# Deep Learning (1470)

**Randall Balestriero**

# Multilayer Perceptrons
## And why we need better



### 3-Layer MLP with ReLU Activations

$$ReLU(z) = max(0, z)$$

Input Layer
(764 units)
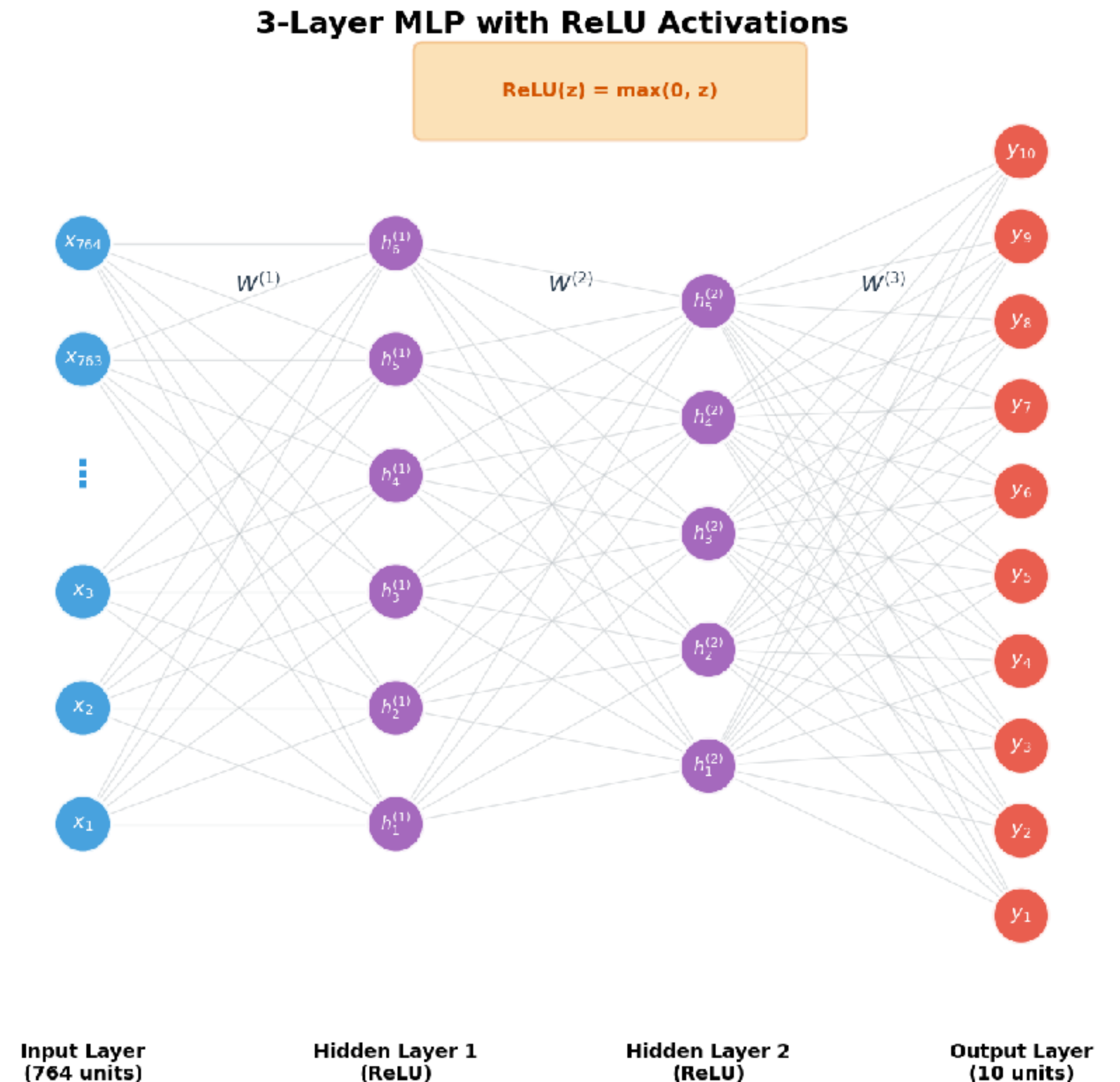
Hidden Layer 1
(ReLU)

Hidden Layer 2
(ReLU)

Output Layer
(10 units)

# Multilayer Perceptrons
## And why we need better

Extracting local features at each layer is enough!



### 3-Layer MLP with ReLU Activations

$ReLU(z) = max(0, z)$

$x_{764}$  $x_{763}$  $x_3$  $x_2$  $x_1$

$W^{(1)}$  $W^{(2)}$  $W^{(3)}$

$h_6^{(1)}$  $h_5^{(1)}$  $h_4^{(1)}$  $h_3^{(1)}$  $h_2^{(1)}$  $h_1^{(1)}$

$h_6^{(2)}$  $h_4^{(2)}$  $h_3^{(2)}$  $h_2^{(2)}$  $h_1^{(2)}$

$y_{10}$  $y_9$  $y_8$  $y_7$  $y_6$  $y_5$  $y_4$  $y_3$  $y_2$  $y_1$

**Input Layer** (764 units)  **Hidden Layer 1** (ReLU)  **Hidden Layer 2** (ReLU)  **Output Layer** (10 units)
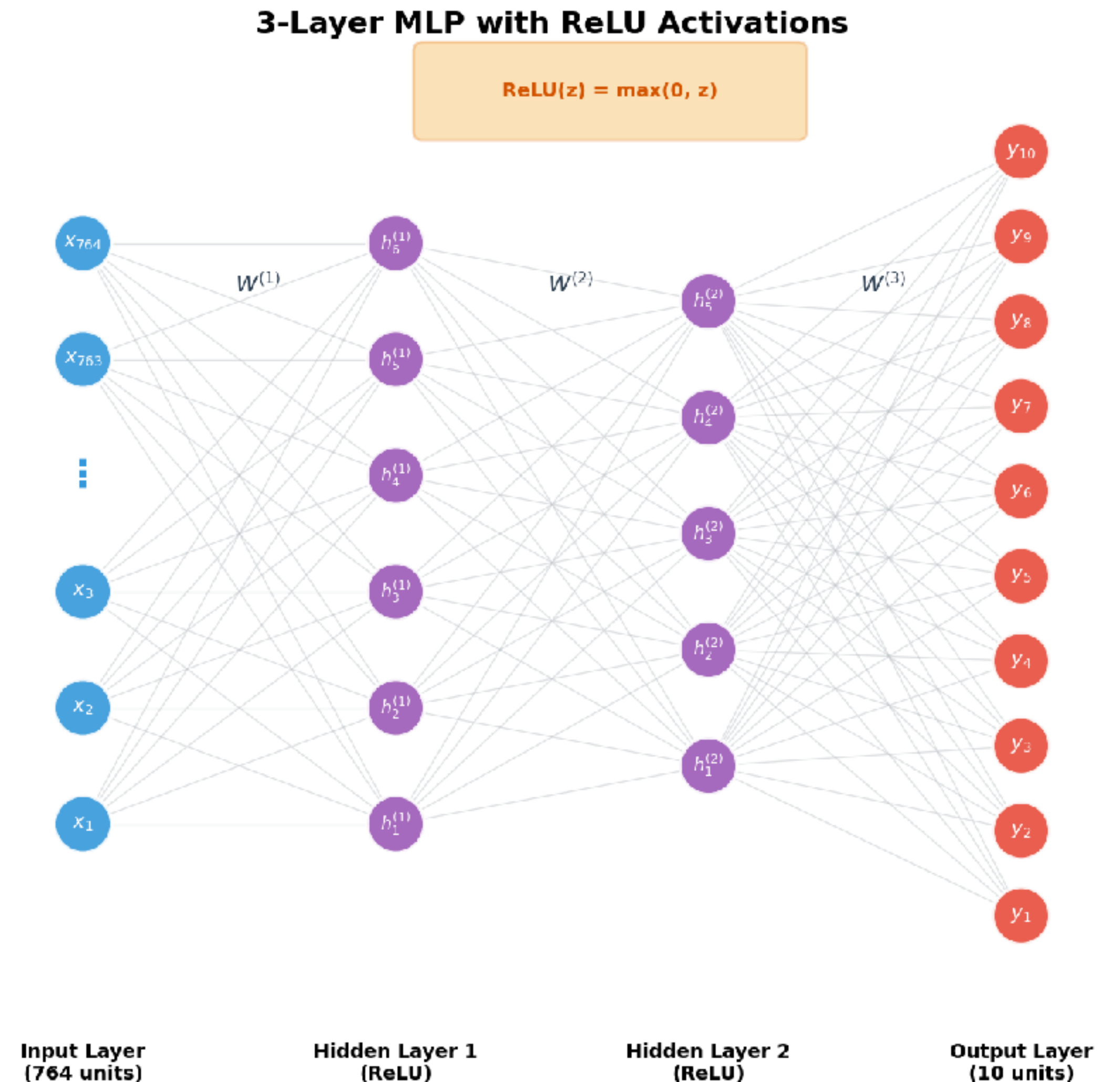
# Multilayer Perceptrons
## And why we need better

Extracting local features
at each layer is enough!



Not only local features… but
translation invariant features!

### 3-Layer MLP with ReLU Activations

$$ReLU(z) = max(0, z)$$



| Input Layer (764 units) | Hidden Layer 1 (ReLU) | Hidden Layer 2 (ReLU) | Output Layer (10 units) |

# The Main Building Block: Convolution

- A convolution is a linear operator

- Convolution is an operation that takes two inputs

(1) An image (2D – B/W)

(2) A filter (also called a kernel)



| 1 | 1 | 1 |
|---|---|---|
| 0 | 0 | 0 |
| −1 | −1 | −1 |

2D array of numbers; could be any values

# What Convolution Does (Visually)

image

| 2 | 0 | 1 | 3 |
|---|---|---|---|
| 7 | 1 | 1 | 0 |
| 0 | 2 | 5 | 0 |
| 0 | 5 | 1 | 4 |

$\otimes$

filter/kernel

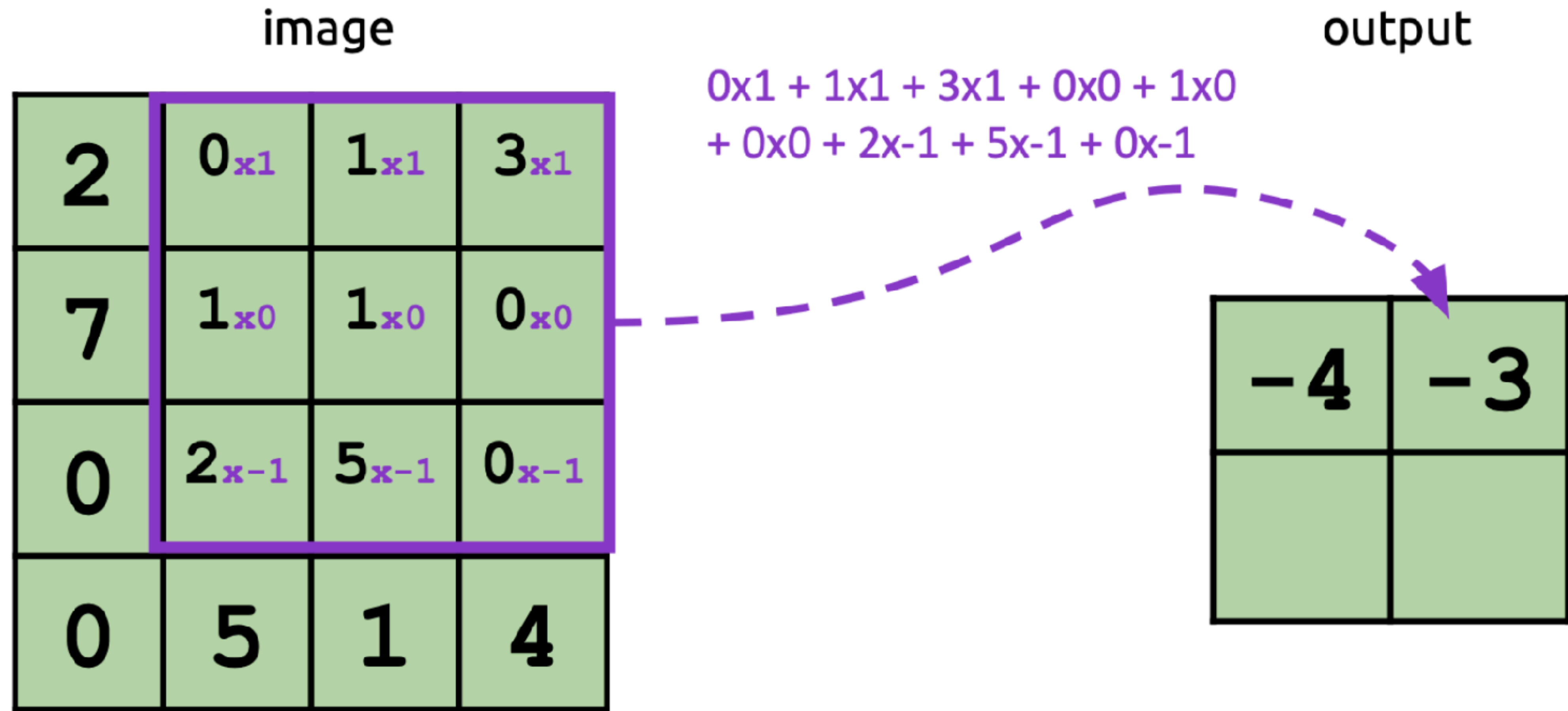| 1 | 1 | 1 |
|---|---|---|
| 0 | 0 | 0 |
| -1 | -1 | -1 |

(We use this symbol for convolution)
(The verb form is "convolve")

# What Convolution Does (Visually)



image

| 2 | 0 | 1 | 3 |
|---|---|---|---|
| 7 | 1 | 1 | 0 |
| 0 | 2 | 5 | 0 |
| 0 | 5 | 1 | 4 |

⊗

filter/kernel

| 1 | 1 | 1 |
|---|---|---|
| 0 | 0 | 0 |
| -1 | -1 | -1 |

=

output

| -4 | -3 |
|----|----|
| 3 | -8 |

# What Convolution Does (Visually)

# What Convolution Does (Visually)



image

$0 \times 1 + 1 \times 1 + 3 \times 1 + 0 \times 0 + 1 \times 0 + 0 \times 0 + 2 \times -1 + 5 \times -1 + 0 \times -1$

output

# Example



CONVOLUTION: Slide, Multiply, Sum

INPUT (8×8)

KERNEL (3×3)
Edge Detector

| +1 | +0 | -1 |
|----|----|----|
| +1 | +0 | -1 |
| +1 | +0 | -1 |

OUTPUT (6×6)

-1

Position [0,0] = -1

# Example



CONVOLUTION: Slide, Multiply, Sum

INPUT (8×8)

KERNEL (3×3)
Edge Detector

| +1 | +0 | -1 |
|----|----|----|
| +1 | +0 | -1 |
| +1 | +0 | -1 |

OUTPUT (6×6)

-1

Position [0,0] = -1

# This was all a lie!

- All those example are actually ``cross-correlation", not ``convolution"

- It is what is implemented in all those frameworks and called ``convolution"

- True convolution is almost that, but you need to flip the filter before applying
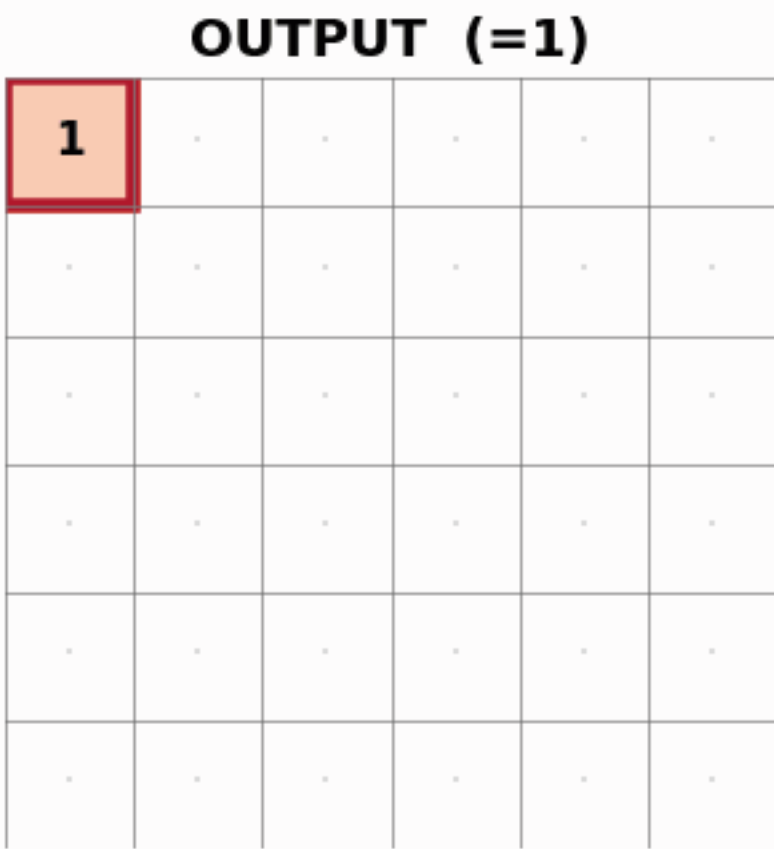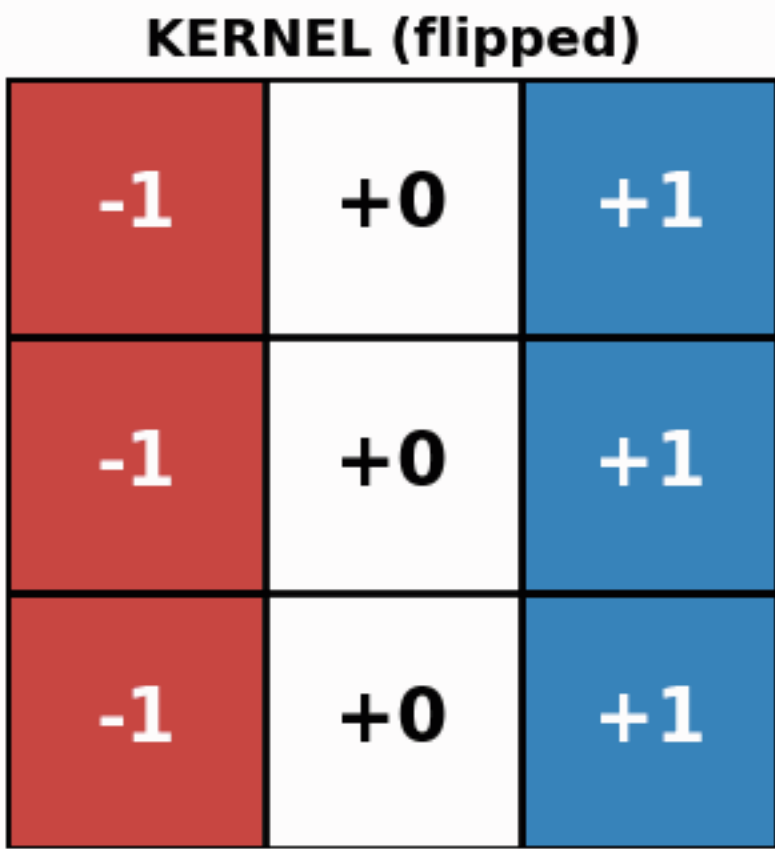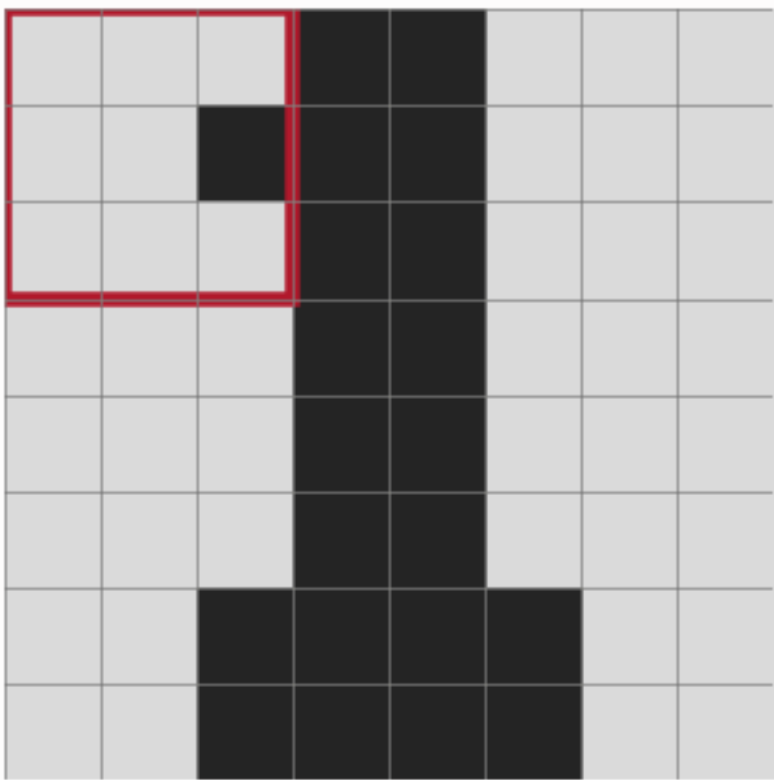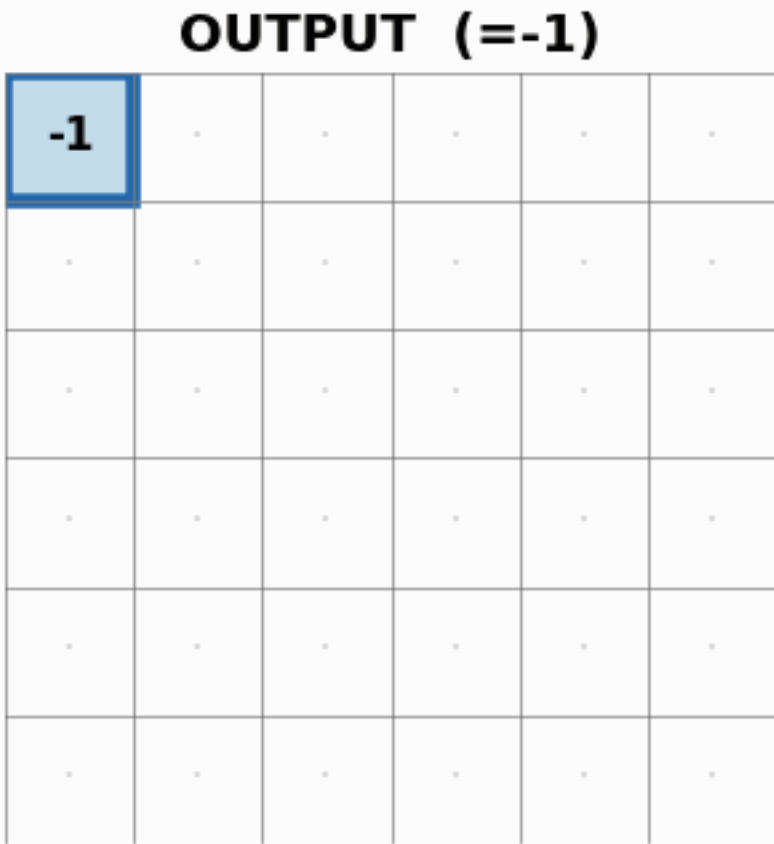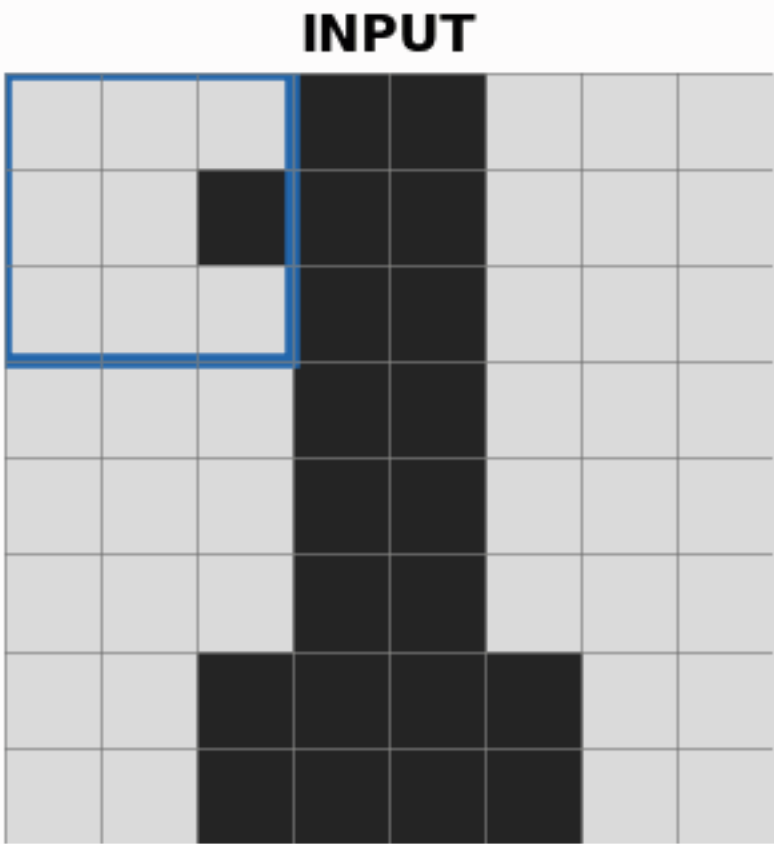
# This was all a lie!

- All those example are actually ``cross-correlation", not ``convolution"

- It is what is implemented in all those frameworks and called ``convolution"

- True convolution is almost that, but you need to flip the filter before applying

Exercise: think why when starting from random init, it doesn't matter!

# Example



CROSS-CORRELATION vs TRUE CONVOLUTION

TRUE CONVOLUTION (kernel flipped 180°) CROSS-CORRELATION (what PyTorch/TF call 'conv')

INPUT

KERNEL

| +1 | +0 | -1 |
| +1 | +0 | -1 |
| +1 | +0 | -1 |

OUTPUT (=-1)

-1

KERNEL (flipped)

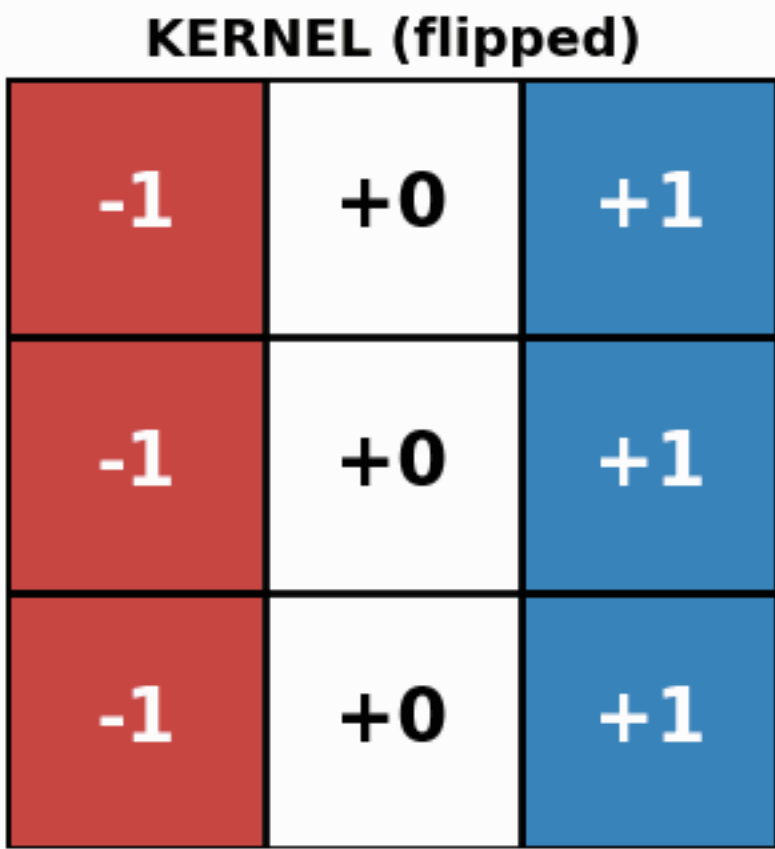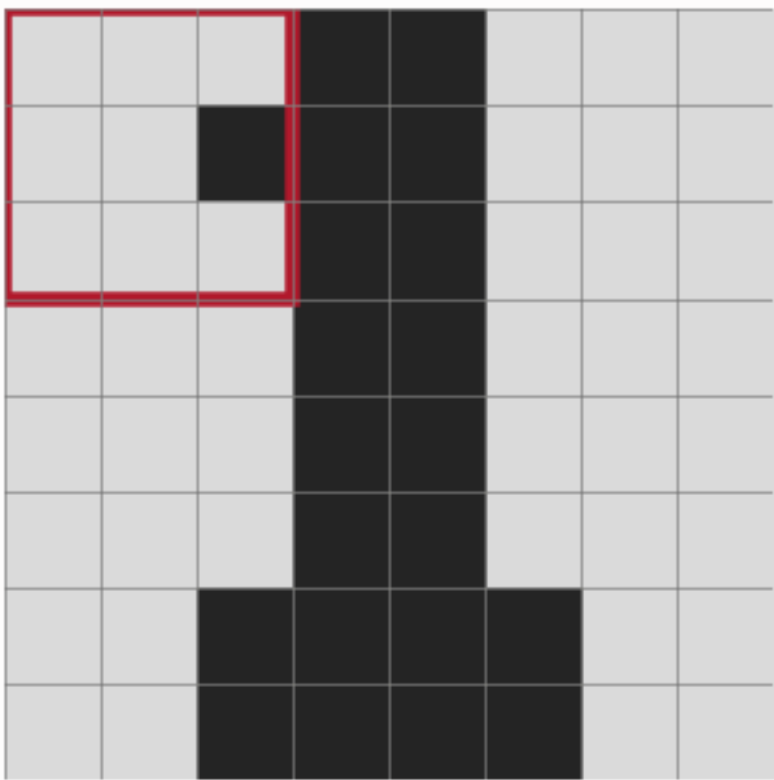| -1 | +0 | +1 |
| -1 | +0 | +1 |
| -1 | +0 | +1 |

OUTPUT (=1)

1

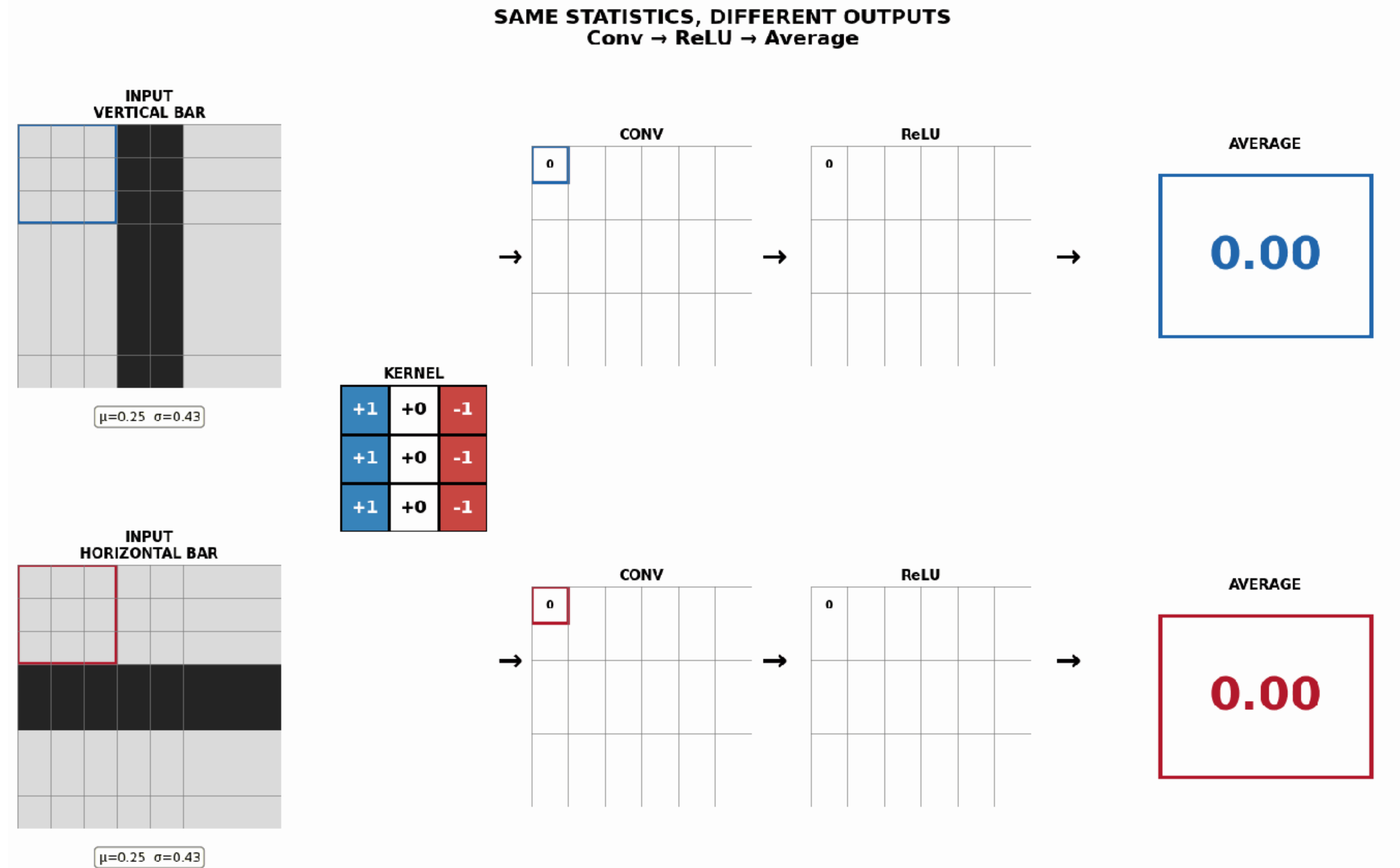Notice: Outputs are NEGATIVES of each other (edges swap polarity)

# Example



CROSS-CORRELATION vs TRUE CONVOLUTION

TRUE CONVOLUTION (kernel flipped 180°) CROSS-CORRELATION (what PyTorch/TF call 'conv')

**INPUT**

**KERNEL**

| +1 | +0 | -1 |
|----|----|----|
| +1 | +0 | -1 |
| +1 | +0 | -1 |

**OUTPUT (=-1)**

-1

**KERNEL (flipped)**

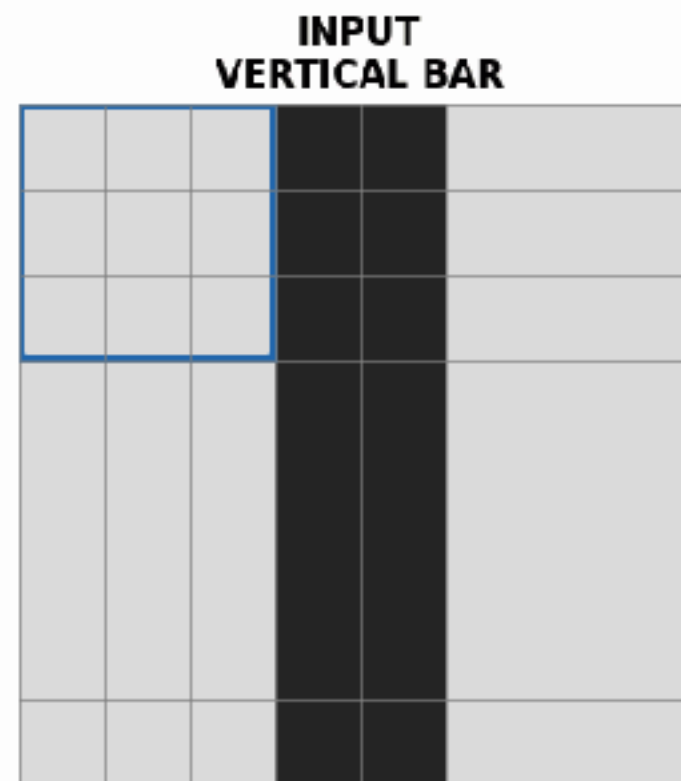| -1 | +0 | +1 |
|----|----|----|
| -1 | +0 | +1 |
| -1 | +0 | +1 |

**OUTPUT (=1)**

1

*Notice: Outputs are NEGATIVES of each other (edges swap polarity)*

# Your First Convolutional Network



SAME STATISTICS, DIFFERENT OUTPUTS
Conv → ReLU → Average

INPUT
VERTICAL BAR

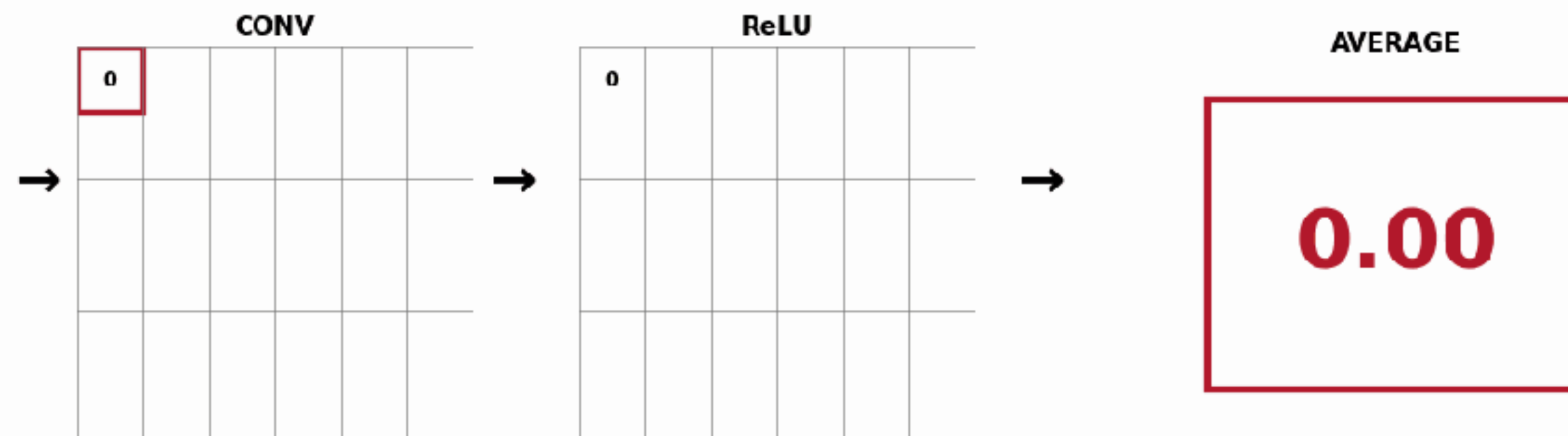μ=0.25 σ=0.43

KERNEL

| +1 | +0 | -1 |
|----|----|----|
| +1 | +0 | -1 |
| +1 | +0 | -1 |

CONV

0

ReLU

0

AVERAGE

0.00

INPUT
HORIZONTAL BAR

μ=0.25 σ=0.43

CONV

0

ReLU

0

AVERAGE

0.00

# Your First Convolutional Network



SAME STATISTICS, DIFFERENT OUTPUTS
Conv → ReLU → Average

INPUT
VERTICAL BAR

μ=0.25  σ=0.43

INPUT
HORIZONTAL BAR

μ=0.25  σ=0.43

KERNEL

| +1 | +0 | -1 |
| +1 | +0 | -1 |
| +1 | +0 | -1 |

CONV

0

ReLU

0

AVERAGE

0.00

CONV

0

ReLU

0

AVERAGE

0.00

# Your First Convolutional Network



VERTICAL
μ=0.11

TRANSLATION INVARIANCE

Conv→ReLU

AVERAGE (invariant!)

0.420

HORIZONTAL
μ=0.11

Conv→ReLU

AVERAGE (invariant!)

0.120

# Your First Convolutional Network



VERTICAL
μ=0.11

TRANSLATION INVARIANCE

Conv→ReLU

AVERAGE (invariant!)

0.420

HORIZONTAL
μ=0.11

Conv→ReLU

AVERAGE (invariant!)

0.120

# Multiple Filters

# Multiple Multi-Channel Filters



Input

Kernel
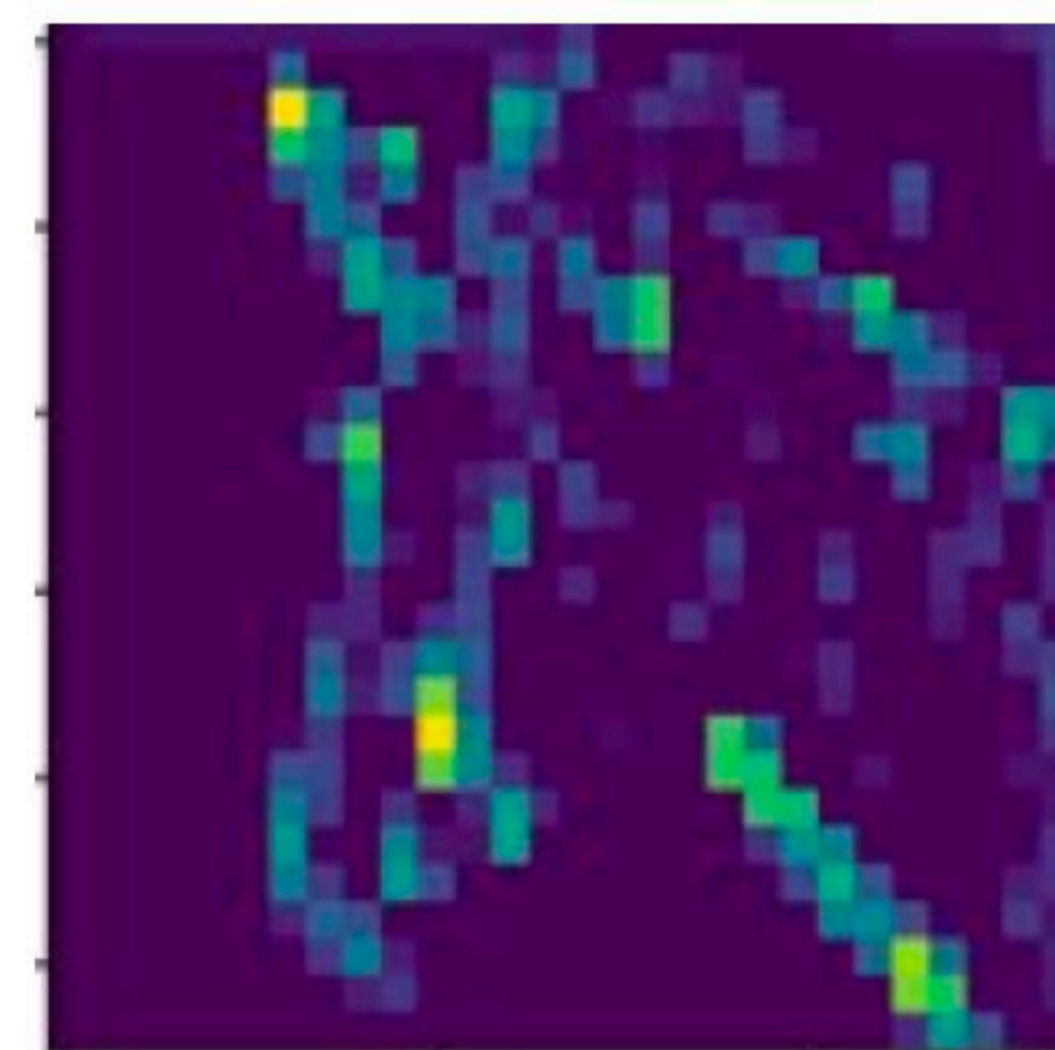
Output

# Multiple Multi-Channel Filters
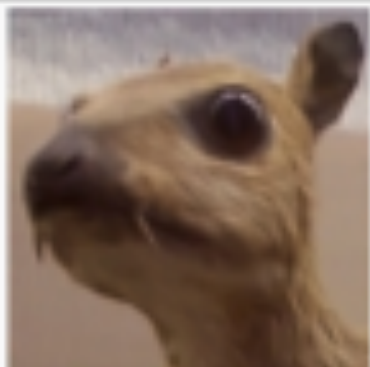
Input

Kernel

Output

# Multiple Multi-Channel Filters



Input image
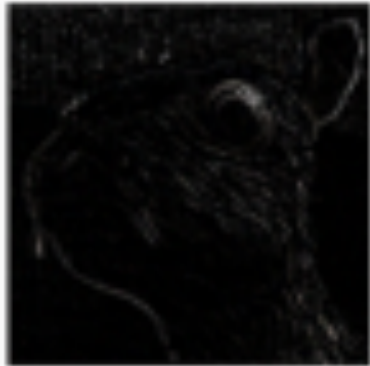
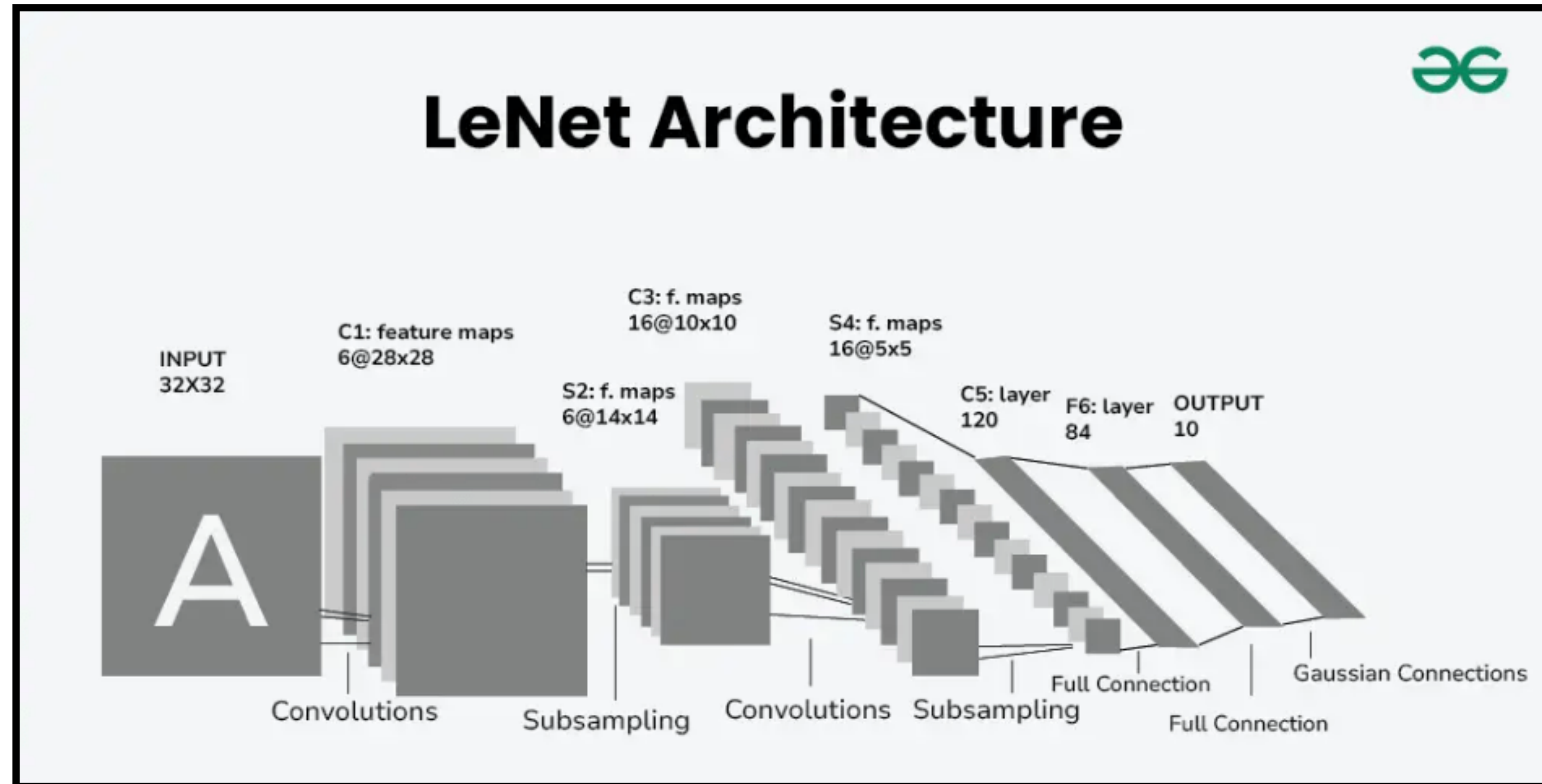Output of filter 1

Output of filter 2

# More Filter Examples



| Operation | Filter | Convolved Image |
|---|---|---|
| **Identity** | $\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$ | |
| **Edge detection** | $\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix}$ | |
| | $\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$ | |
| | $\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$ | |

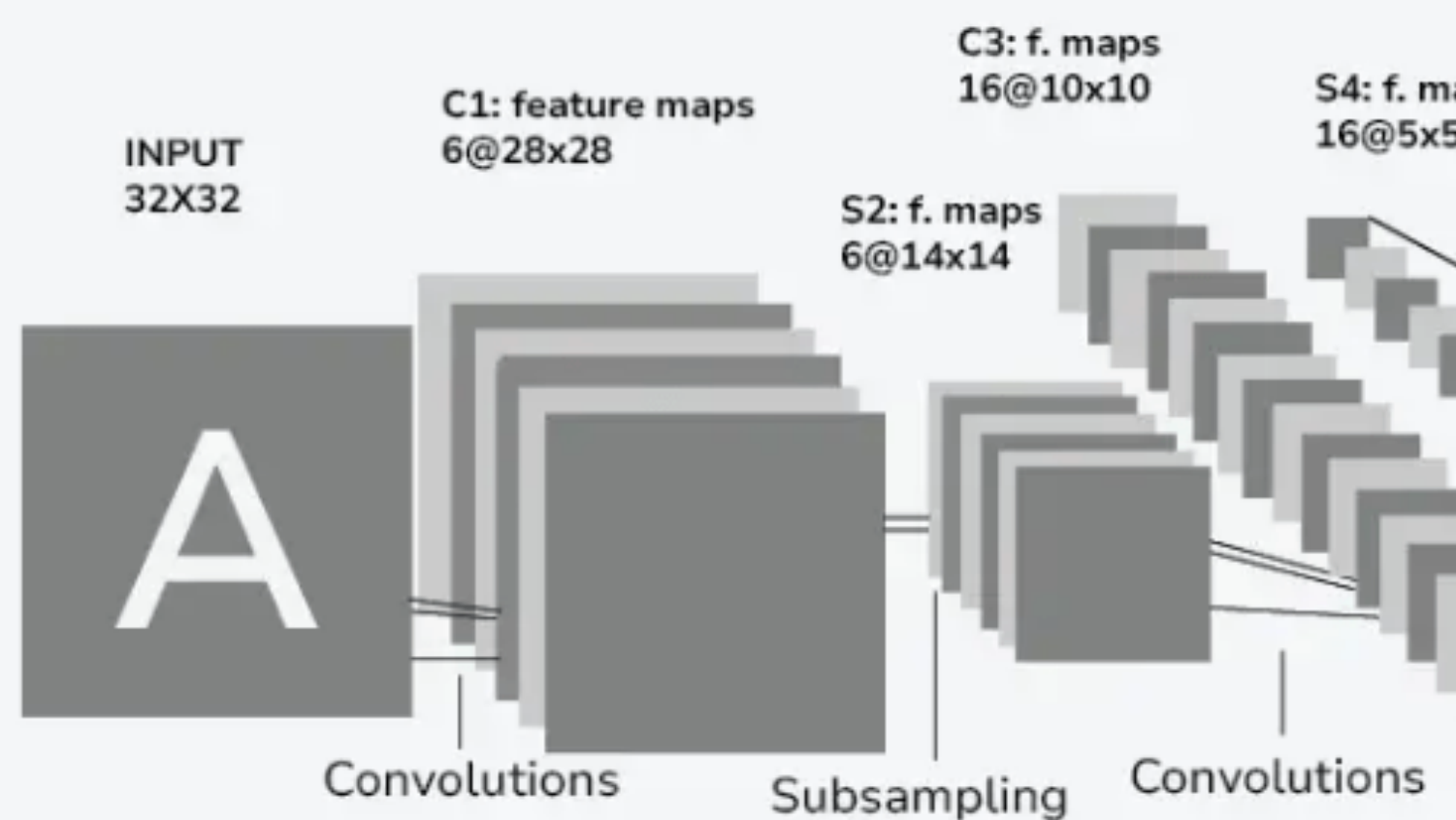| | | |
|---|---|---|
| **Sharpen** | $\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$ | |
| **Box blur** (normalized) | $\frac{1}{9}\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$ | |
| **Gaussian blur** (approximation) | $\frac{1}{16}\begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$ | |

# Deep Convolutional Networks
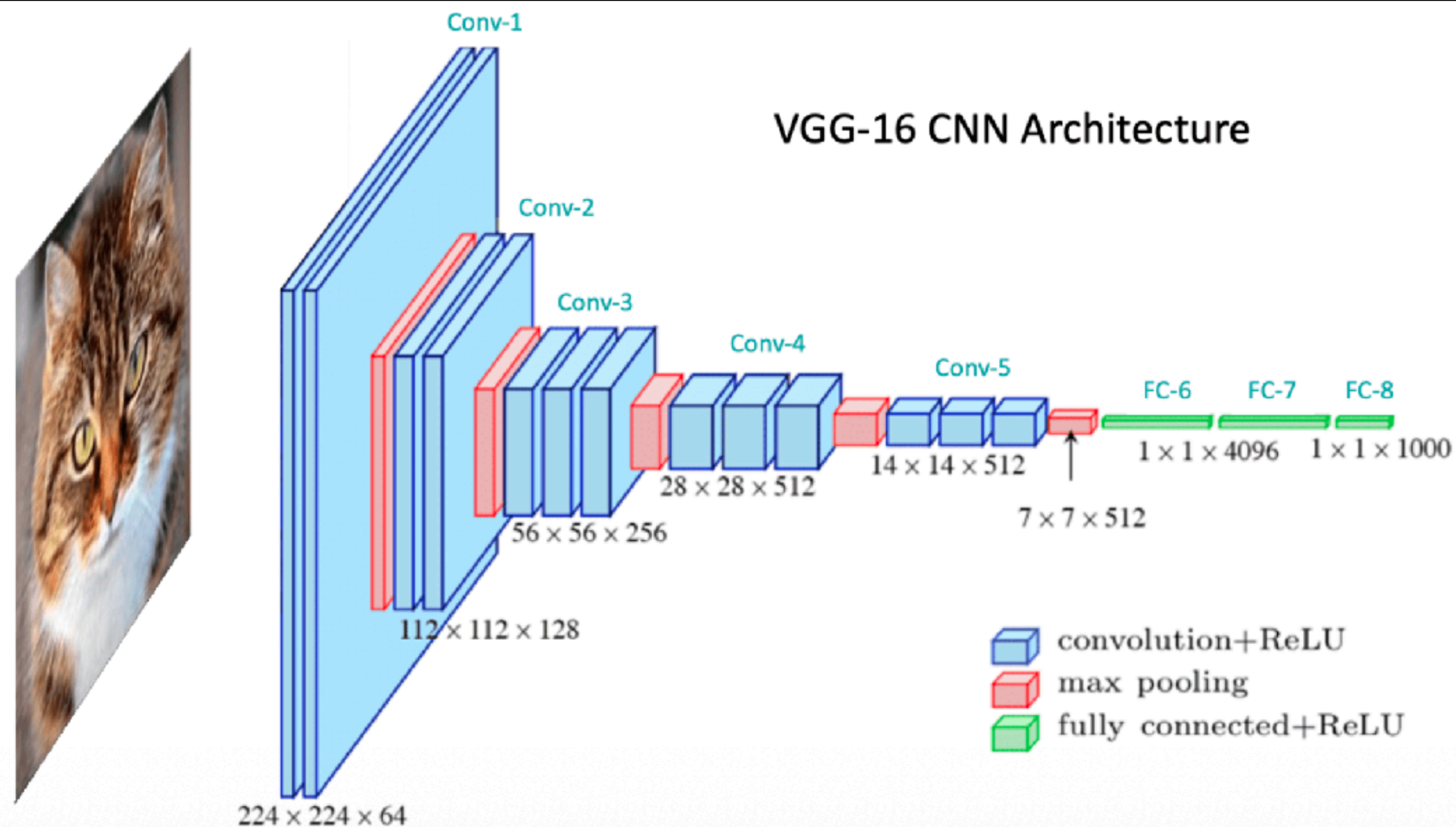
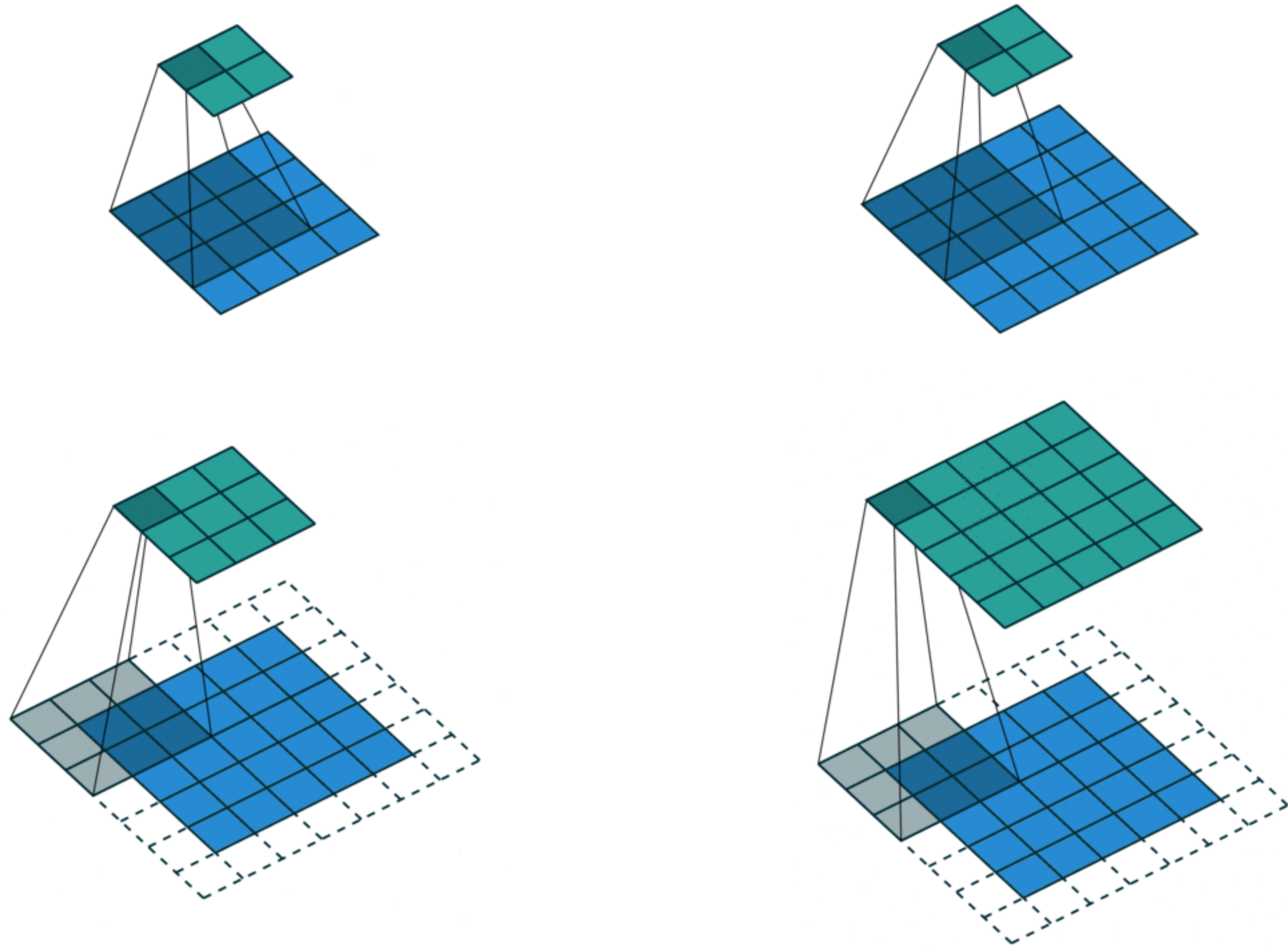# Deep Convolutional Networks

# Deep Convolutional Networks
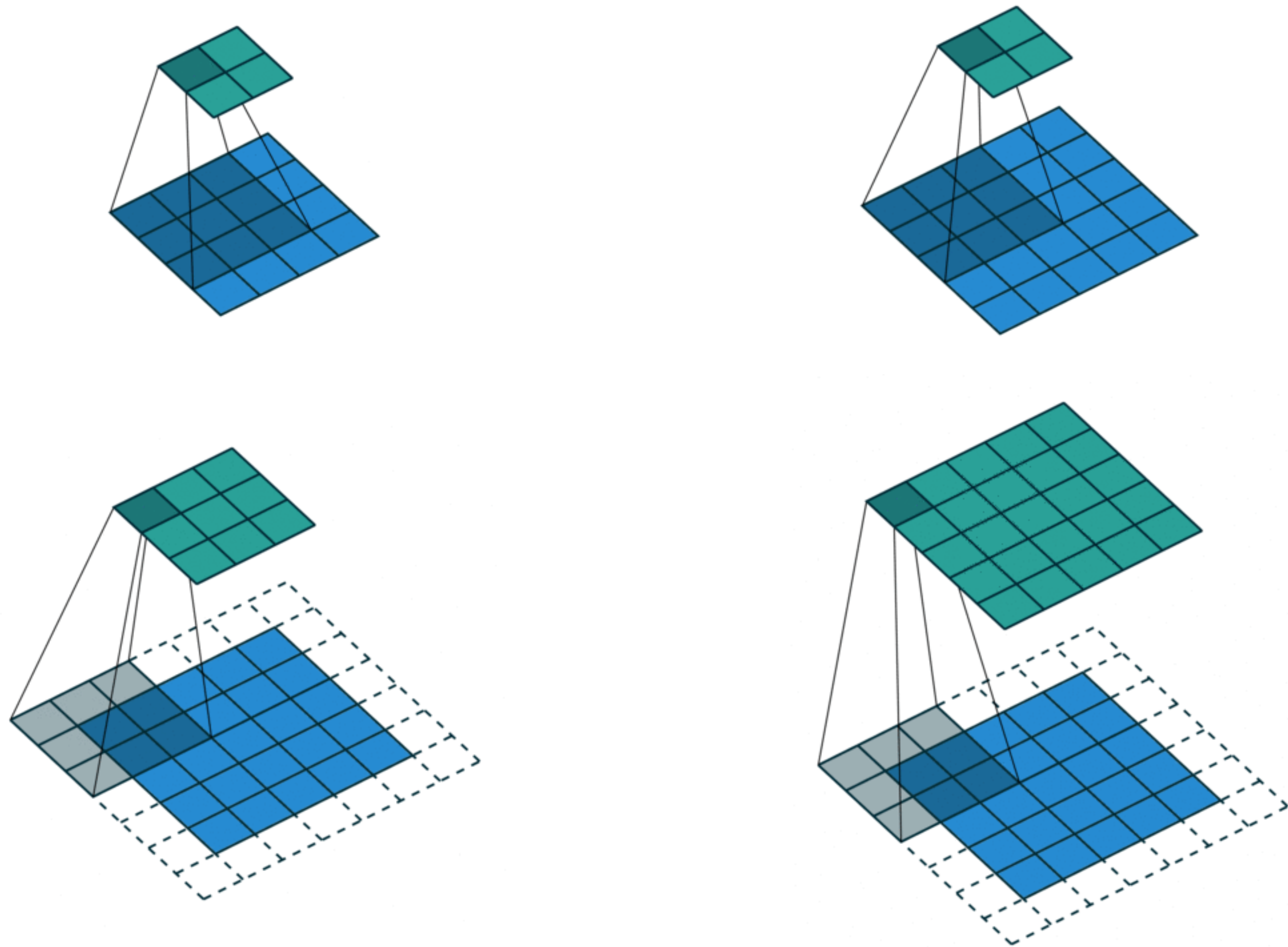
# Other Convolution Parameters

**Stride and padding!**



https://github.com/vdumoulin/conv_arithmetic

# Other Convolution Parameters

**Stride and padding!**



https://github.com/vdumoulin/conv_arithmetic

# Provable Benefit of Convolution

- You will encounter the same "patterns" at different translations (position)

- Holds for 1d (audio, music), 2d (images), 3d (videos)

- This is what the brain does!

# Provable Benefit of Convolution

- You will encounter the same "patterns" at different translations (position)

- Holds for 1d (audio, music), 2d (images), 3d (videos)

- This is what the brain does!

Exercise: can you think of inputs for which we DO NOT want convolutions?