



성균관대학교  
SUNGKYUNKWAN UNIVERSITY

# Deep Learning

## - Convolutional Neural Networks 1 -

**Eunbyung Park**

Assistant Professor

School of Electronic and Electrical Engineering

[Eunbyung Park \(silverbottlep.github.io\)](https://github.com/silverbottlep)

# Convolution

# 1D Convolution

- Convolution is a mathematical operation on two functions ( $f$ ,  $g$ ) that produces a third function  $f * g$

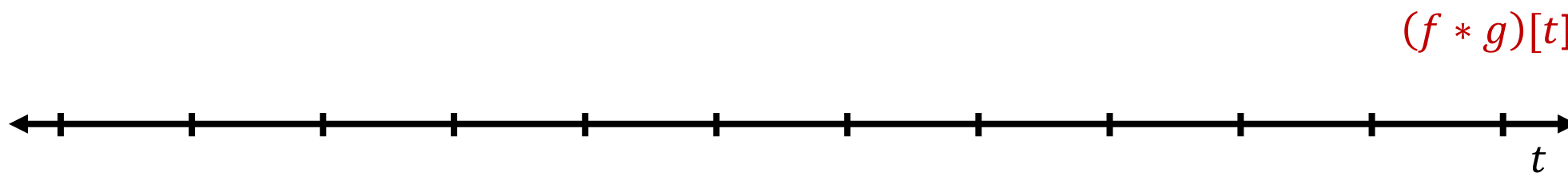
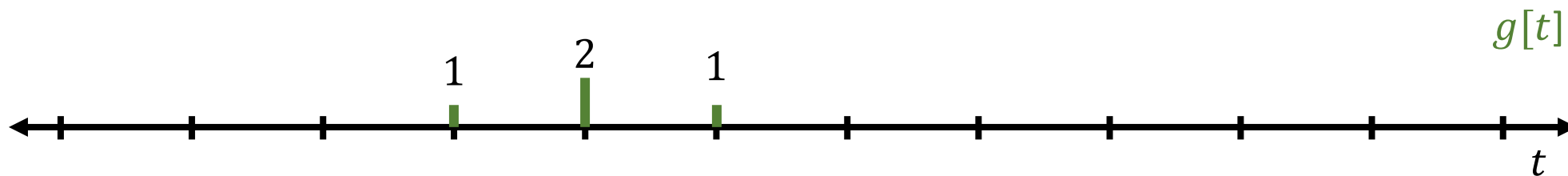
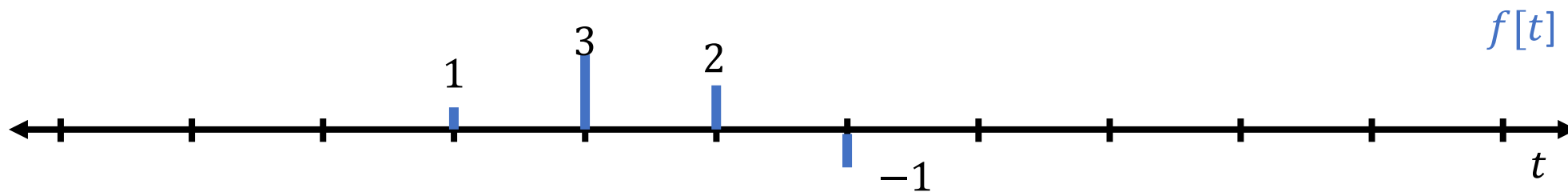
$$(f * g)(t) := \int_{-\infty}^{\infty} f(t - \tau)g(\tau)d\tau$$

$$(f * g)[t] := \sum_{\tau} f[t - \tau]g[\tau]$$

# 1D Convolution

$$(f * g)[t] := \sum_{\tau} f[t - \tau]g[\tau]$$

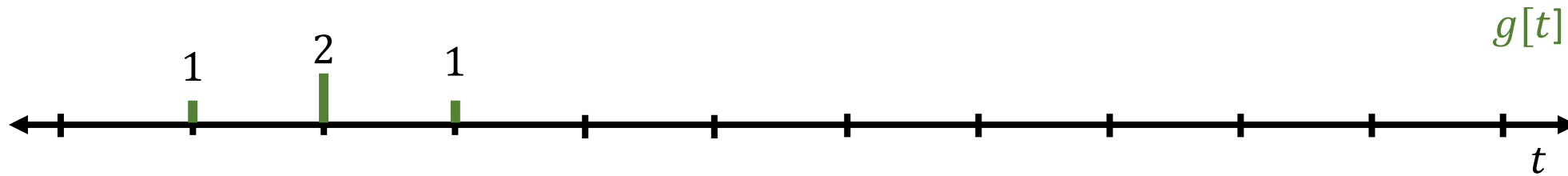
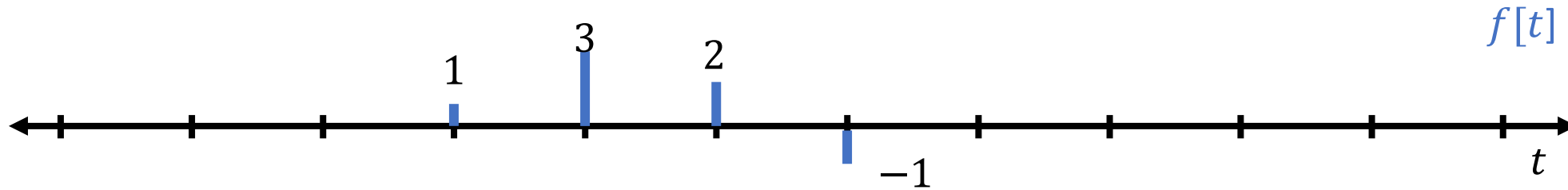
- Flip the filter and sliding



# 1D Convolution

$$(f * g)[t] := \sum_{\tau} f[t - \tau]g[\tau]$$

- Flip the filter and sliding



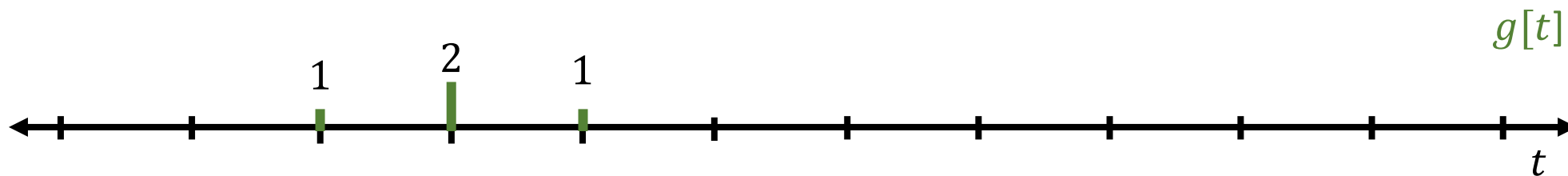
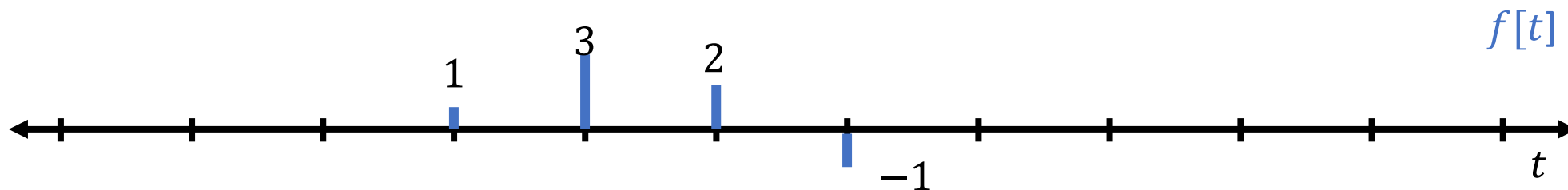
$$0 \cdot 1 + 0 \cdot 2 + 1 \cdot 1 = 1$$



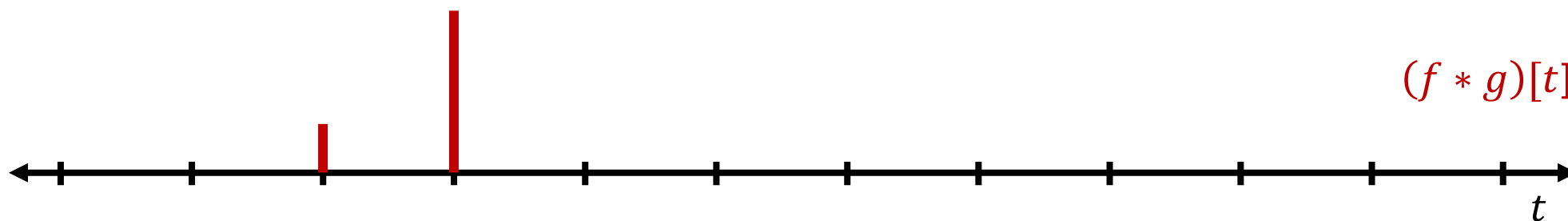
# 1D Convolution

$$(f * g)[t] := \sum_{\tau} f[t - \tau]g[\tau]$$

- Flip the filter and sliding



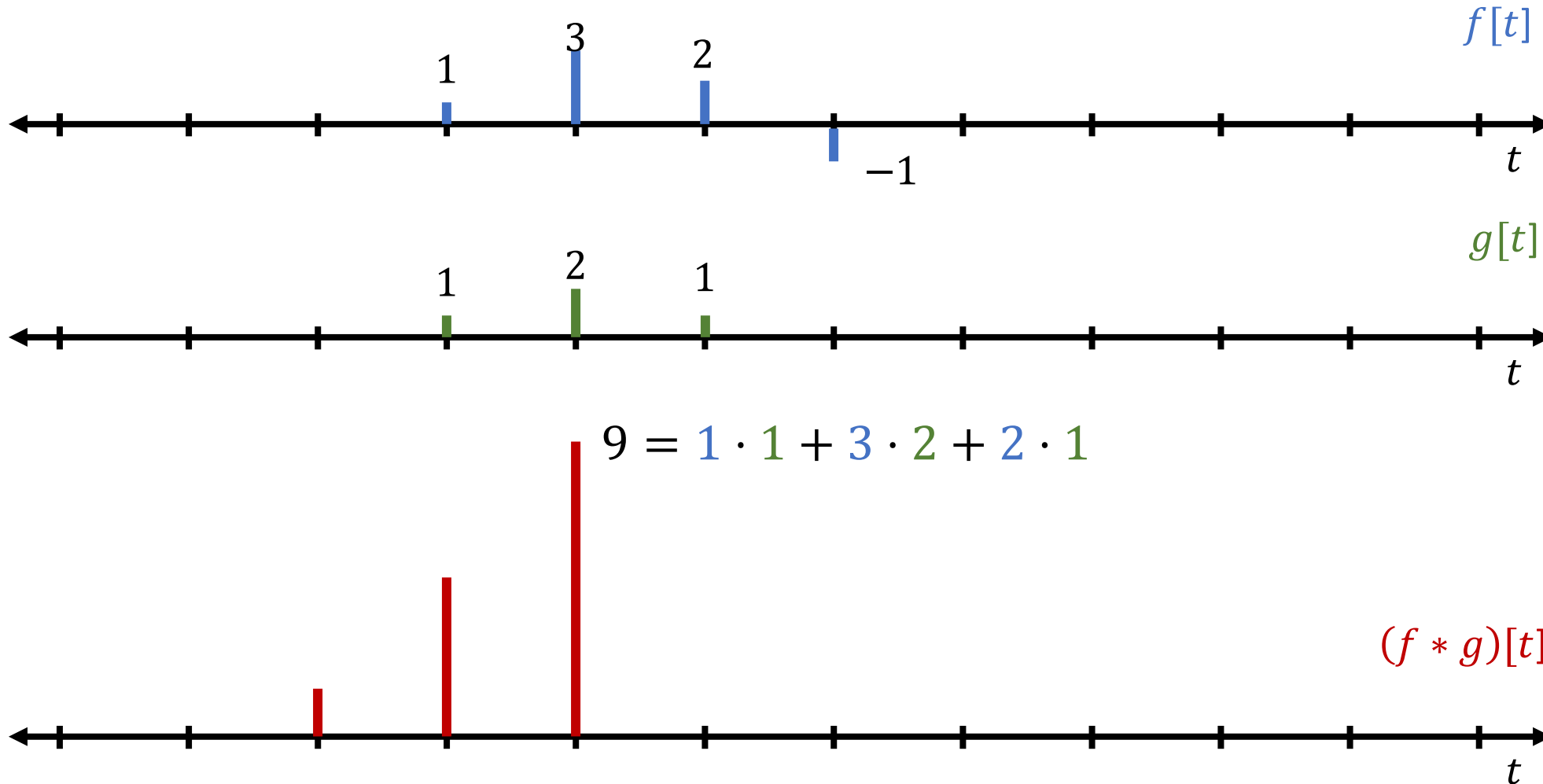
$$0 \cdot 1 + 1 \cdot 2 + 3 \cdot 1 = 5$$



# 1D Convolution

$$(f * g)[t] := \sum_{\tau} f[t - \tau]g[\tau]$$

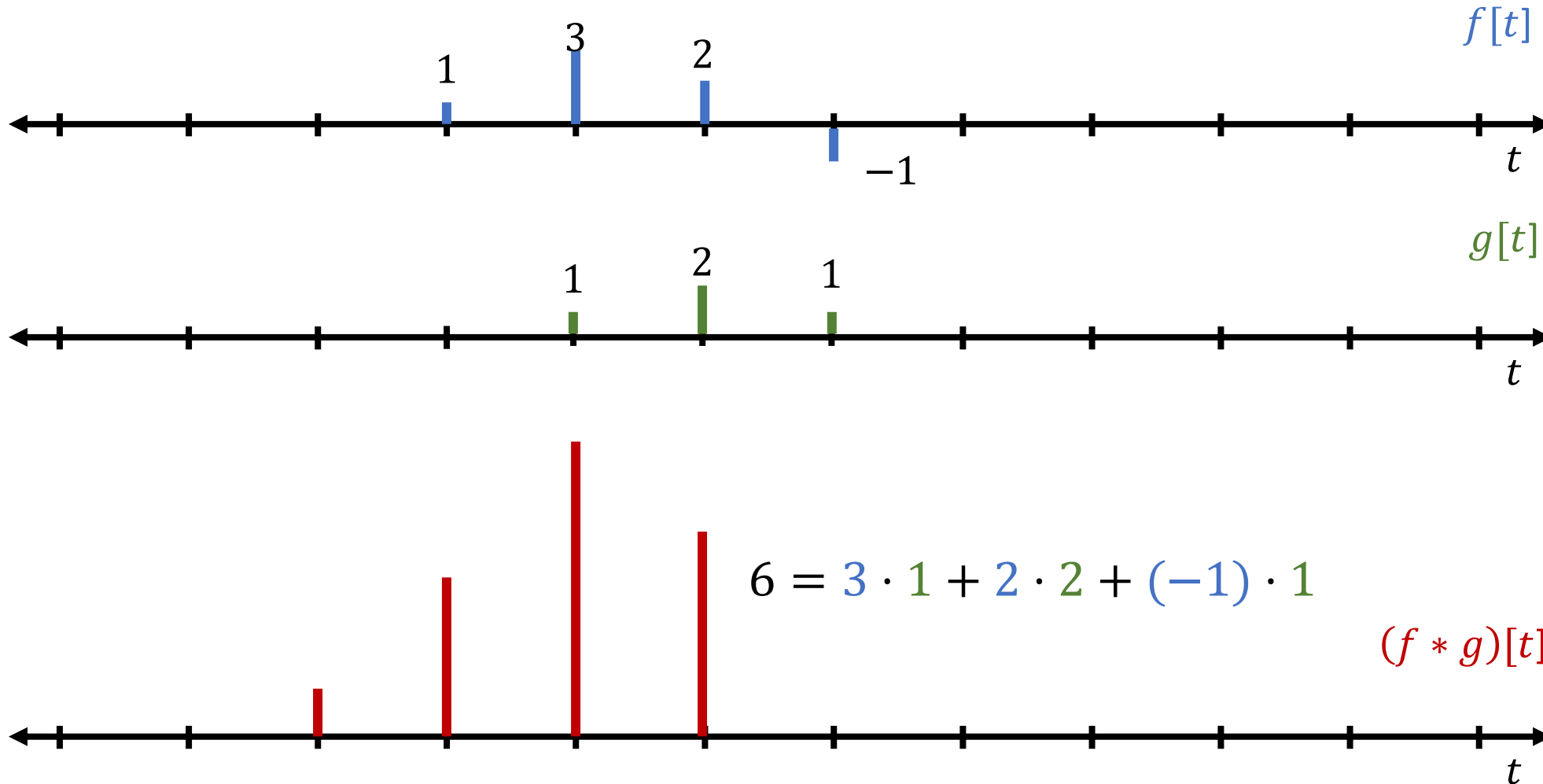
- Flip the filter and sliding



# 1D Convolution

$$(f * g)[t] := \sum_{\tau} f[t - \tau]g[\tau]$$

- Flip the filter and sliding

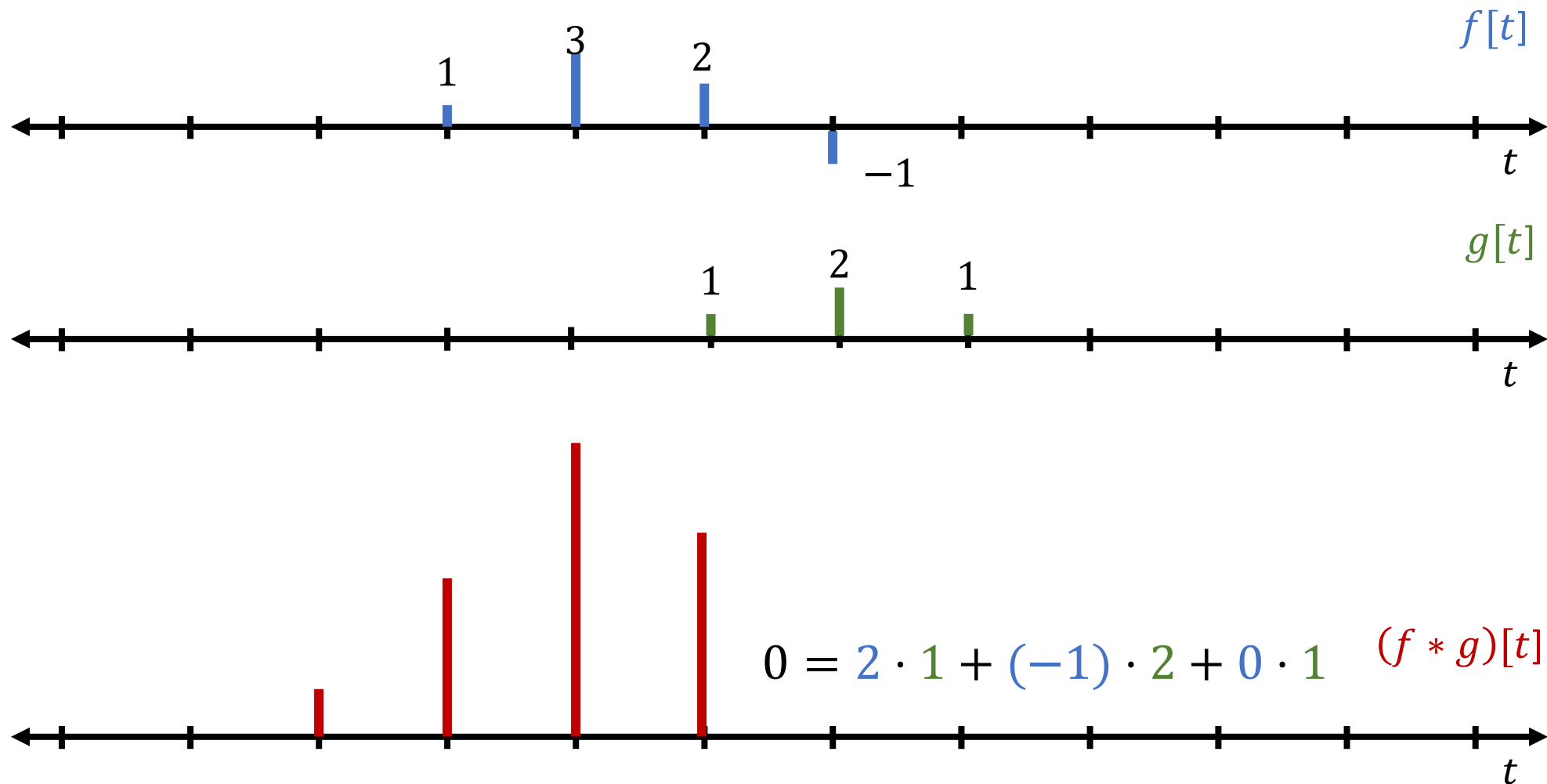




# 1D Convolution

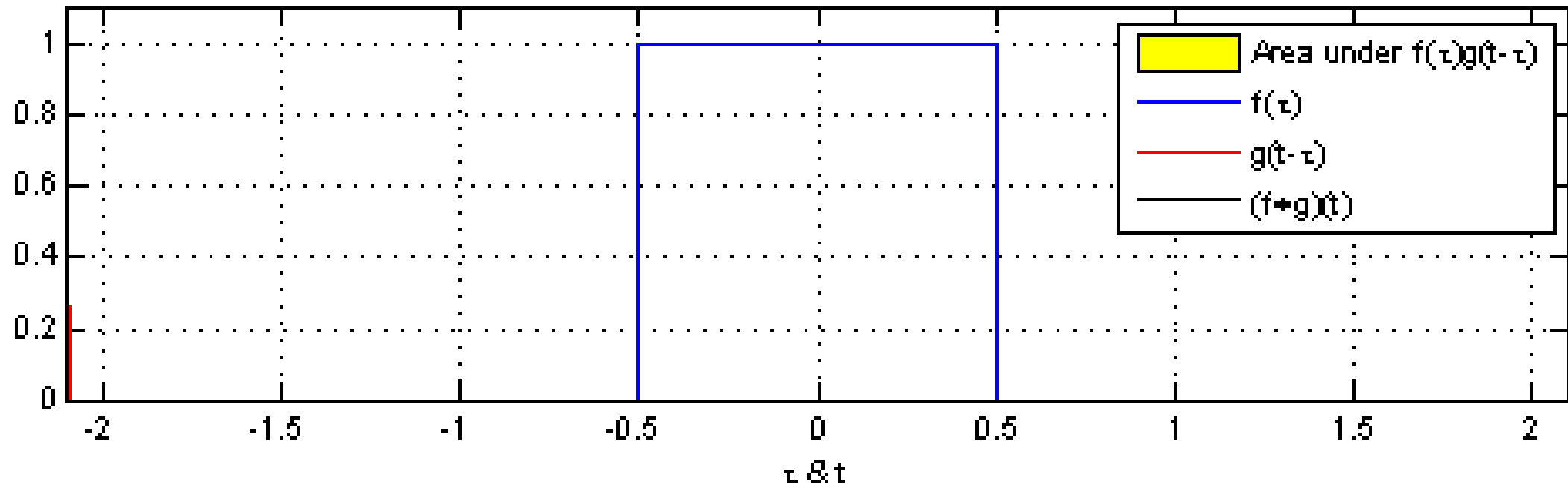
$$(f * g)[t] := \sum_{\tau} f[t - \tau]g[\tau]$$

- Flip the filter and sliding



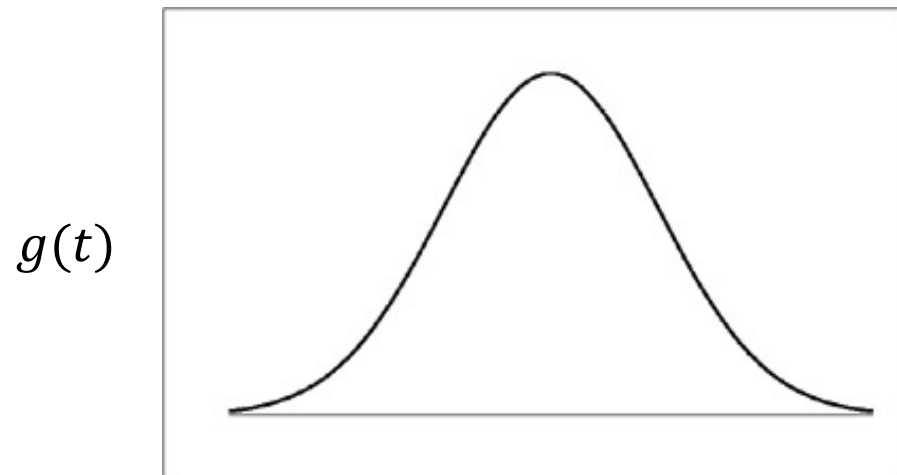
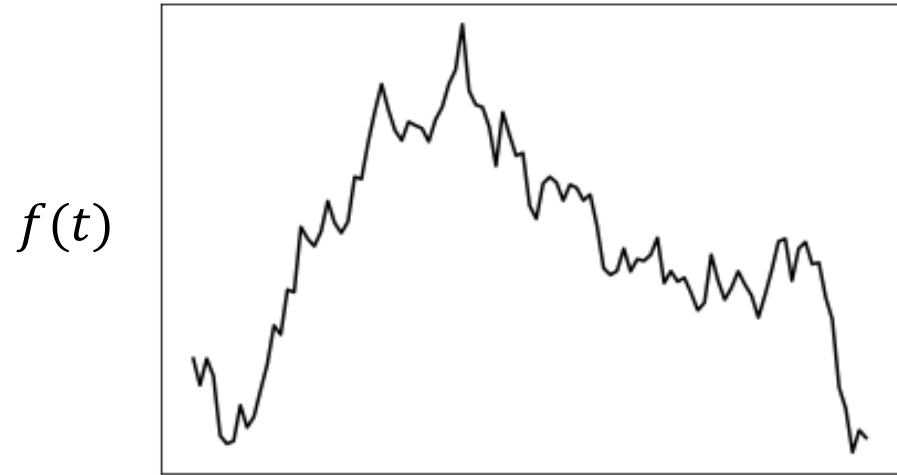
# 1D Convolution

- Example

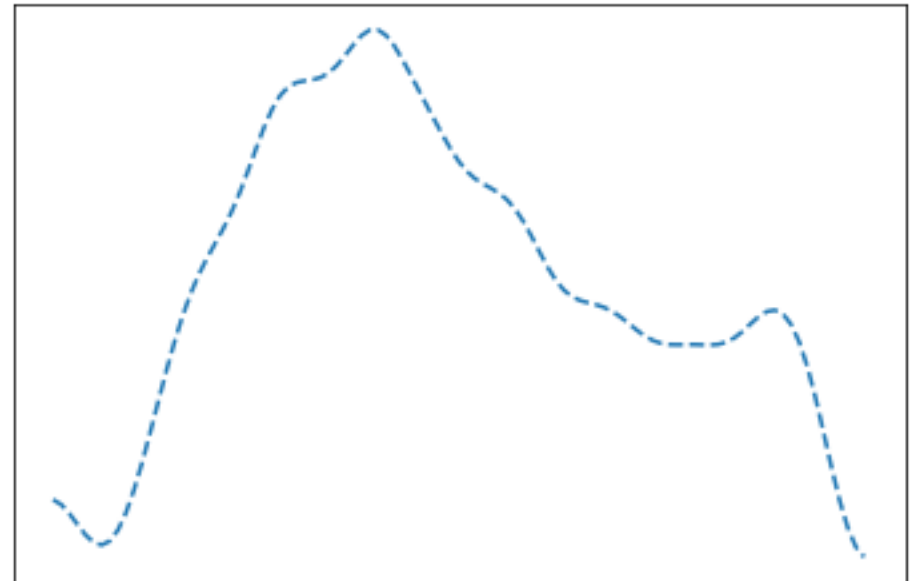


# 1D Convolution

- Gaussian filter



$$f(t) * g(t)$$



# 2D Convolution

$$(f * g)(s, t) := \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(s - \tau_1, t - \tau_2) g(\tau_1, \tau_2) d\tau_1 d\tau_2$$

$$(f * g)[s, t] := \sum_{\tau_1} \sum_{\tau_2} f[s - \tau_1, t - \tau_2] g[\tau_1, \tau_2]$$

# 2D Convolution

- One input channel, e.g. gray color image
  - Padding=1, stride=1

1 <sub>0</sub>	3 <sub>0</sub>	2 <sub>0</sub>	0	0	0	0
1 <sub>0</sub>	3 <sub>1</sub>	3 <sub>3</sub>	2	3	3	0
3 <sub>0</sub>	1 <sub>3</sub>	1 <sub>1</sub>	2	1	1	0
0	3	3	3	1	2	0
0	2	2	1	2	1	0
0	2	3	2	1	2	0
0	0	0	0	0	0	0

16				

# 2D Convolution

- One input channel, e.g. gray color image
  - Padding=1, stride=1

0	1	3	2	0	0	0
0	1	3	3	3	3	0
0	3	1	1	1	1	0
0	3	3	3	1	2	0
0	2	2	1	2	1	0
0	2	3	2	1	2	0
0	0	0	0	0	0	0

16	28			

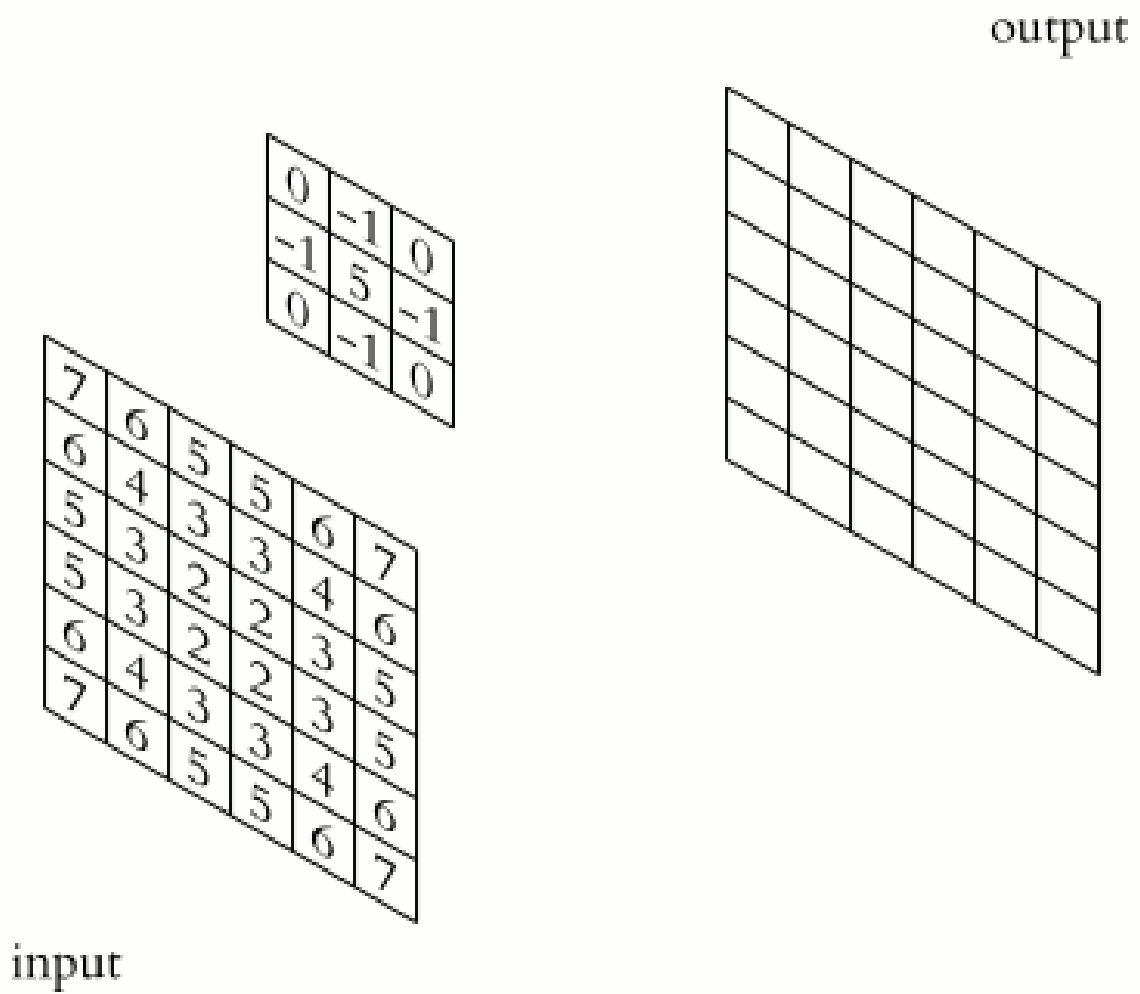
# 2D Convolution

- One input channel, e.g. gray color image
  - Padding=1, stride=1

0	0	1 <sub>0</sub>	3 <sub>0</sub>	2 <sub>0</sub>	0	0
0	1	1 <sub>1</sub>	3 <sub>2</sub>	3 <sub>3</sub>	3	0
0	3	3 <sub>1</sub>	1 <sub>2</sub>	1 <sub>1</sub>	1	0
0	3	3	3	1	2	0
0	2	2	1	2	1	0
0	2	3	2	1	2	0
0	0	0	0	0	0	0

16	28	24		

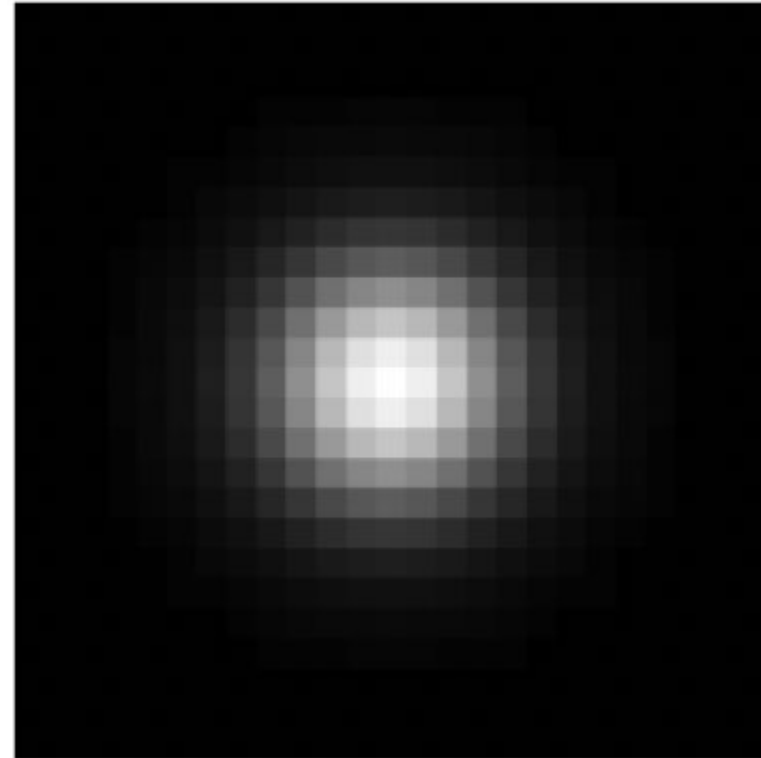
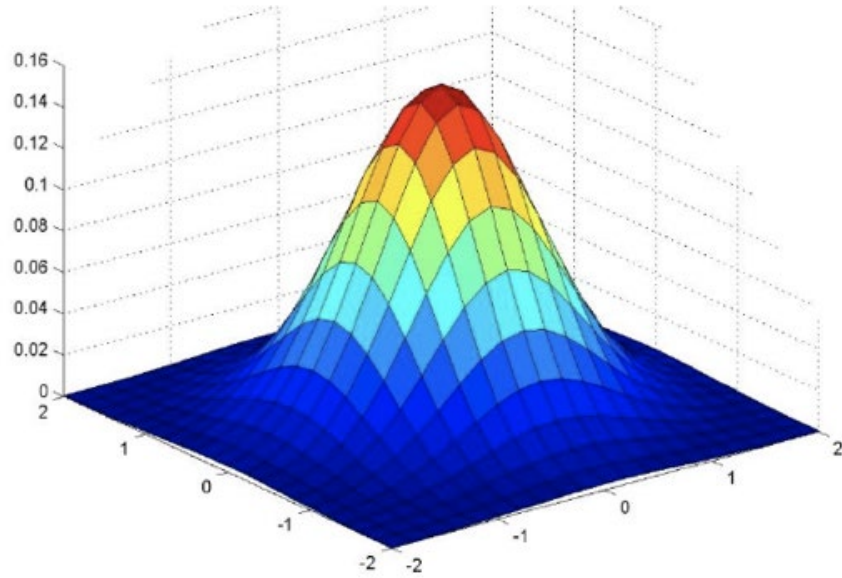
# 2D Convolution





# 2D Convolution

- 2D gaussian filter



# 2D Convolution

- 2D gaussian filter

Original Image (Left) Vs. Gaussian Filtered Image (Right)



# Properties of Convolution

$$f * g = g * f$$

Commutative

$$(f * g) * h = f * (g * h)$$

Associative

$$(f + g) * h = (f * h) + (g * h)$$

Linear

$$cf * h = c(f * h)$$

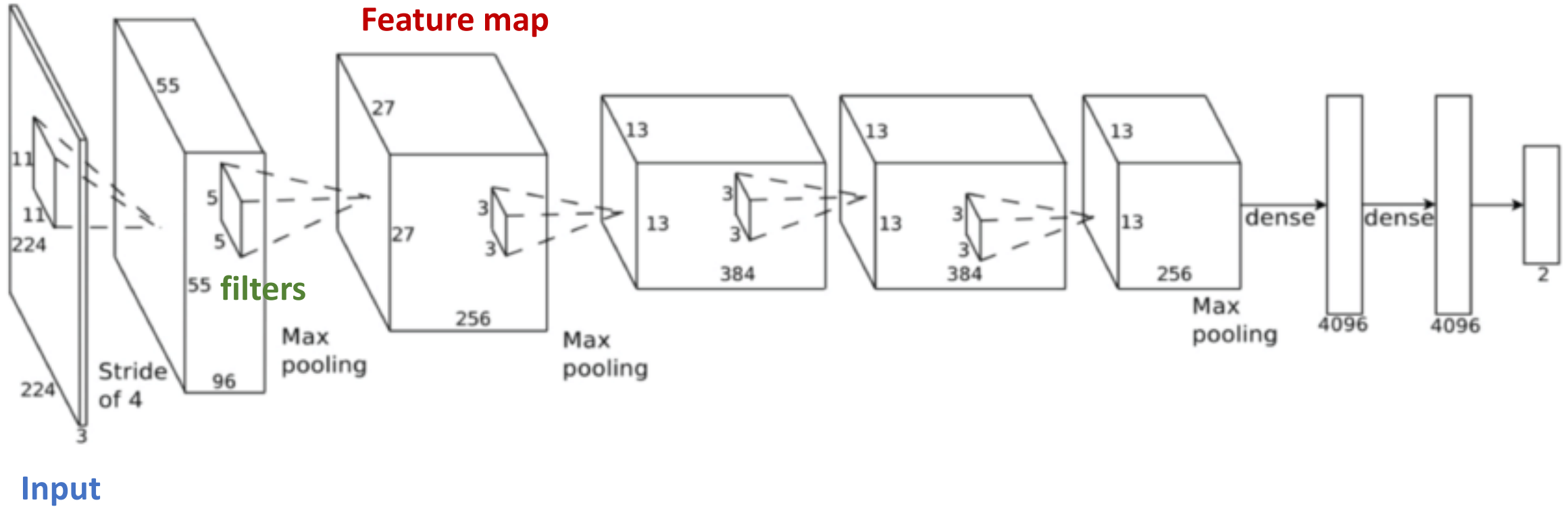
$$(L_t f) * h = L_t(f * h)$$

Translation equivariance

Translate  $f$  by  $t$

# Convolutional Neural Networks

# Convolutional Neural Network



# 2D Convolution

- Input\_channel=1, output\_channel=1, padding=1, stride=1

1 <sub>0</sub>	3 <sub>0</sub>	2 <sub>0</sub>	0	0	0	0
1 <sub>0</sub>	3 <sub>1</sub>	3 <sub>3</sub>	2	3	3	0
3 <sub>0</sub>	1 <sub>3</sub>	1 <sub>1</sub>	2	1	1	0
0	3	3	3	1	2	0
0	2	2	1	2	1	0
0	2	3	2	1	2	0
0	0	0	0	0	0	0

16				

# 2D Convolution

- Input\_channel=1, output\_channel=1, padding=1, stride=1

0	1	3	2	0	0	0
0	1	3	3	3	3	0
0	3	1	1	1	1	0
0	3	3	3	1	2	0
0	2	2	1	2	1	0
0	2	3	2	1	2	0
0	0	0	0	0	0	0

16	28			

# 2D Convolution

- Input\_channel=1, output\_channel=1, padding=1, stride=1

0	0	1 <sub>0</sub>	3 <sub>0</sub>	2 <sub>0</sub>	0	0
0	1	1 <sub>1</sub>	3 <sub>1</sub>	3 <sub>1</sub>	3	0
0	3	3 <sub>1</sub>	1 <sub>2</sub>	1 <sub>1</sub>	1	0
0	3	3	3	1	2	0
0	2	2	1	2	1	0
0	2	3	2	1	2	0
0	0	0	0	0	0	0

16	28	24		



# 2D Convolution

- Input\_channel=1, output\_channel=1, padding=1, stride=1

0	0	0	0	0	0	0
0	1	3	2	3	3	0
0	3	1	2	1	1	0
0	3	3	3	1	2	0
0	2	2	1	2	1	0
0	2	3	2	1	2	0
0	0	0	0	0	0	0

Input

\*

1	3	2
1	3	3
3	1	1

filter

=

16	28	24	28	16
27	41	38	37	21
33	40	33	25	18
32	40	37	29	17
25	27	21	20	12

Output

# 2D Convolution

- Input\_channel=1, output\_channel=1, padding=1, stride=2

1	3	2	0	0	0	0
1	3	3	2	3	3	0
3	1	1	2	1	1	0
0	3	3	3	1	2	0
0	2	2	1	2	1	0
0	2	3	2	1	2	0
0	0	0	0	0	0	0

16		

# 2D Convolution

- Input\_channel=1, output\_channel=1, padding=1, stride=2

0	0	1 <sub>0</sub>	3 <sub>0</sub>	2 <sub>0</sub>	0	0
0	1	3 <sub>3</sub>	3 <sub>2</sub>	3 <sub>3</sub>	3	0
0	3	3 <sub>1</sub>	1 <sub>2</sub>	1 <sub>1</sub>	1	0
0	3	3	3	1	2	0
0	2	2	1	2	1	0
0	2	3	2	1	2	0
0	0	0	0	0	0	0

16	24	

# 2D Convolution

- Input\_channel=1, output\_channel=1, padding=1, stride=2

0	0	0	0	1	3	2
0	1	3	2	1	3	3
0	3	1	2	3	1	1
0	3	3	3	1	2	0
0	2	2	1	2	1	0
0	2	3	2	1	2	0
0	0	0	0	0	0	0

16	24	16

# 2D Convolution

- Input\_channel=1, output\_channel=1, padding=1, stride=2

0	0	0	0	0	0	0
0	1	3	2	3	3	0
1	3	2	2	1	1	0
1	3	3	3	1	2	0
3	1	1	1	2	1	0
0	2	3	2	1	2	0
0	0	0	0	0	0	0

16	24	16
33		

# 2D Convolution

- Input\_channel=1, output\_channel=1, padding=1, stride=2

0	0	0	0	0	0	0
0	1	3	2	3	3	0
0	3	1	2	1	1	0
0	3	3	3	1	2	0
0	2	2	1	2	1	0
0	2	3	2	1	2	0
0	0	0	0	0	0	0

Input

\*

1	3	2
1	3	3
3	1	1

filter

=

16	24	16
33	33	18
25	21	12

Output

# 2D Convolution

- Input\_channel=1, output\_channel=1, padding=2, stride=1

0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0
0	0	1	3	2	3	3	0	0
0	0	3	1	2	1	1	0	0
0	0	3	3	3	1	2	0	0
0	0	2	2	1	2	1	0	0
0	0	2	3	2	1	2	0	0
0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0

\*

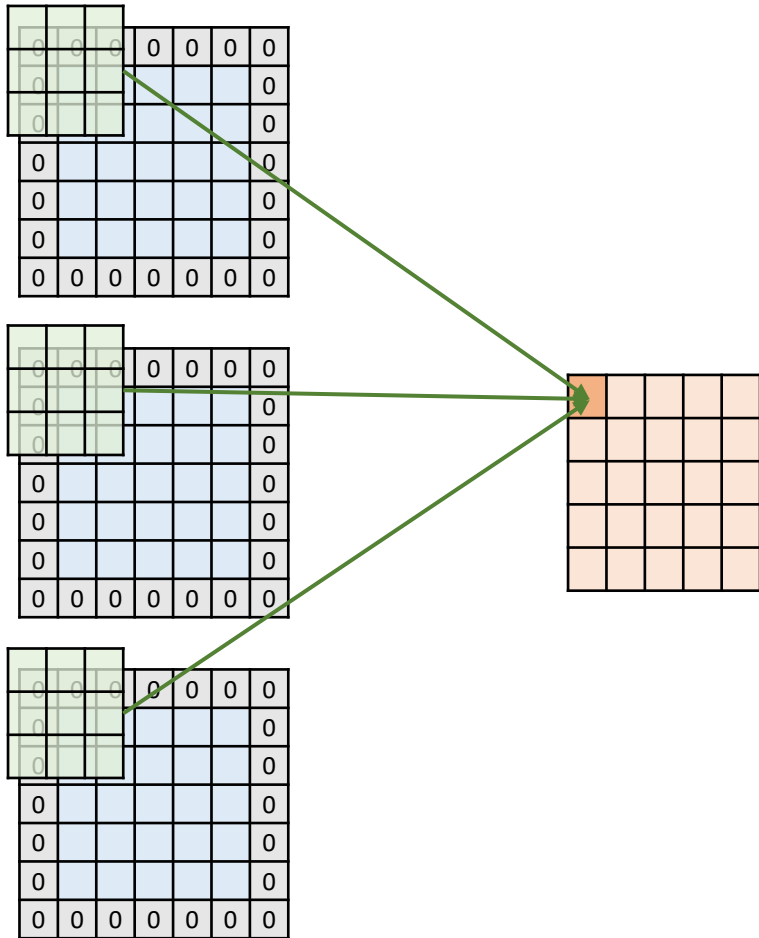
1	3	2
1	3	3
3	1	1

=

1	4	8	14	12	12	9
6	16	28	24	28	16	6
14	27	41	38	37	21	10
17	33	40	33	25	18	6
14	32	40	37	29	17	9
10	25	27	21	20	12	3
4	12	15	11	9	7	2

# 2D Convolution

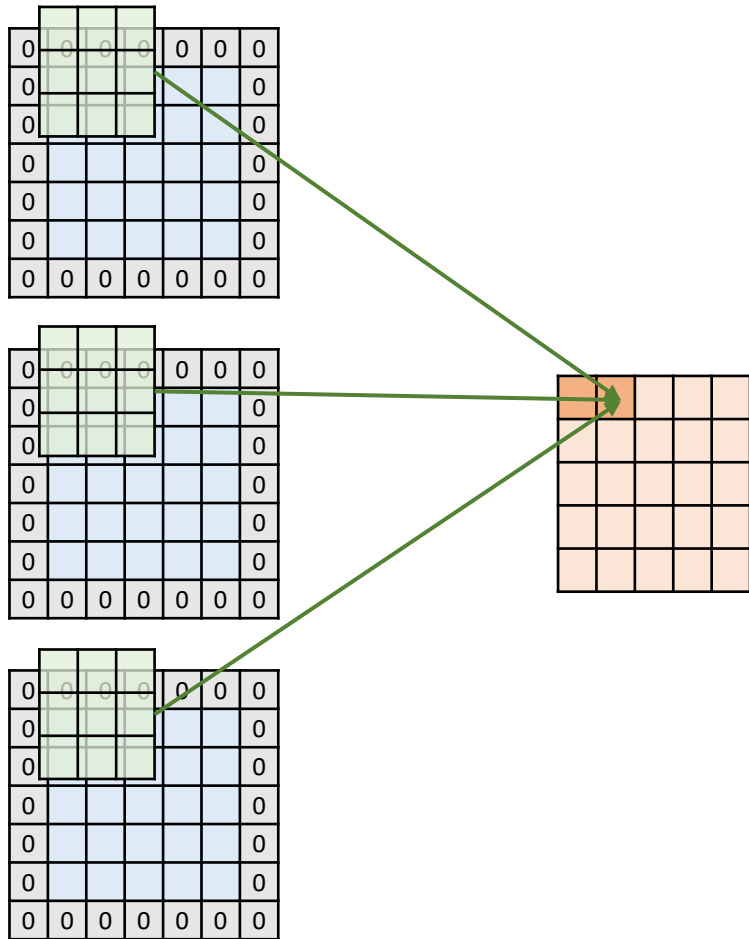
- **Input\_channel=3**, output\_channel=1, padding=1, stride=1





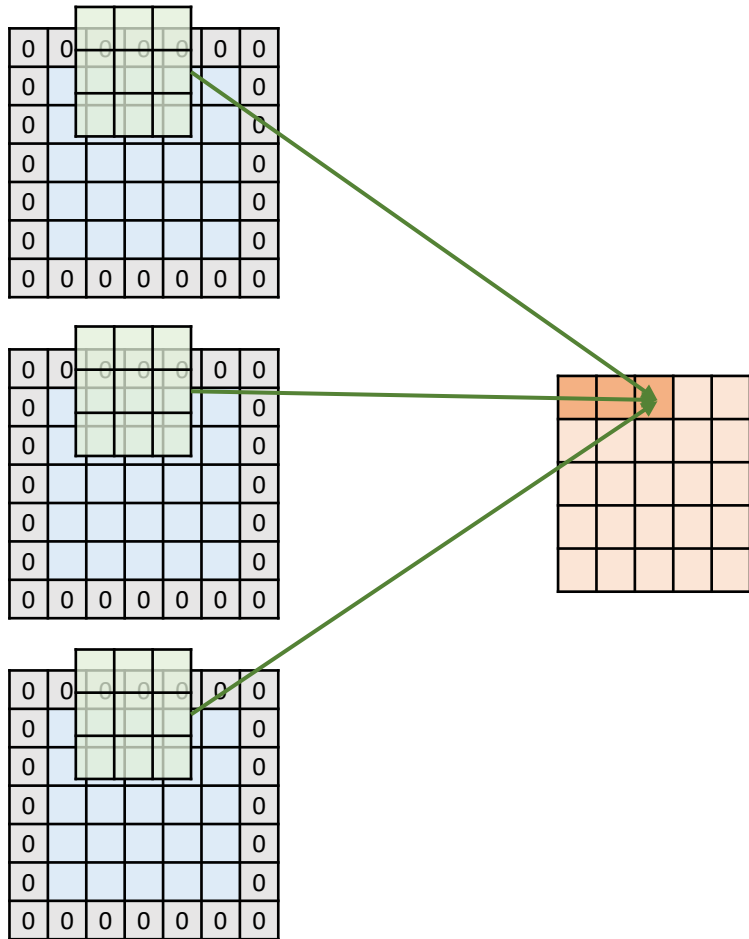
# 2D Convolution

- **Input\_channel=3**, output\_channel=1, padding=1, stride=1



# 2D Convolution

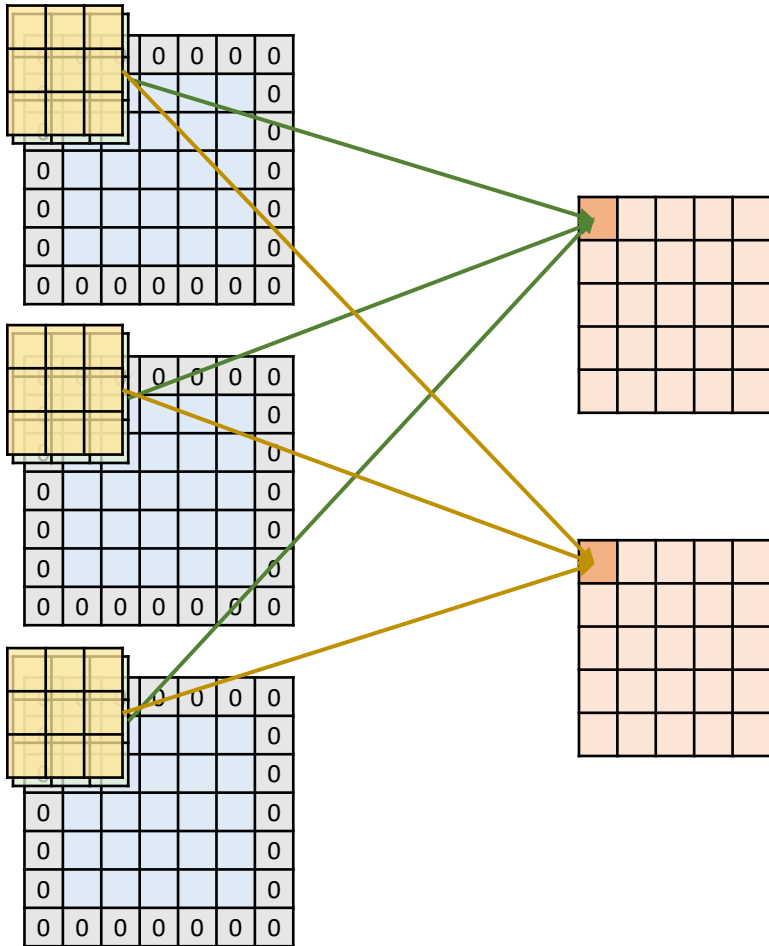
- **Input\_channel=3**, output\_channel=1, padding=1, stride=1





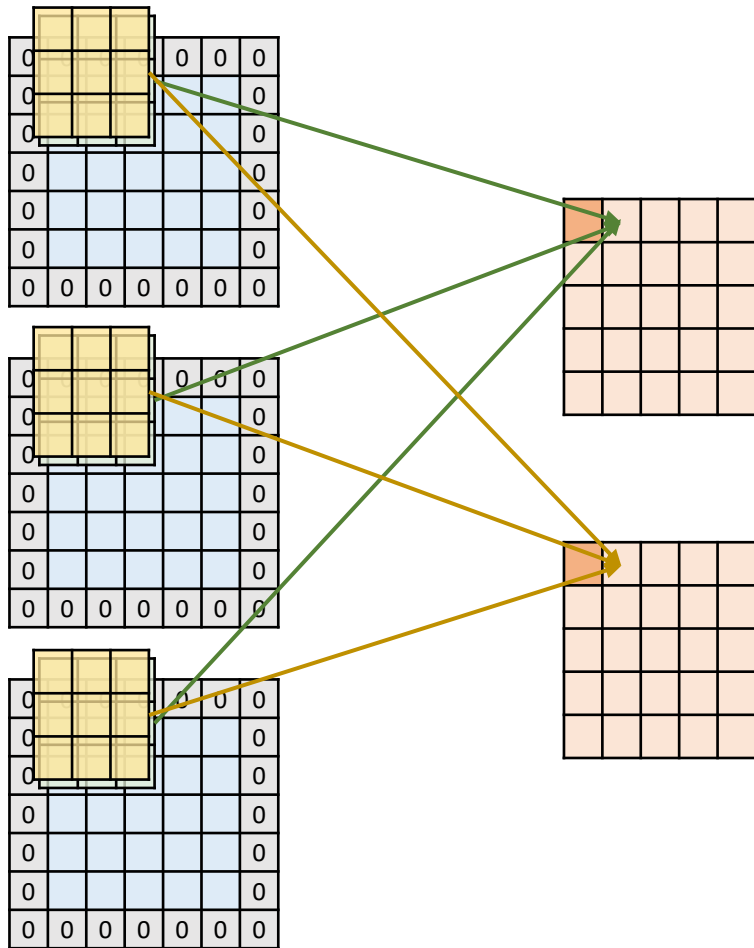
# 2D Convolution

- Input\_channel=3, output\_channel=2, padding=1, stride=1



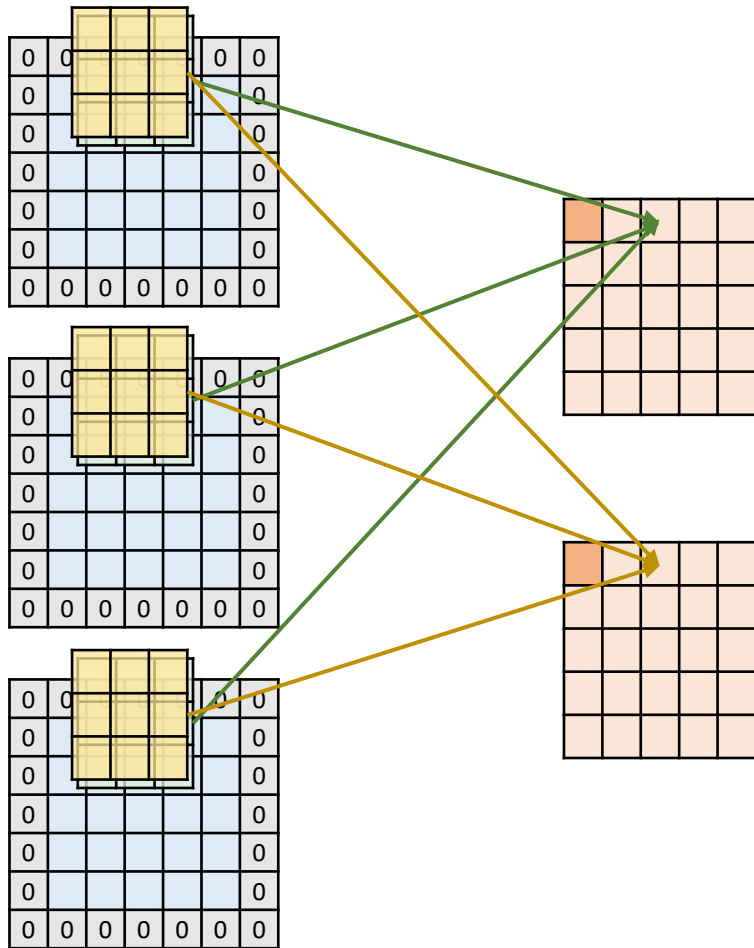
# 2D Convolution

- $\text{Input\_channel}=3$ ,  $\text{output\_channel}=2$ ,  $\text{padding}=1$ ,  $\text{stride}=1$



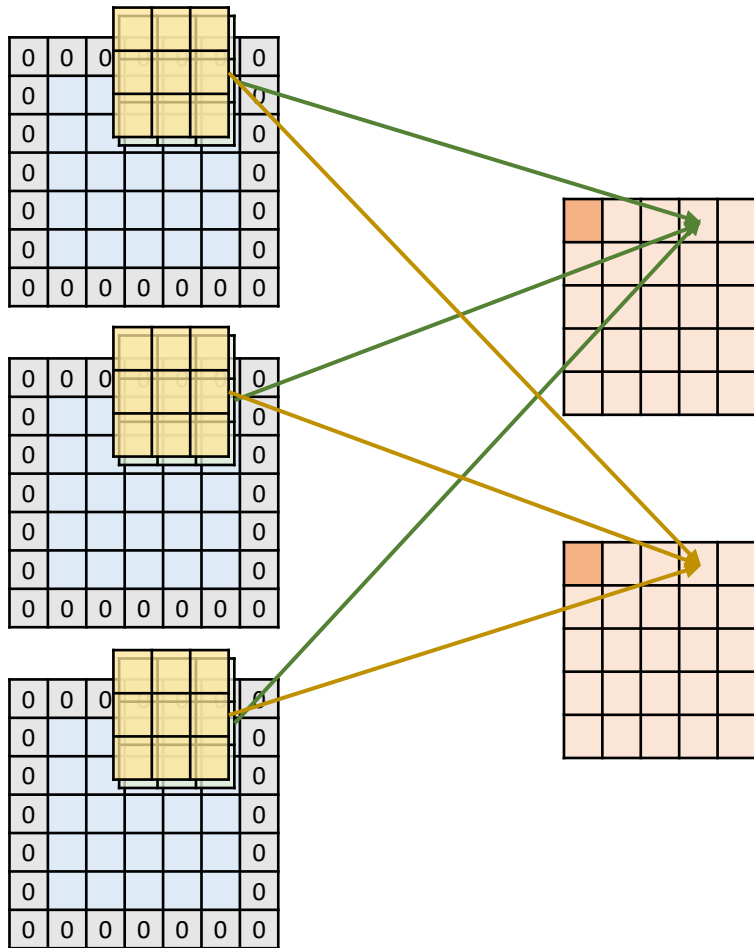
# 2D Convolution

- **Input\_channel=3, output\_channel=2, padding=1, stride=1**



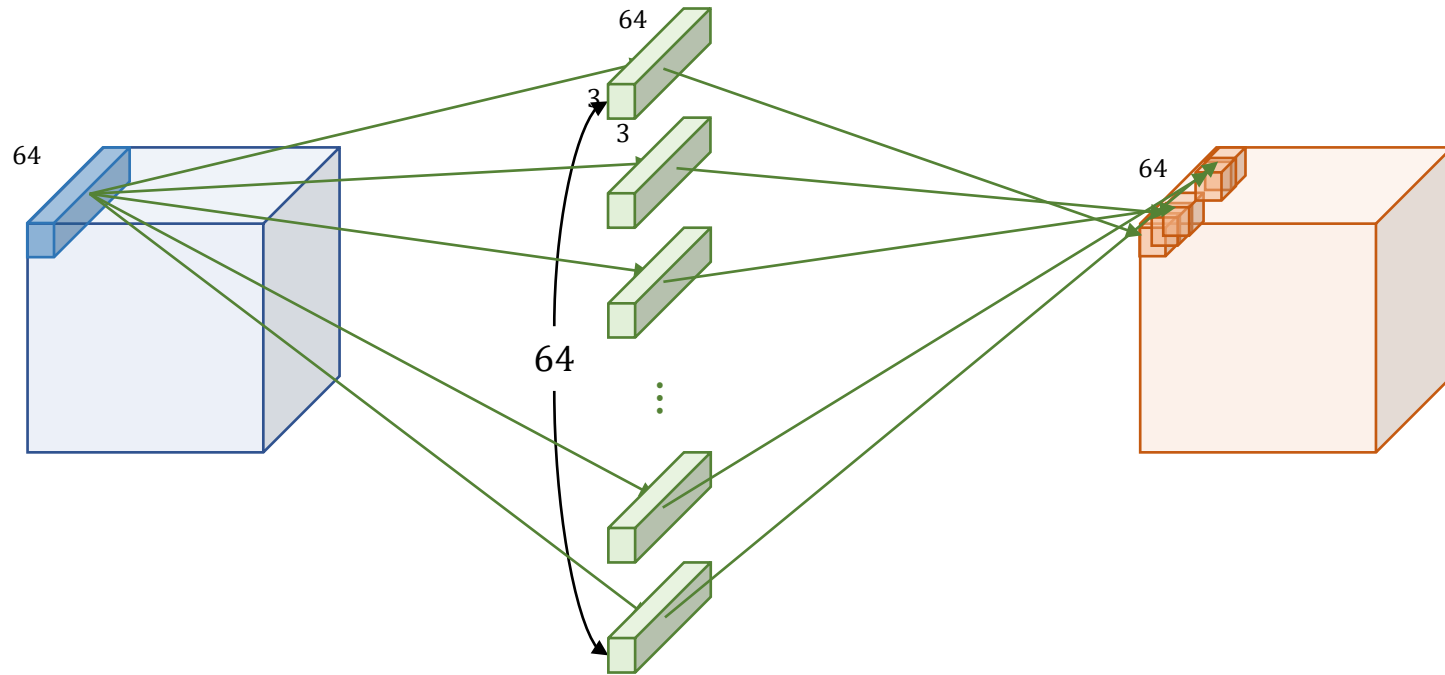
# 2D Convolution

- $\text{Input\_channel}=3$ ,  $\text{output\_channel}=2$ ,  $\text{padding}=1$ ,  $\text{stride}=1$



# 2D Convolution

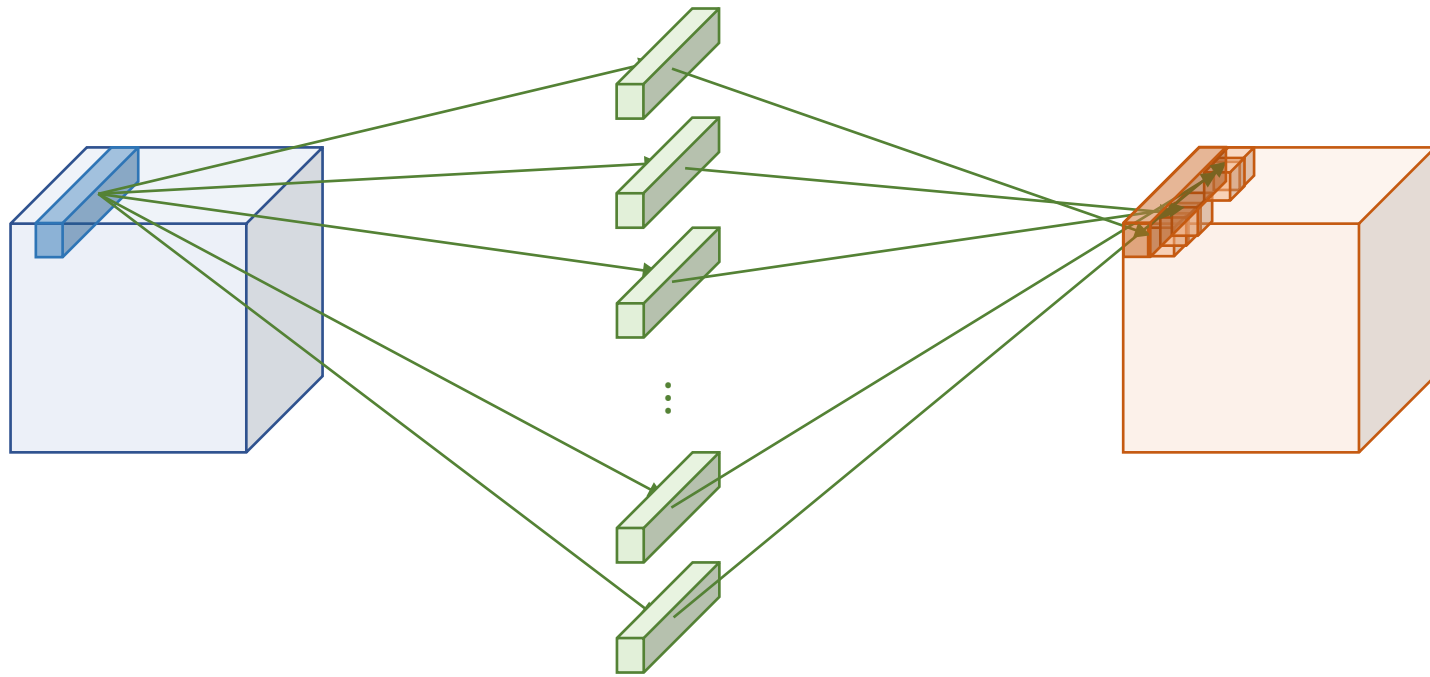
- **Input\_channel=64, output\_channel=64, kernel\_size=3, padding=1, stride=1**





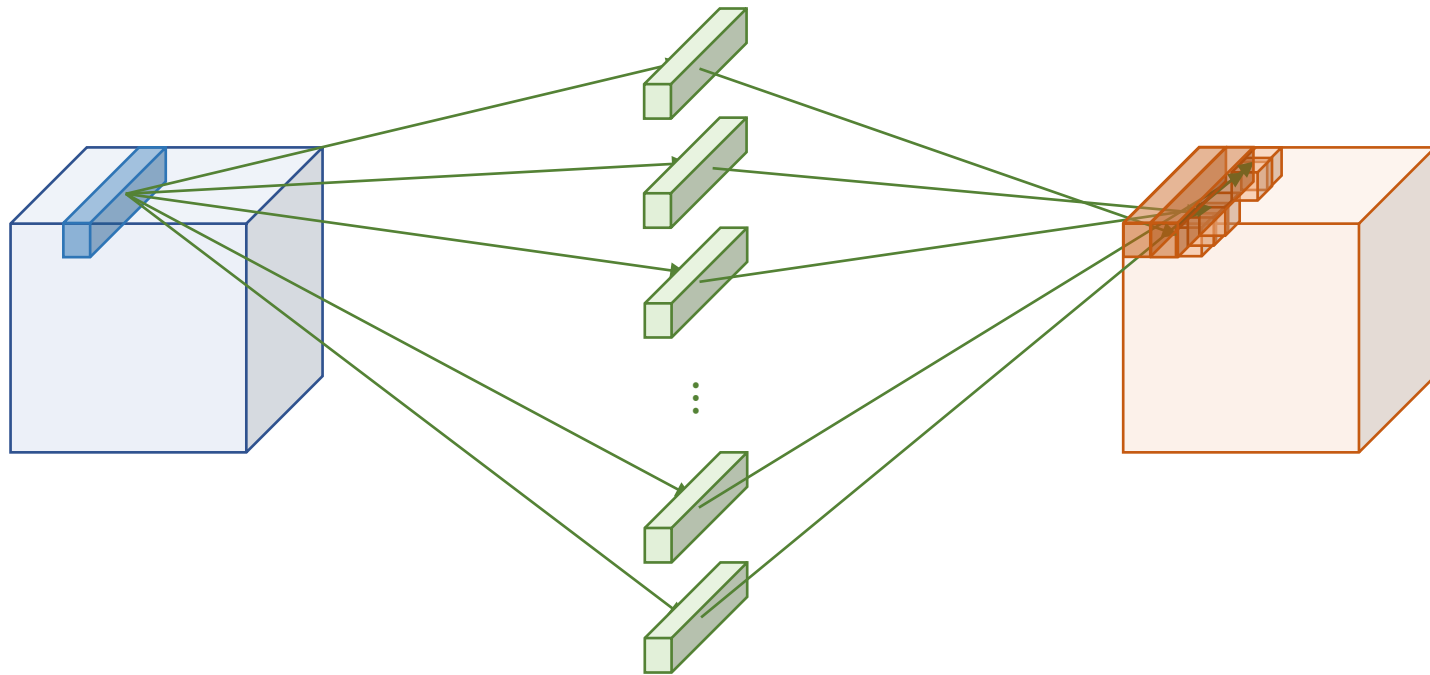
# 2D Convolution

- $\text{Input\_channel}=64$ ,  $\text{output\_channel}=64$ ,  $\text{kernel\_size}=3$ ,  $\text{padding}=1$ ,  $\text{stride}=1$



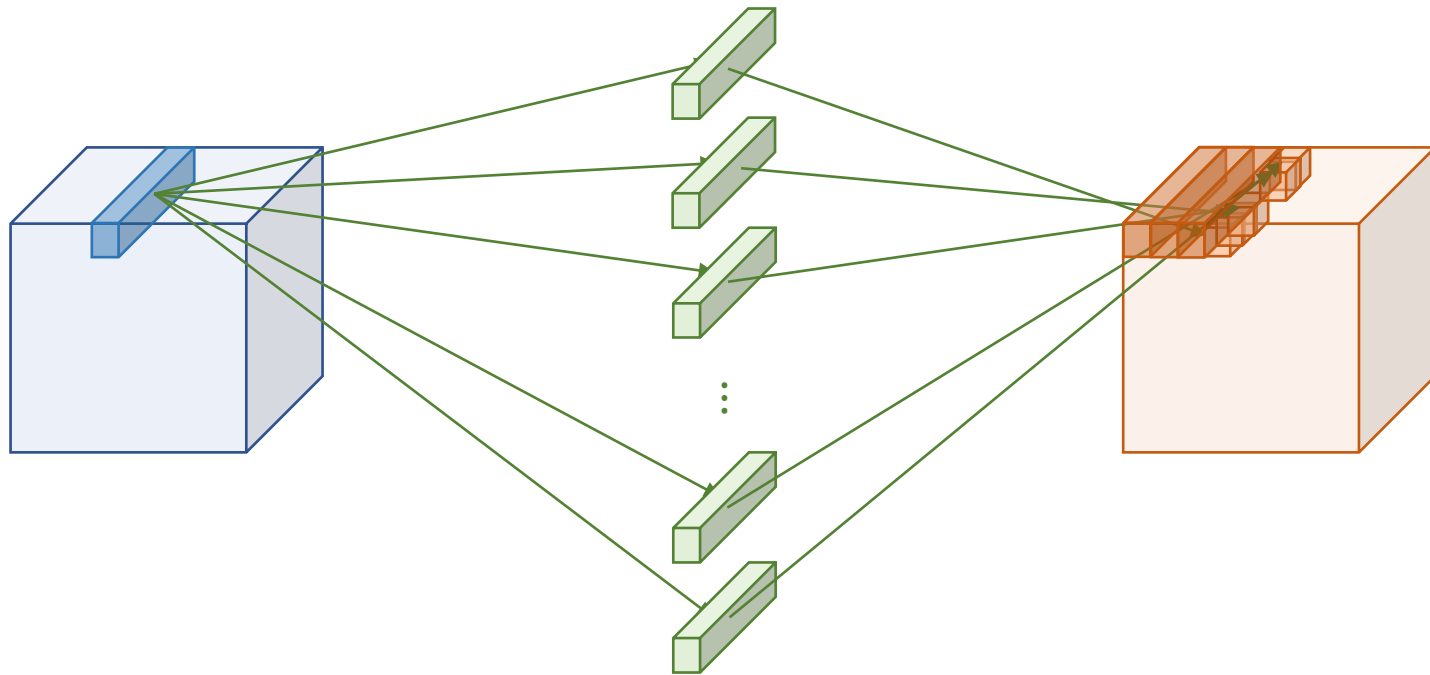
# 2D Convolution

- **Input\_channel=64, output\_channel=64, kernel\_size=3, padding=1, stride=1**



# 2D Convolution

- **Input\_channel=64, output\_channel=64, kernel\_size=3, padding=1, stride=1**



# Convolutions in PyTorch

---

**CLASS** `torch.nn.Conv1d(in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1, groups=1, bias=True, padding_mode='zeros', device=None, dtype=None)` [\[SOURCE\]](#)

---

**CLASS** `torch.nn.Conv2d(in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1, groups=1, bias=True, padding_mode='zeros', device=None, dtype=None)` [\[SOURCE\]](#)

---

**CLASS** `torch.nn.Conv3d(in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1, groups=1, bias=True, padding_mode='zeros', device=None, dtype=None)` [\[SOURCE\]](#)

# Max Pooling

- Pooling a maximum value given the window
- Used to reduce the size of feature maps
- Example) stride=2, padding=1

0	0	0	0	0	0	0	0
0	1	3	2	3	3	0	
0	3	1	2	1	1	0	
0	3	3	3	1	2	0	
0	2	2	1	2	1	0	
0	2	3	2	1	2	0	
0	0	0	0	0	0	0	

3		

# Max Pooling

- Pooling a maximum value given the window
- Used to reduce the size of feature maps
- Example) stride=2, padding=1

0	0	0	0	0	0	0
0	1	3	2	3	3	0
0	3	1	2	1	1	0
0	3	3	3	1	2	0
0	2	2	1	2	1	0
0	2	3	2	1	2	0
0	0	0	0	0	0	0

3	3	

# Max Pooling

- Pooling a maximum value given the window
- Used to reduce the size of feature maps
- Example) stride=2, padding=1

0	0	0	0	0	0	0
0	1	3	2	3	3	0
0	3	1	2	1	1	0
0	3	3	3	1	2	0
0	2	2	1	2	1	0
0	2	3	2	1	2	0
0	0	0	0	0	0	0

3	3	3

# Max Pooling in PyTorch

---

**CLASS** `torch.nn.MaxPool1d(kernel_size, stride=None, padding=0, dilation=1, return_indices=False, ceil_mode=False)`

[\[SOURCE\]](#)

---

**CLASS** `torch.nn.MaxPool2d(kernel_size, stride=None, padding=0, dilation=1, return_indices=False, ceil_mode=False)`

[\[SOURCE\]](#)

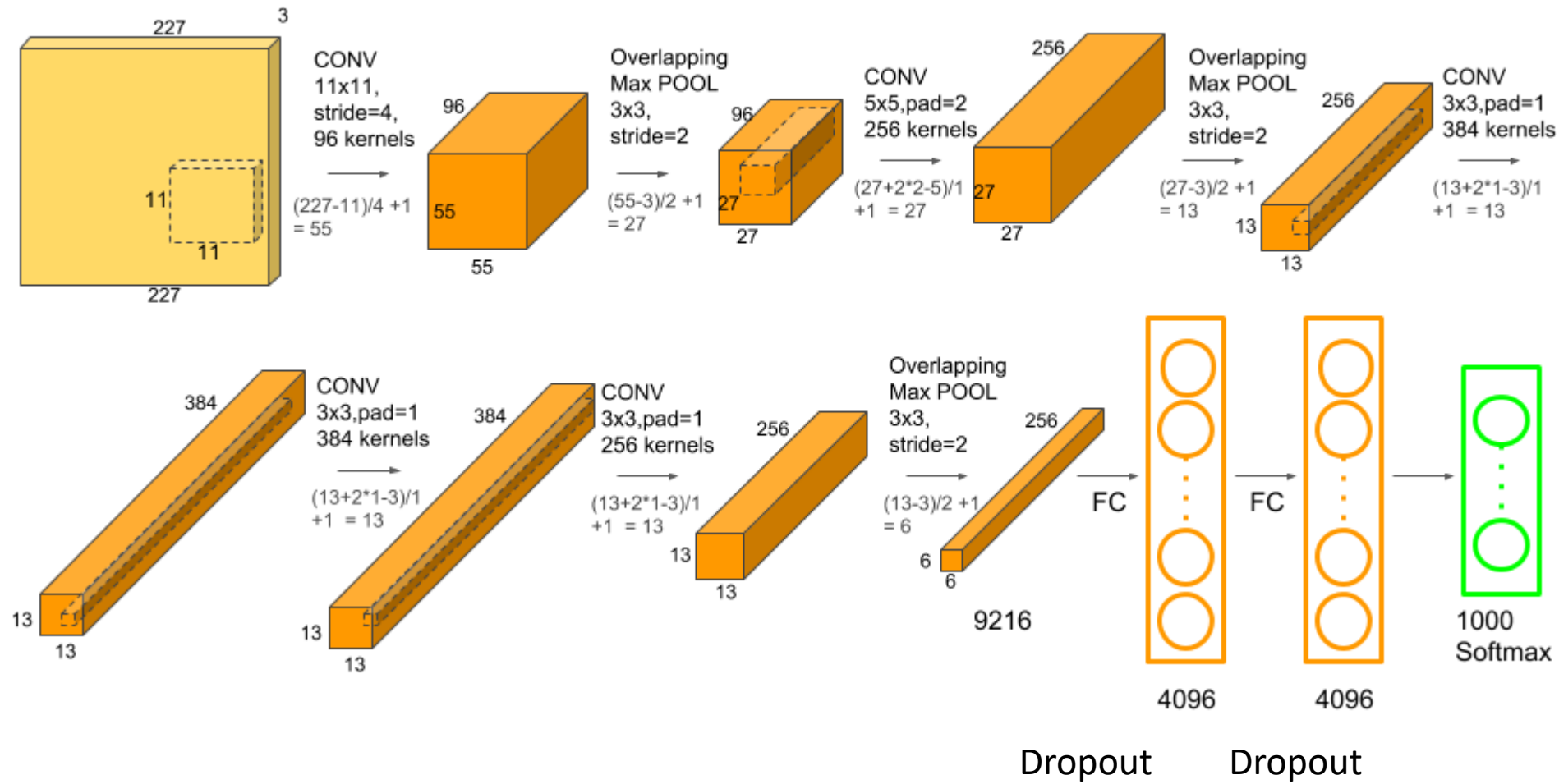
---

**CLASS** `torch.nn.MaxPool3d(kernel_size, stride=None, padding=0, dilation=1, return_indices=False, ceil_mode=False)`

[\[SOURCE\]](#)

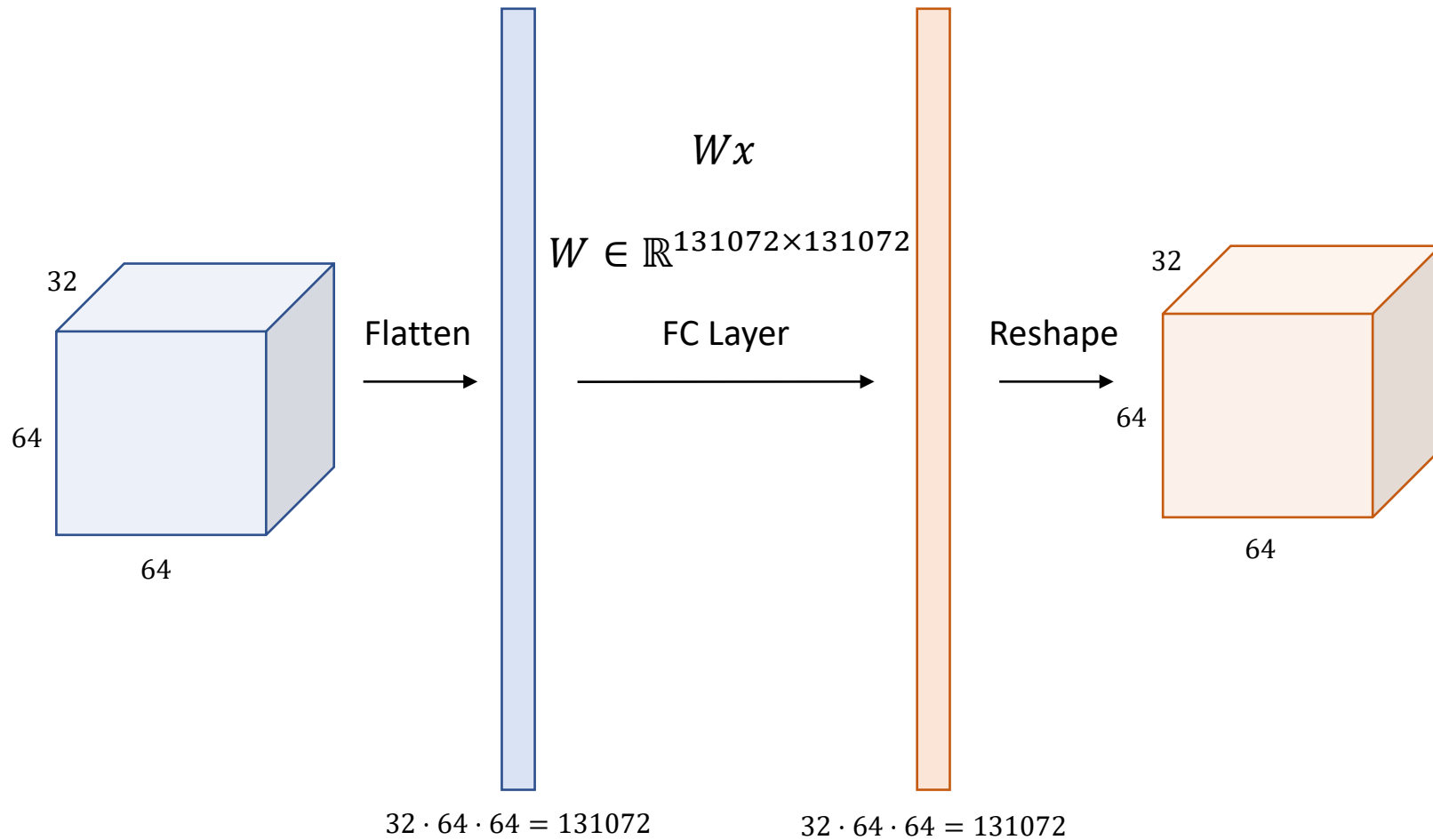


# AlexNet



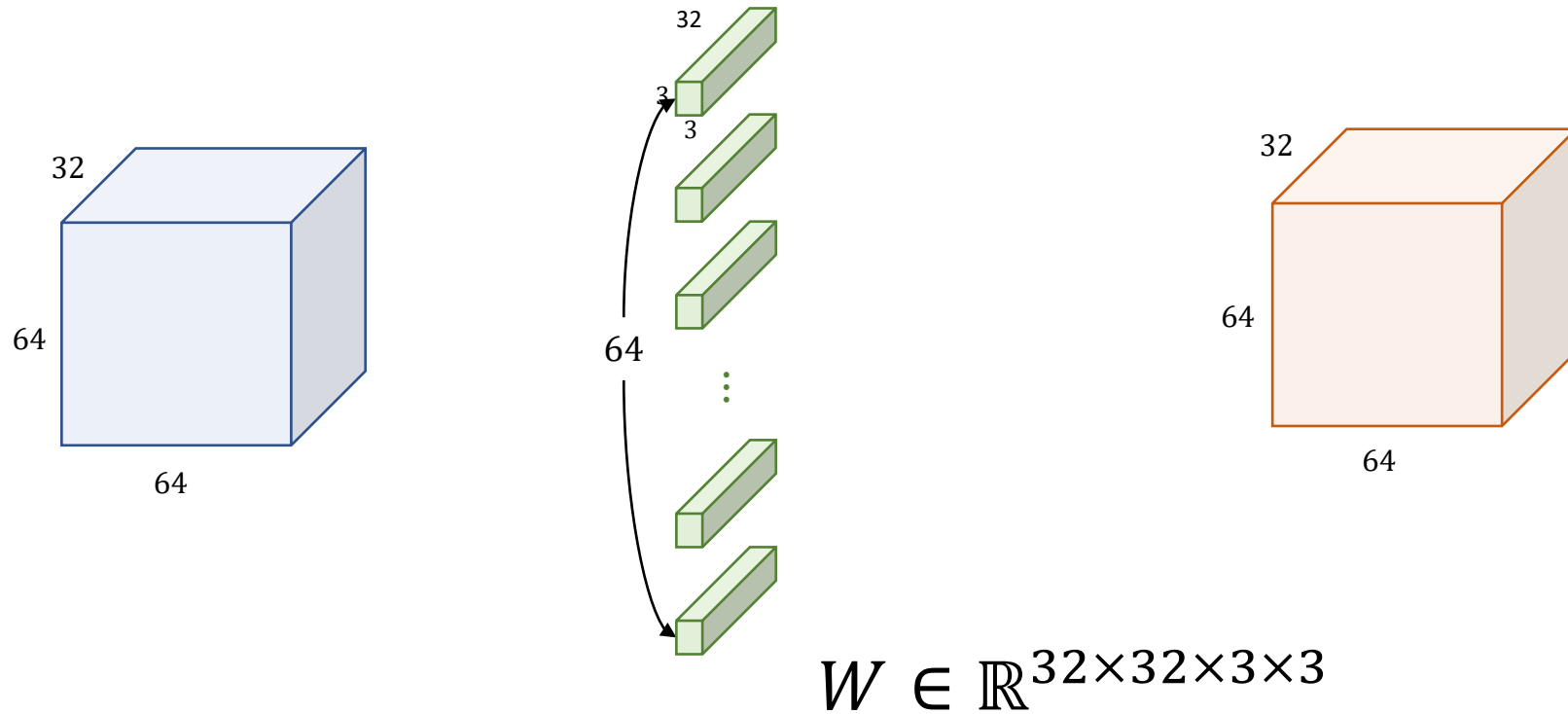
# Fully Connected Layer vs Convolutional Layer

- Translation equivariance and parameter sharing



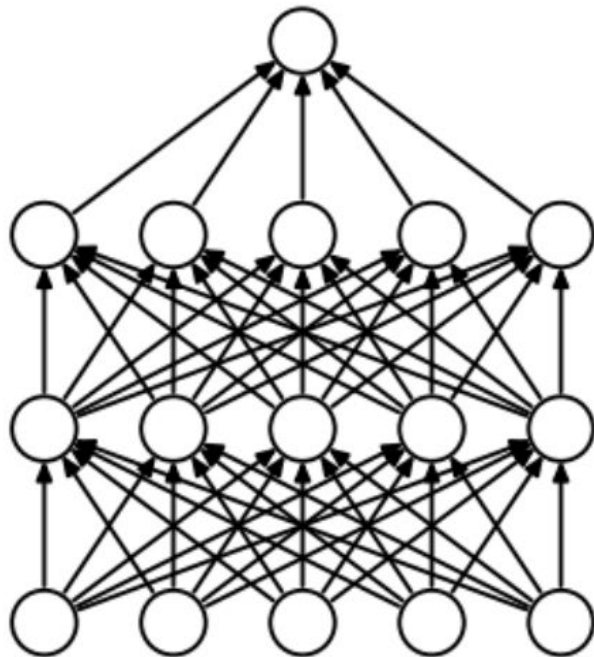
# Fully Connected Layer vs Convolutional Layer

- Translation equivariance and parameter sharing

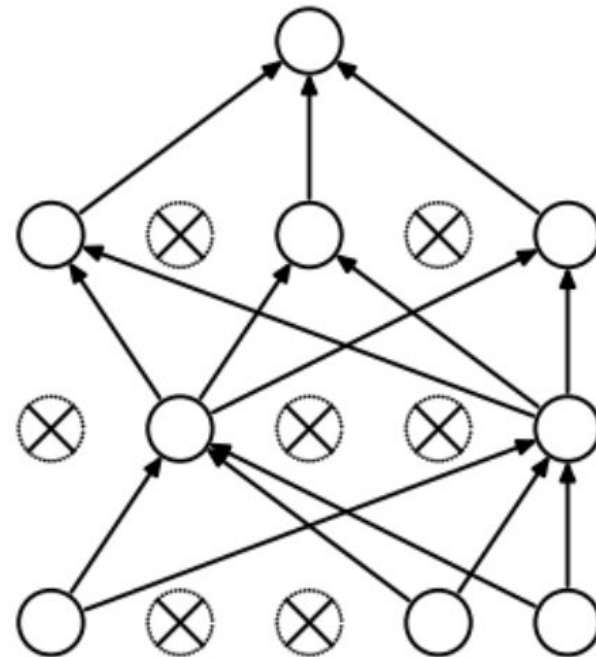


# Dropout

- Turning off neurons w/ given probability (e.g. 0.5)
- Every iterations, new network architectures emerge



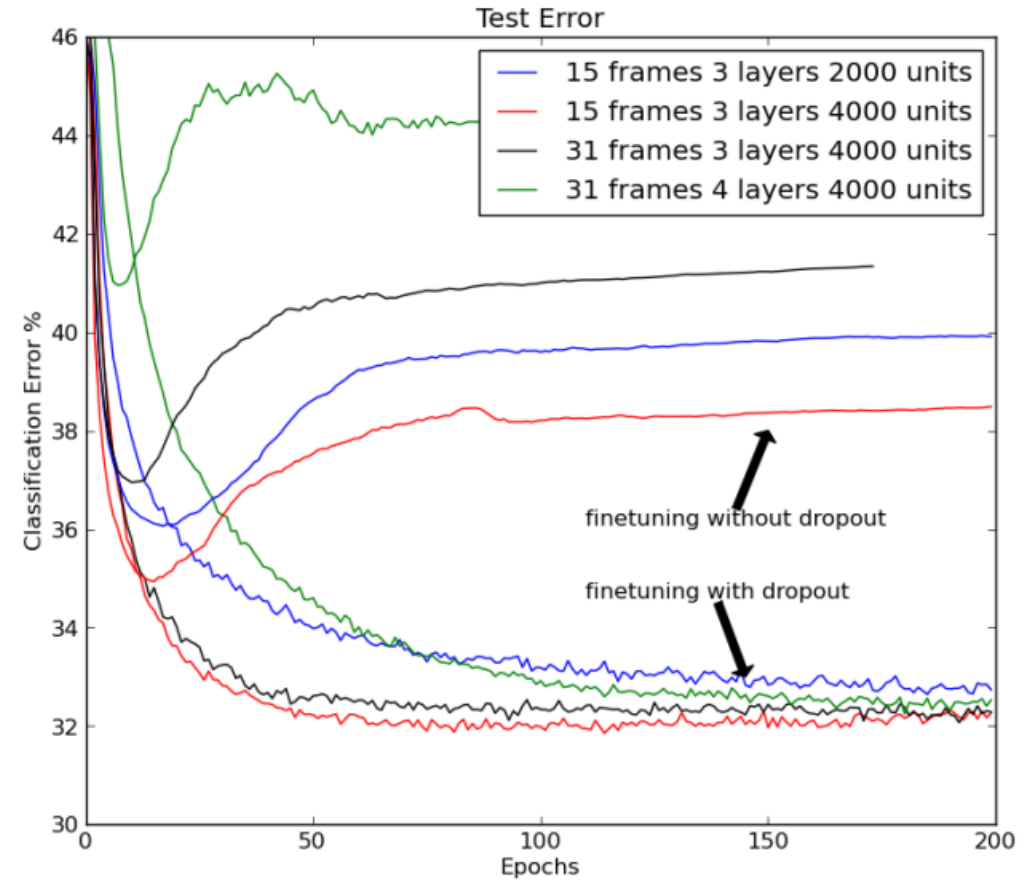
(a) Standard Neural Net



(b) After applying dropout.

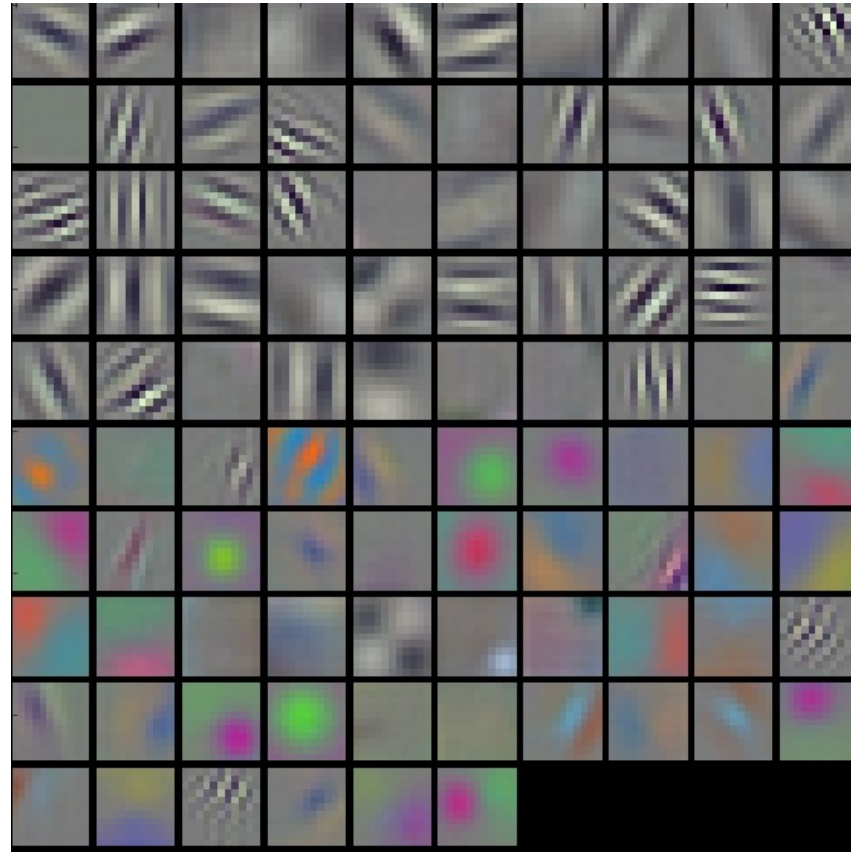
# Dropout

- A simple way to train deep neural networks for improving generalization performance
- Avoiding co-adaptations: a hidden unit cannot rely on other hidden units being present
- Model averaging

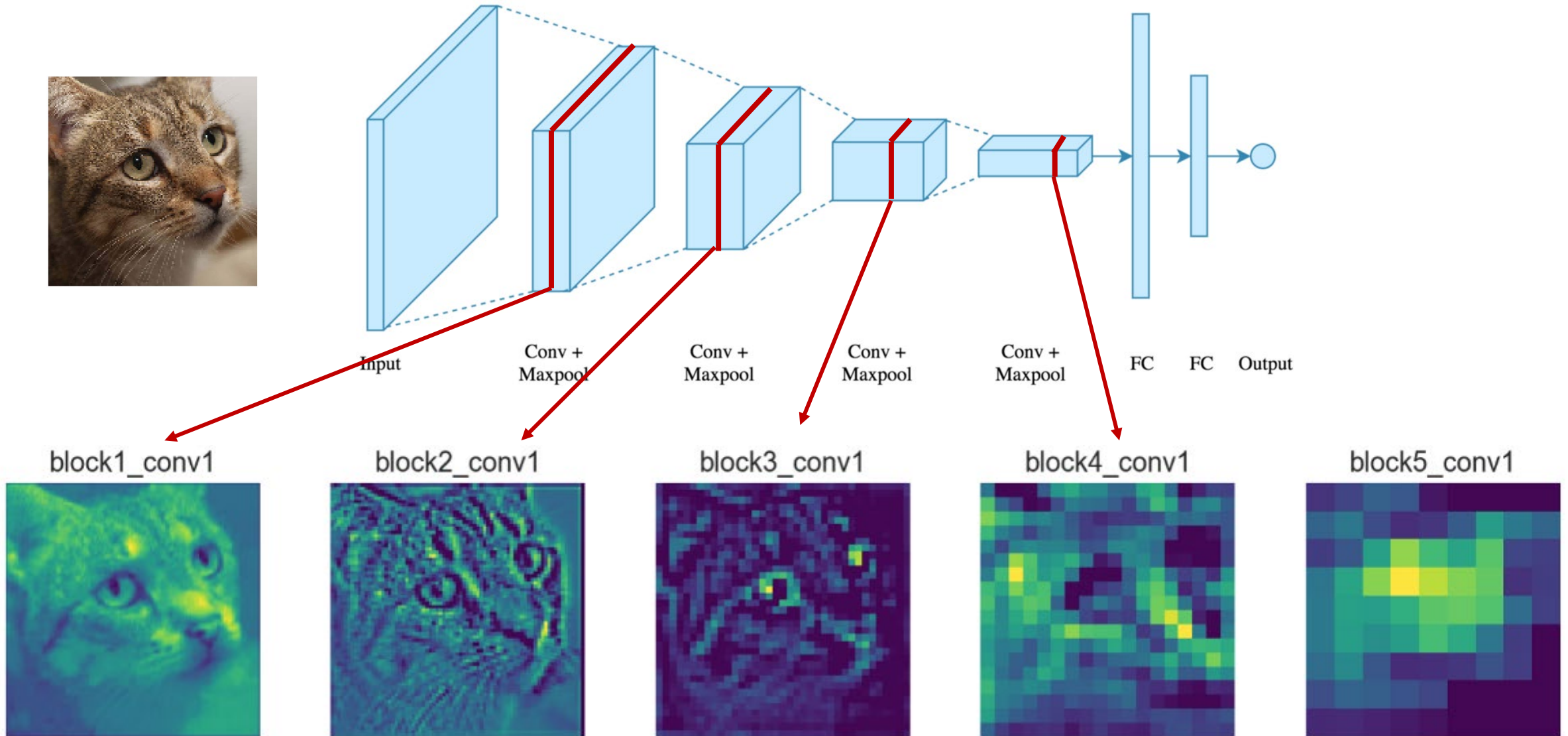


# Visualization of Learned Filter

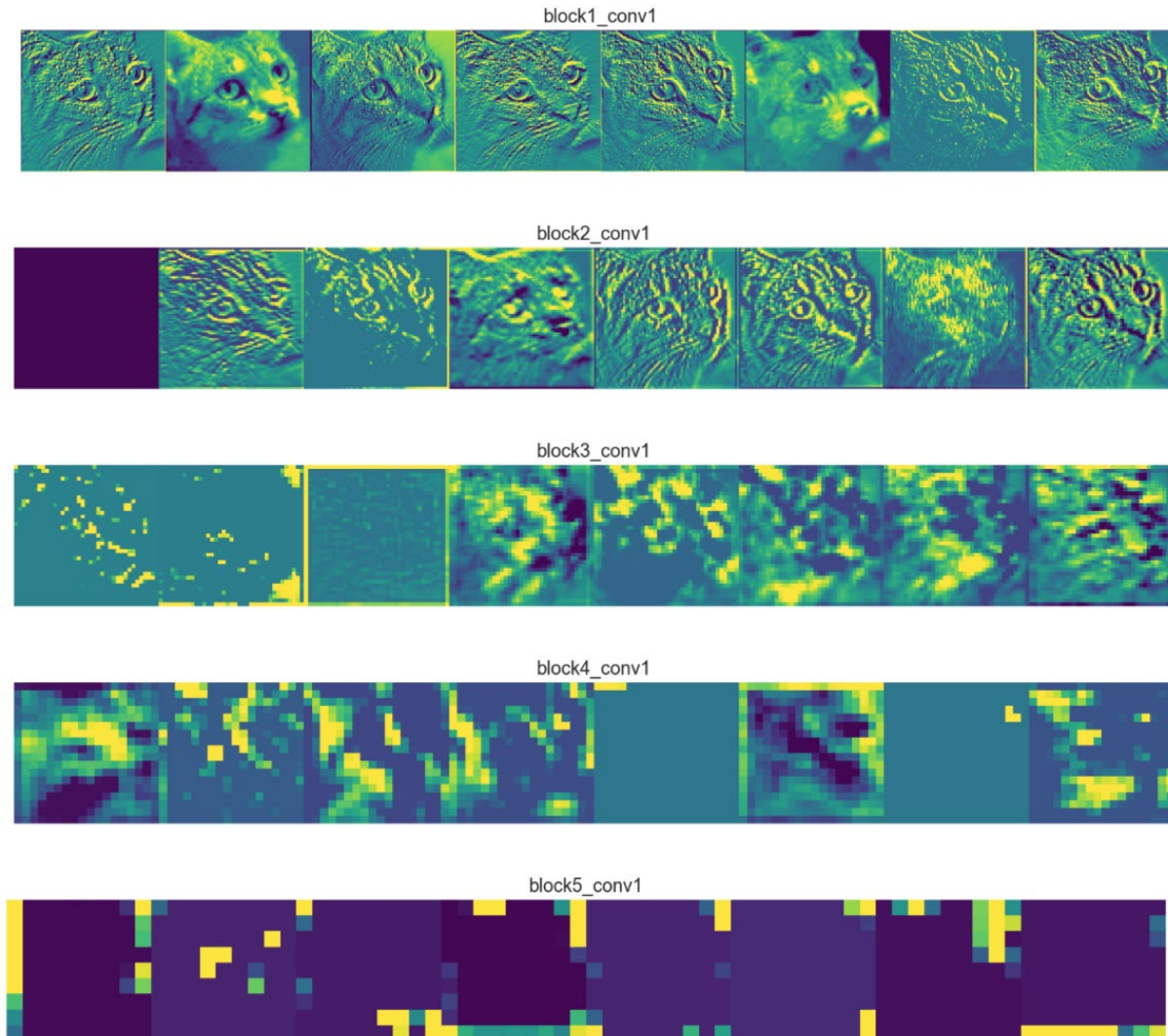
- First layer conv filters



# Visualization of Learned Feature Maps



# Visualization of Learned Feature Maps

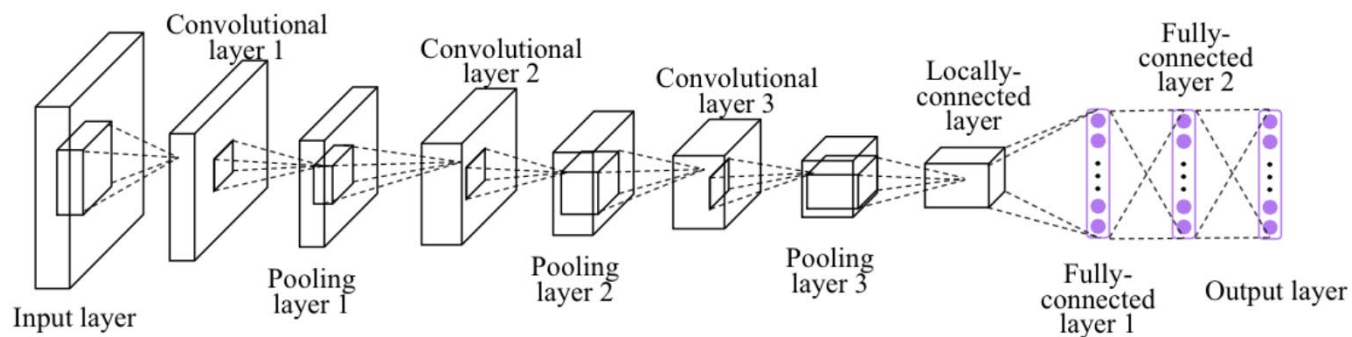




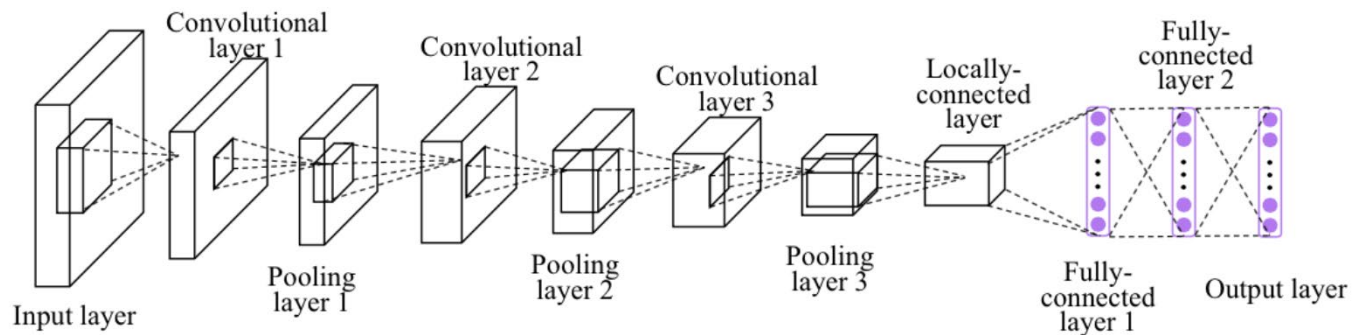
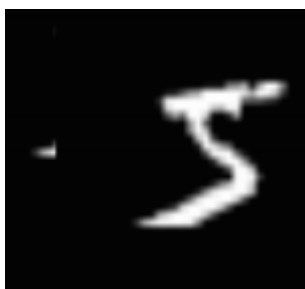
# Translation Equivariance vs Invariance

- Convolutional layers are translation equivariant
  - If you shift 1 pixel of input image, then the outputs will be 1 pixel shifted
- Pooling layers are translation invariant upto small shifting
  - If you shift 1 pixel of input image, then the outputs will be same
- We want network's prediction to be translation invariant
  - If you shift, you still want network to classify same as before

# Translation Equivariance vs Invariance

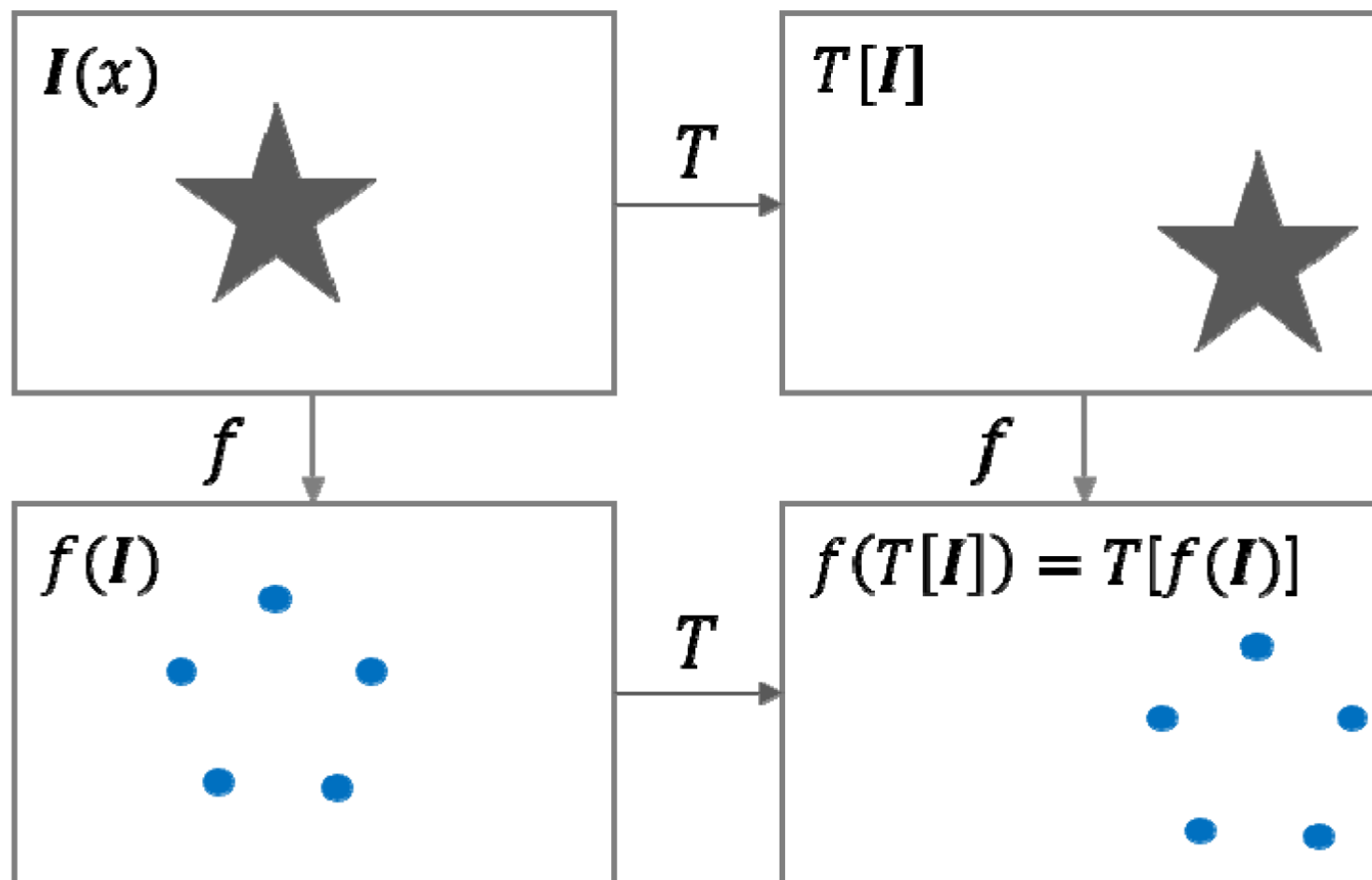


→  
→ '5'



→ '5'

# Translation Equivariance vs Invariance



# Cross correlation vs Convolution

- Cross-correlation: sliding a kernel across an image
- Convolution: sliding a flipped kernel across an image
  
- Most of deep learning libraries are actually doing 'cross-correlation'
  - But, it doesn't change the results since the same weight values would be learned in flipped manner

# CNNExplainer

- [CNN Explainer \(poloclub.github.io\)](https://poloclub.github.io/cnn-explainer/)

# Convolutions Behind the Scene

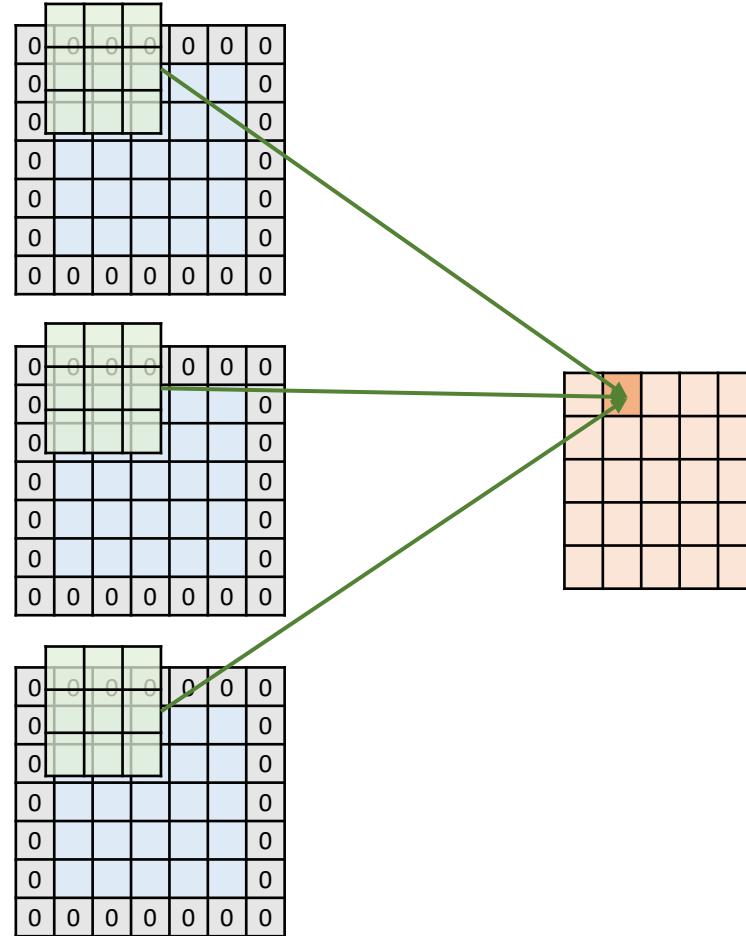
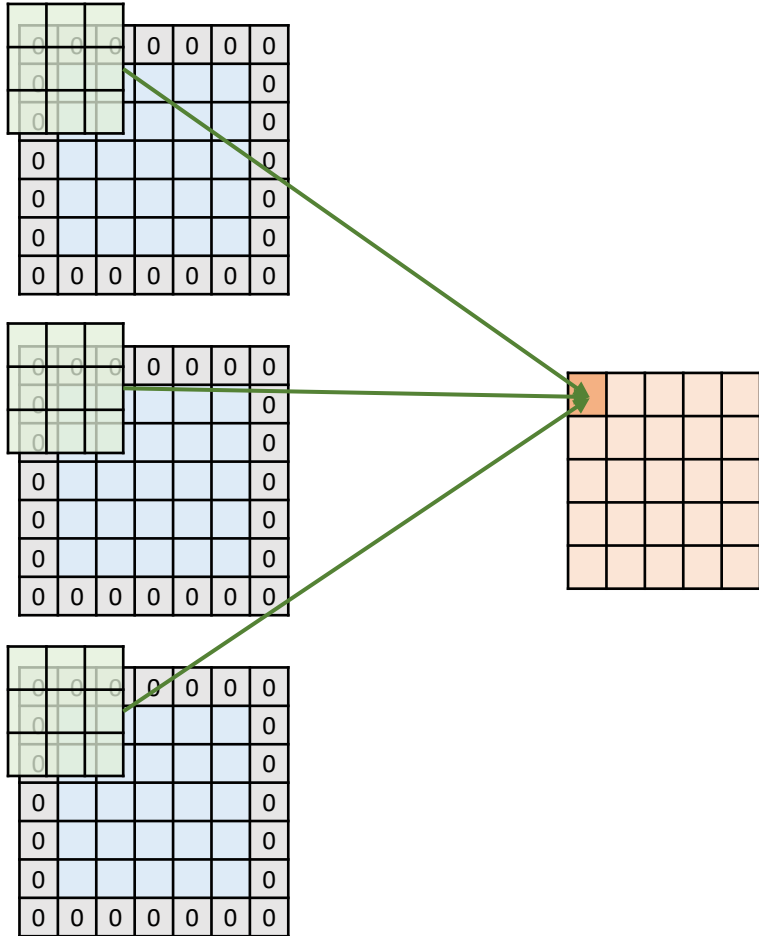
# Sequential Implementation (Conv2D)

- Sliding filters via 'for loops'
  - Slow, and we never do this

```
for row in range(x.shape[0] - 1):  
    for col in range(x.shape[1] - 1):  
        window = x[row: row + kernel_shape, col: col + kernel_shape]  
        result[row, col] = np.sum(np.multiply(kernel, window))
```

# Parallel Implementation

- Convolutions are parallelizable

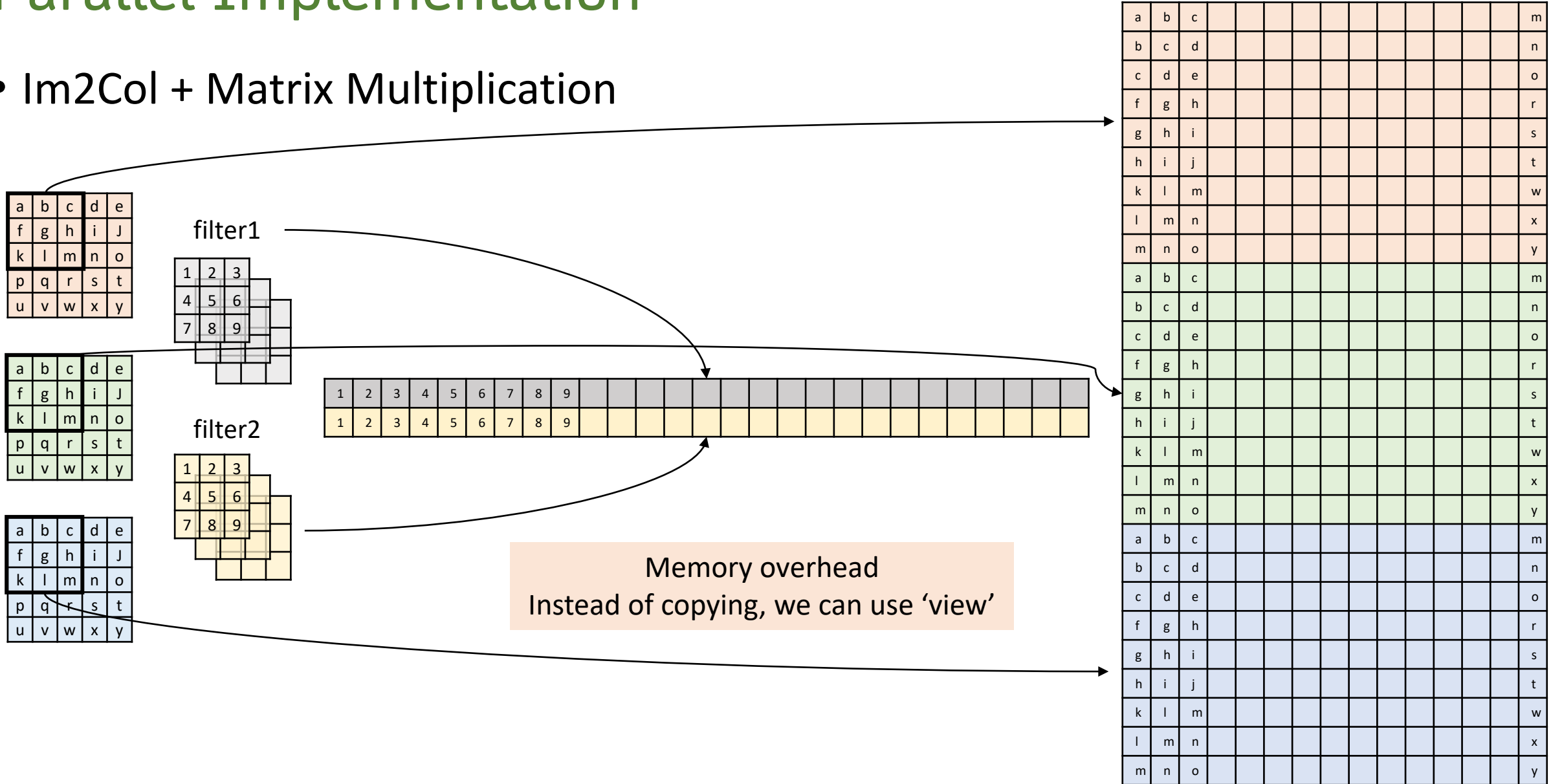






# Parallel Implementation

- Im2Col + Matrix Multiplication



# Further Optimizations

- GEMM optimization
- Winograd fast convolution
  - [\[1509.09308\] Fast Algorithms for Convolutional Neural Networks \(arxiv.org\)](#)
- FFT
  - For larger kernels

# Backprop in Convolutional Layers

# Transposed Convolution

- Opposite of normal convolution
- Can be used to up-sampling

1	0	1		
1	1	0	1	
0	1	2	1	
	2	1	2	

1	0	1		
1	1	0		
0	1	2		

# Transposed Convolution

- Opposite of normal convolution
- Can be used to up-sampling

1	0	1
<u>1</u>	<u>1</u>	<u>0</u>
<u>0</u>	<u>1</u>	<u>2</u>
2	1	2

1	0	1		
1	1	0		
0	1	2		

+

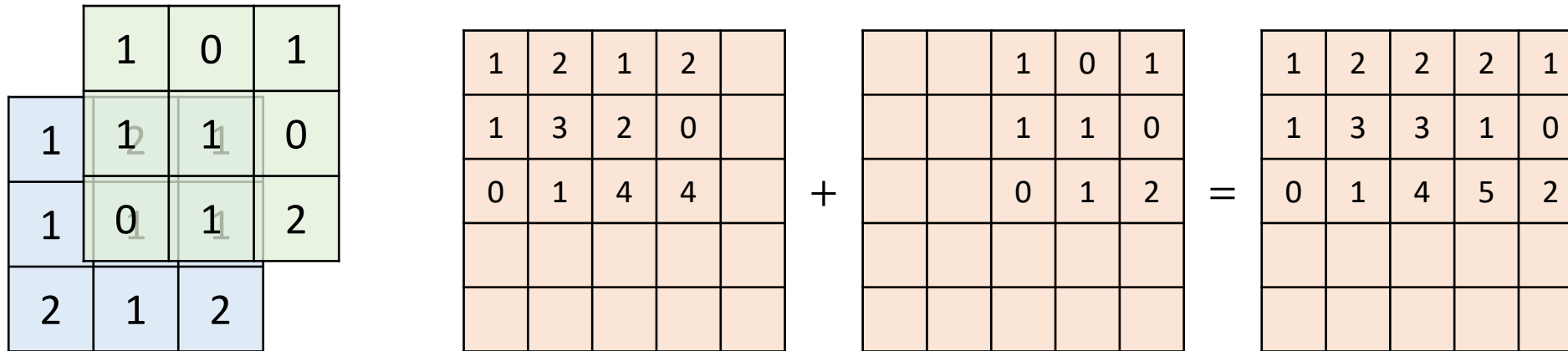
	2	0	2	
	2	2	0	
	0	2	4	

=

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

# Transposed Convolution

- Opposite of normal convolution
- Can be used to up-sampling



# Transposed Convolution

- Opposite of normal convolution
- Can be used to up-sampling

1	0	1	1
1	1	0	1
0	1	2	2

1	2	2	2	1
1	3	3	1	0
0	1	4	5	2

+

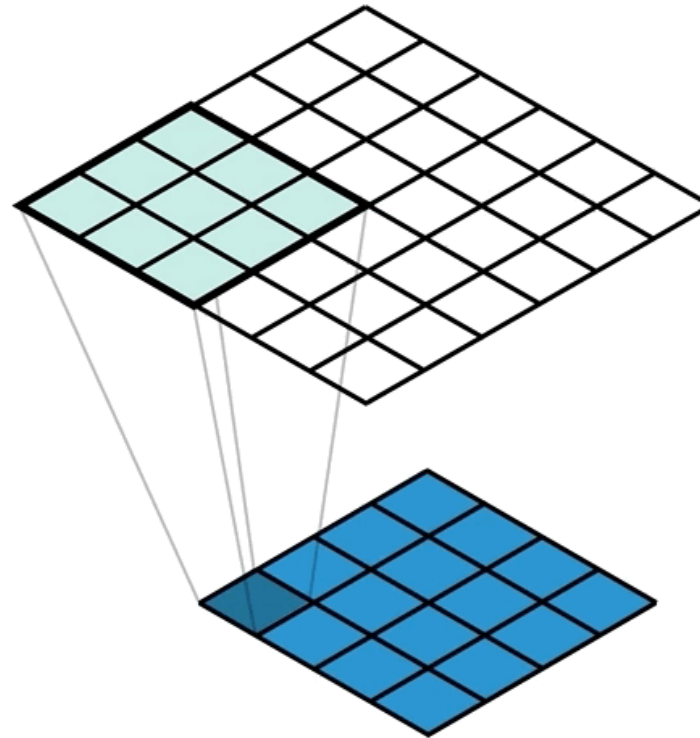
1	0	1		
1	1	0		
0	1	2		

=

1	2	2	2	1
2	3	4	1	0
1	2	4	5	2
0	1	2		

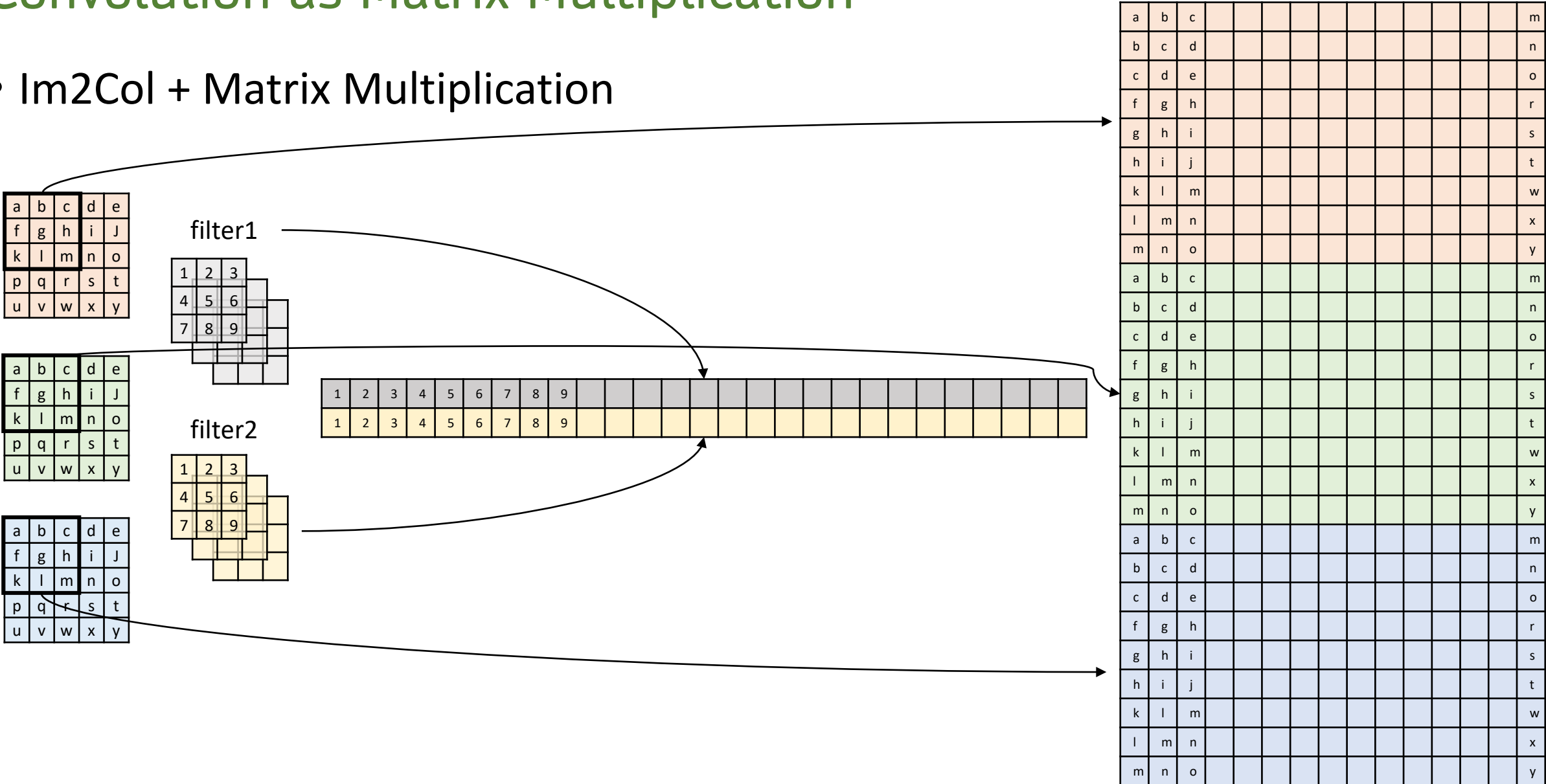


# Transposed Convolution



# Convolution as Matrix Multiplication

- Im2Col + Matrix Multiplication



# Convolution as Matrix Multiplication

$$X \in \mathbb{R}^{n \times d}, Y \in \mathbb{R}^{m \times d}, z \in \mathbb{R}, W \in \mathbb{R}^{m \times n}$$

$$Y = WX$$

$$\mathbb{R}^{n \times d} \quad \mathbb{R}^{n \times m} \quad \mathbb{R}^{m \times d}$$

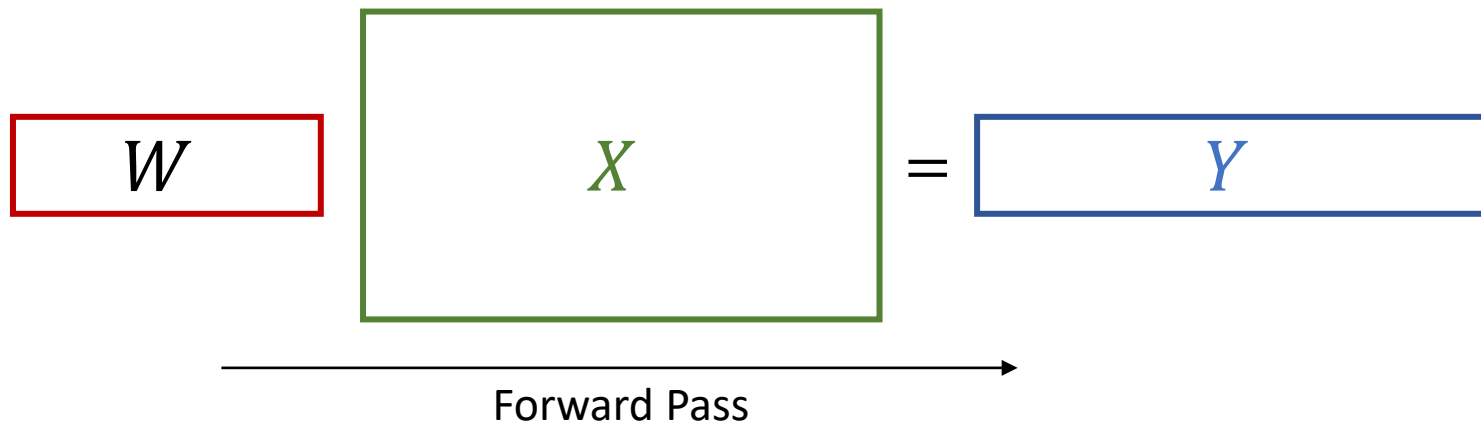
$$\frac{dz}{dX} = \frac{dz}{dY} \frac{dY}{dX} = W^T \frac{dz}{dY}$$

$$\mathbb{R}^{m \times d} \quad \mathbb{R}^{m \times d \times n \times d}$$

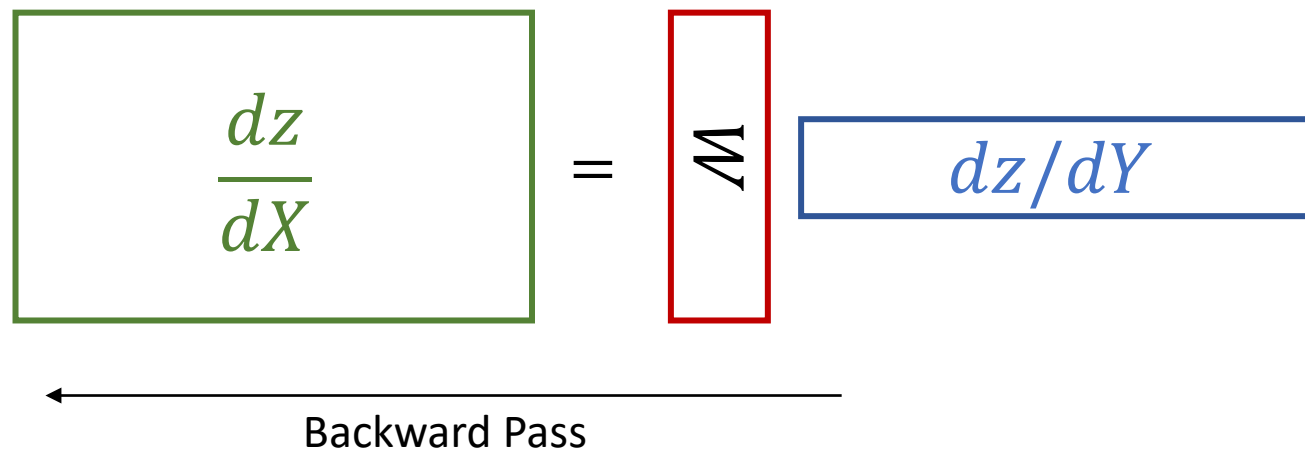
# Forward and Backward

$$X \in \mathbb{R}^{n \times d}, Y \in \mathbb{R}^{m \times d}, z \in \mathbb{R}, W \in \mathbb{R}^{m \times n}$$

$$Y = WX$$

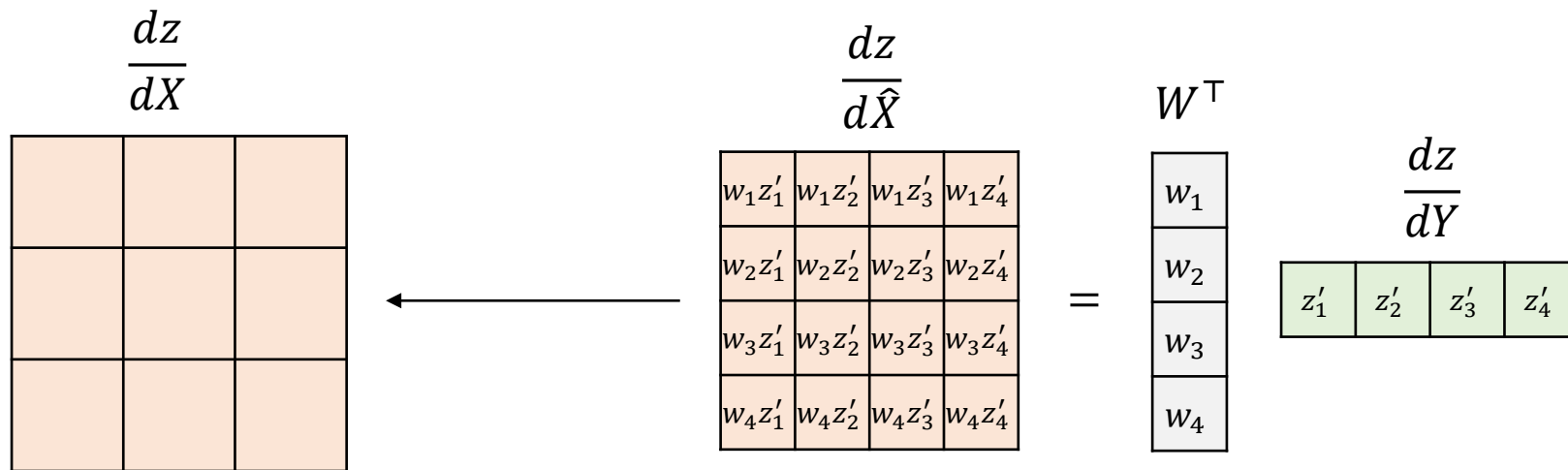


$$\frac{dz}{dX} = W^T \frac{dz}{dY}$$



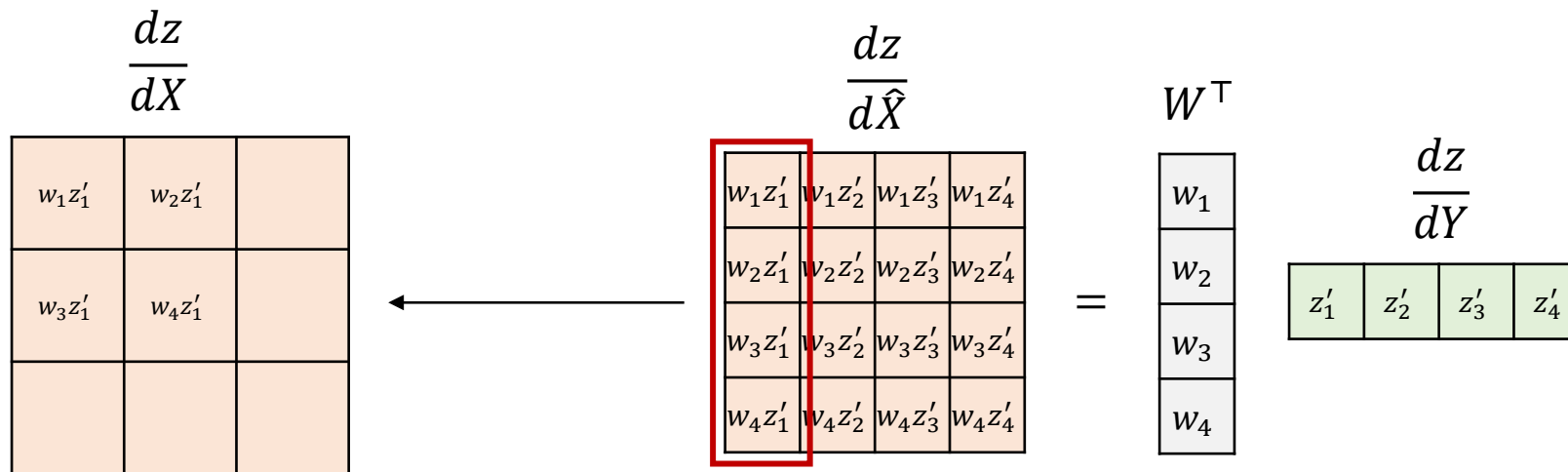
# Backward Pass

- Matrix multiplication view



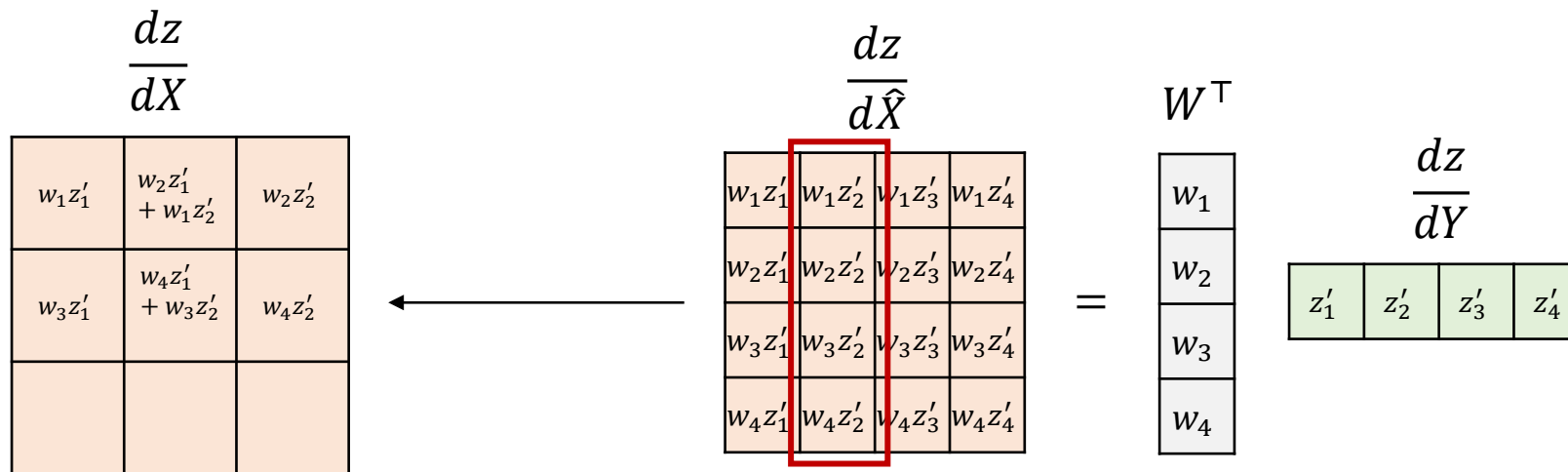
# Backward Pass

- Matrix multiplication view



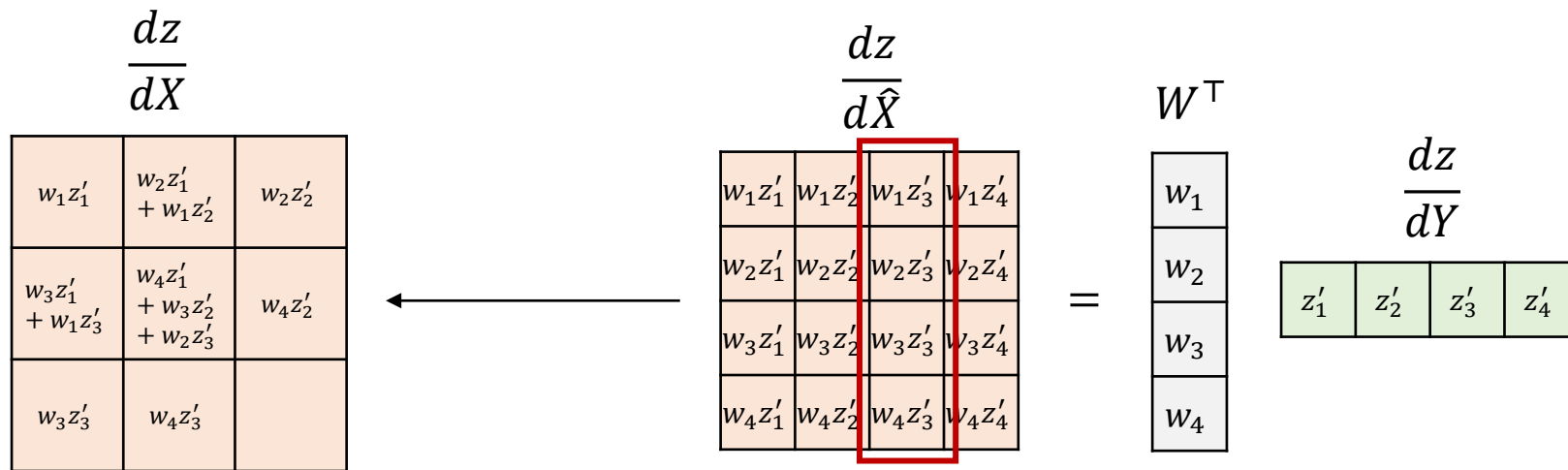
# Backward Pass

- Matrix multiplication view



# Backward Pass

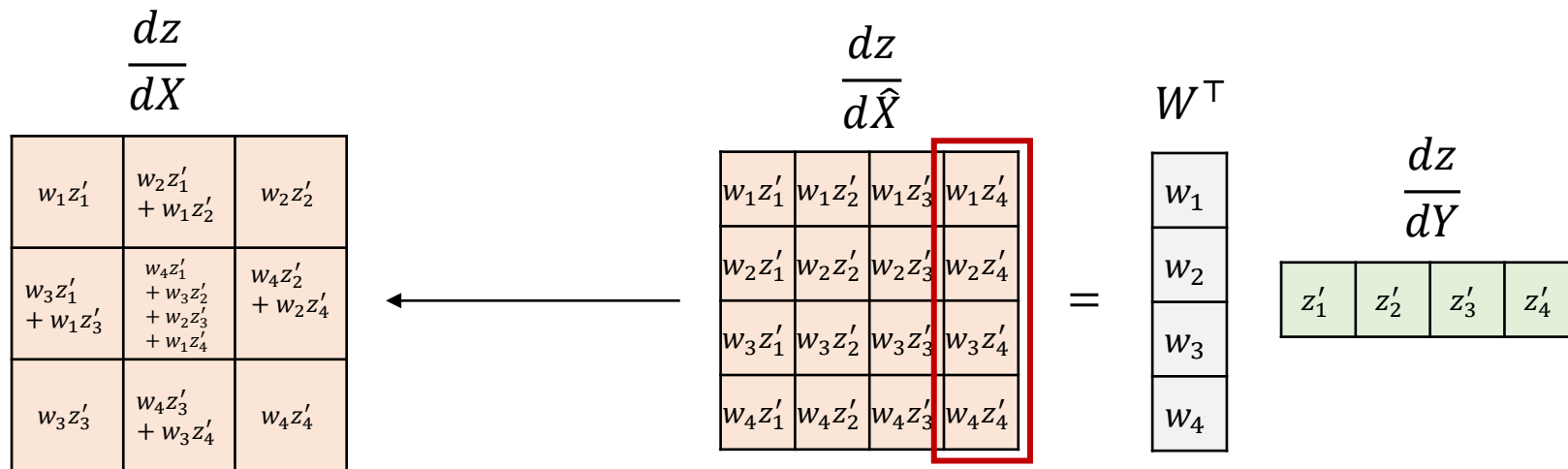
- Matrix multiplication view





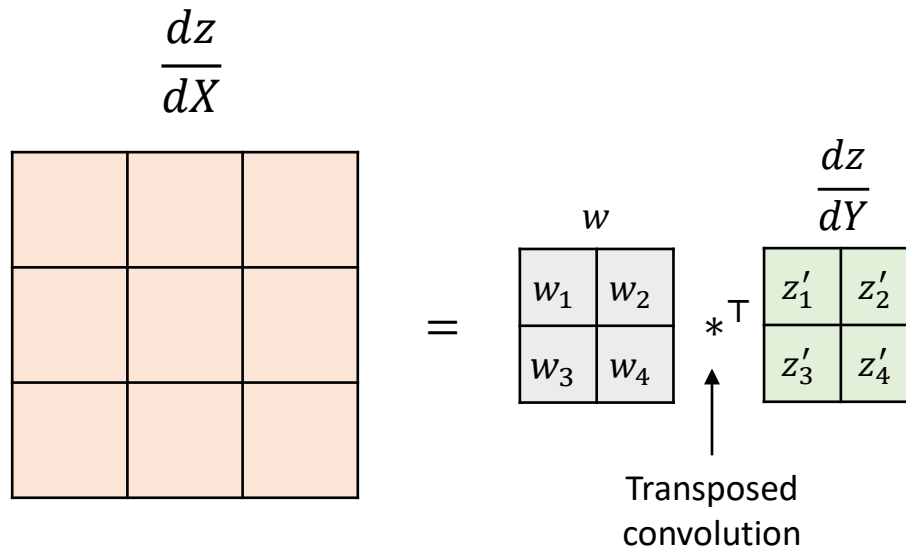
# Backward Pass

- Matrix multiplication view



# Backward Pass

- Matrix multiplication view



# Backward Pass

- Matrix multiplication view

$$\frac{dz}{dX} = \begin{matrix} & \frac{dz}{dY} \\ \begin{matrix} w_1 z'_1 & w_1 z'_2 \\ w_1 z'_3 & w_1 z'_4 \\ & & & \end{matrix} & \begin{matrix} w \\ \begin{matrix} w_1 & w_2 \\ w_3 & w_4 \end{matrix} \end{matrix} *^T \begin{matrix} \begin{matrix} z'_1 & z'_2 \\ z'_3 & z'_4 \end{matrix} \\ \frac{dz}{dY} \end{matrix}$$

↑  
Transposed convolution

# Backward Pass

- Matrix multiplication view

$$\frac{dz}{dX} = \begin{bmatrix} w_1 z'_1 & w_1 z'_2 + w_2 z'_1 & w_2 z'_2 \\ w_1 z'_3 & w_1 z'_4 + w_2 z'_3 & w_2 z'_4 \\ & & \end{bmatrix} = \begin{bmatrix} w_1 & w_2 \\ w_3 & w_4 \end{bmatrix} *^T \begin{bmatrix} z'_1 & z'_2 \\ z'_3 & z'_4 \end{bmatrix} \frac{dz}{dY}$$

↑  
Transposed convolution

# Backward Pass

- Matrix multiplication view

$$\frac{dz}{dX}$$

$w_1 z'_1$	$w_1 z'_2 + w_2 z'_1$	$w_2 z'_2$
$w_1 z'_3 + w_3 z'_1$	$w_1 z'_4 + w_2 z'_3 + w_3 z'_2$	$w_2 z'_4$
$w_3 z'_3$	$w_3 z'_4$	

=

$w$		$\frac{dz}{dY}$	
$w_1$	$w_2$	$z'_1$	$z'_2$
$w_3$	$w_4$	$z'_3$	$z'_4$

$\uparrow$   
\*<sup>T</sup>  
Transposed convolution

# Backward Pass

- Matrix multiplication view

$$\begin{array}{|c|c|c|} \hline w_1 z'_1 & w_2 z'_1 + w_1 z'_2 & w_2 z'_2 \\ \hline w_3 z'_1 + w_1 z'_3 & w_4 z'_1 + w_3 z'_2 + w_2 z'_3 + w_1 z'_4 & w_4 z'_2 + w_2 z'_4 \\ \hline w_3 z'_3 & w_4 z'_3 + w_3 z'_4 & w_4 z'_4 \\ \hline \end{array} = \frac{dz}{dX}$$

$$\begin{array}{|c|c|c|} \hline w_1 z'_1 & w_1 z'_2 + w_2 z'_1 & w_2 z'_2 \\ \hline w_1 z'_3 + w_3 z'_1 & w_1 z'_4 + w_2 z'_3 + w_3 z'_2 + w_4 z'_1 & w_2 z'_4 + w_4 z'_2 \\ \hline w_3 z'_3 & w_3 z'_4 + w_4 z'_3 & w_4 z'_4 \\ \hline \end{array}$$

$$= \begin{array}{|c|c|} \hline w_1 & w_2 \\ \hline w_3 & w_4 \\ \hline \end{array} *^T \begin{array}{|c|c|} \hline z'_1 & z'_2 \\ \hline z'_3 & z'_4 \\ \hline \end{array} \frac{dz}{dY}$$

↑  
Transposed convolution