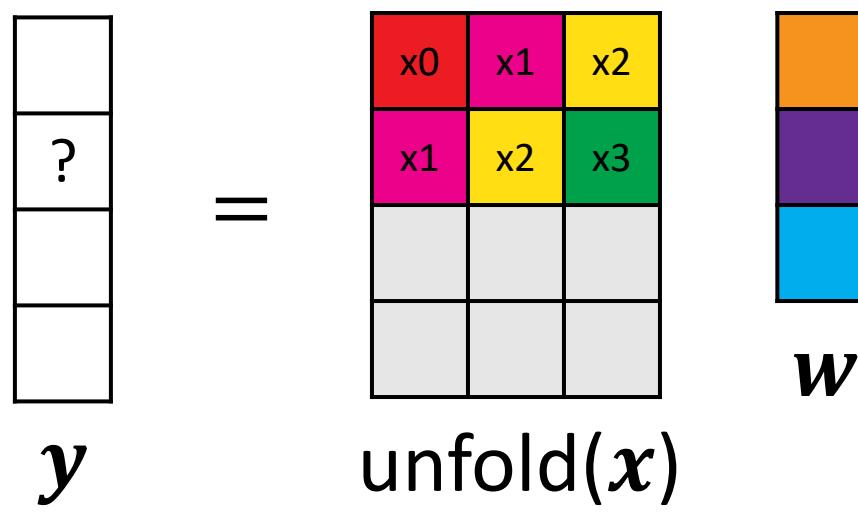
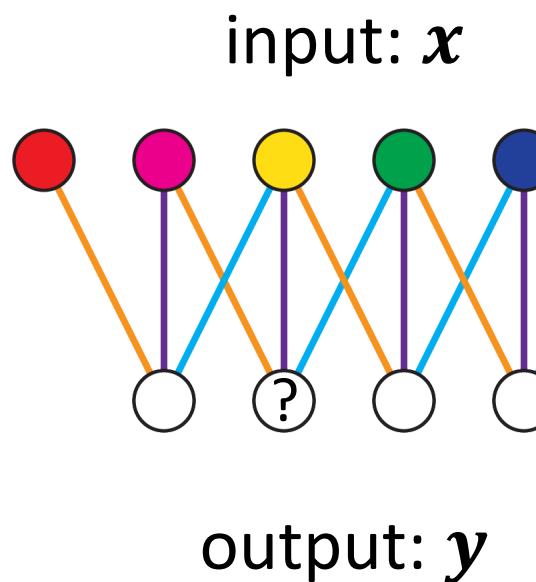
Convolution: unfold / im2col



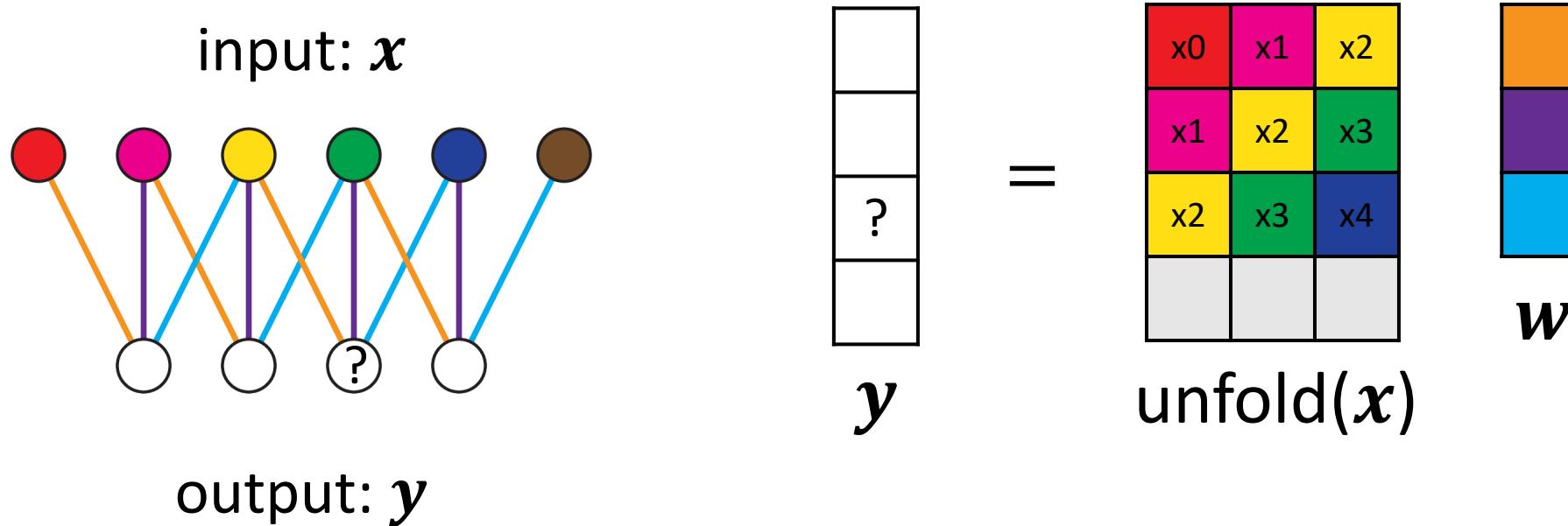
Convolution: unfold / im2col



Convolution: unfold / im2col

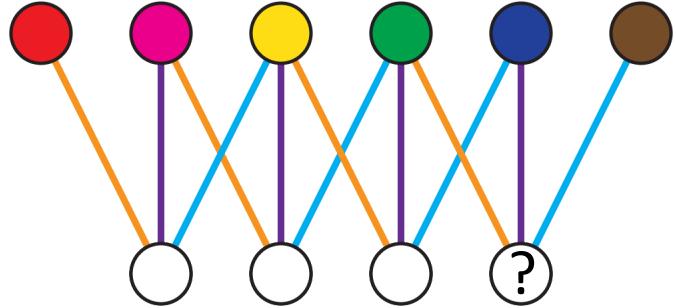


Convolution: unfold / im2col

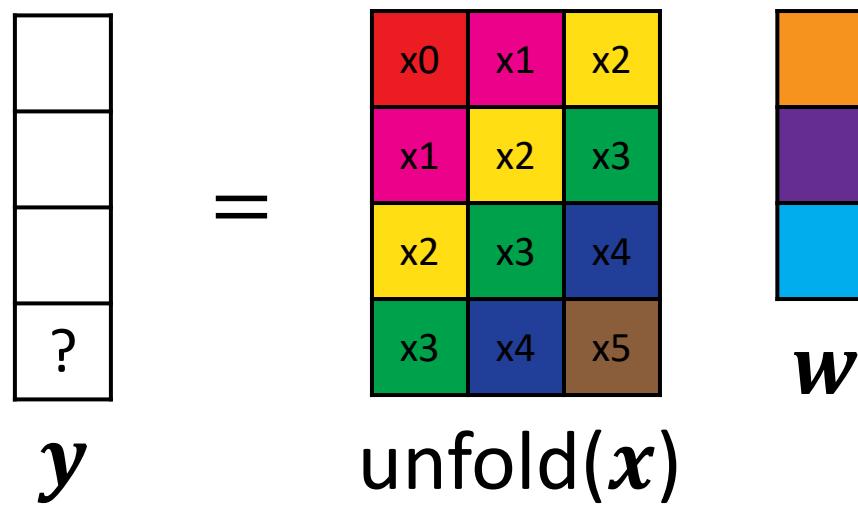


Convolution: unfold / im2col

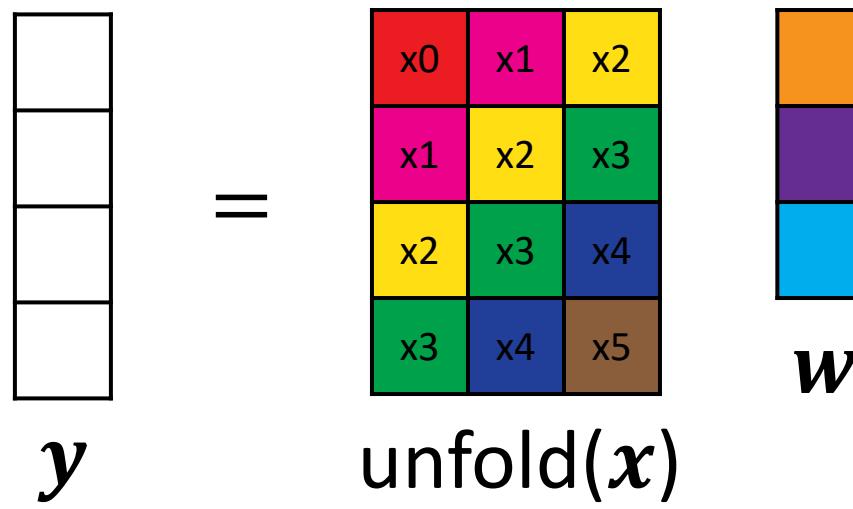
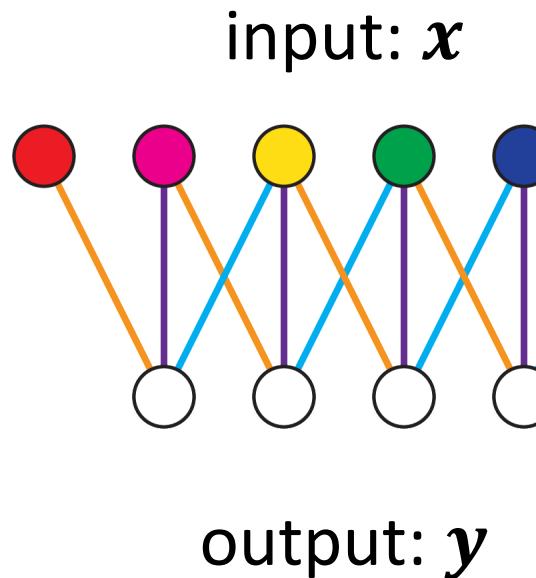
input: x



output: y



Convolution: unfold / im2col

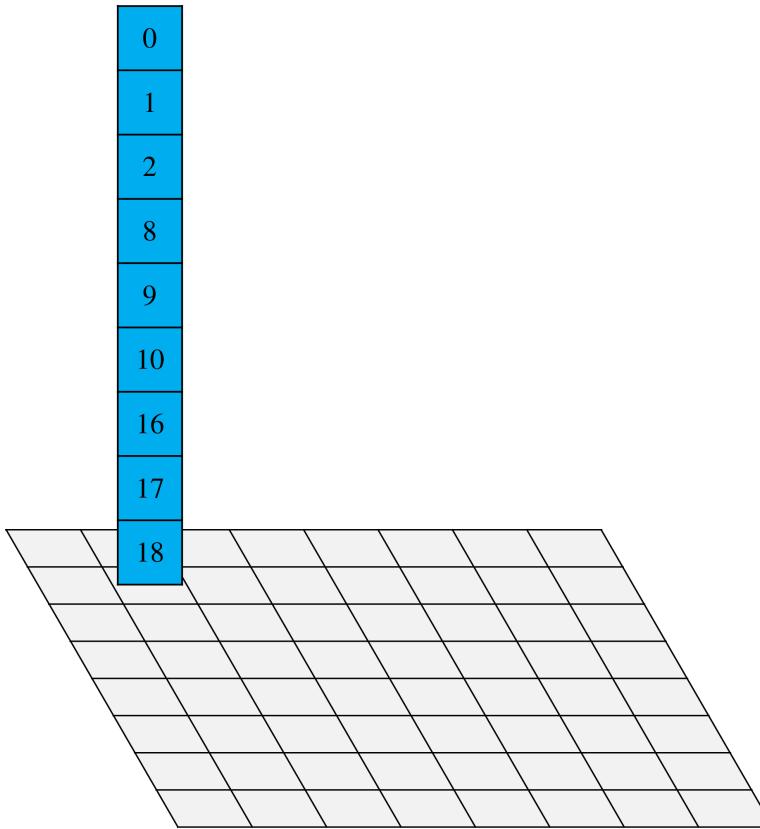


- this transform is called “unfold” (PyTorch)
- in 2D, it’s also called “im2col” (Matlab)
- w is dense; can easily do multiple filters
- “unfold” in itself is convolution (*Exercise: why?*)

Convolution: unfold / im2col

0	1	2	3	4	5	6	7
8	9	10	11	12	13	14	15
16	17	18	19	20	21	22	23
24	25	26	27	28	29	30	31
32	33	34	35	36	37	38	39
40	41	42	43	44	45	46	47
48	49	50	51	52	53	54	55
56	57	58	59	60	61	62	63

2D im

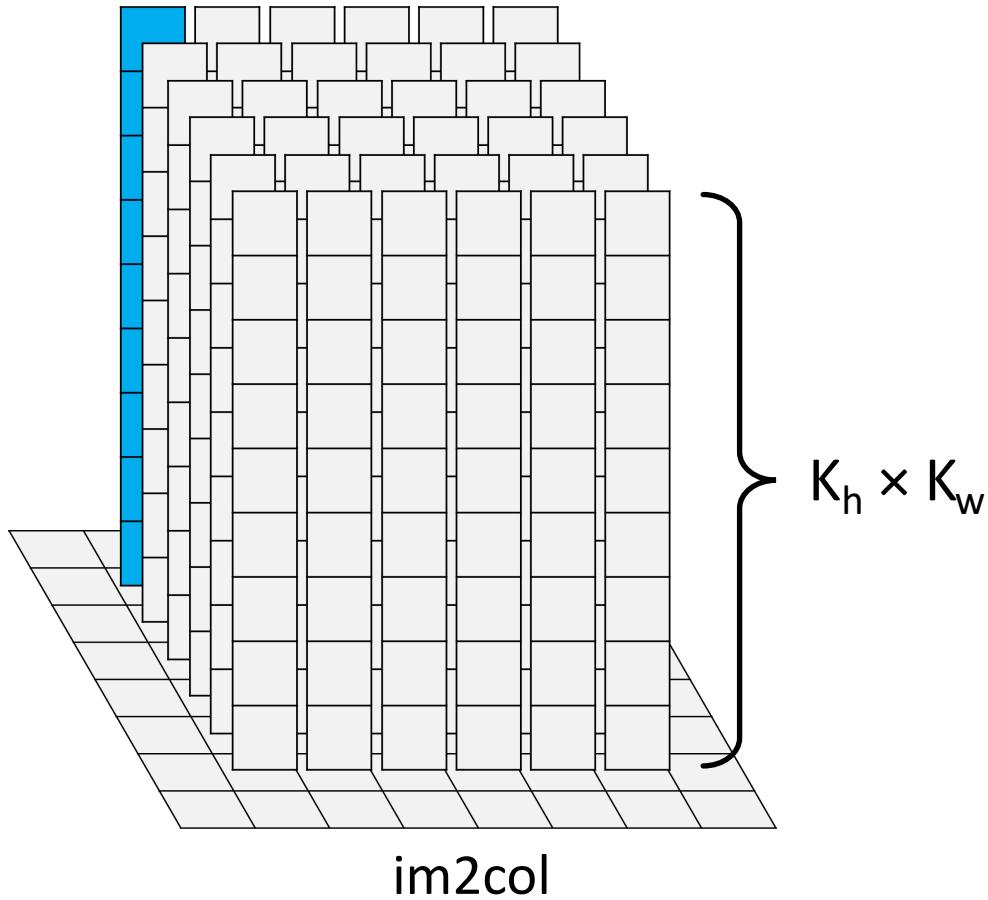


im2col

Convolution: unfold / im2col

0	1	2	3	4	5	6	7
8	9	10	11	12	13	14	15
16	17	18	19	20	21	22	23
24	25	26	27	28	29	30	31
32	33	34	35	36	37	38	39
40	41	42	43	44	45	46	47
48	49	50	51	52	53	54	55
56	57	58	59	60	61	62	63

2D im



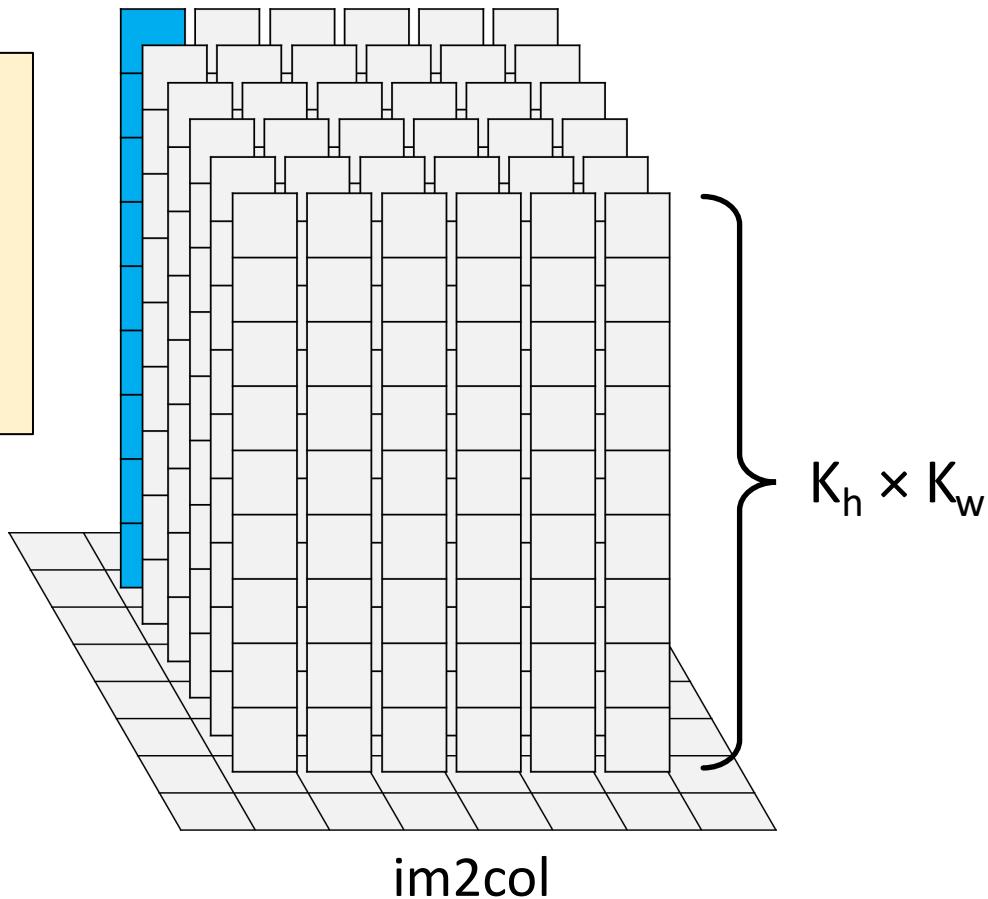
Convolution: unfold / im2col

im2col

- gather local information
- flatten into a “col” ($K_h \times K_w$ channels)
- followed by 1×1 conv

0	1	2	3	4	5	6	7
8	9	10	11	12	13	14	15
16	17	18	19	20	21	22	23
24	25	26	27	28	29	30	31
32	33	34	35	36	37	38	39
40	41	42	43	44	45	46	47
48	49	50	51	52	53	54	55
56	57	58	59	60	61	62	63

2D im

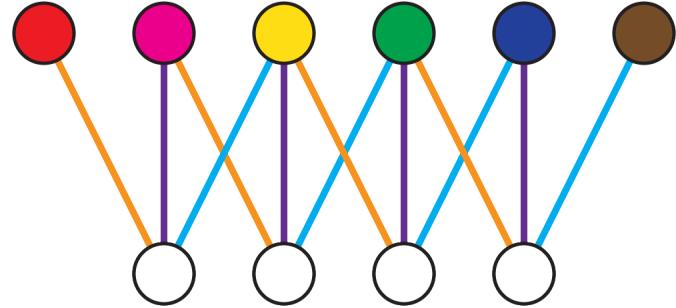


im2col

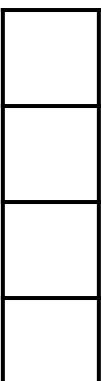
$$K_h \times K_w$$

Two views of convolution

input: x

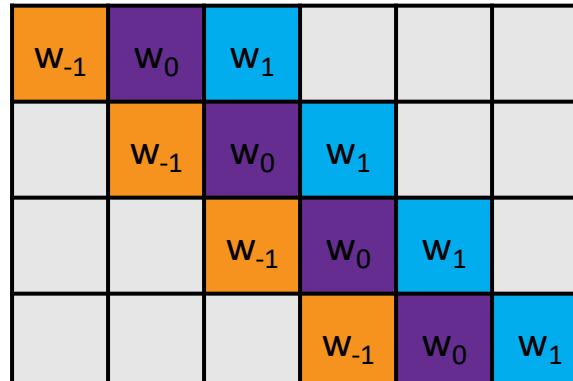


output: y



y

=



sliding w

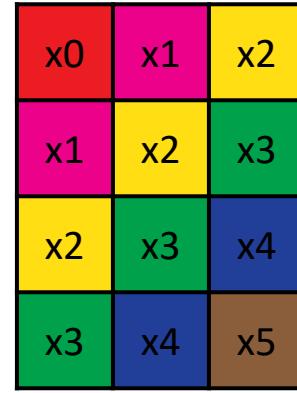


x



y

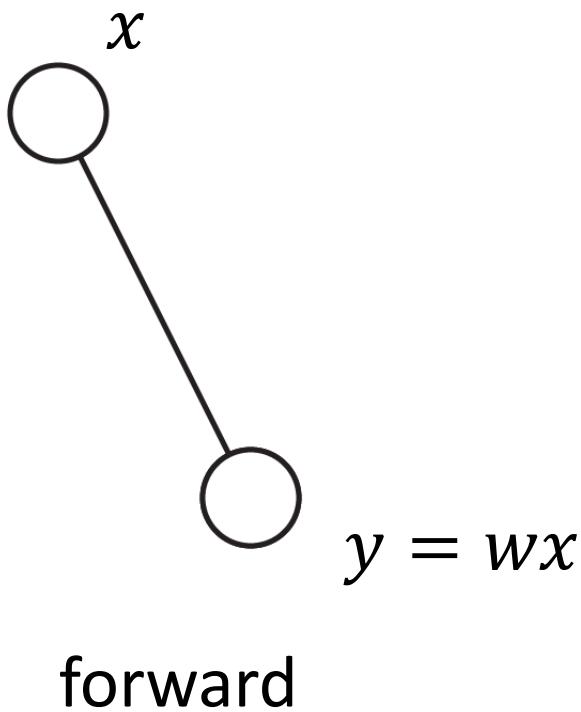
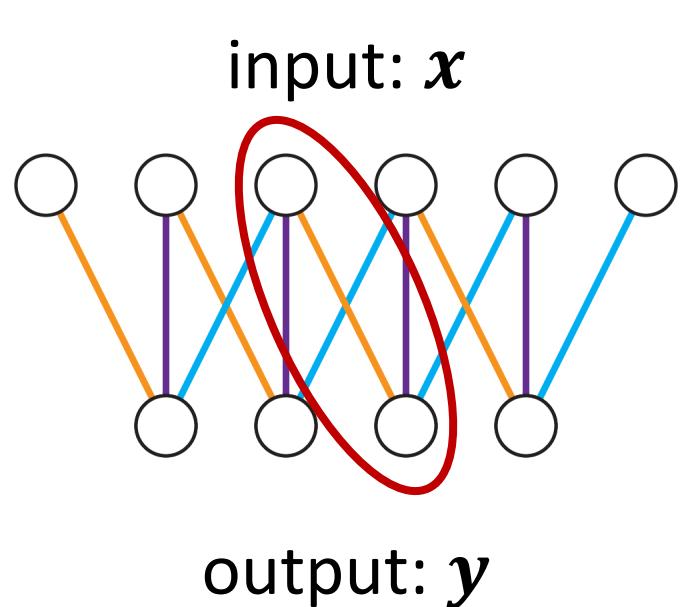
=



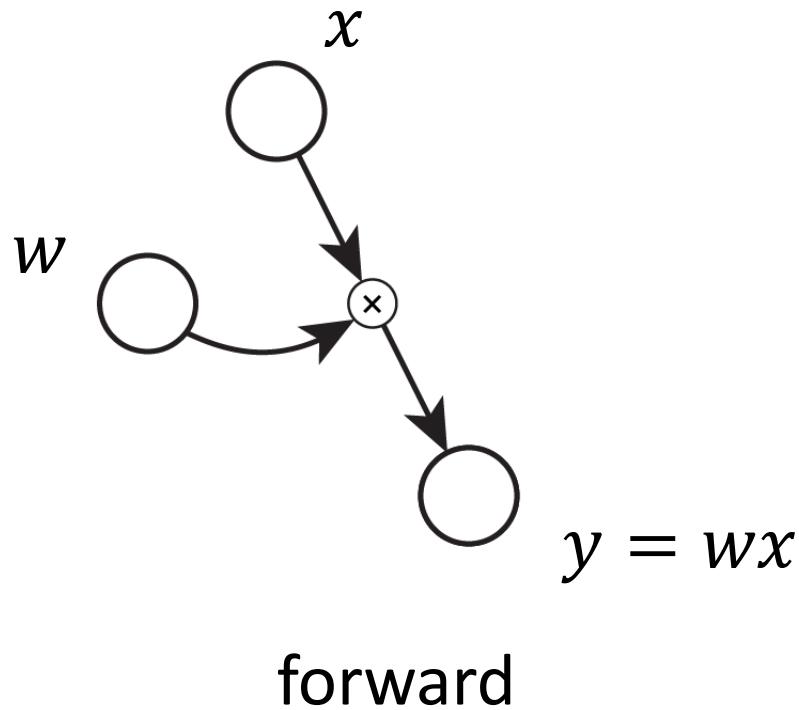
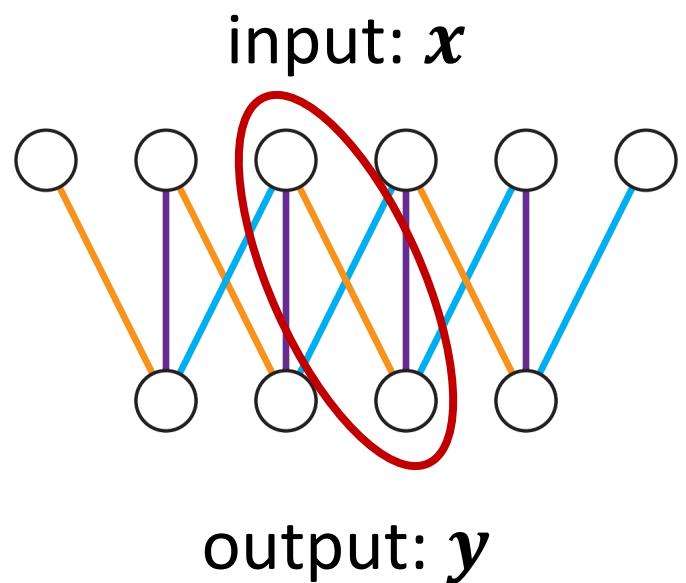
sliding x

w

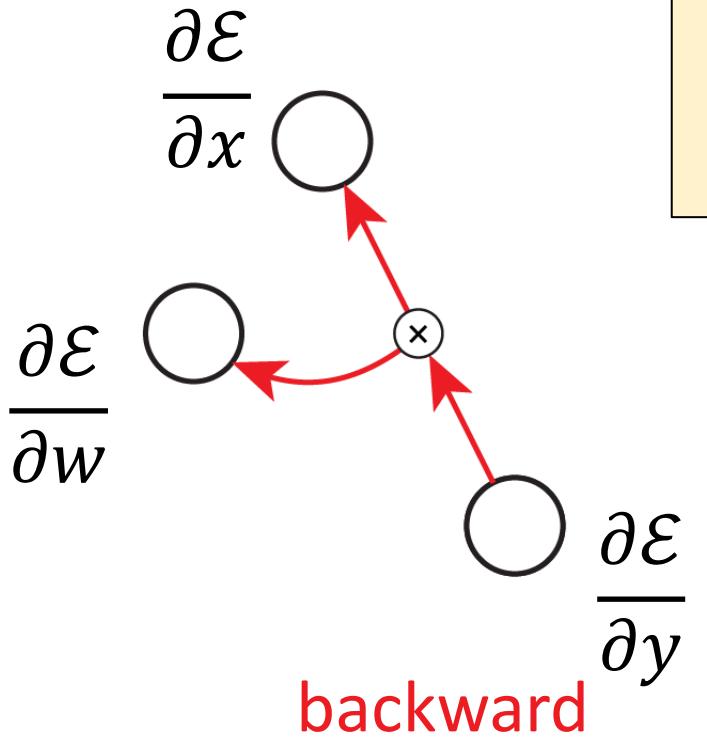
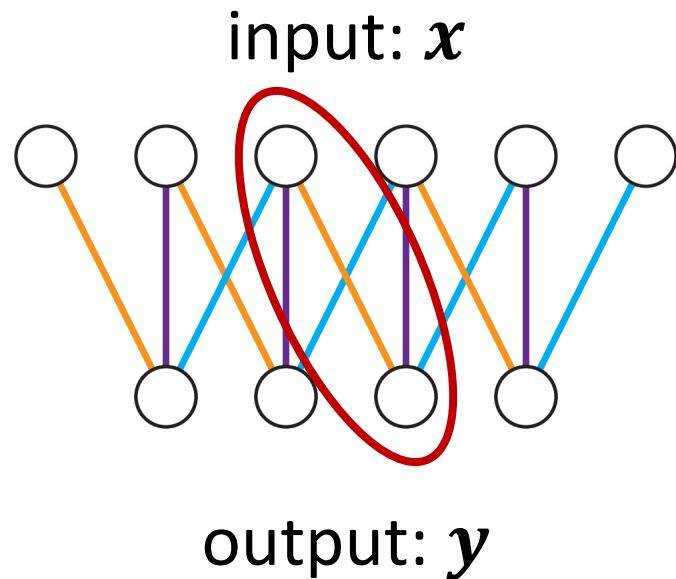
Convolution: BackProp



Convolution: BackProp



Convolution: BackProp

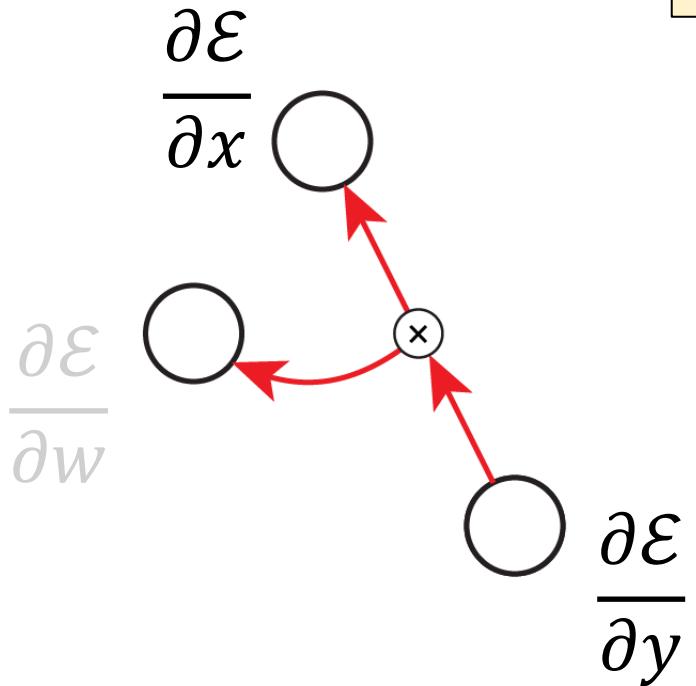
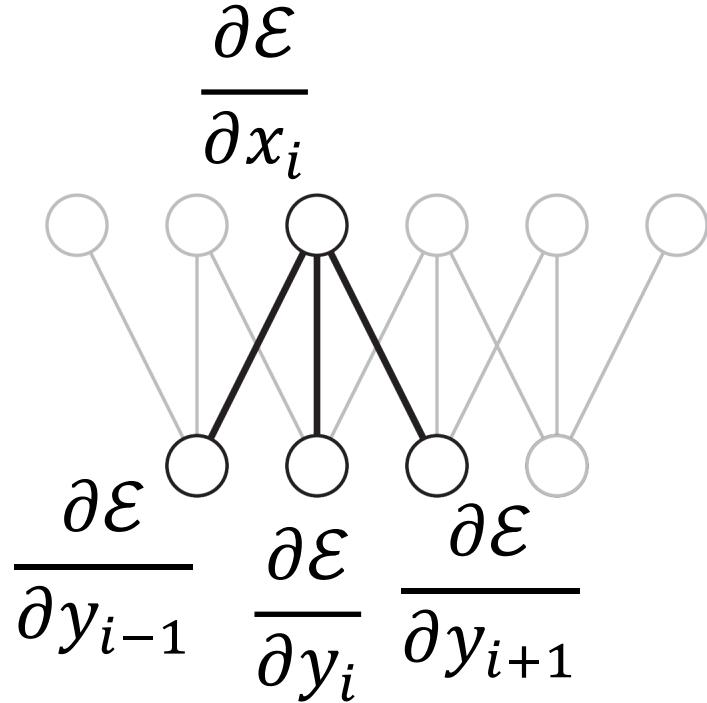


chain rule (scalars):

$$\frac{\partial \mathcal{E}}{\partial x} = \frac{\partial \mathcal{E}}{\partial y} \frac{\partial y}{\partial x} = \frac{\partial \mathcal{E}}{\partial y} w$$

$$\frac{\partial \mathcal{E}}{\partial w} = \frac{\partial \mathcal{E}}{\partial y} \frac{\partial y}{\partial w} = \frac{\partial \mathcal{E}}{\partial y} x$$

Convolution: BackProp



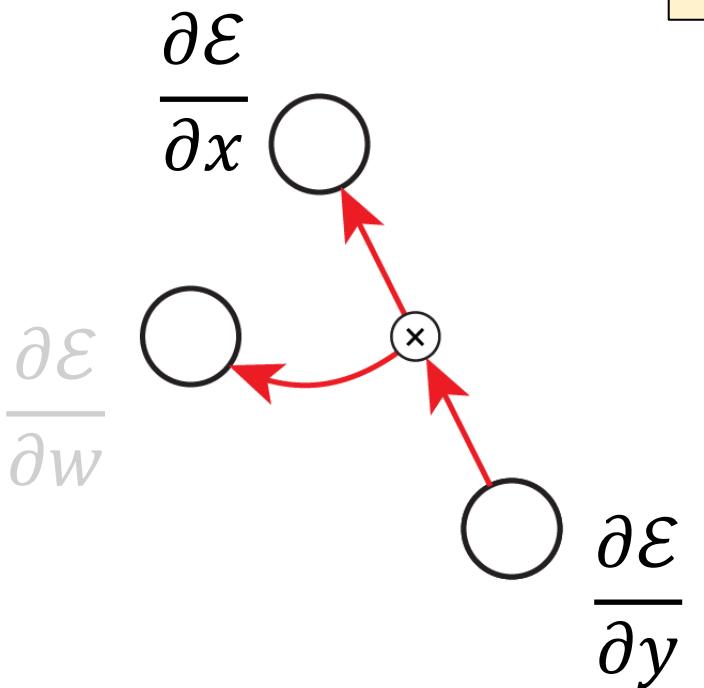
chain rule (scalars):

$$\frac{\partial \mathcal{E}}{\partial x} = \frac{\partial \mathcal{E}}{\partial y} \frac{\partial y}{\partial x} = \frac{\partial \mathcal{E}}{\partial y} w$$

Convolution: BackProp

$$\frac{\partial \mathcal{E}}{\partial x_i} = \sum_{j=-1}^1 \frac{\partial \mathcal{E}}{\partial y_{i+j}} w_{-j}$$

$\frac{\partial \mathcal{E}}{\partial y_{i-1}}$ $\frac{\partial \mathcal{E}}{\partial y_i}$ $\frac{\partial \mathcal{E}}{\partial y_{i+1}}$

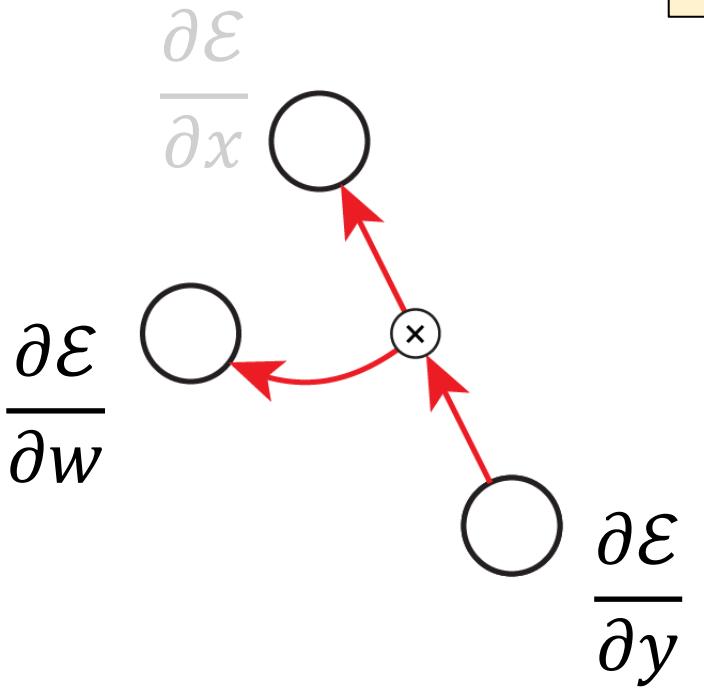
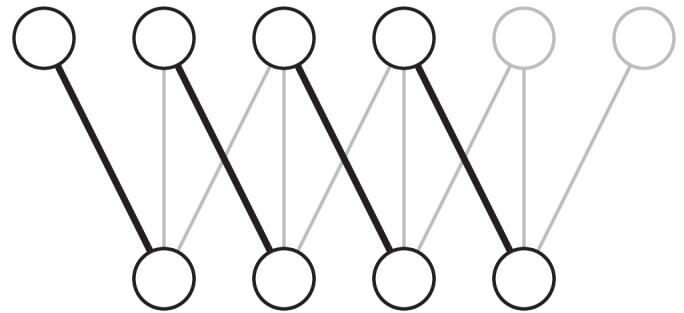


backprop wrt x is conv

chain rule (scalars):

$$\frac{\partial \mathcal{E}}{\partial x} = \frac{\partial \mathcal{E}}{\partial y} \frac{\partial y}{\partial x} = \frac{\partial \mathcal{E}}{\partial y} w$$

Convolution: BackProp

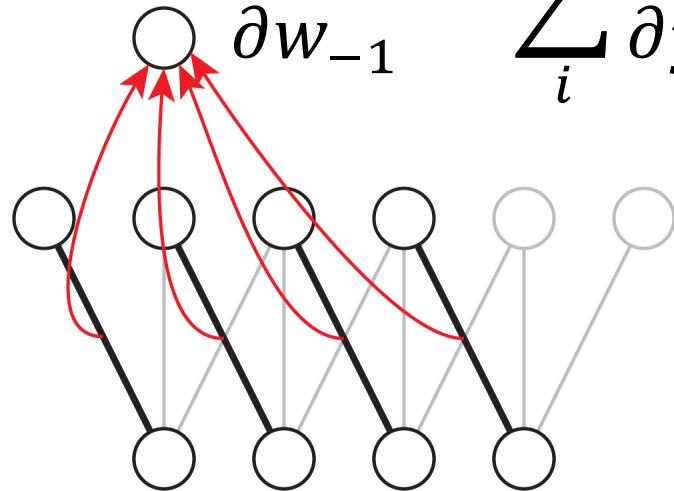


chain rule (scalars):

$$\frac{\partial \mathcal{E}}{\partial w} = \frac{\partial \mathcal{E}}{\partial y} \frac{\partial y}{\partial w} = \frac{\partial \mathcal{E}}{\partial y} x$$

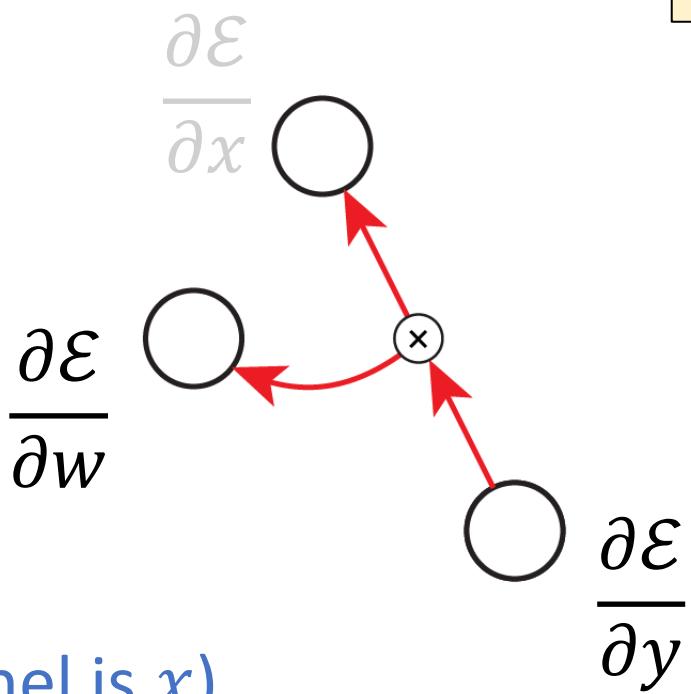
Convolution: BackProp

$$\frac{\partial \mathcal{E}}{\partial w_{-1}} = \sum_i \frac{\partial \mathcal{E}}{\partial y_i} x_{i-1}$$



chain rule (scalars):

$$\frac{\partial \mathcal{E}}{\partial w} = \frac{\partial \mathcal{E}}{\partial y} \frac{\partial y}{\partial w} = \frac{\partial \mathcal{E}}{\partial y} x$$



backprop wrt w is also conv (kernel is x)

Exercise: derive the backprop in matrix forms

Takeaways about Convolution

- sliding window
- local connection + weight-sharing
- linear transform
- translation-equivariant

Convolutional Neural Networks

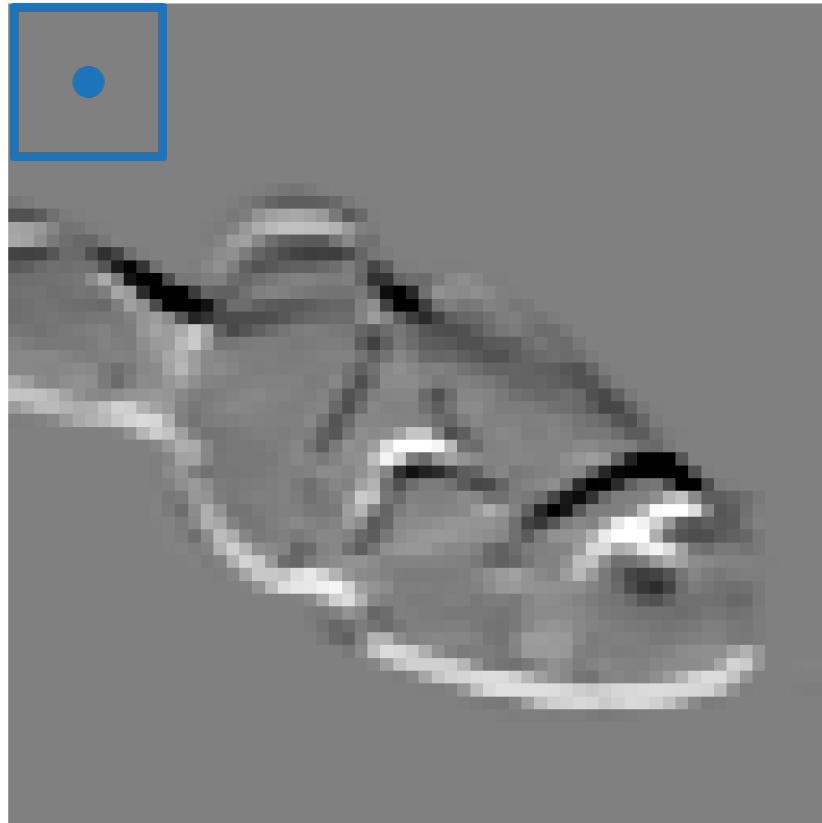
Convolutional Neural Network (ConvNet, CNN)

- Core idea of deep learning:
Composing basic operations into complex functions
- Common operations of ConvNet
 - convolution
 - activation function
 - pooling
 - fully-connected layer
 - softmax



Composing basic operations

output of prev. layer



output of this layer



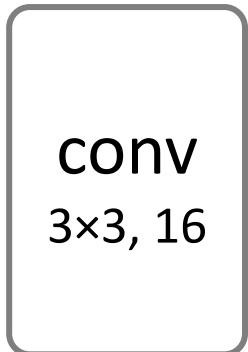
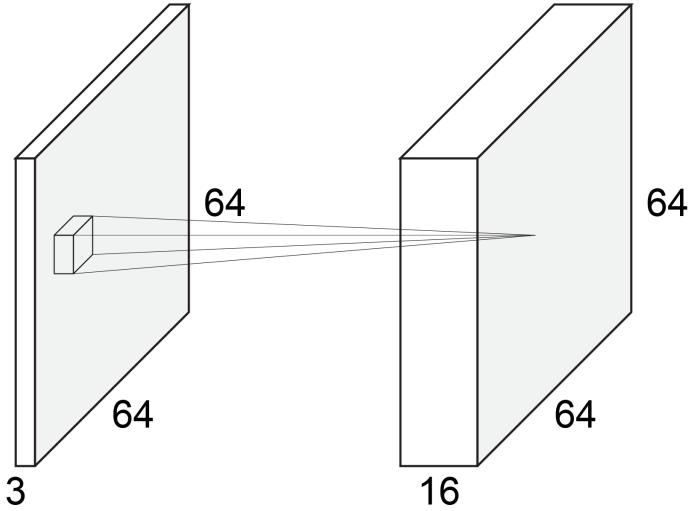
filter

$$\begin{array}{|c|c|c|} \hline & & \\ \hline \end{array}$$

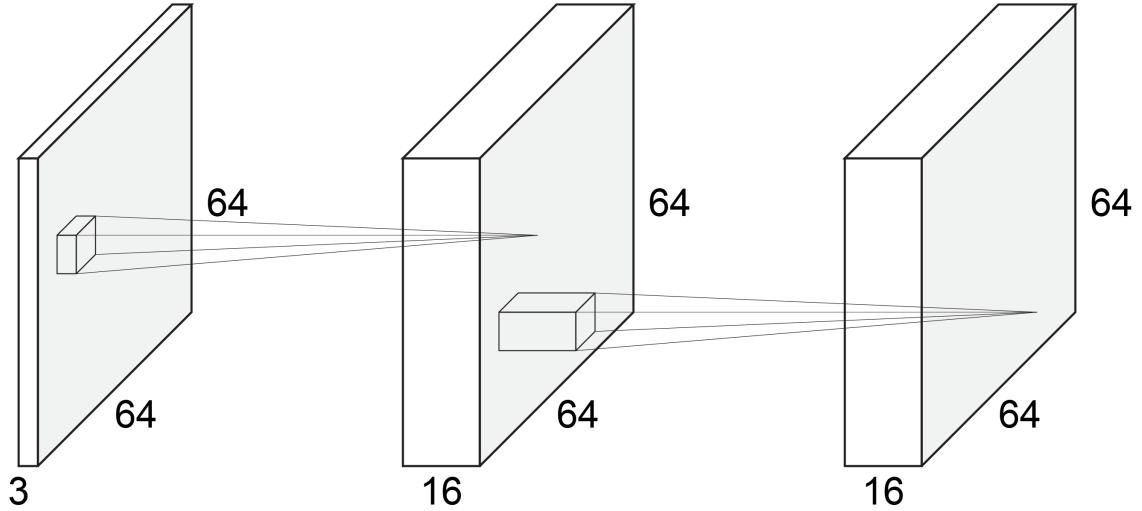
*

=

Composing basic operations



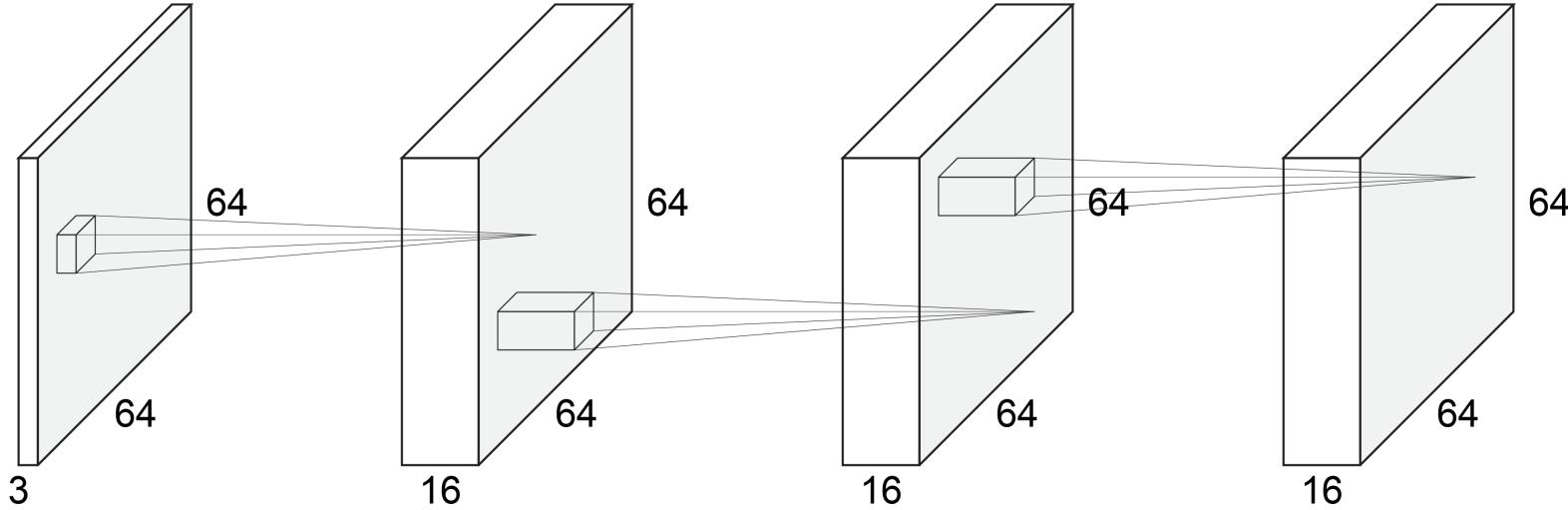
Composing basic operations



conv
3x3, 16

conv
3x3, 16

Composing basic operations

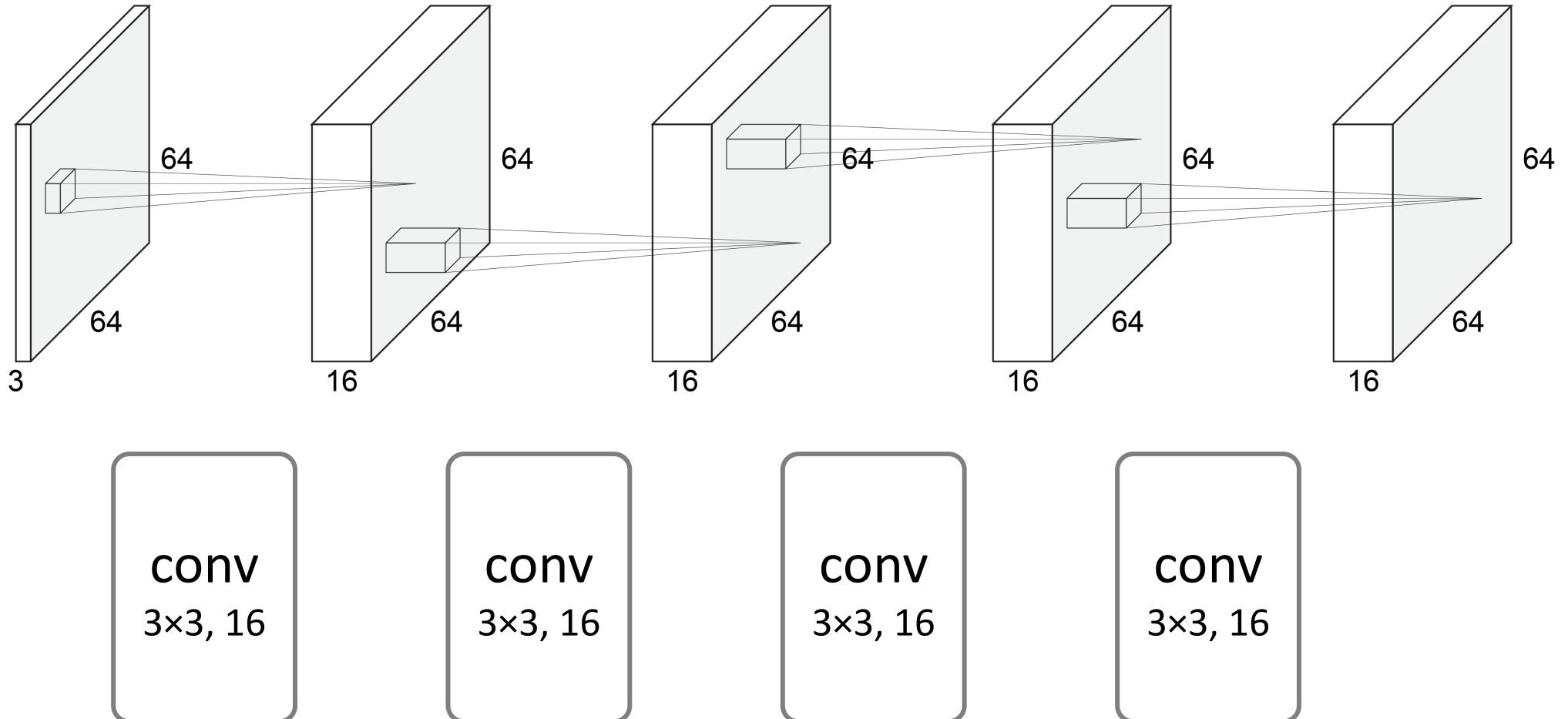


conv
3x3, 16

conv
3x3, 16

conv
3x3, 16

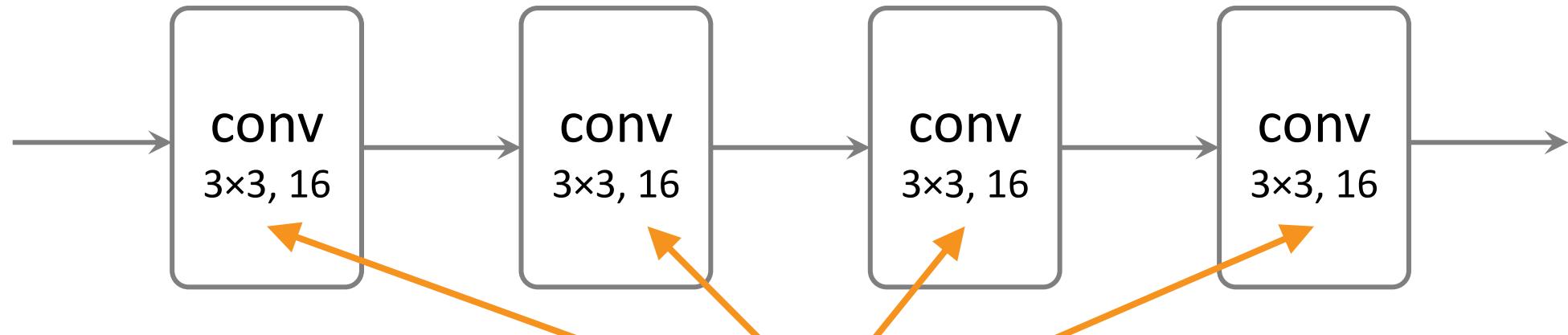
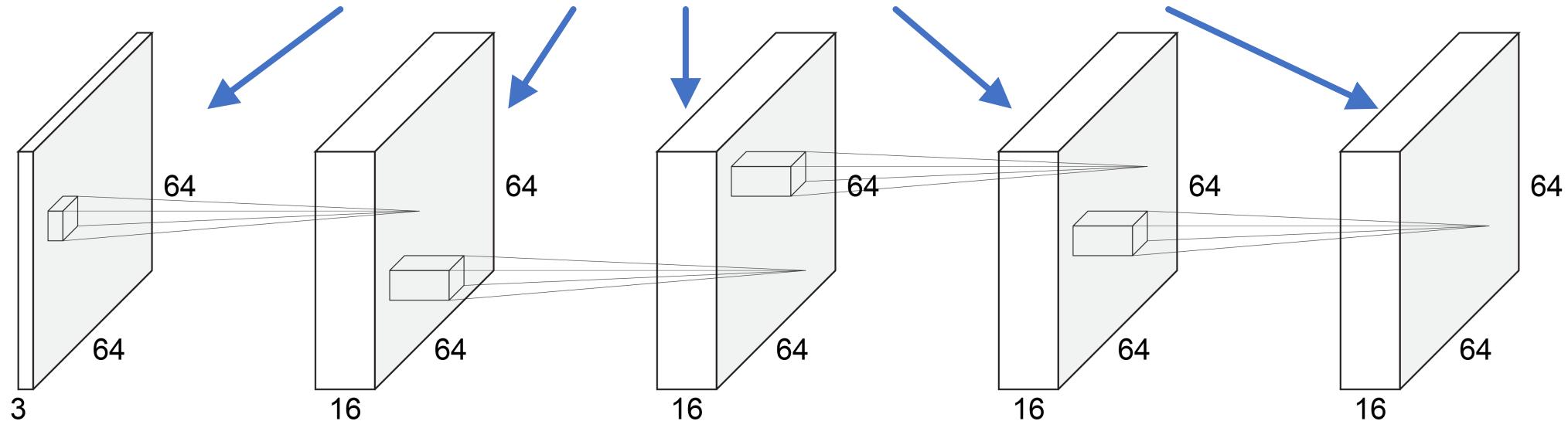
Composing basic operations



two ways of showing
neural nets

Composing basic operations

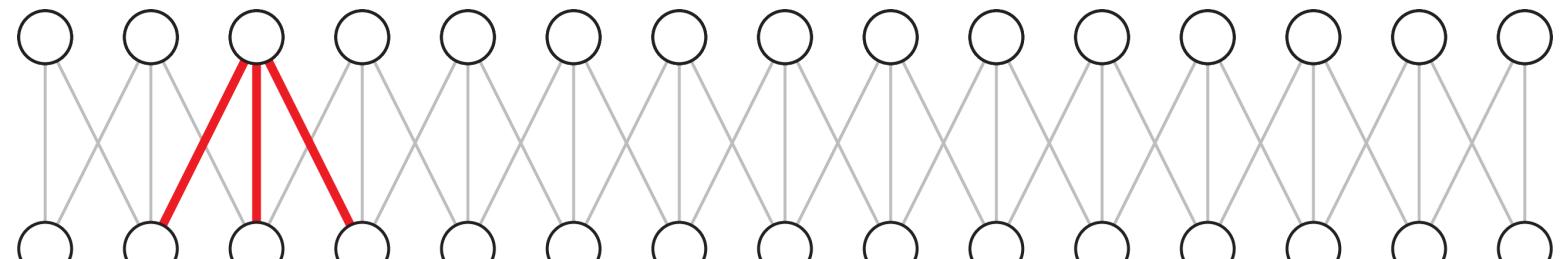
these are activations (features, embeddings, tensors ...)



these are operations (functions, transforms, layers ...)

Receptive Field

kernel = 3

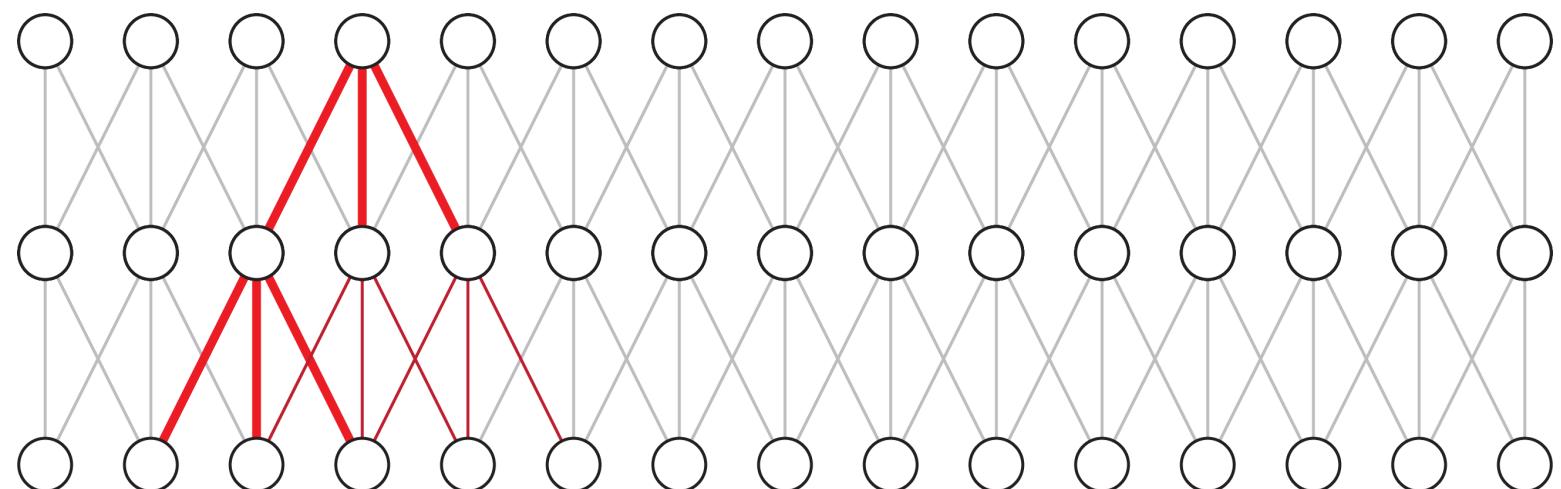


receptive field = 3

Receptive Field

kernel = 3

kernel = 3



receptive field = 5

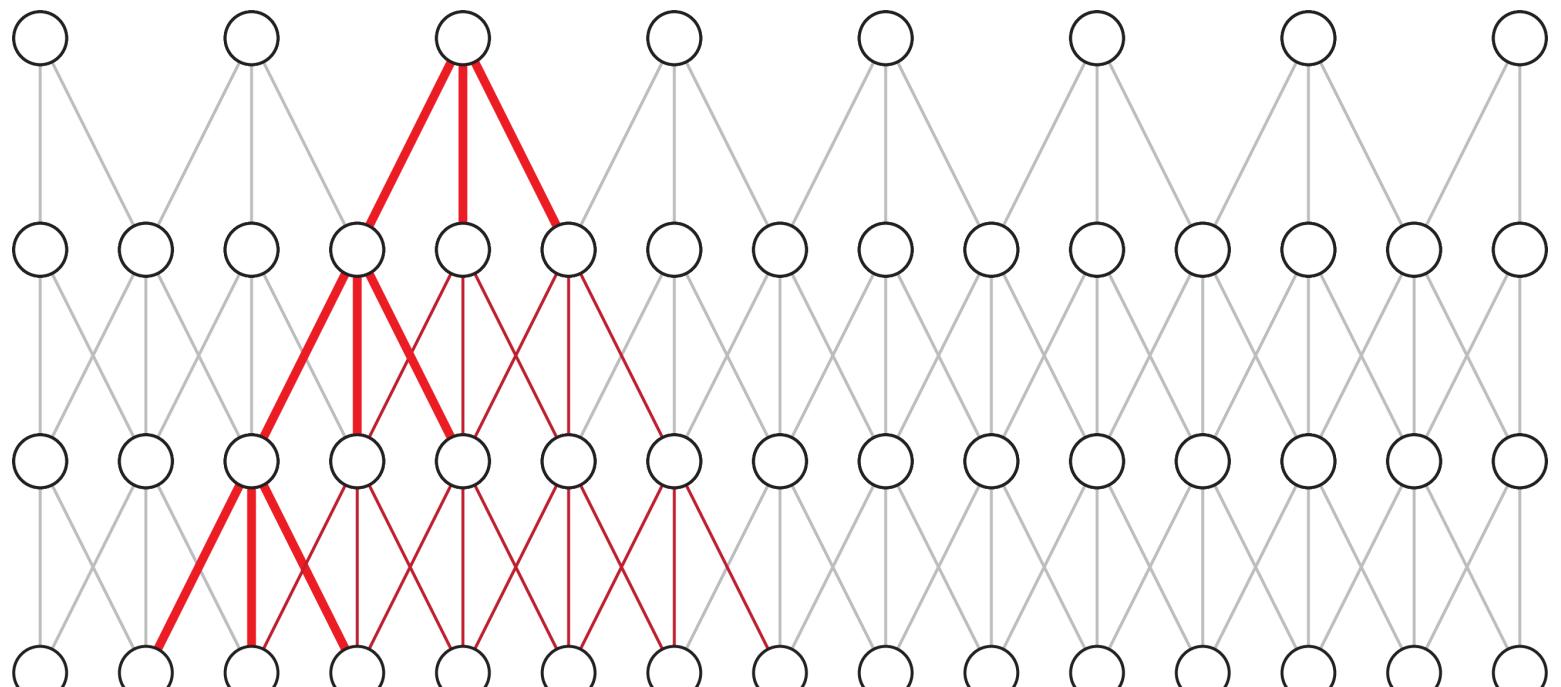
Receptive Field

kernel = 3, str = 2

kernel = 3

kernel = 3

receptive field = 7



Receptive Field

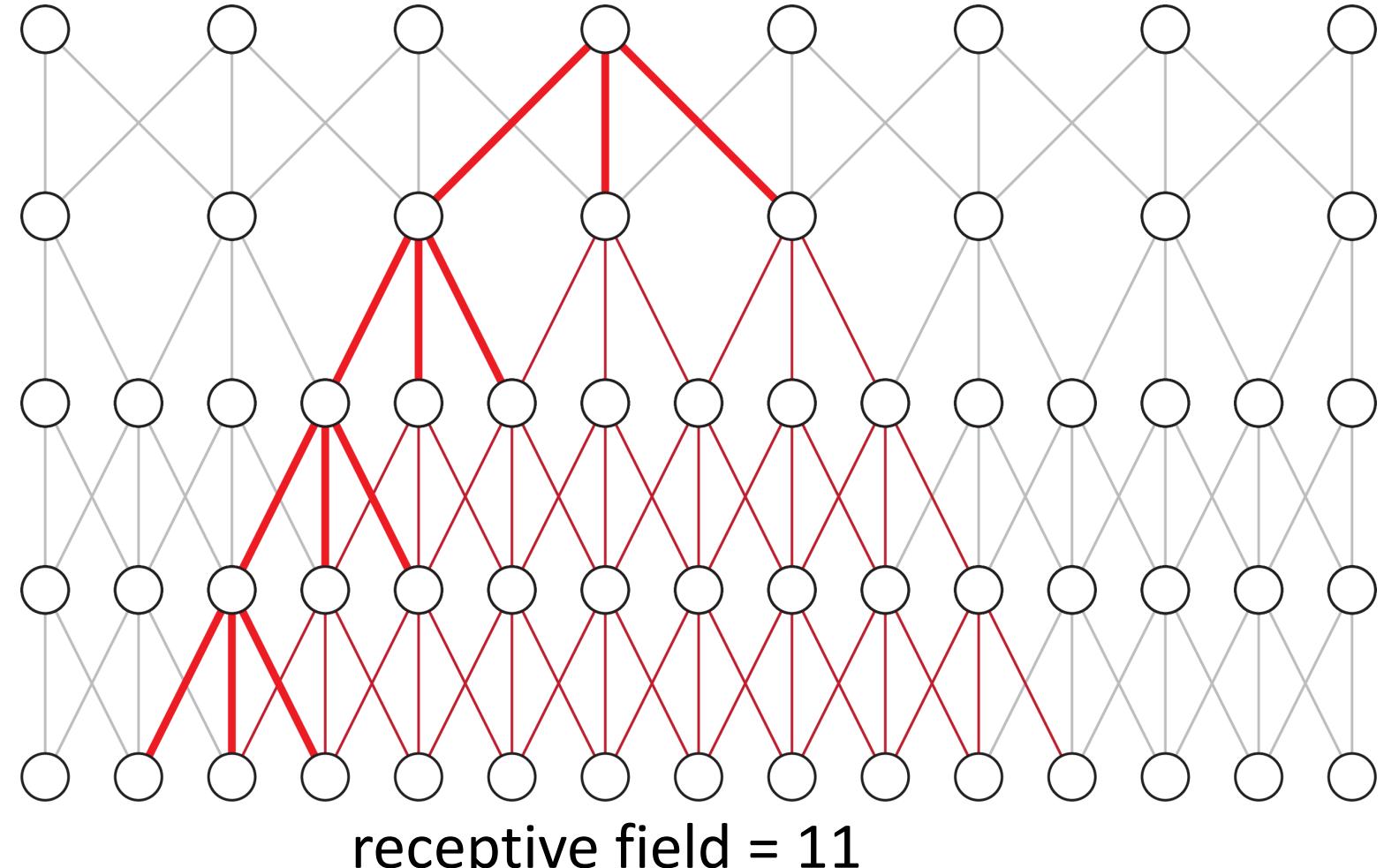
Composing many conv layers can create a large receptive field.

kernel = 3

kernel = 3, str = 2

kernel = 3

kernel = 3



Convolutional Neural Network (ConvNet, CNN)

- Core idea of deep learning:
Composing basic operations into complex functions
- Common operations of ConvNet
 - convolution
 - activation function
 - pooling
 - fully-connected layer
 - softmax



Activation function

- Activation functions introduce non-linearity

x

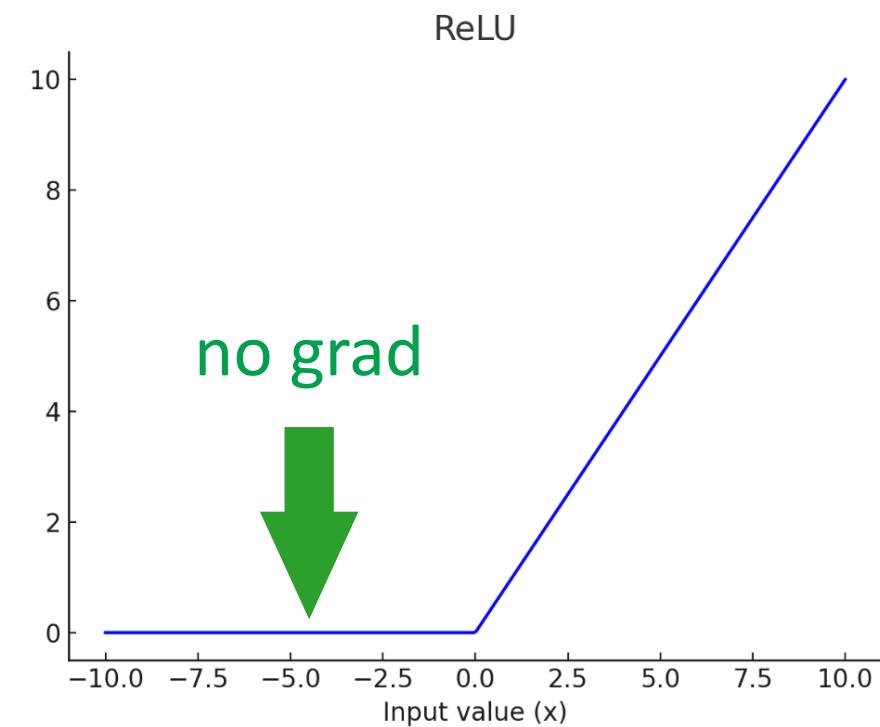
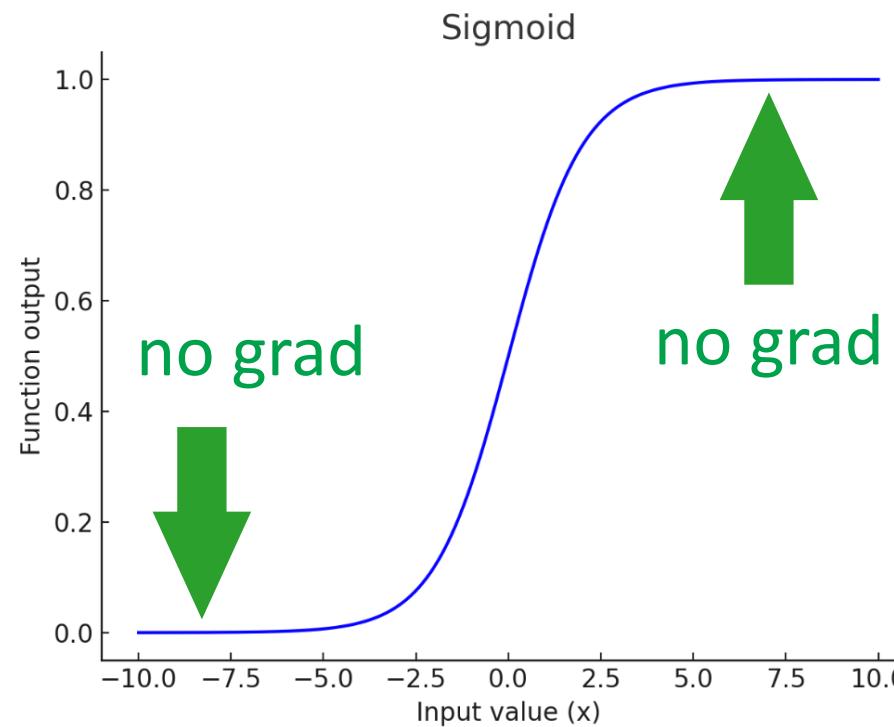


$$\text{ReLU}(x) = \max(x, 0)$$



Activation function

- Rectified Linear Unit (ReLU) helps gradient propagation



Pooling

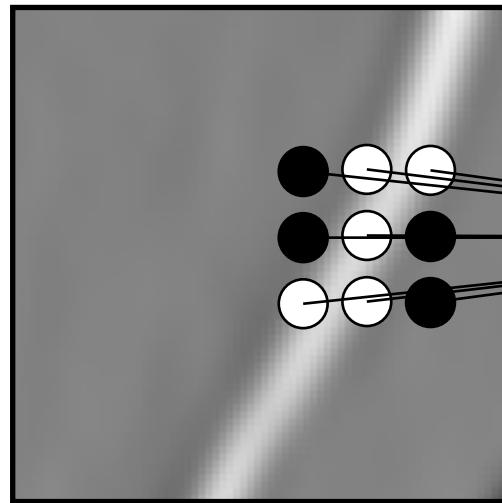
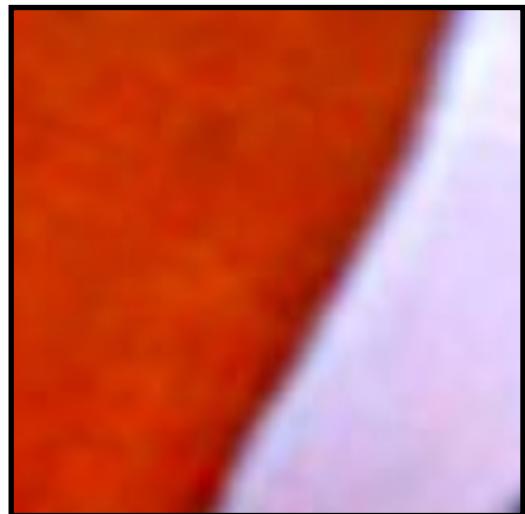
- Pooling is similar to convolution:
 - sliding window + aggregation
 - often stride > 1
 - local invariance w.r.t. translation

max pooling:

$$y[i] = \max_{j \in \mathcal{N}(i)} x[j]$$

average pooling:

$$y[i] = \frac{1}{|\mathcal{N}|} \sum_{j \in \mathcal{N}(i)} x[j]$$



“There is an edge!”

large activation

Pooling

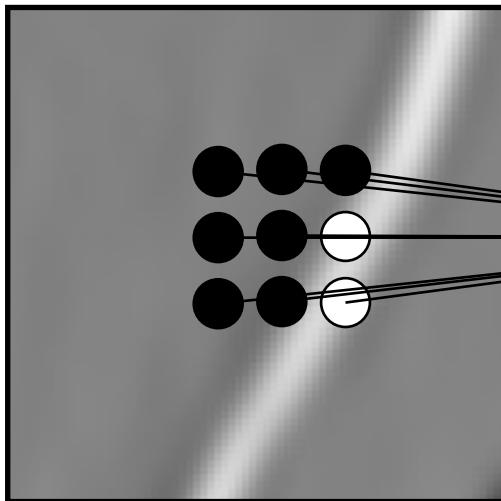
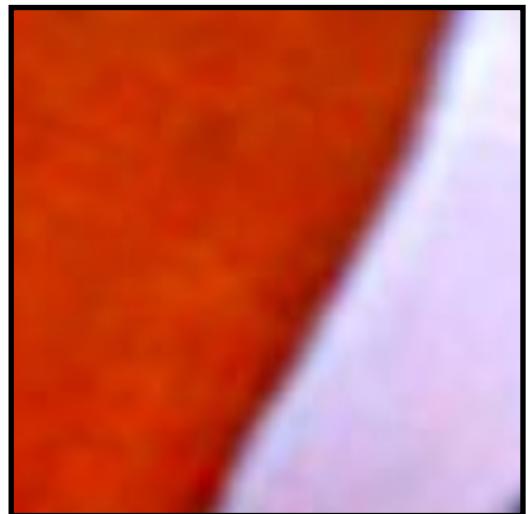
- Pooling is similar to convolution:
 - sliding window + aggregation
 - often stride > 1
 - local invariance w.r.t. translation

max pooling:

$$y[i] = \max_{j \in \mathcal{N}(i)} x[j]$$

average pooling:

$$y[i] = \frac{1}{|\mathcal{N}|} \sum_{j \in \mathcal{N}(i)} x[j]$$



max

“There is an edge!”

large activation

Pooling

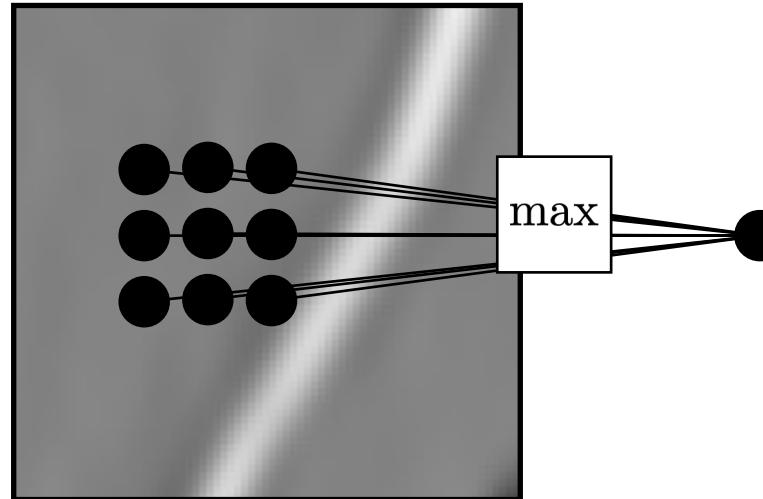
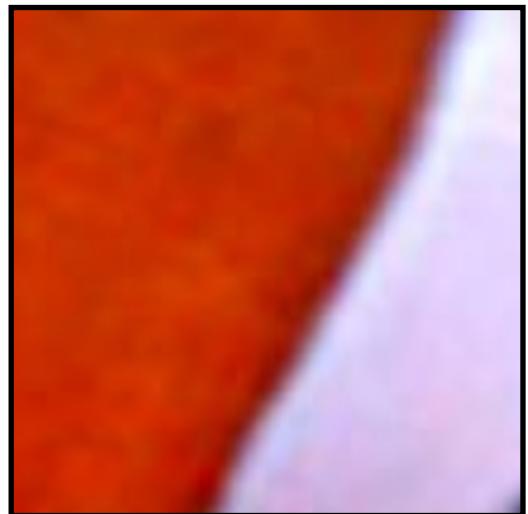
- Pooling is similar to convolution:
 - sliding window + aggregation
 - often stride > 1
 - local invariance w.r.t. translation

max pooling:

$$y[i] = \max_{j \in \mathcal{N}(i)} x[j]$$

average pooling:

$$y[i] = \frac{1}{|\mathcal{N}|} \sum_{j \in \mathcal{N}(i)} x[j]$$



“No edge!”

no activation

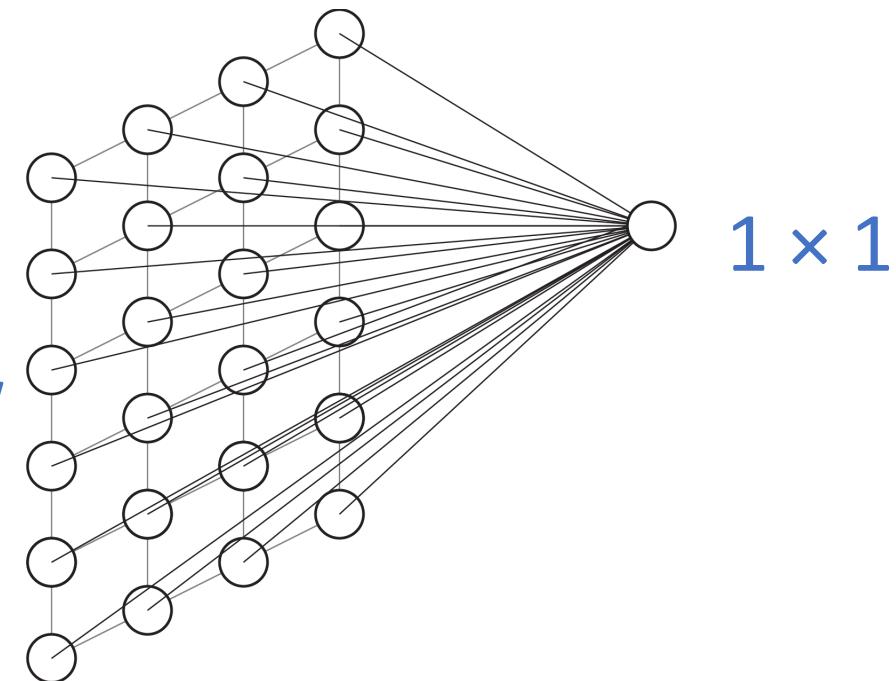
Fully-connected layer in ConvNet

- every output is connected to every input

or alternatively:

- convolution with $\text{kernel_size} = H \times W$ (w/o padding)
- output map size = 1×1
- useful for “fully conv net”

$H \times W$
each node represent
a multi-channel vector

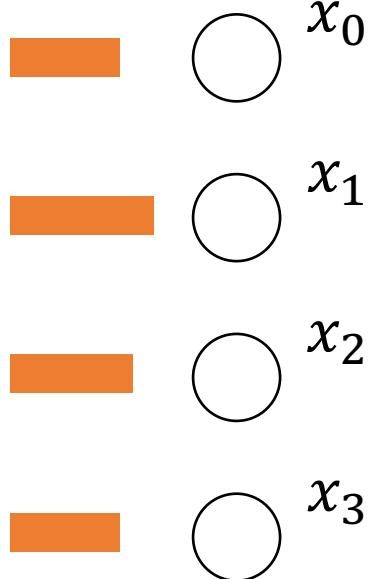


Softmax

$$\text{softmax}(\mathbf{x})_i = e^{x_i} / \sum e^{x_j}$$

- softmax = exponential activation & normalization

logits



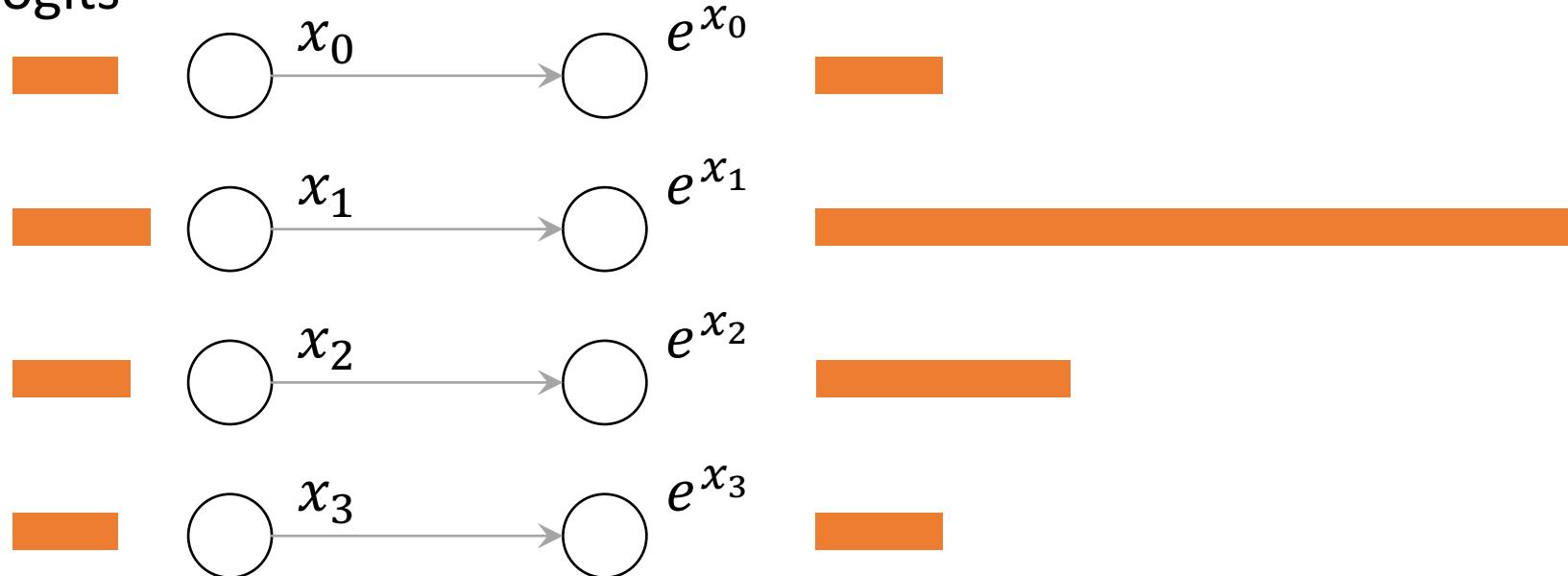
this example: 4 channels for 4-class classification

Softmax

$$\text{softmax}(x)_i = e^{x_i} / \sum e^{x_j}$$

- softmax = exponential activation & normalization

logits



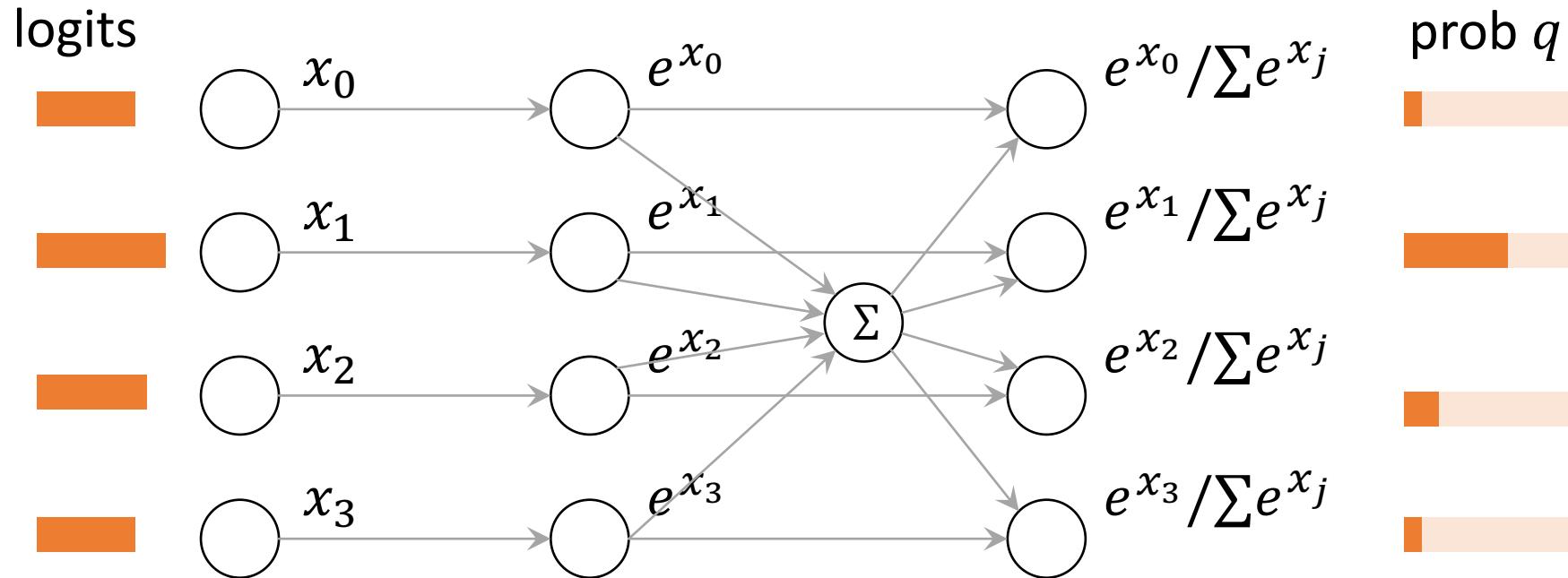
exponential:

- bigger activation gets even bigger
- suppresses non-maximum activation

Softmax

$$\text{softmax}(\mathbf{x})_i = e^{x_i} / \sum e^{x_j}$$

- softmax = exponential activation & normalization



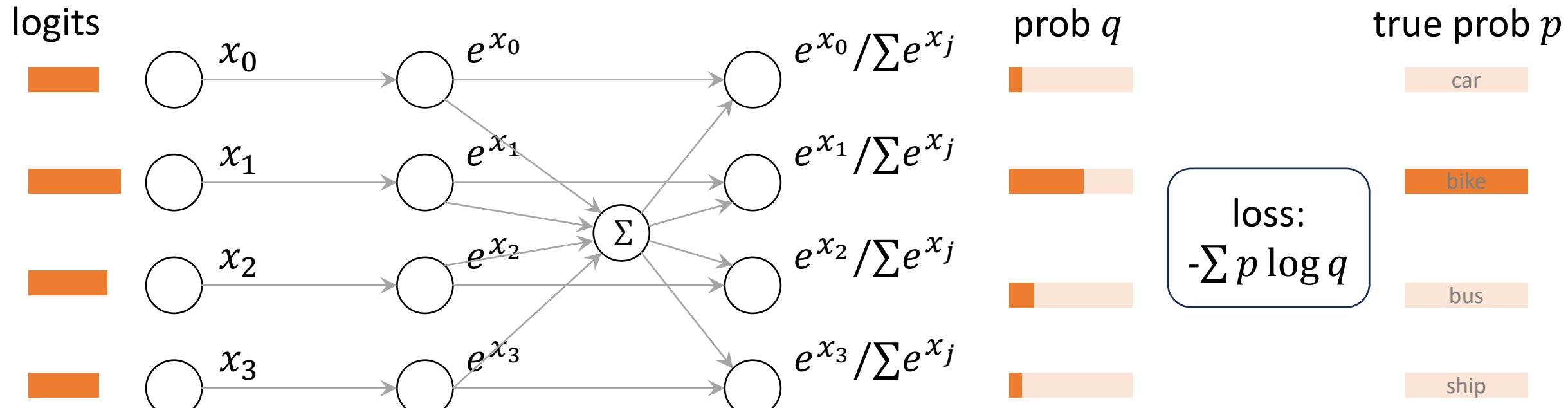
normalization:

- turn into a probability distribution

Softmax

$$\text{softmax}(\mathbf{x})_i = e^{x_i} / \sum e^{x_j}$$

- softmax = exponential activation & normalization



often implemented as `SoftmaxCrossEntropyLoss`
for numerical stability

Convolutional Neural Network (ConvNet, CNN)

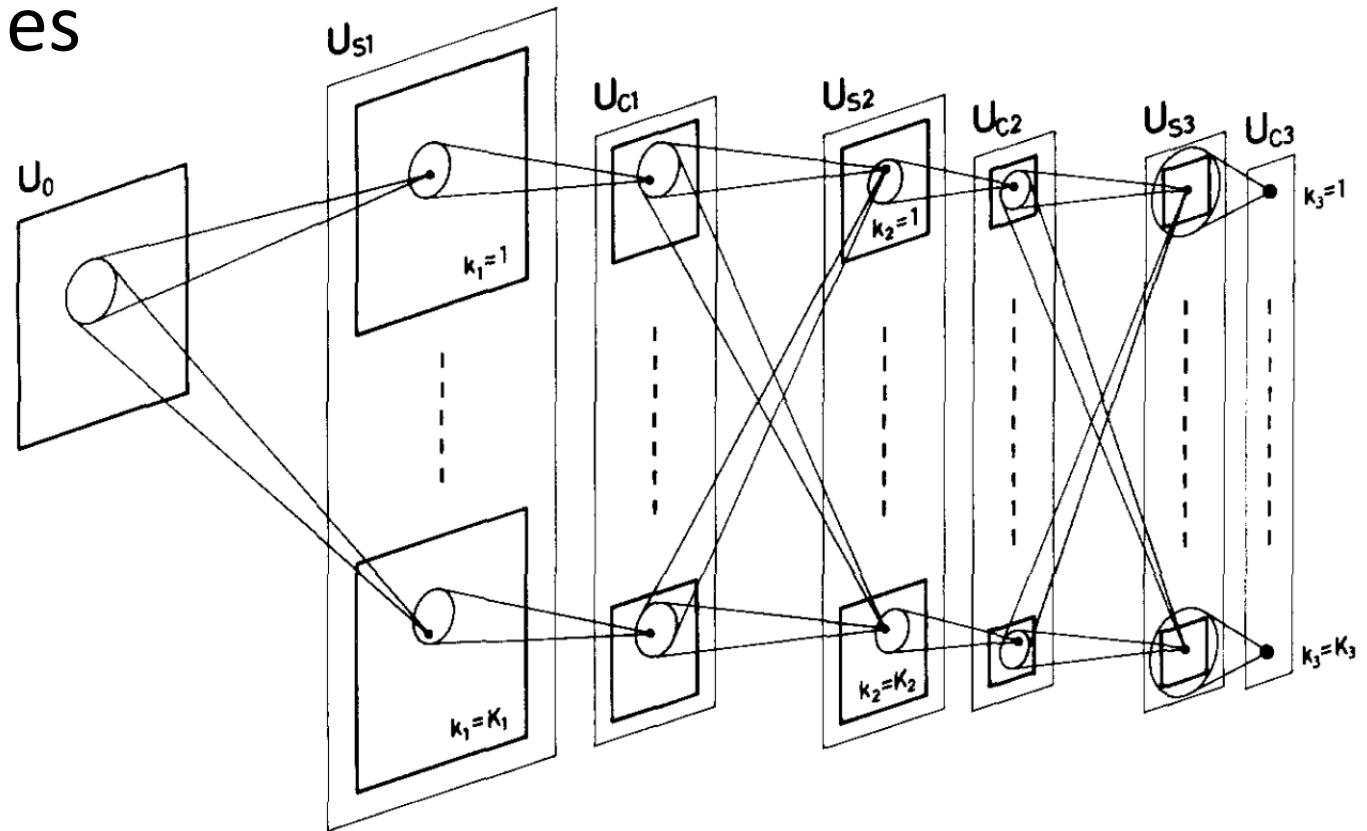
- Core idea of deep learning:
Composing basic operations into complex functions
- Common operations of ConvNet
 - convolution
 - activation function
 - pooling
 - fully-connected layer
 - softmax
- **Train with BackProp!**



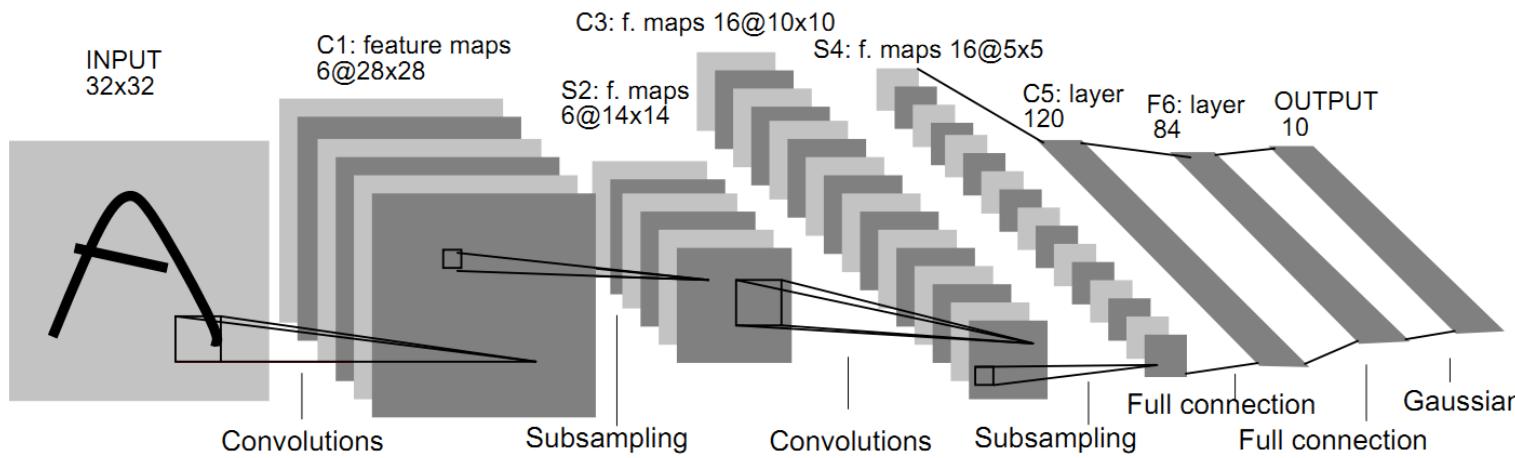
Convolutional Neural Networks: Case Study

Neocognitron, 1979

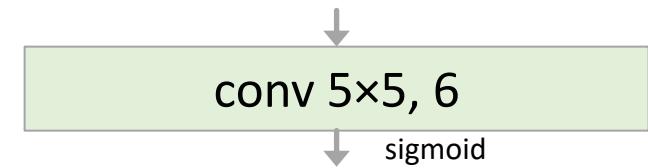
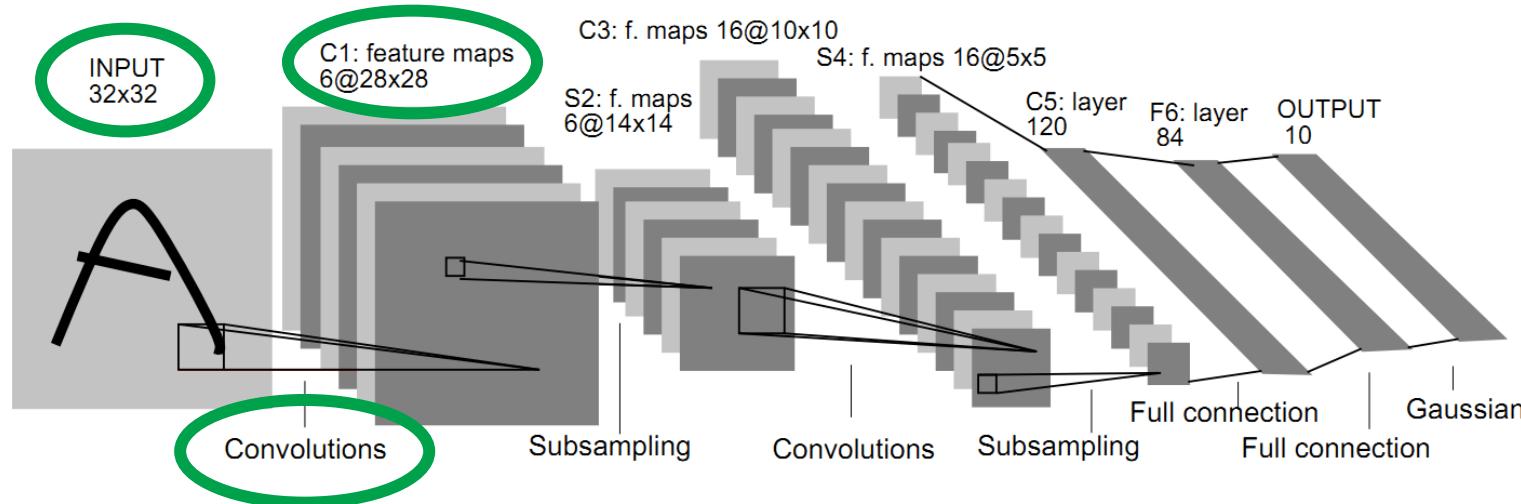
- root of ConvNet
- sliding-window features
- composing basic modules
- no backprop



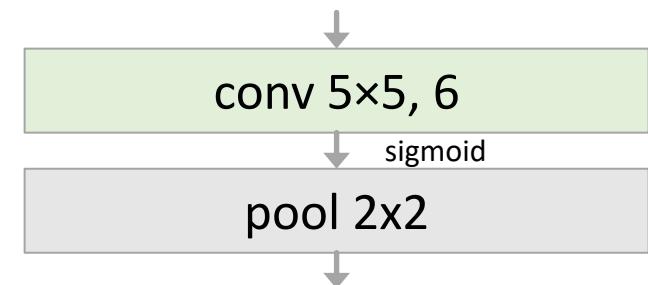
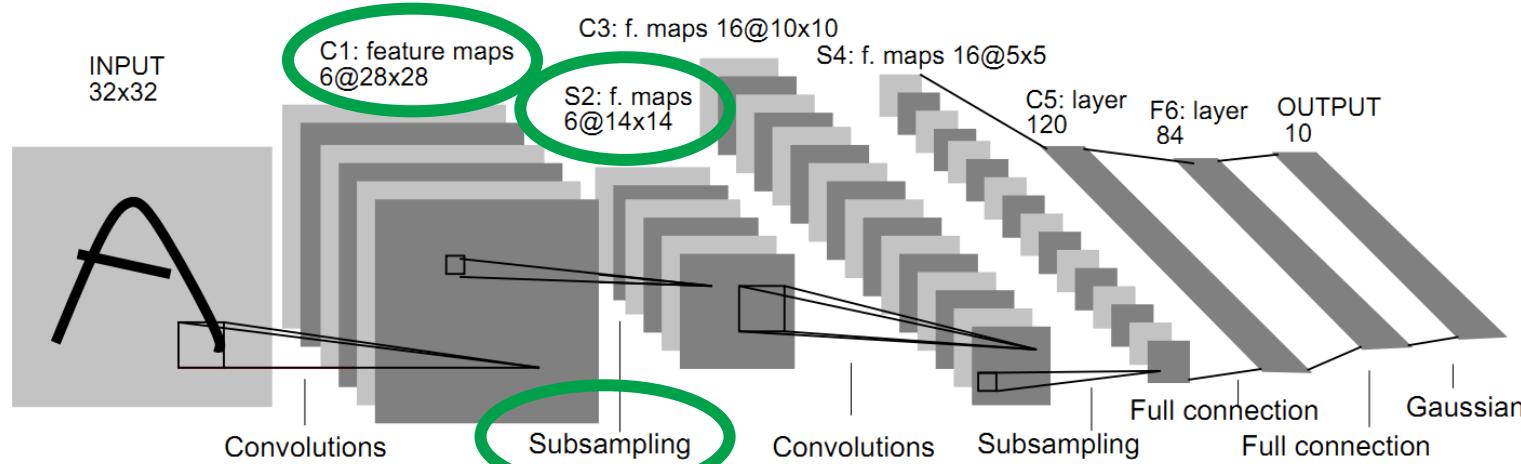
LeNet, 1989



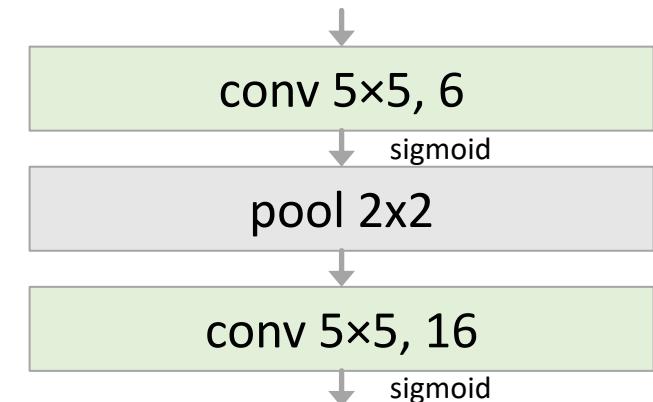
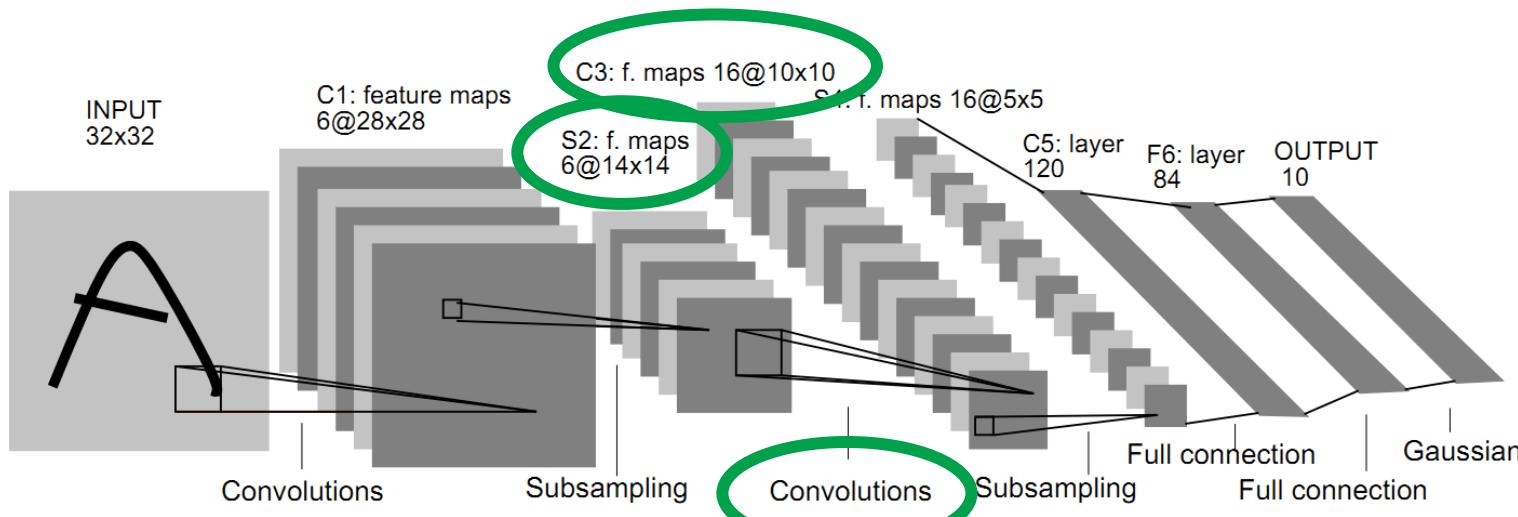
LeNet, 1989



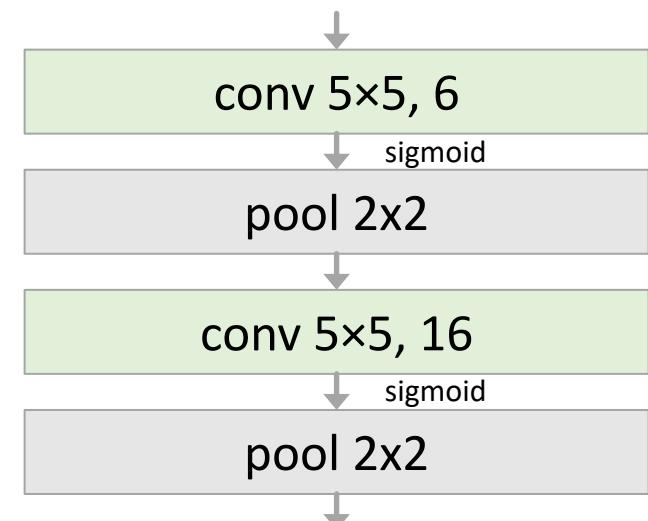
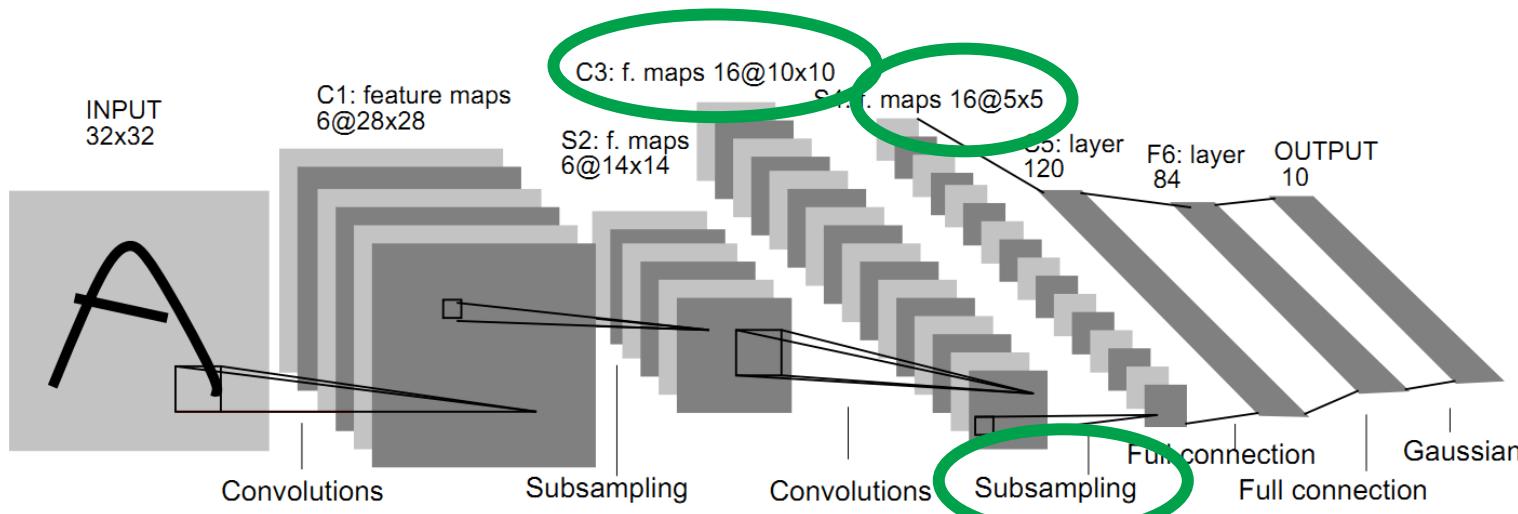
LeNet, 1989



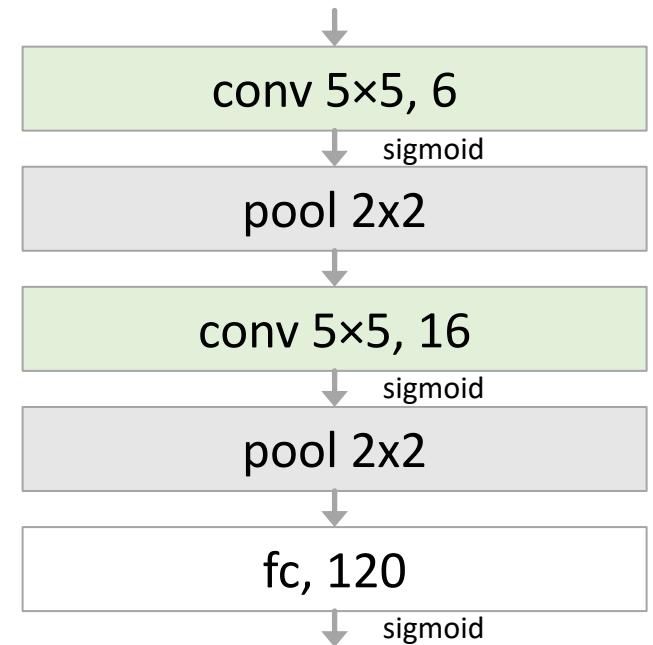
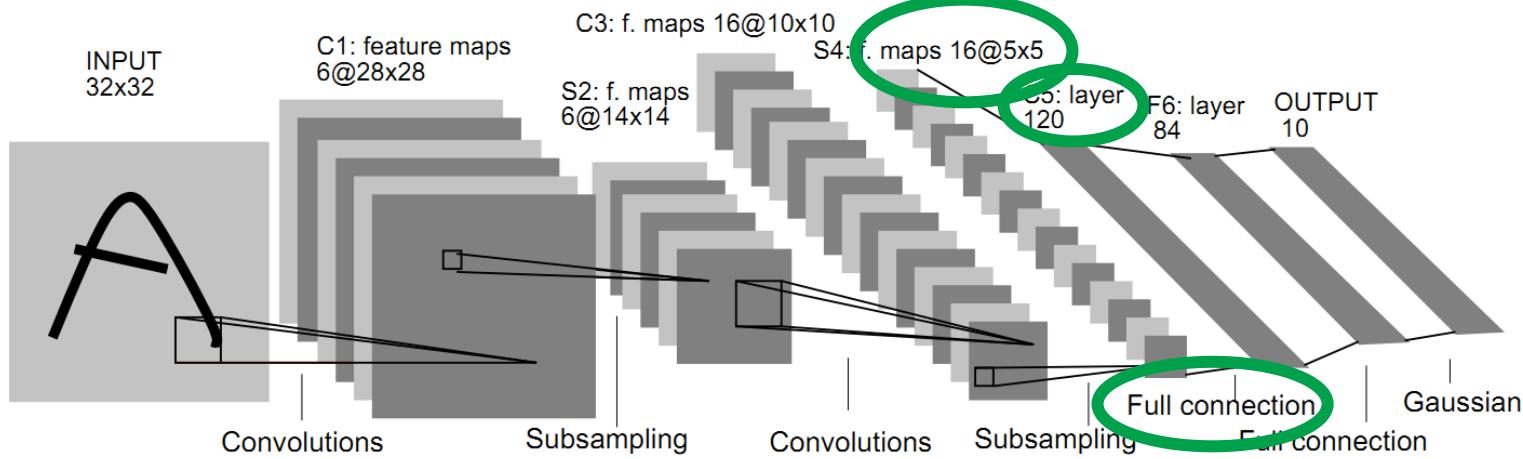
LeNet, 1989



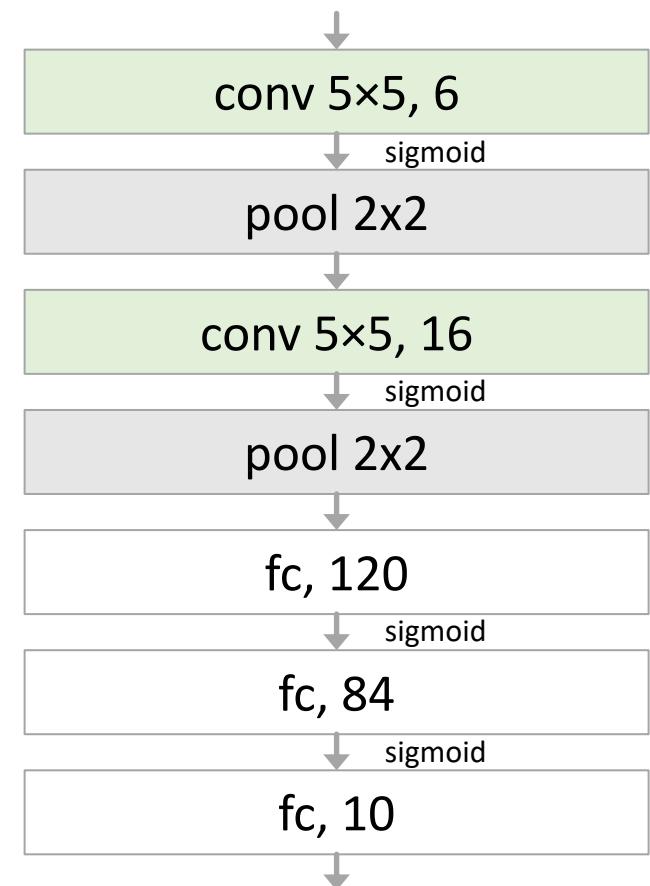
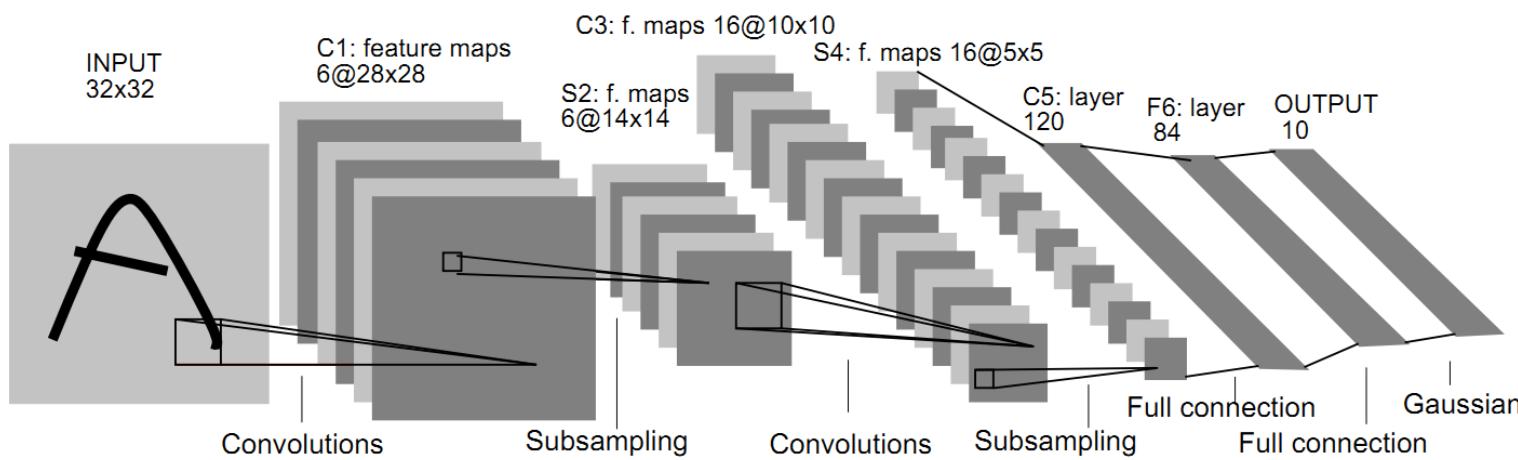
LeNet, 1989



LeNet, 1989



LeNet, 1989

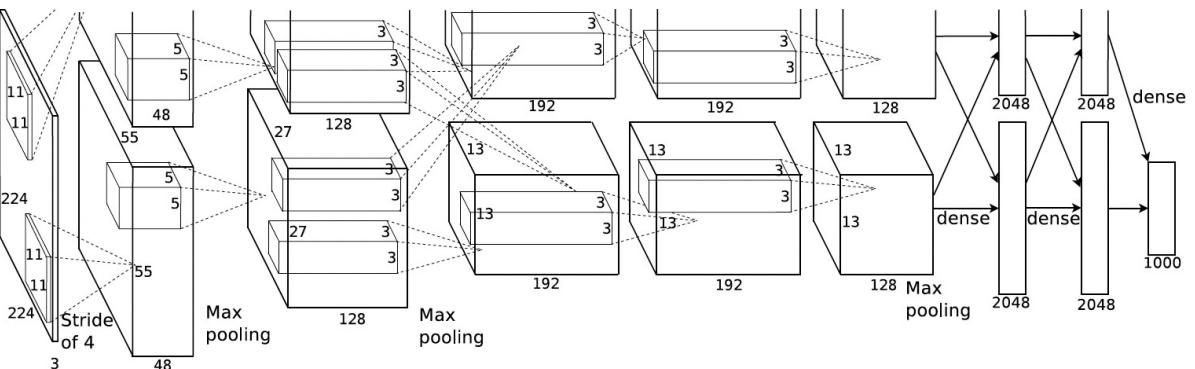


AlexNet, 2012

The Deep Learning Revolution

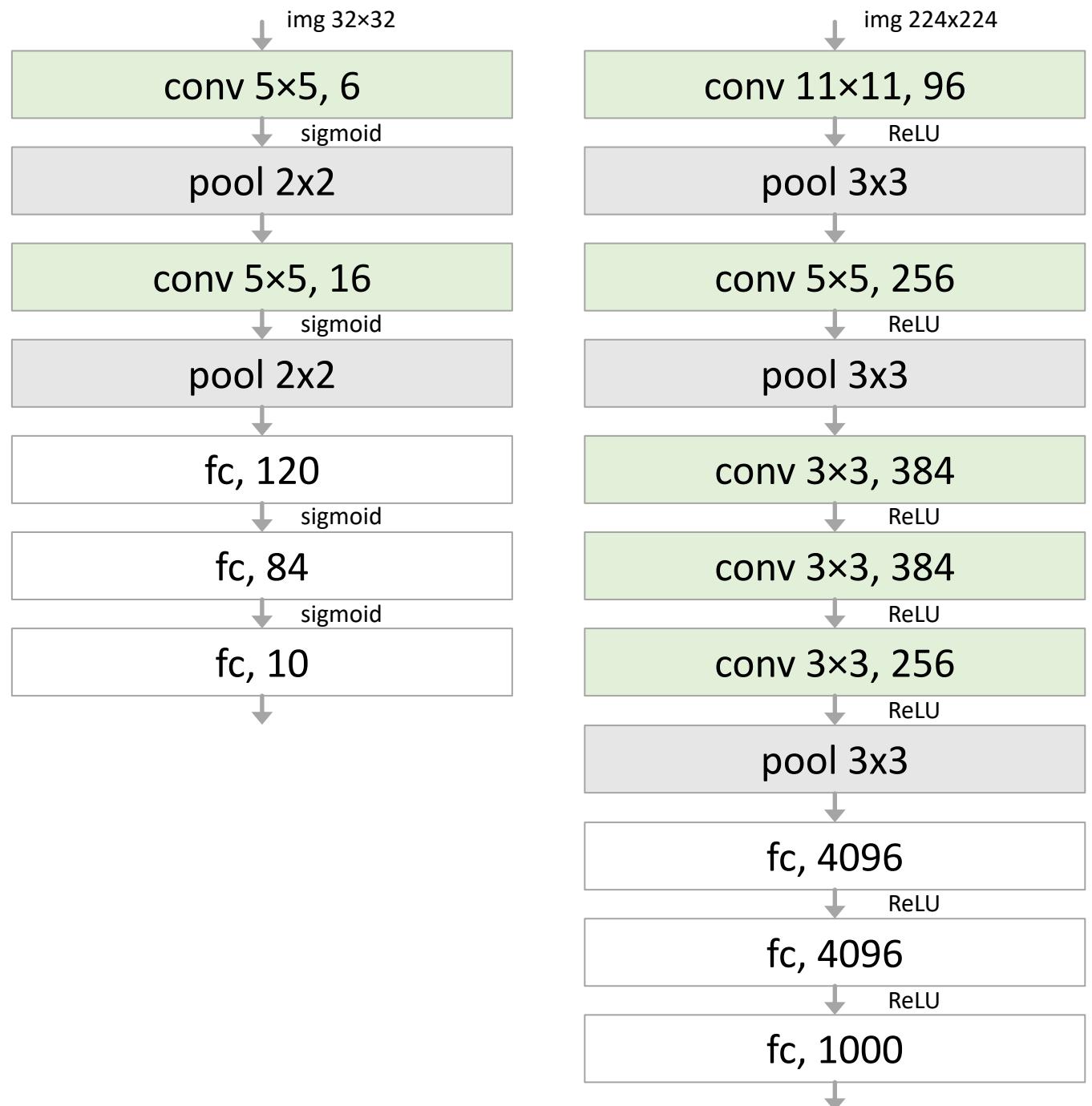
Scaling up ConvNet

- data scaling
 - ImageNet, 1.28 million images, 1000 classes
- model scaling
 - “60 million parameters and 650,000 neurons”
- reducing overfitting
 - data augmentation
 - dropout
- GPU implementation



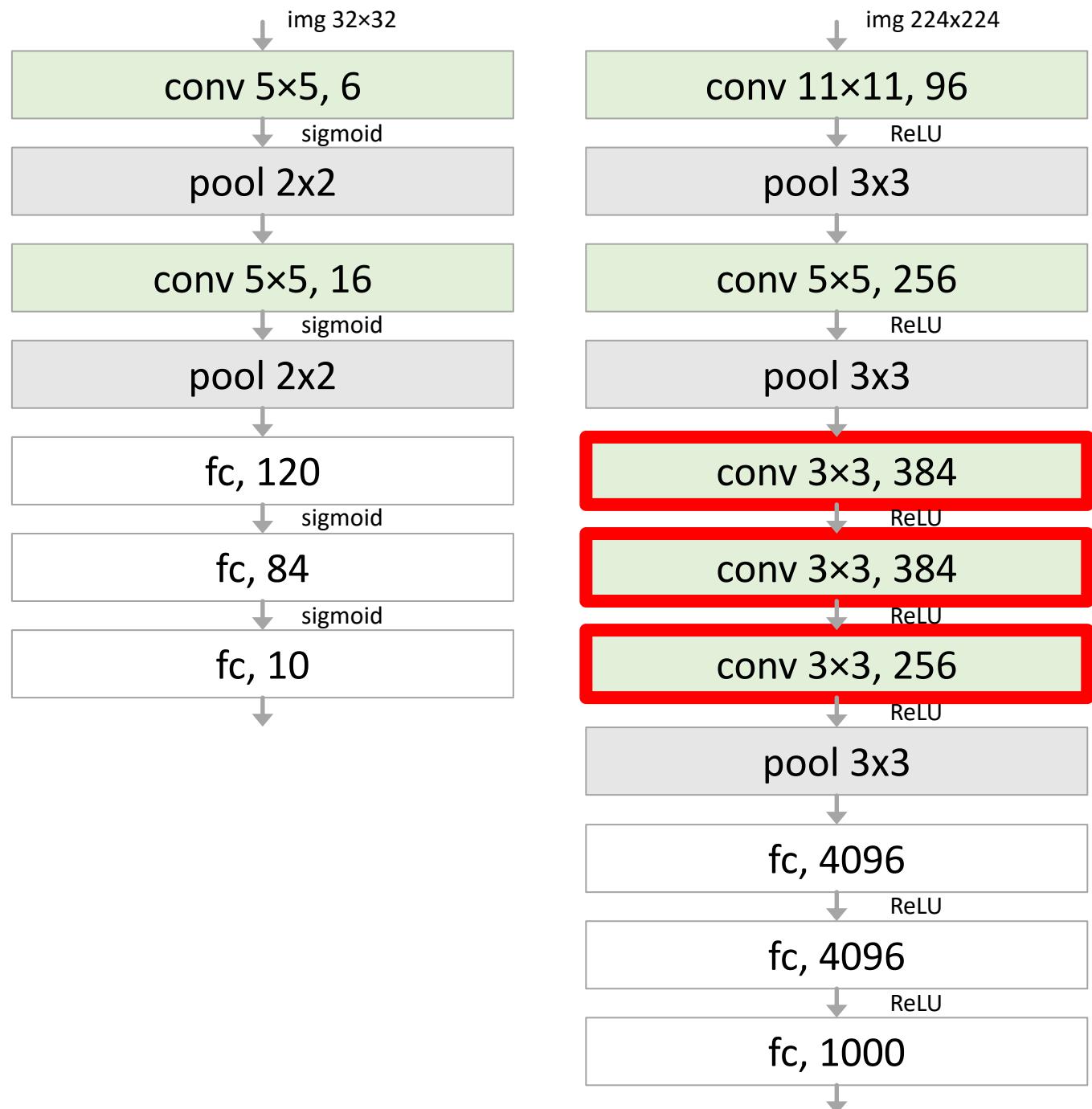
LeNet vs. AlexNet

- deeper (more layers)
- wider (more channels)
- ReLU (not sigmoid)



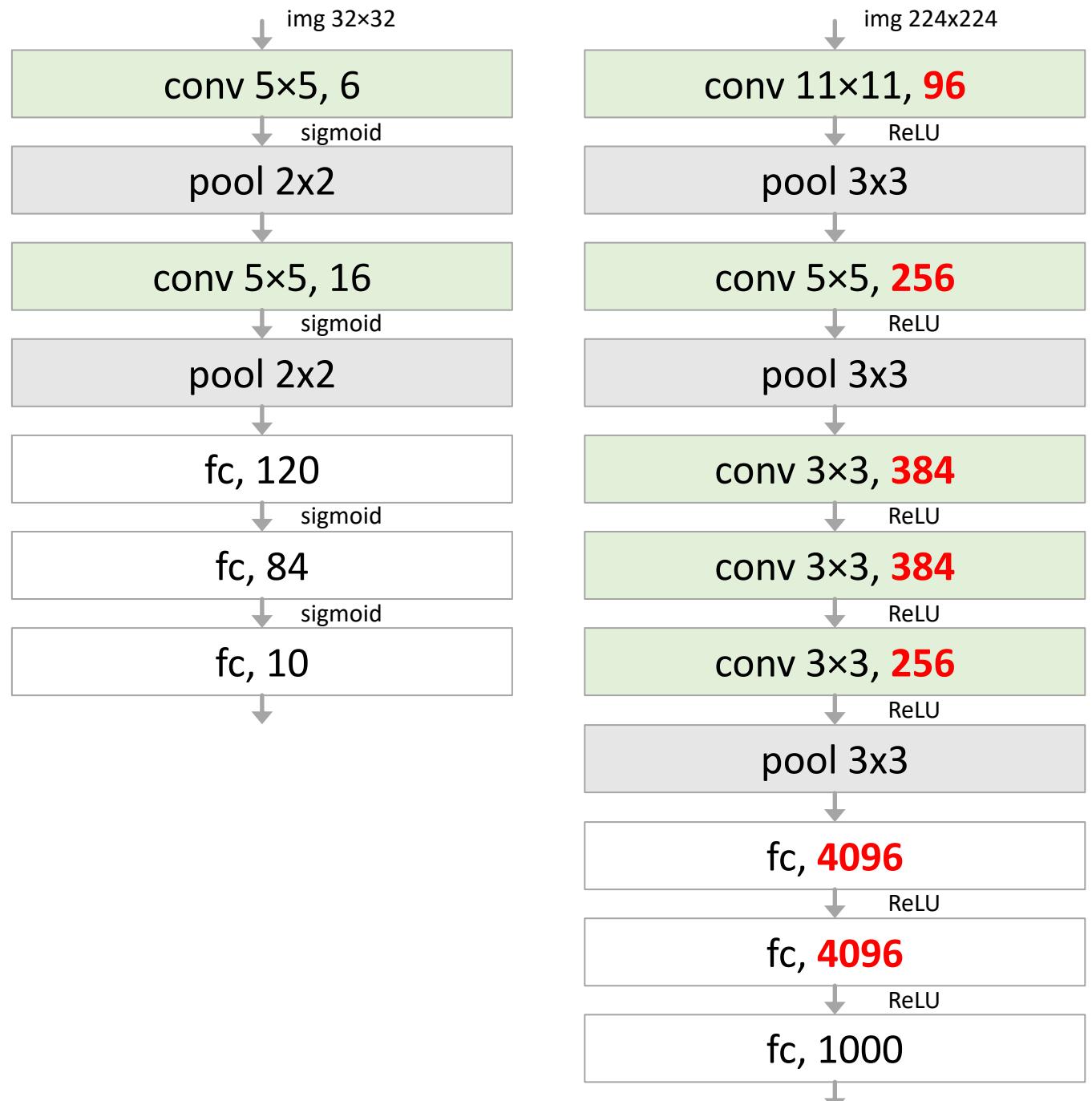
LeNet vs. AlexNet

- deeper (more layers)
- wider (more channels)
- ReLU (not sigmoid)



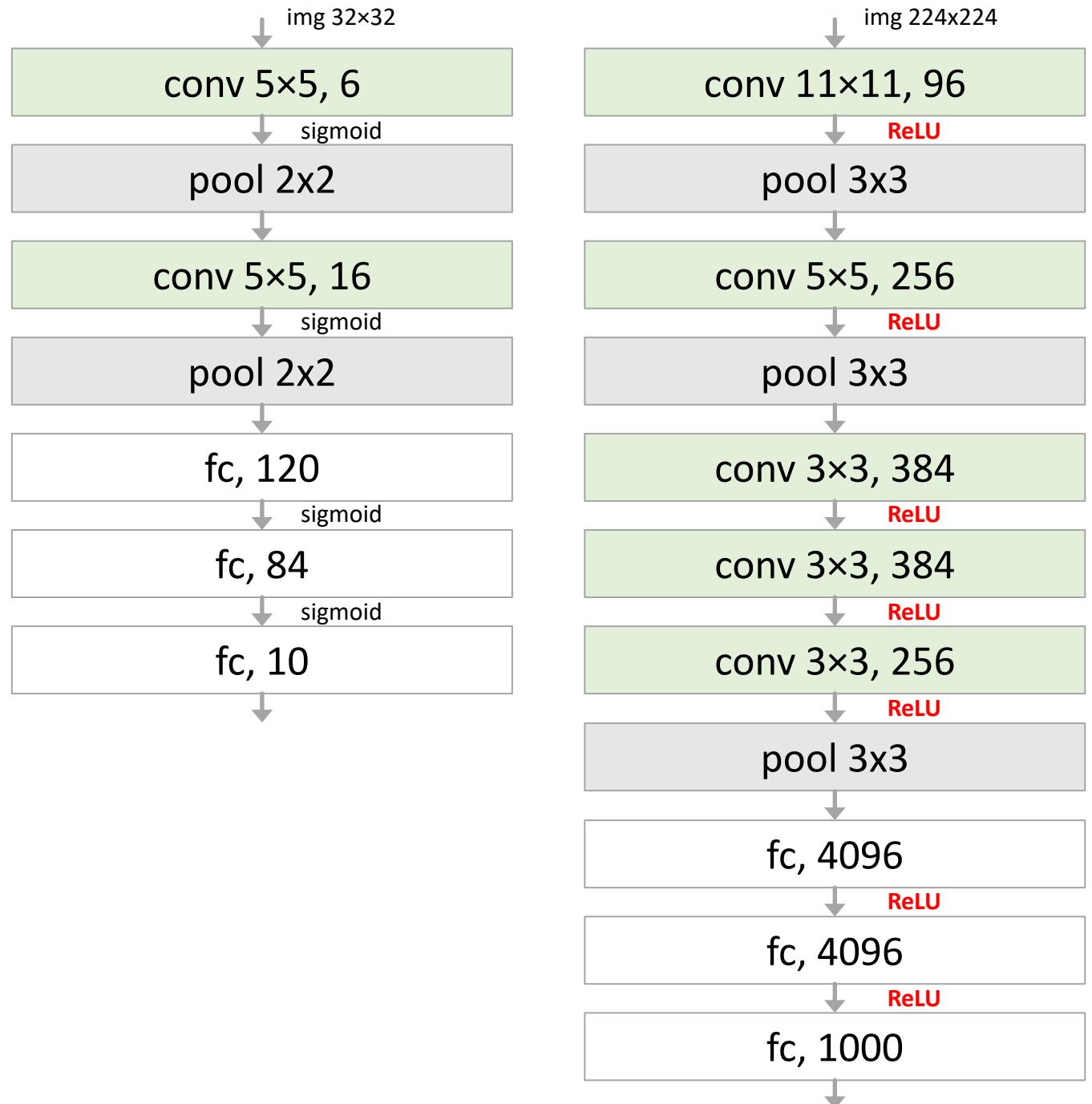
LeNet vs. AlexNet

- deeper (more layers)
- wider (more channels)
- ReLU (not sigmoid)



LeNet vs. AlexNet

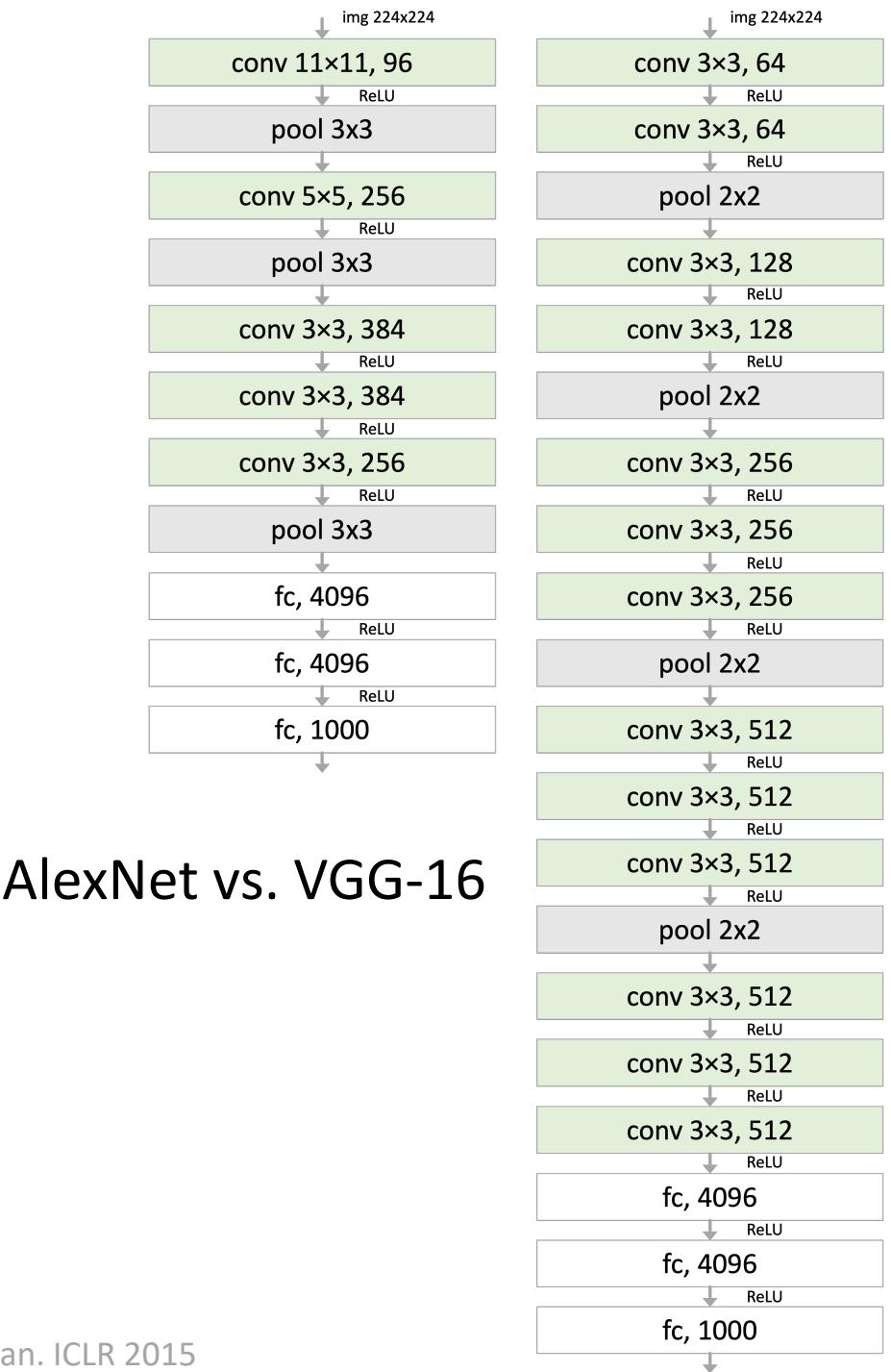
- deeper (more layers)
- wider (more channels)
- **ReLU (not sigmoid)**



VGG nets, 2014

Very deep ConvNets

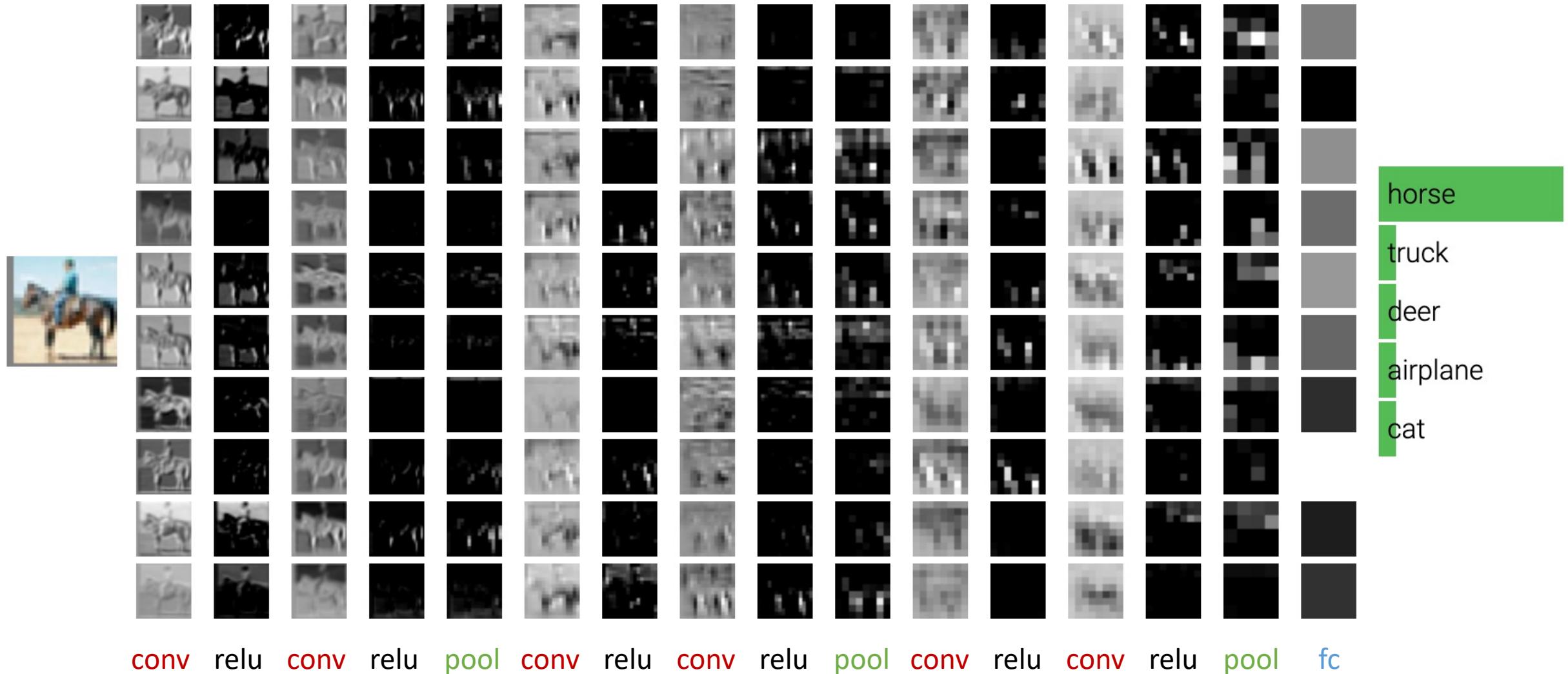
- stack the same modules
 - conv: all 3x3
- very deep
 - 16 or 19 layers
- clear evidence: deeper is better
- **Composing basic operations into complex functions**



AlexNet vs. VGG-16

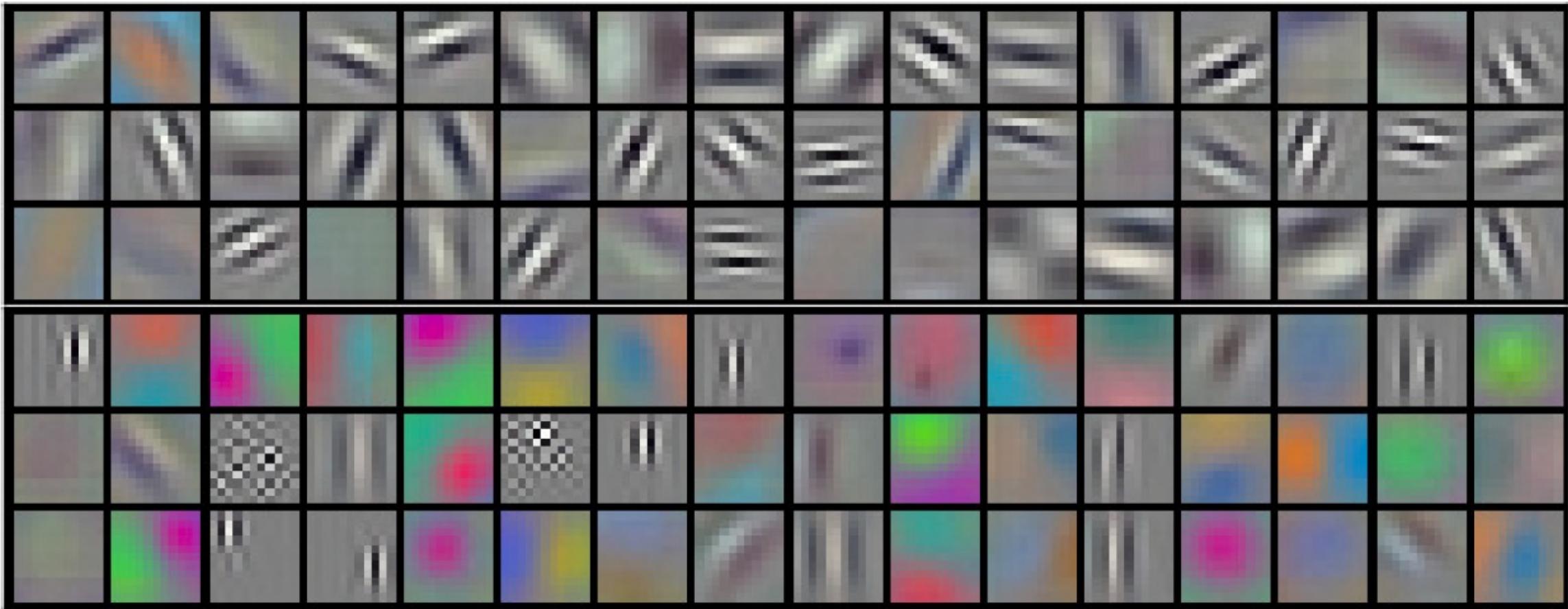
Visualizing Convolutional Neural Networks

Visualizing feature maps



see more examples in <http://cs231n.stanford.edu/>

Visualizing filters

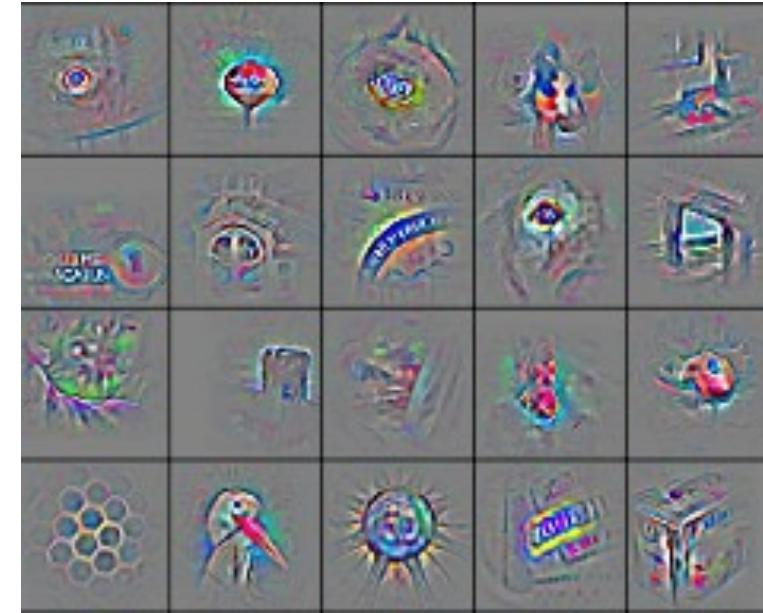


AlexNet's 1st conv layer (96 filters):
edge or color detectors

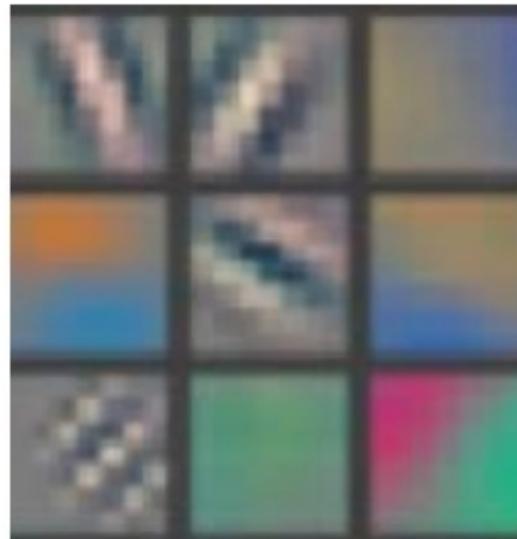
Visualizing 1-hot feature maps by BackProp

To find “what input can produce a specific feature”

- set a one-hot feature map
- back-prop to pixels



Visualizing 1-hot feature maps by BackProp



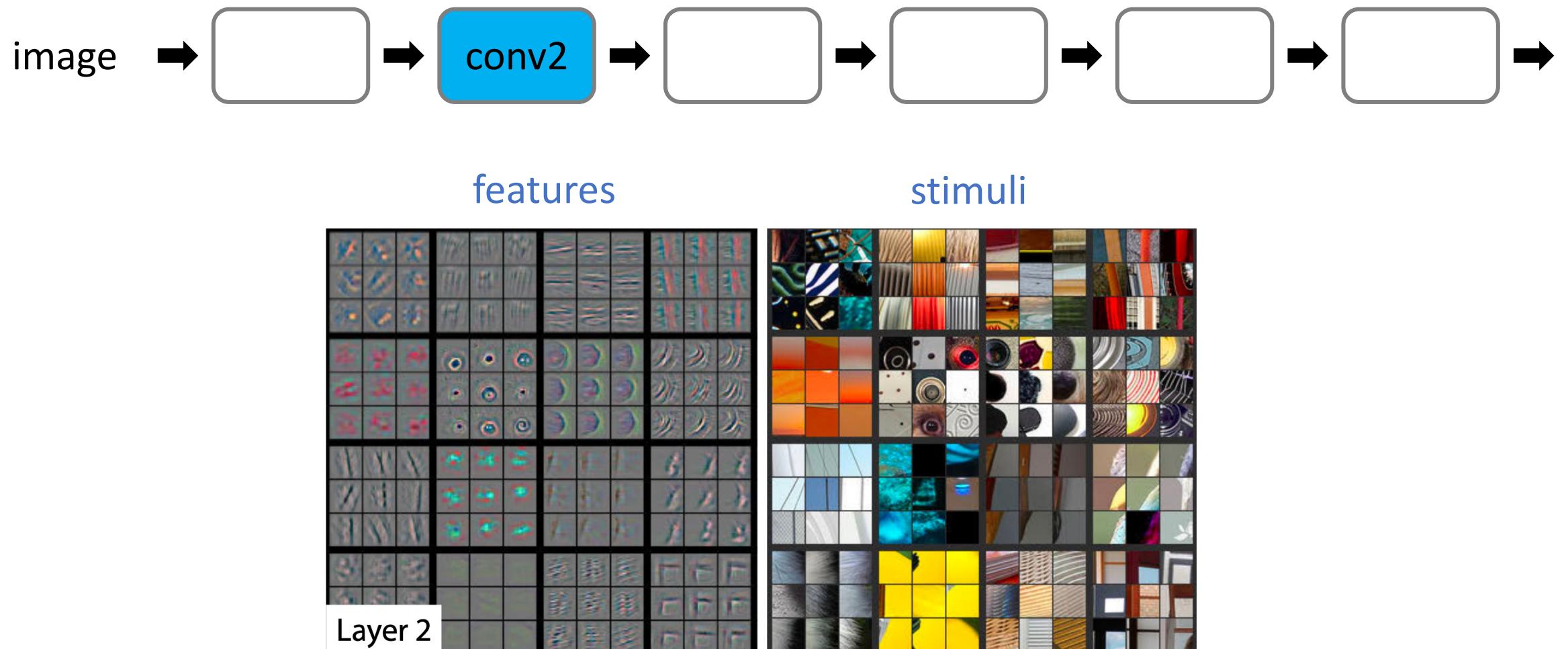
features



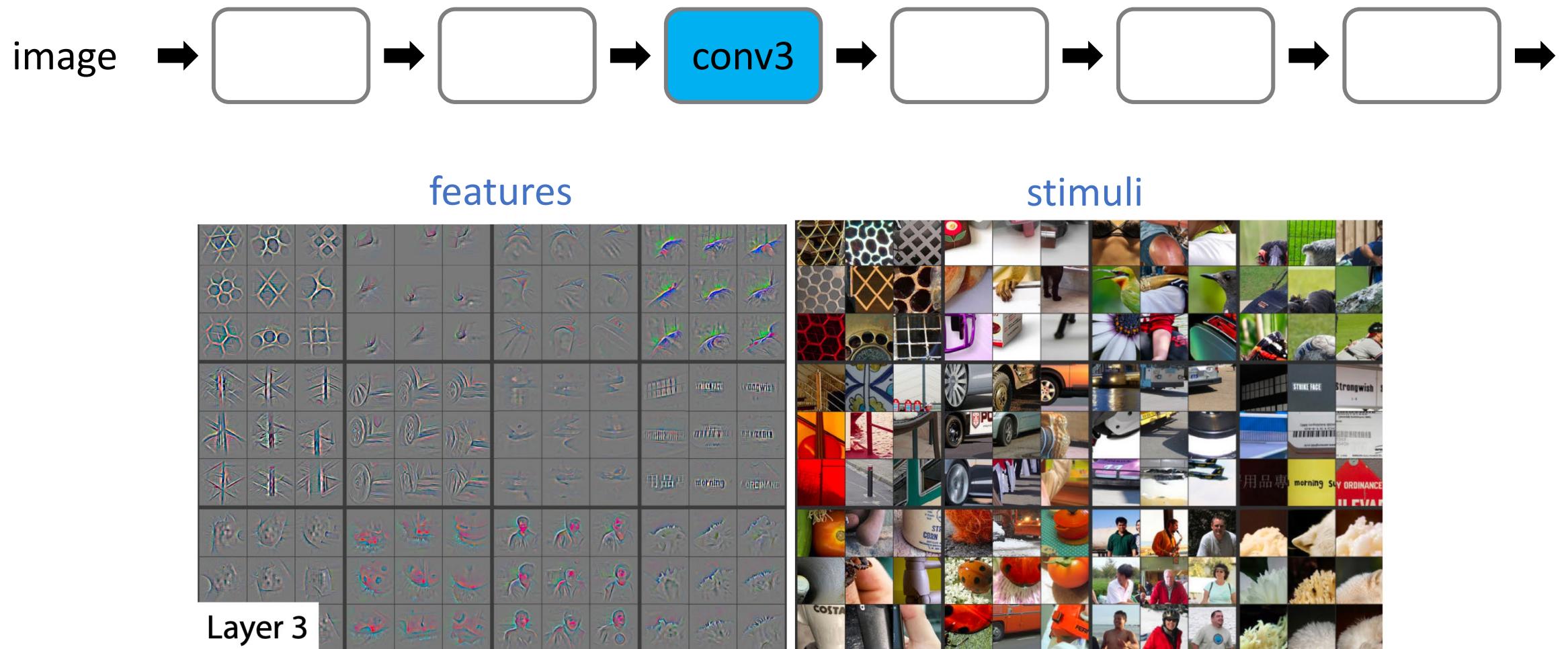
stimuli

(patches with the highest
1-hot activations)

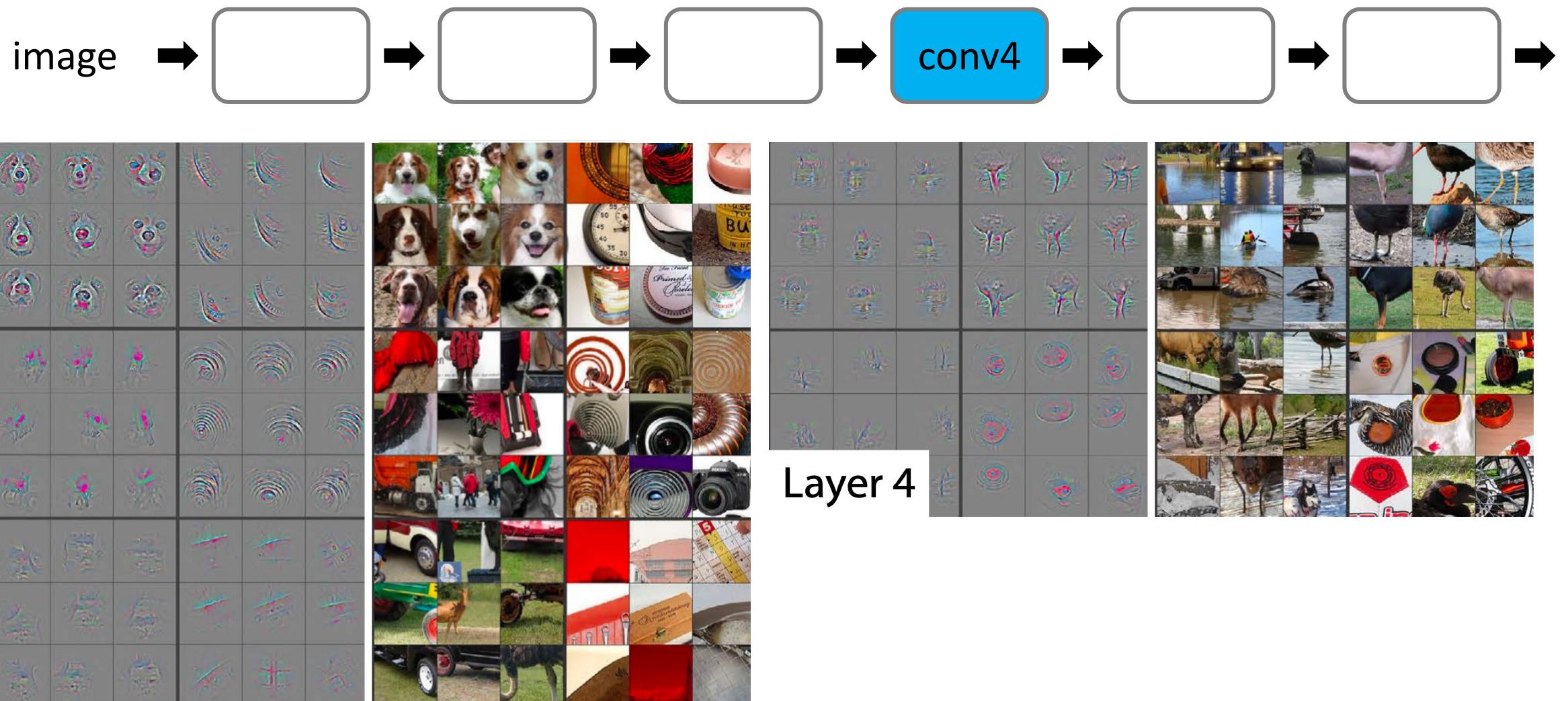
Visualizing 1-hot feature maps by BackProp



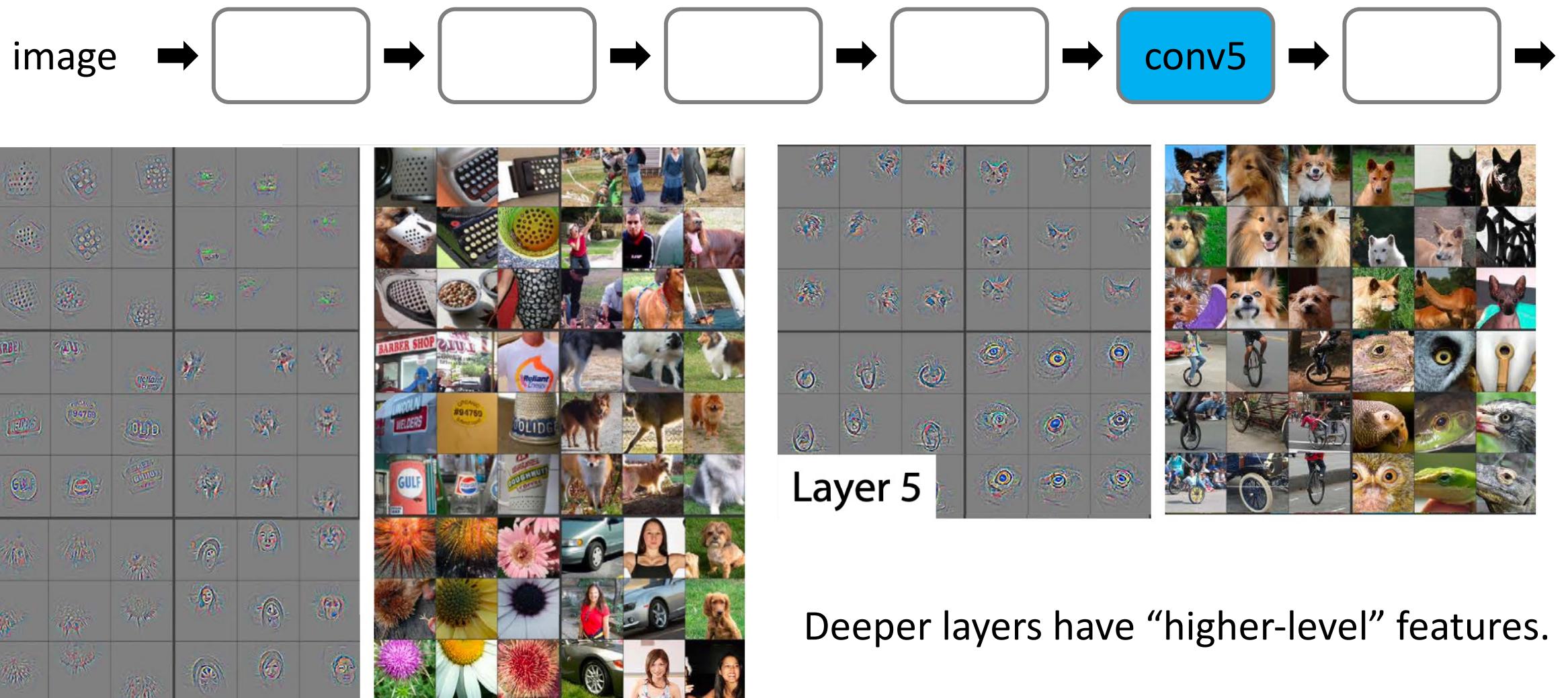
Visualizing 1-hot feature maps by BackProp



Visualizing 1-hot feature maps by BackProp



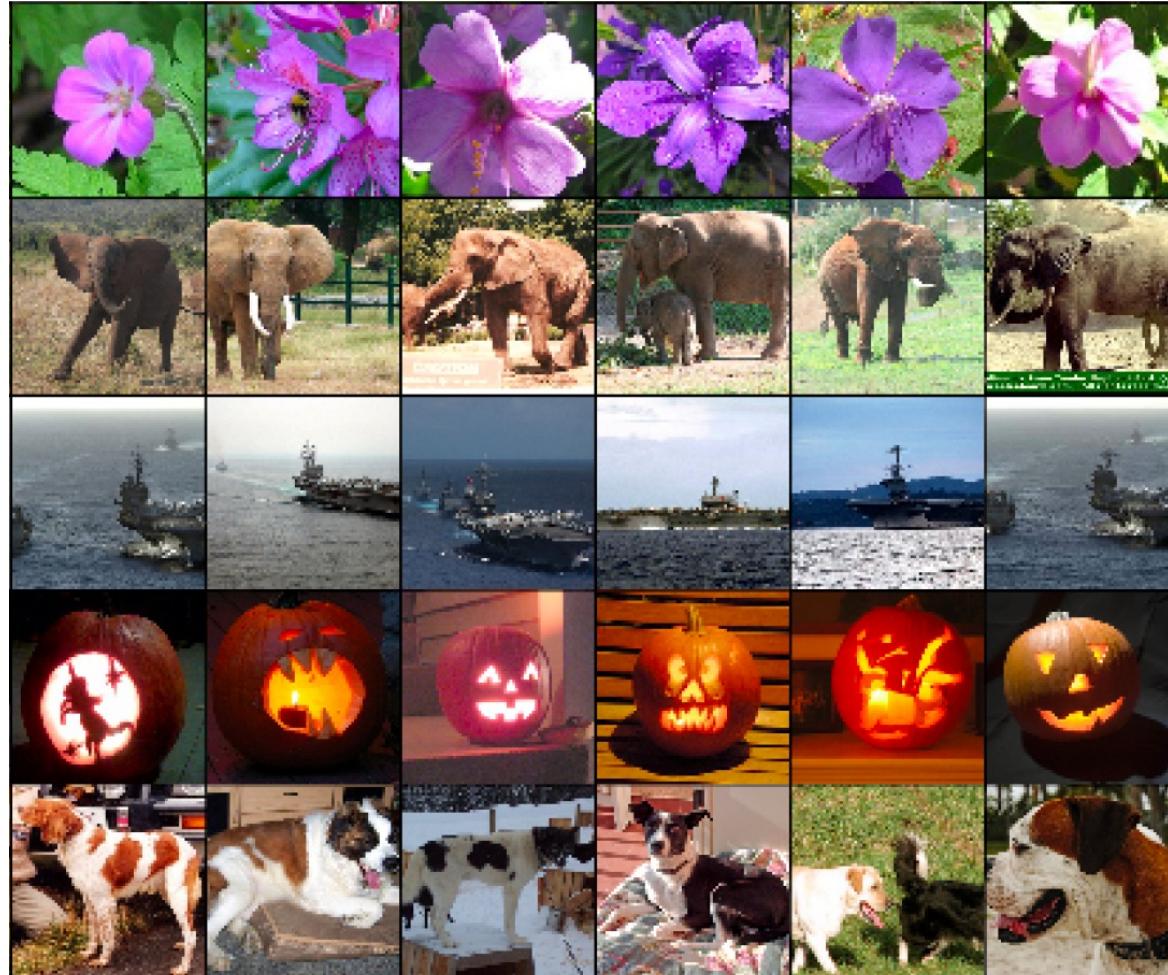
Visualizing 1-hot feature maps by BackProp



Visualizing by retrieval

nearest neighbors in feature space of
AlexNet's last layer

query



Lecture 9: Convolutional Neural Networks

- Convolution
- Convolutional Neural Networks
- Case Study: LeNet, AlexNet, VGG
- Visualization