

# CS189 Summer 2018

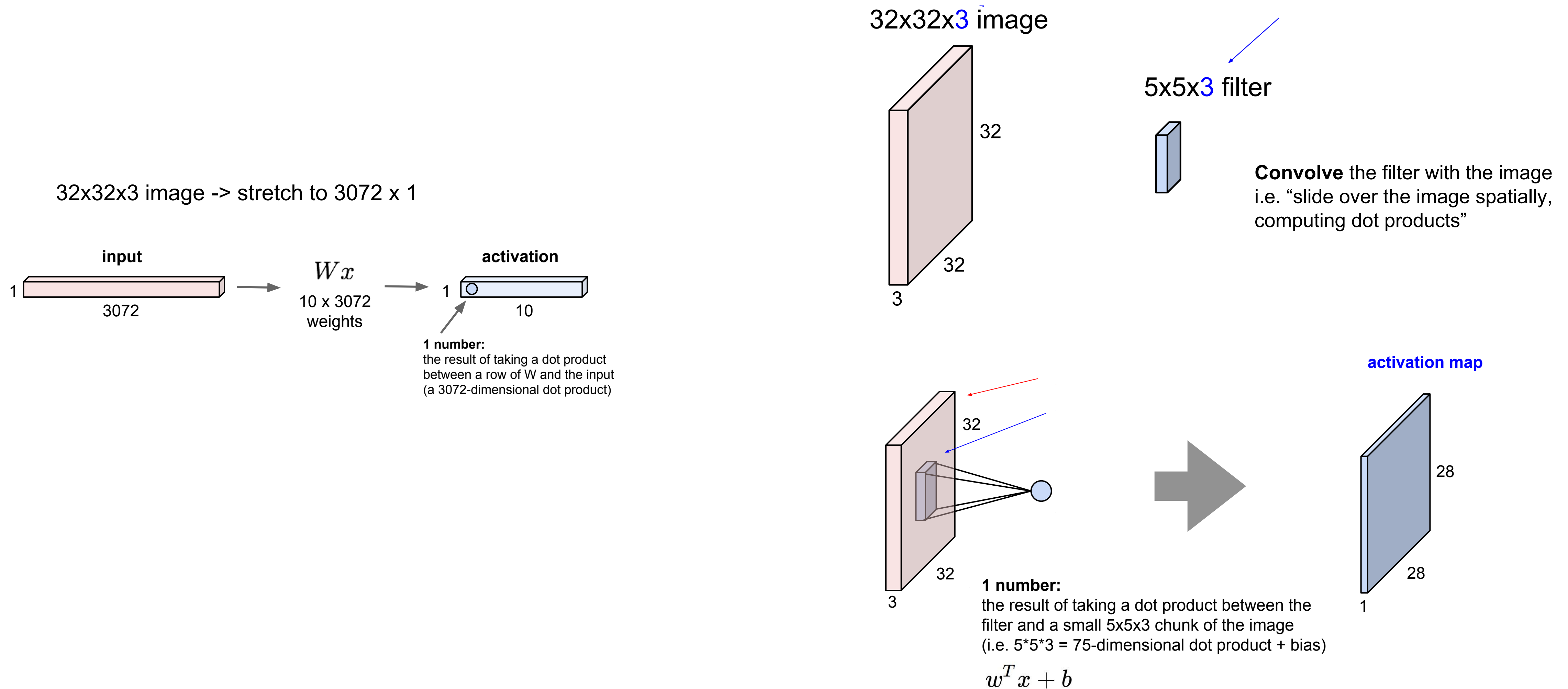
Convolutional Neural Networks

Josh Tobin

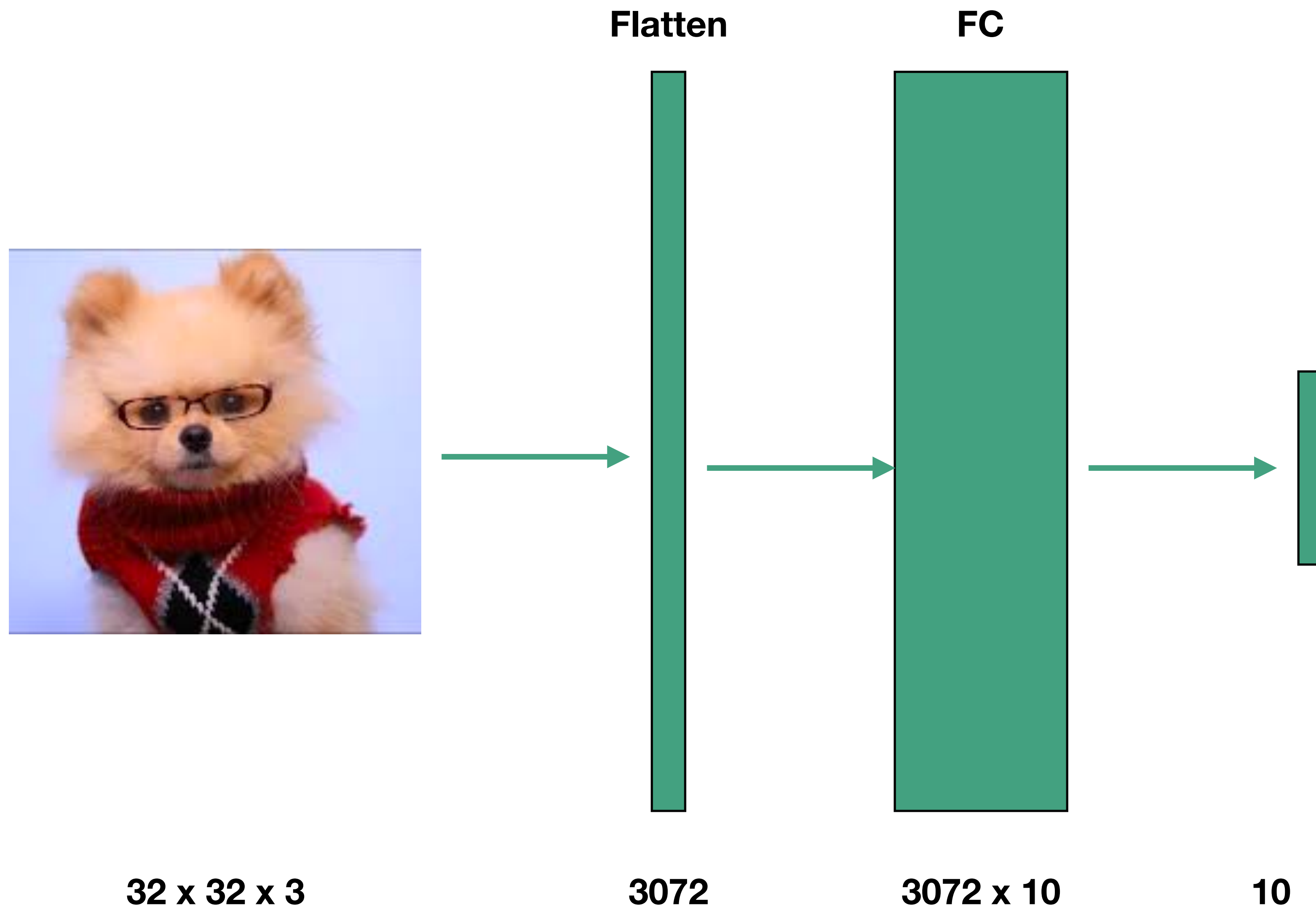
# Agenda

- 1. Overview of convolutions**
2. Other ConvNet operations
3. Basic ConvNet architectures

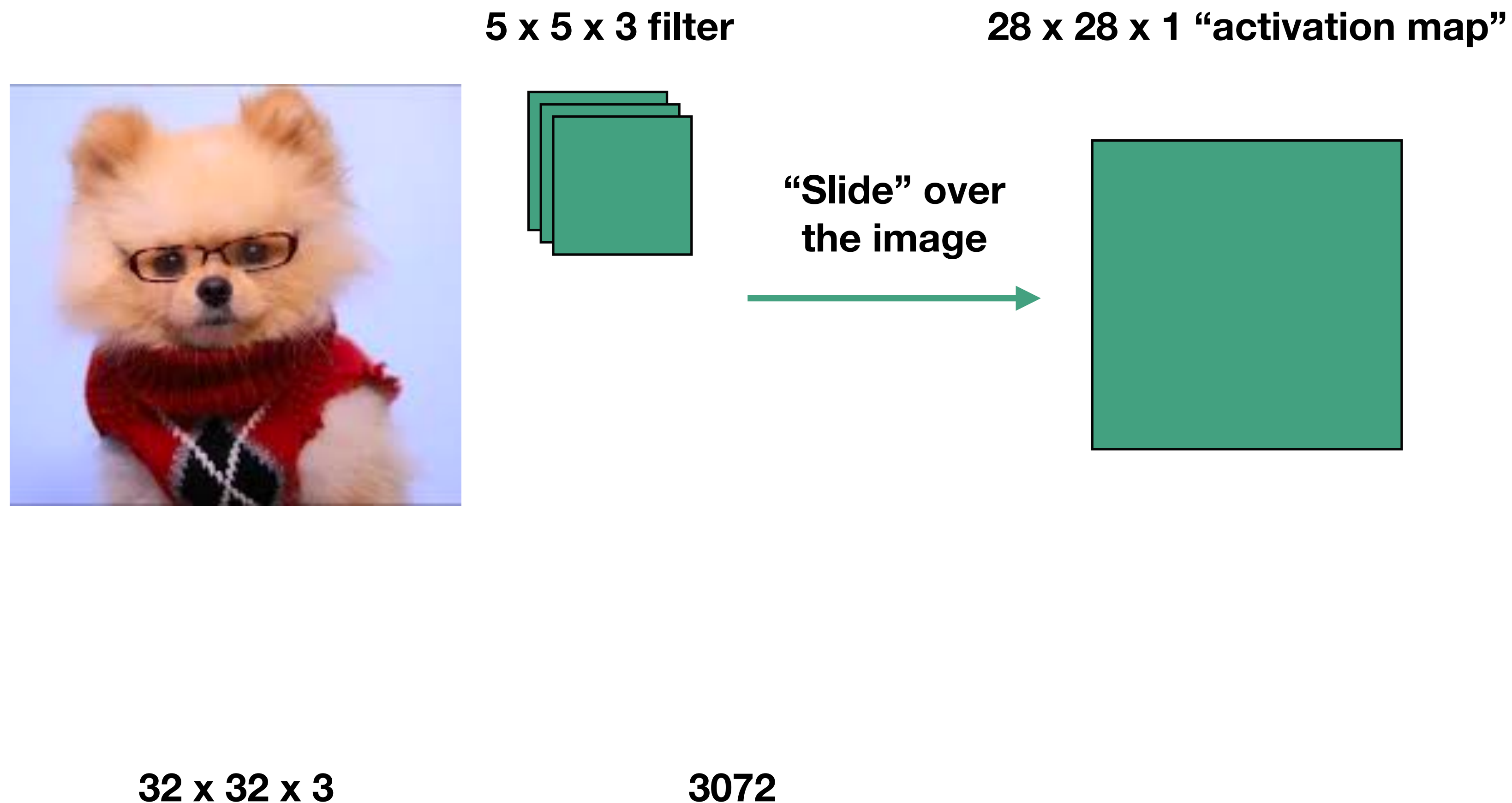
# Why convolutions?



# Fully connected vs Conv



# Fully connected vs Conv



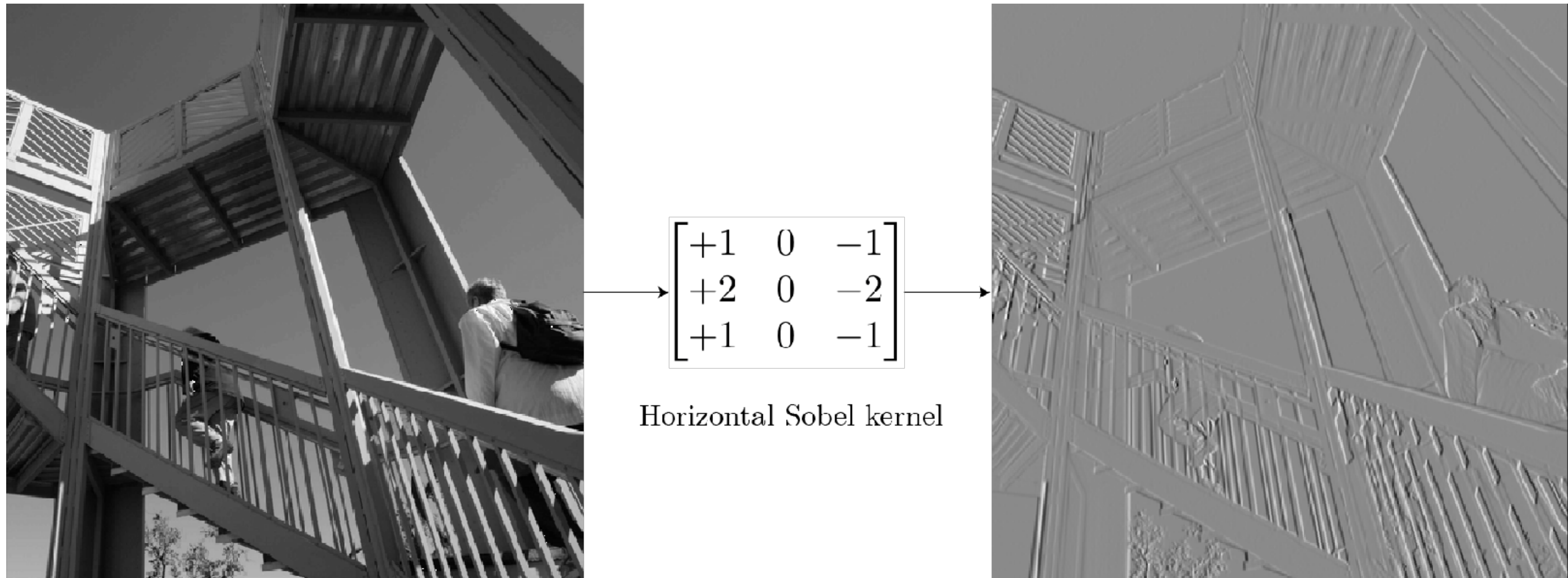
# The convolution operation

$3_0$	$3_1$	$2_2$	1	0
$0_2$	$0_2$	$1_0$	3	1
$3_0$	$1_1$	$2_2$	2	3
2	0	0	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0



# What can a conv filter do?



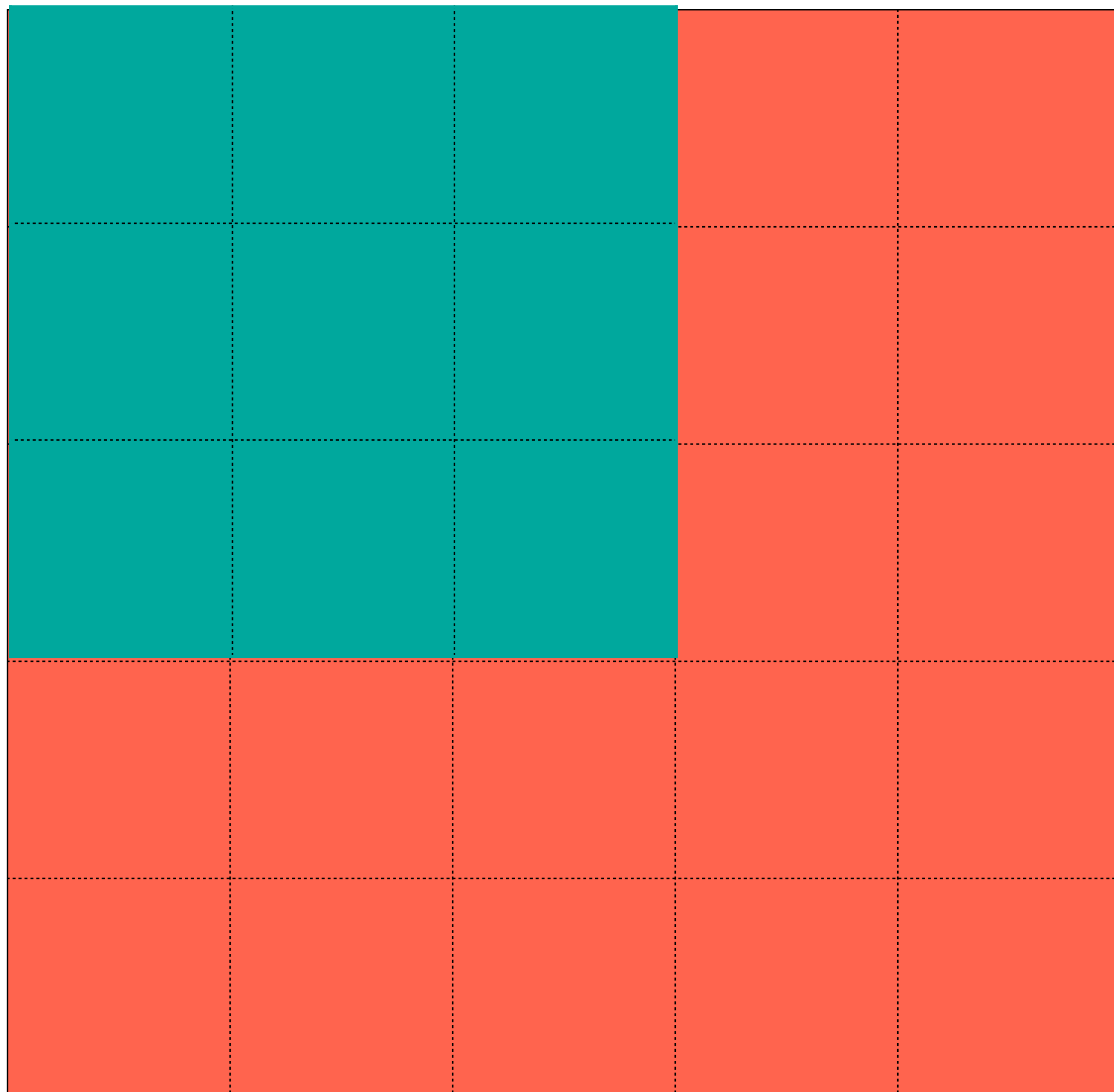
**Instead of hard-coding the weights,  
we can learn them!**

# Strides

- Convolutions can subsample the image by jumping across some locations — this is called ‘stride’



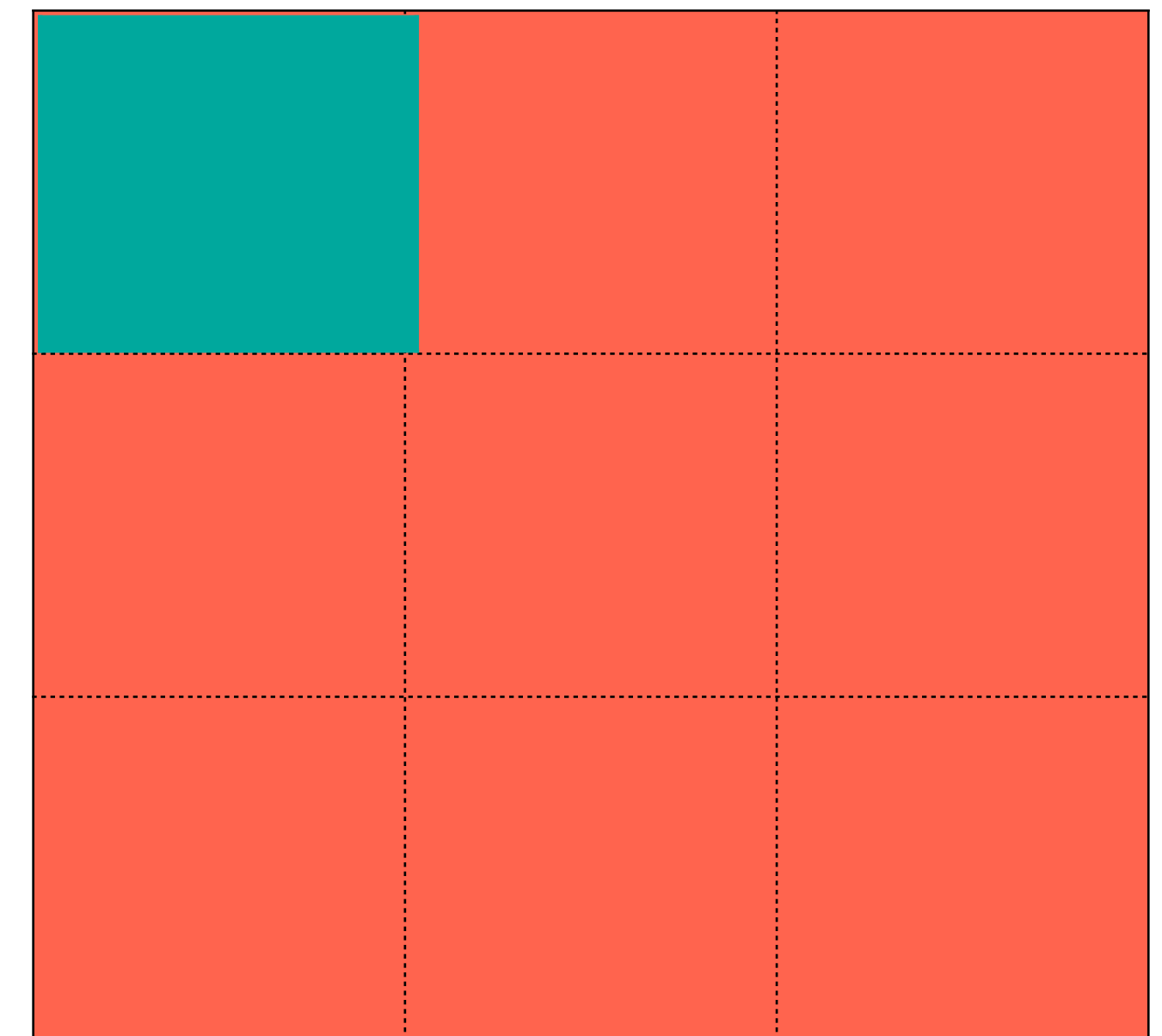
# Strides



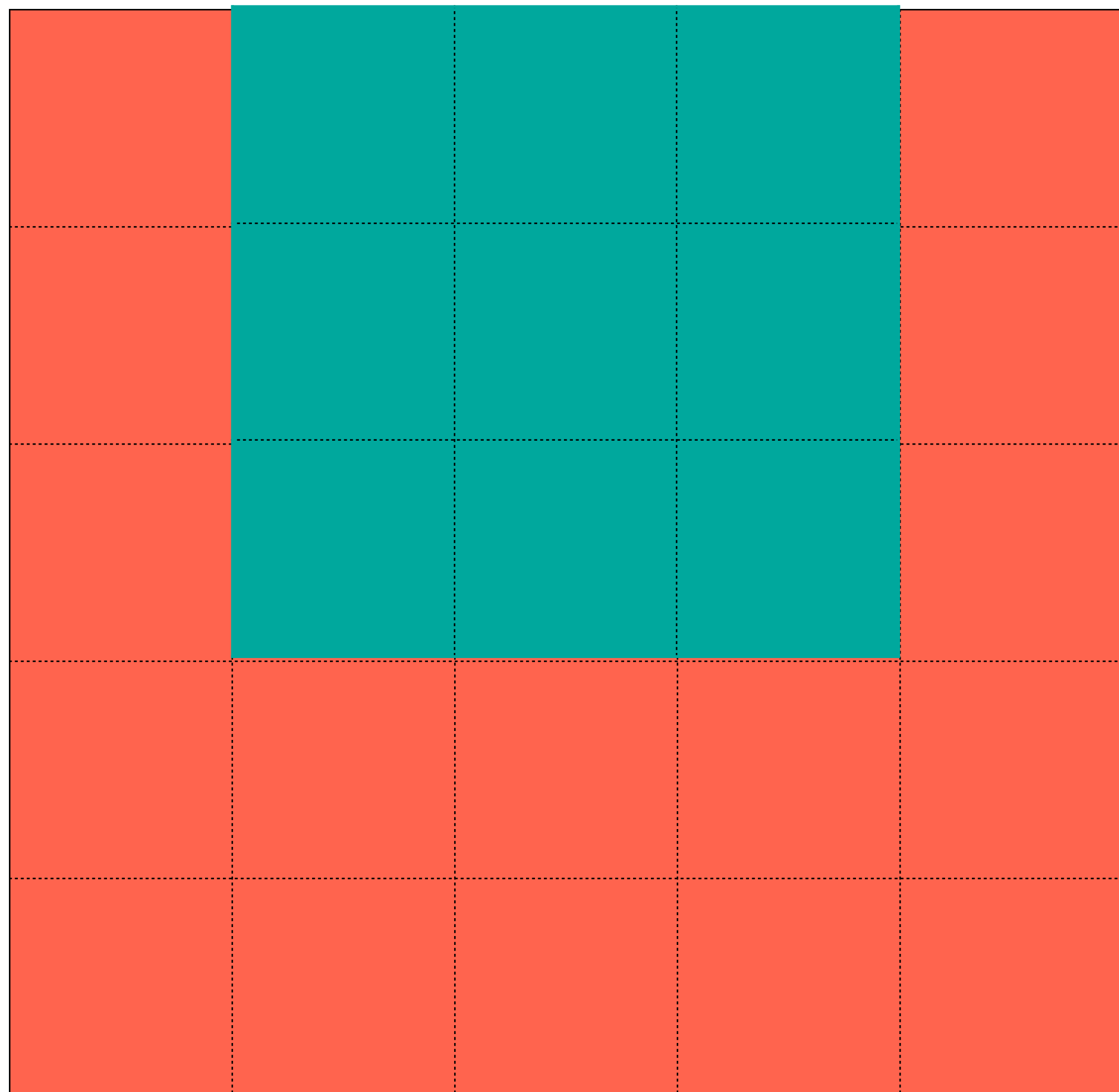
Conv2D

Filter = (3, 3)

Stride = (1, 1)



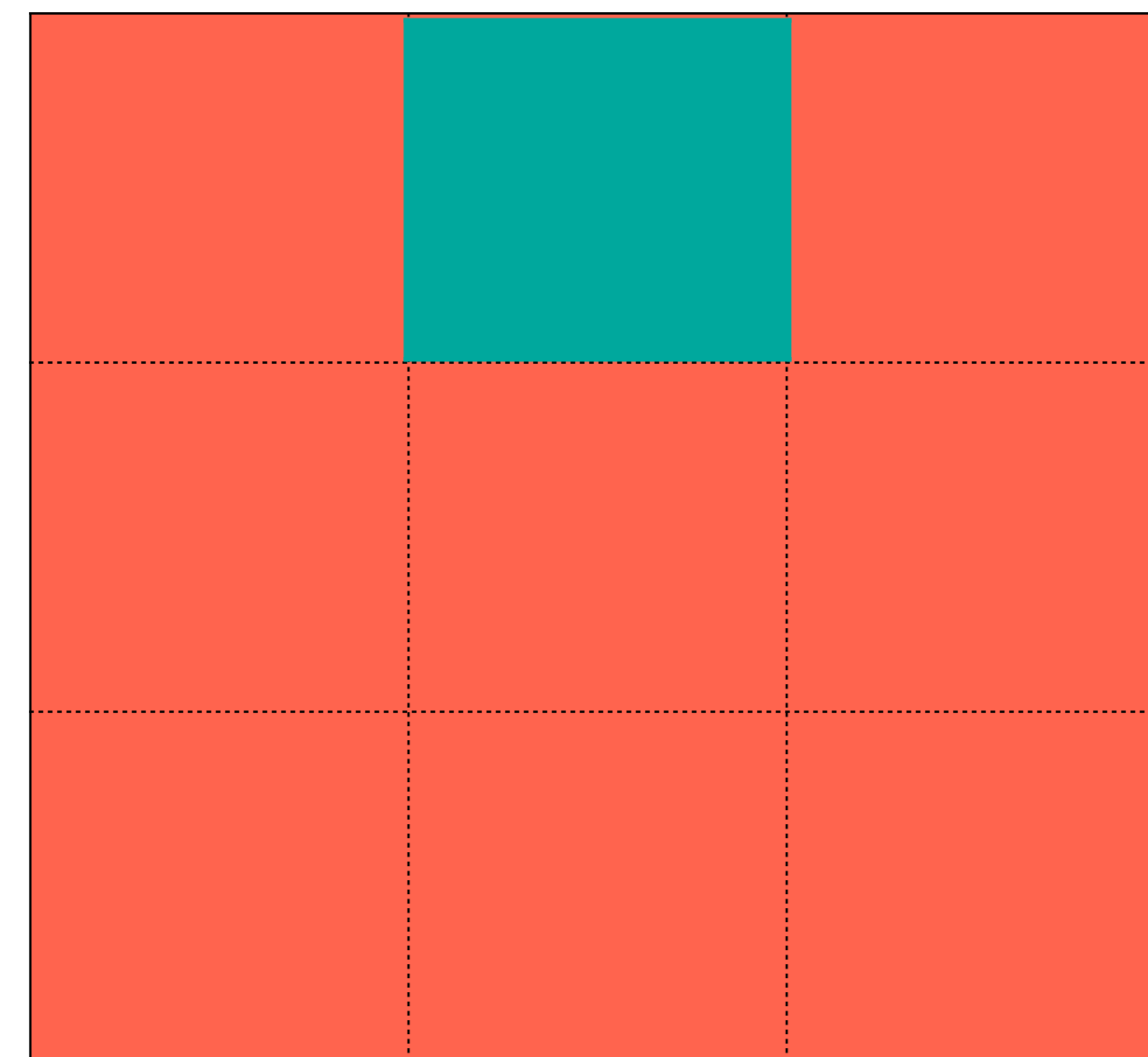
# Strides



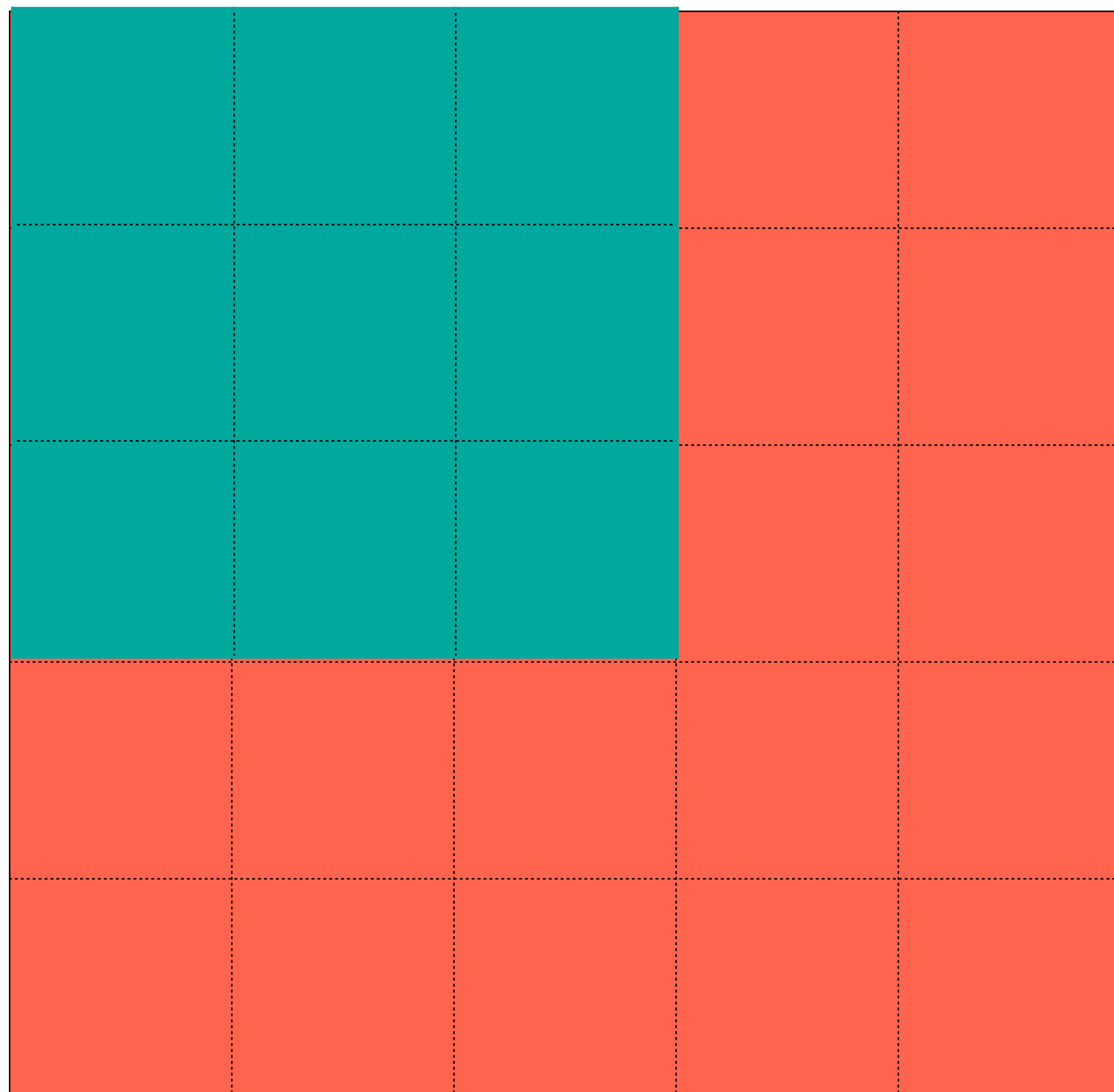
Conv2D

Filter = (3, 3)

Stride = (1, 1)



# Strides



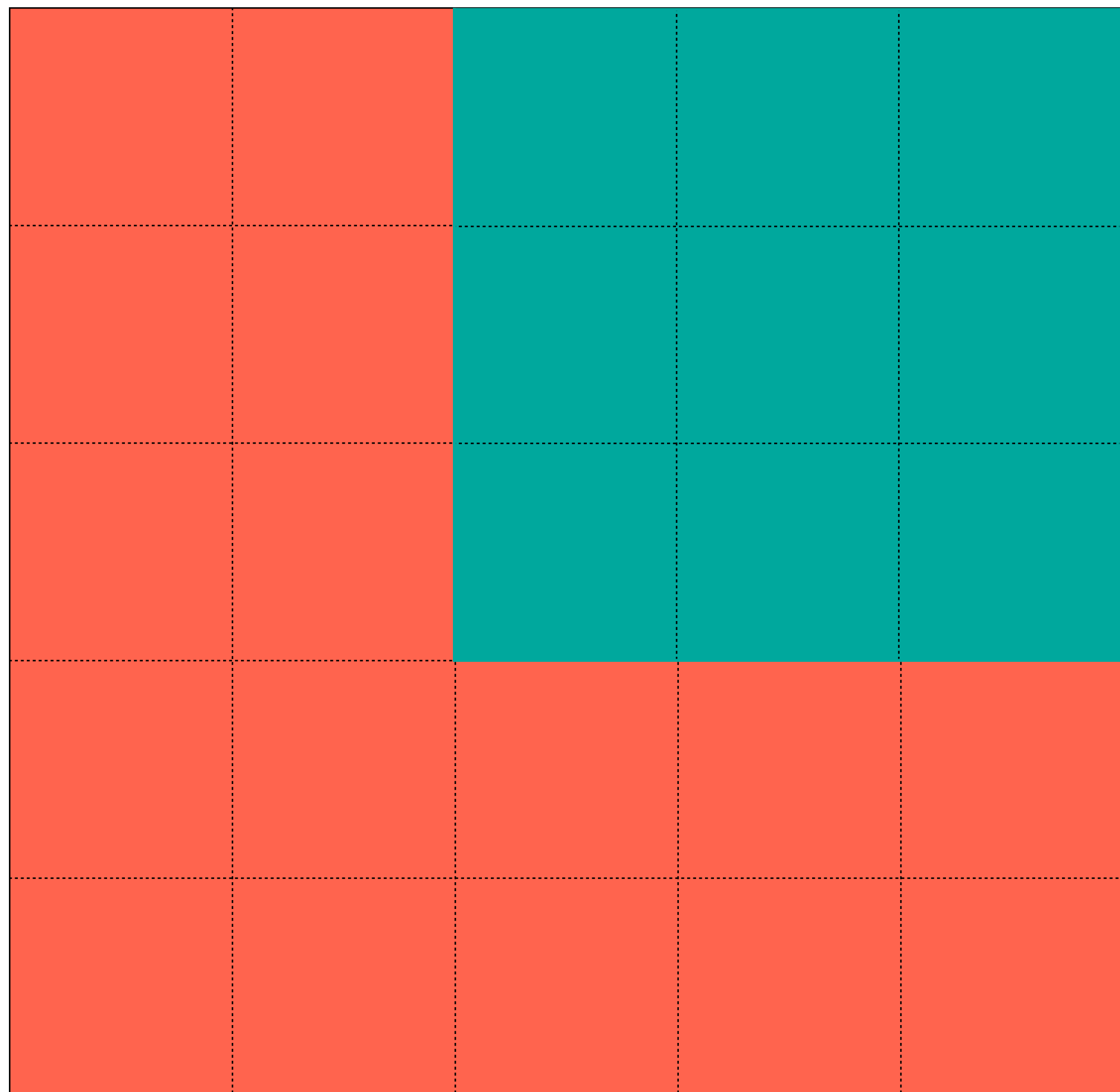
Conv2D

Filter = (3, 3)

**Stride = (2, 2)**



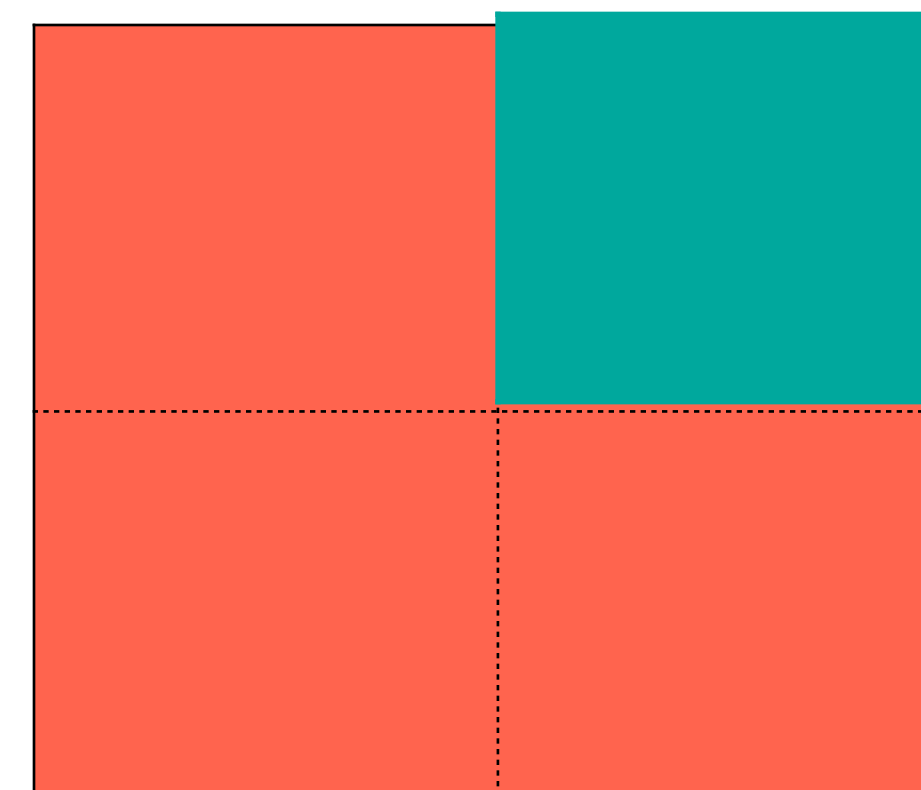
# Strides



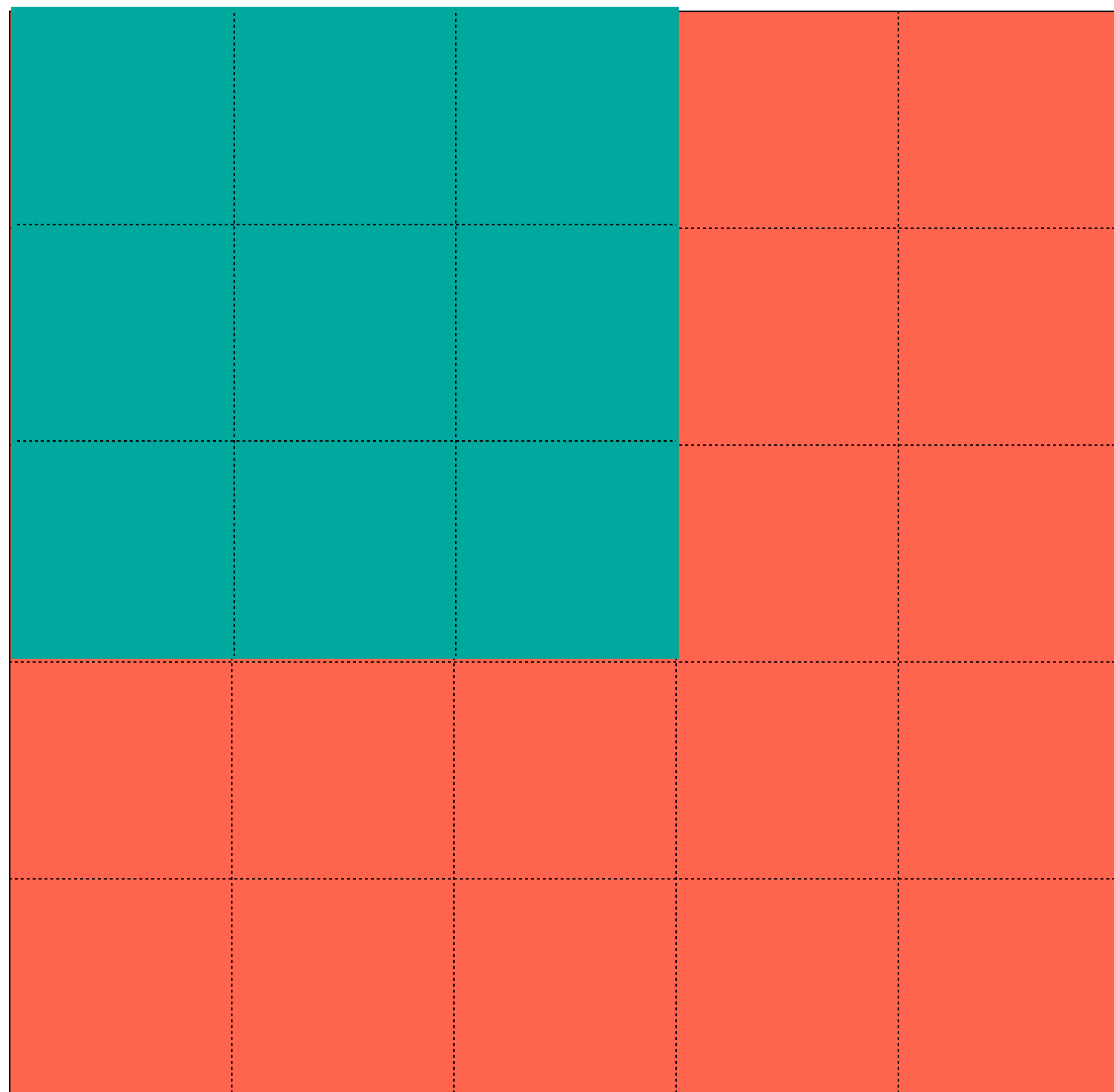
Conv2D

Filter = (3, 3)

**Stride = (2, 2)**



# Strides



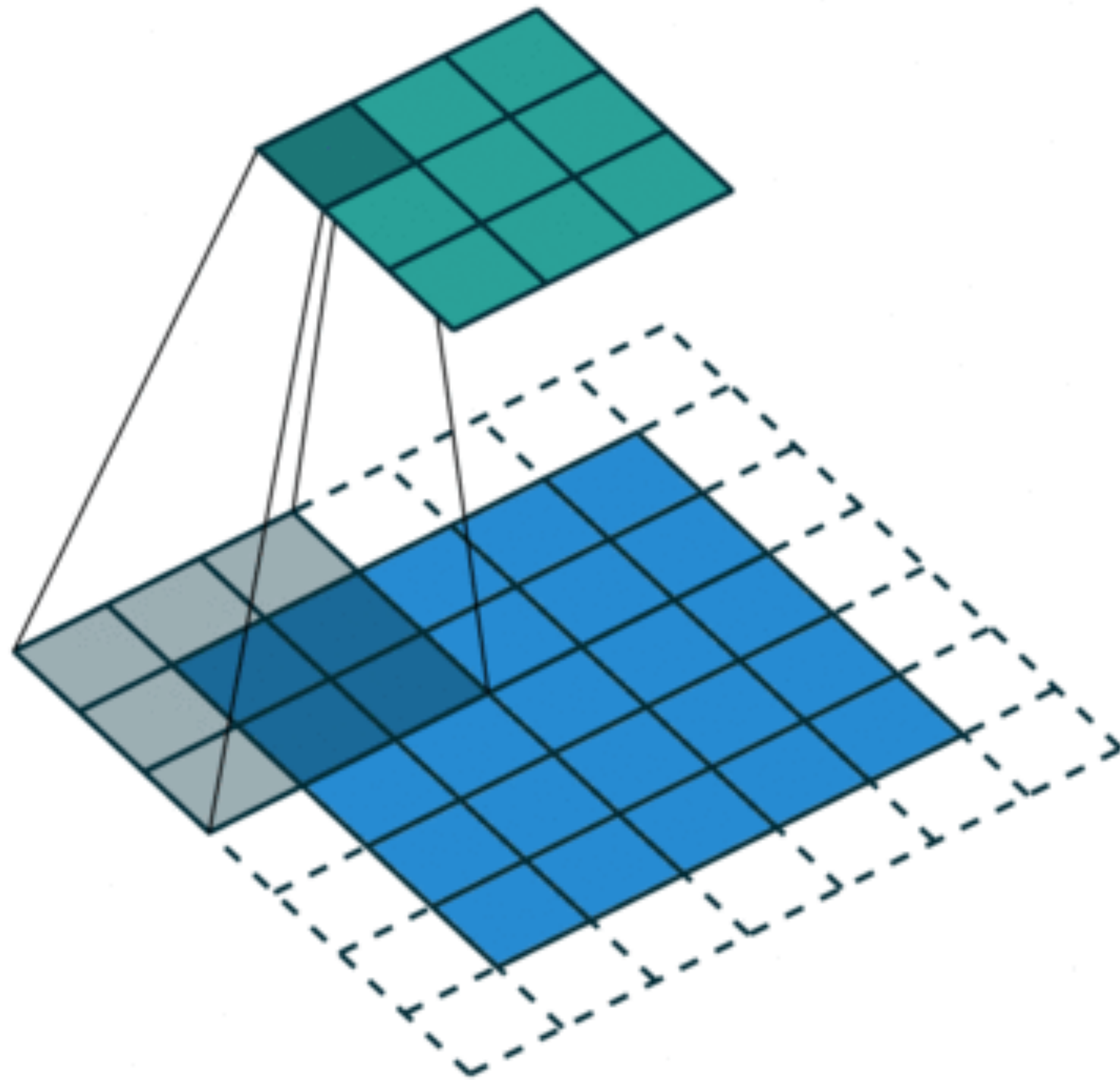
Conv2D

Filter = (3, 3)

**Stride = (3, 3)**

?

# Padding




- Padding solves the problem of filters running out of image
- Done by adding extra rows/cols to the input (usually set to 0)
- ‘SAME’ padding is illustrated here for filter=(3,3) with stride=(2,2)
- Not padding is called ‘VALID’ padding




# Conv2D Math

- Input:  $W \times H \times D$  volume
- Parameters:
  - $K$  is the number of filters, each one of size  $(F_w, F_h)$ ...
  - ...moving at stride  $(S_w, S_h)$
  - ...with the input padded symmetrically with  $P$  all-zero rows/cols
- Output:  $W' \times H' \times K$  volume
  - $W' = (W - F_w + 2P) / S_w + 1$
  - $H' = (H - F_h + 2P) / S_h + 1$
- Each filter has  $(F_w * F_h * D)$  parameters, for  $K * (F_w * F_h * D)$  total in the layer

# Conv2D Math

- Input: WxHxD volume
- Parameters:
  - K is the number of filters... 
  - ...each one of size (F\_w, F\_h)
  - ...moving at stride (S\_w, S\_h)
  - ...with the input padded symmetrically with P all-zero rows/cols
- **Commonly set to powers of 2 (e.g. 32, 64, 128)**
- Output: W'xH'xK volume
  - $W' = (W - F_w + 2P) / S_w + 1$
  - $H' = (H - F_h + 2P) / S_h + 1$
- Each filter has (F\_w \* F\_h \* D) parameters, for K \* (F\_w \* F\_h \* D) total in the layer

# Conv2D Math

- Input: WxHxD volume
  - Parameters:
    - K is the number of filters...
    - ...each one of size (F\_w, F\_h)
    - ...moving at stride (S\_w, S\_h)
    - ...with the input padded symmetrically with P all-zero rows/cols
  - Output: W'xH'xK volume
    - $W' = (W - F_w + 2P) / S_w + 1$
    - $H' = (H - F_h + 2P) / S_h + 1$
  - Each filter has (F\_w \* F\_h \* D) parameters, for K \* (F\_w \* F\_h \* D) total in the layer
- 
- **Commonly (5, 5), (3, 3), (2, 2), (1, 1)**

# Conv2D Math

- Input:  $W \times H \times D$  volume
  - Parameters:
    - $K$  is the number of filters...
    - ...each one of size  $(F_w, F_h)$
    - ...moving at stride  $(S_w, S_h)$
    - ...with the input padded symmetrically with  $P$  all-zero rows/cols
  - Output:  $W' \times H' \times K$  volume
    - $W' = (W - F_w + 2P) / S_w + 1$
    - $H' = (H - F_h + 2P) / S_h + 1$
  - Each filter has  $(F_w * F_h * D)$  parameters, for  $K * (F_w * F_h * D)$  total in the layer
- • **‘SAME’ sets it automatically**

# A guide to convolution arithmetic for deep learning

Vincent Dumoulin<sup>1★</sup> and Francesco Visin<sup>2★†</sup>

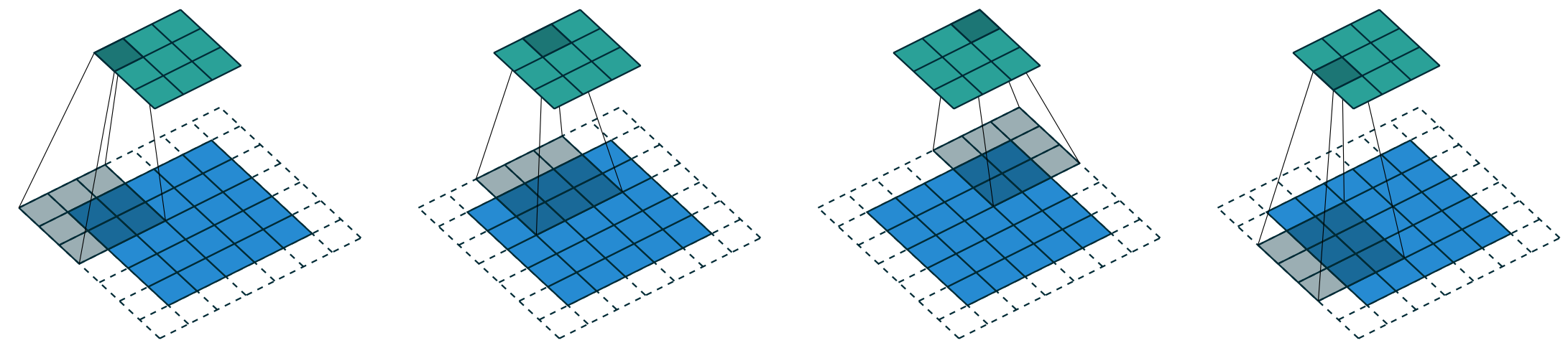


Figure 2.6: (Arbitrary padding and strides) Convolving a  $3 \times 3$  kernel over a  $5 \times 5$  input padded with a  $1 \times 1$  border of zeros using  $2 \times 2$  strides (i.e.,  $i = 5$ ,  $k = 3$ ,  $s = 2$  and  $p = 1$ ).

- Lots of cool visualizations and comforting equations

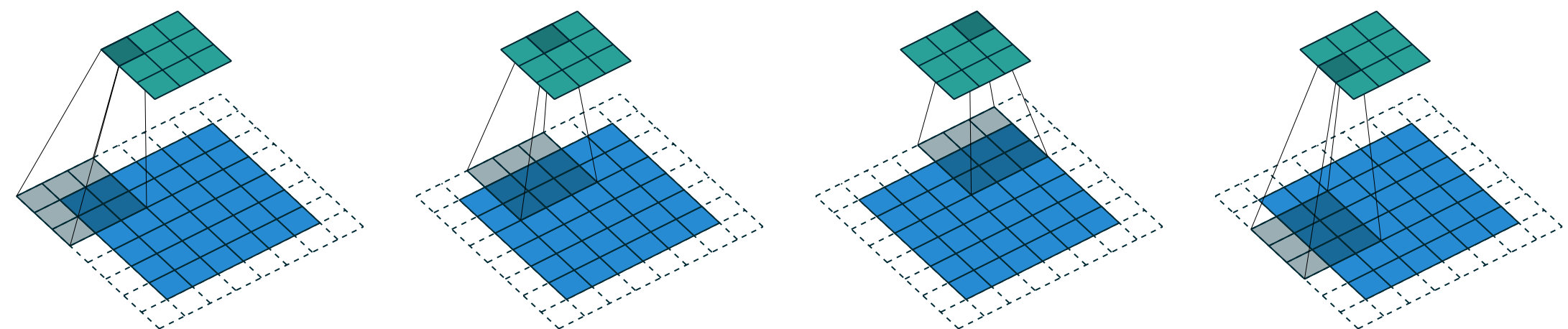
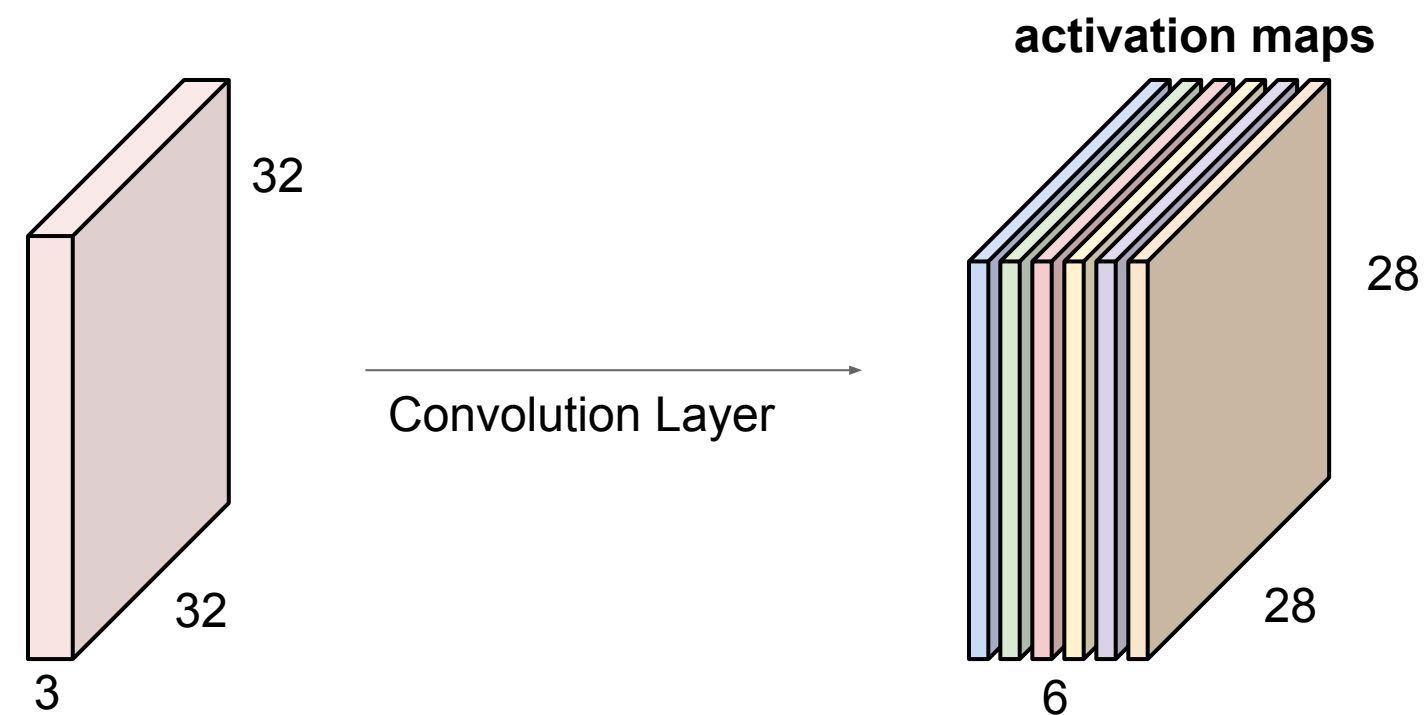


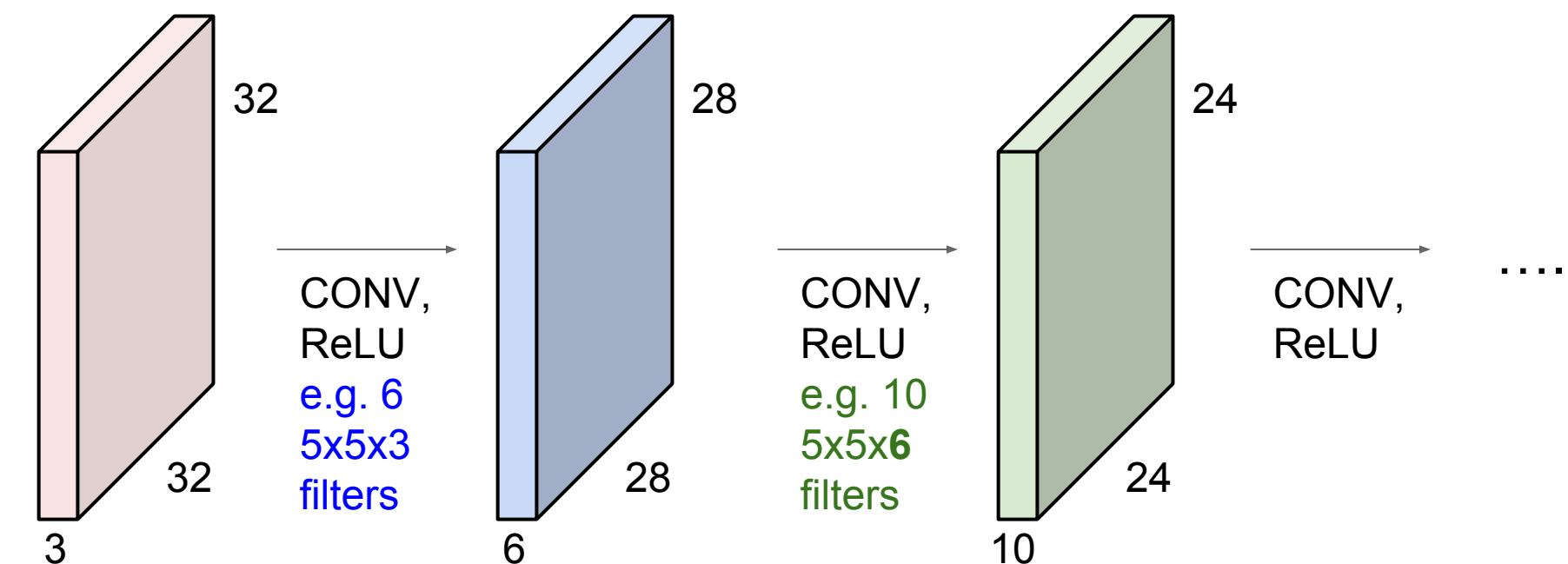
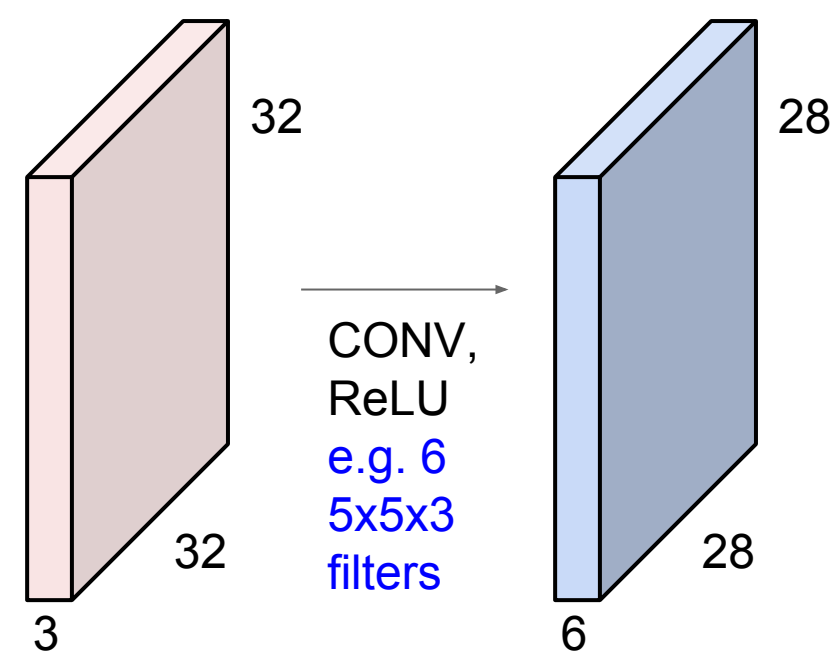
Figure 2.7: (Arbitrary padding and strides) Convolving a  $3 \times 3$  kernel over a  $6 \times 6$  input padded with a  $1 \times 1$  border of zeros using  $2 \times 2$  strides (i.e.,  $i = 6$ ,  $k = 3$ ,  $s = 2$  and  $p = 1$ ). In this case, the bottom row and right column of the zero padded input are not covered by the kernel.

# Conv2D output is another “image”

For example, if we had 6 5x5 filters, we'll get 6 separate activation maps



We stack these up to get a “new image” of size 28x28x6!

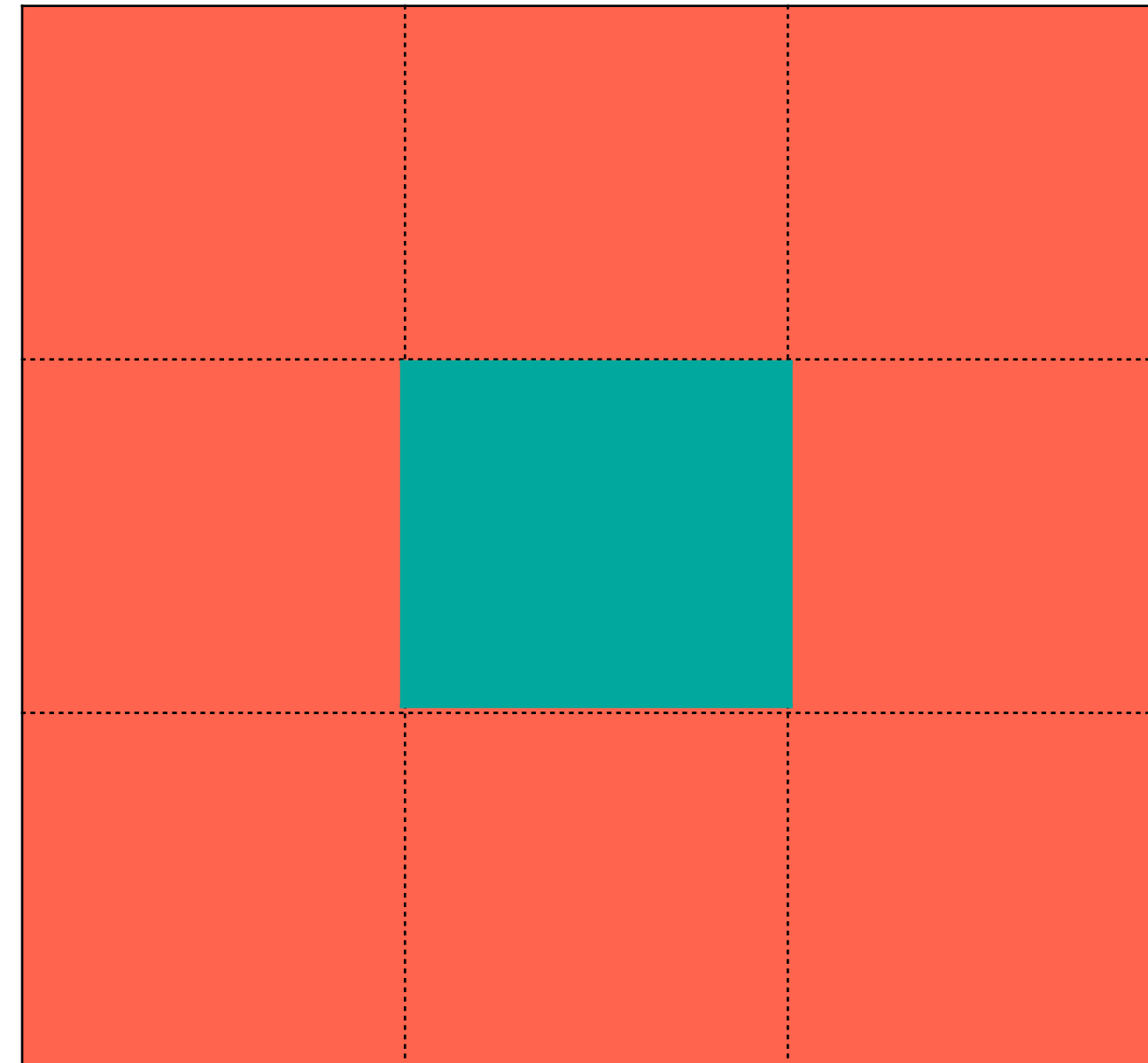
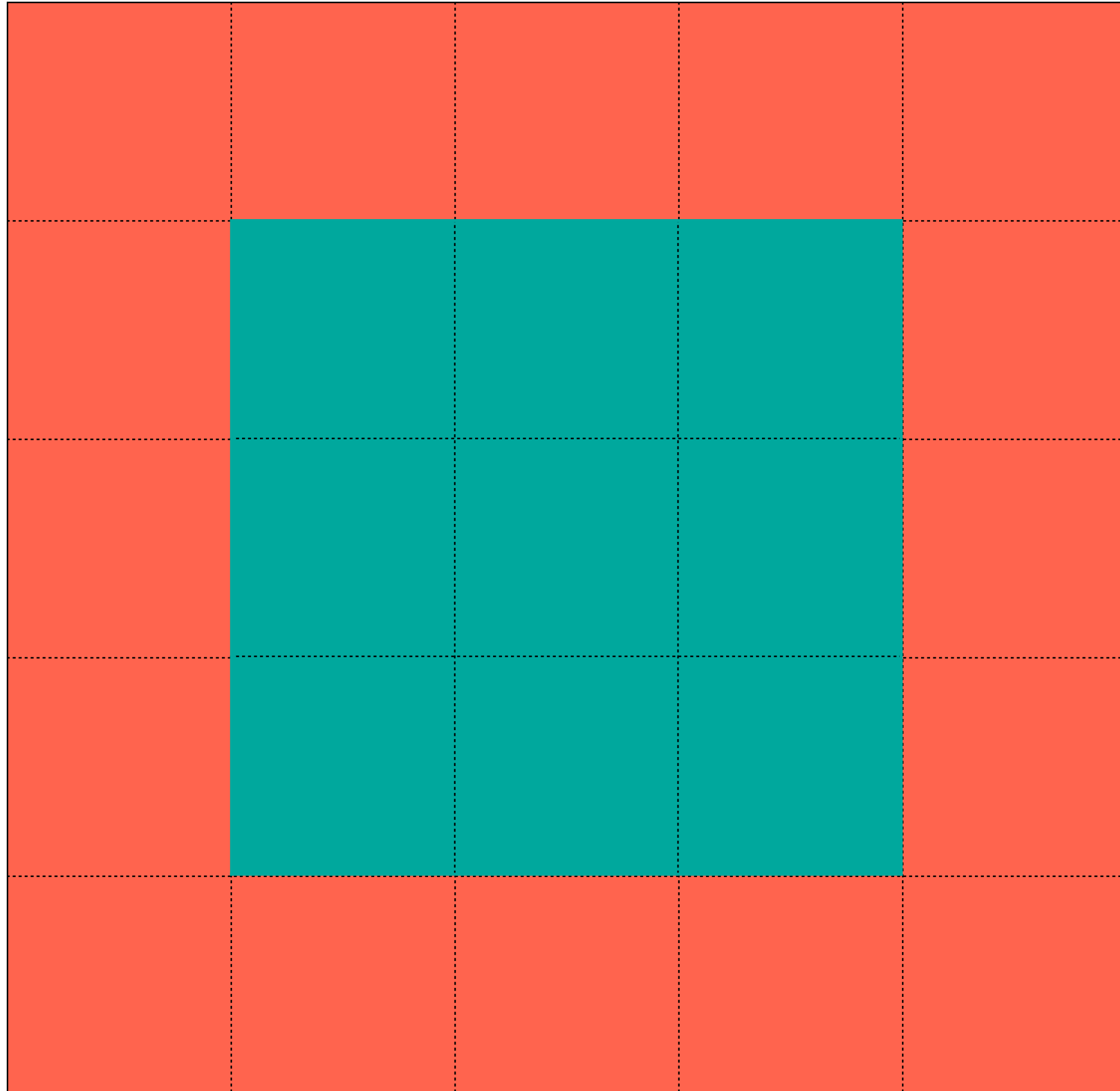


**Build a network by repeating the conv operation,  
with activation function in between**



# Receptive fields

Receptive field: 3x3

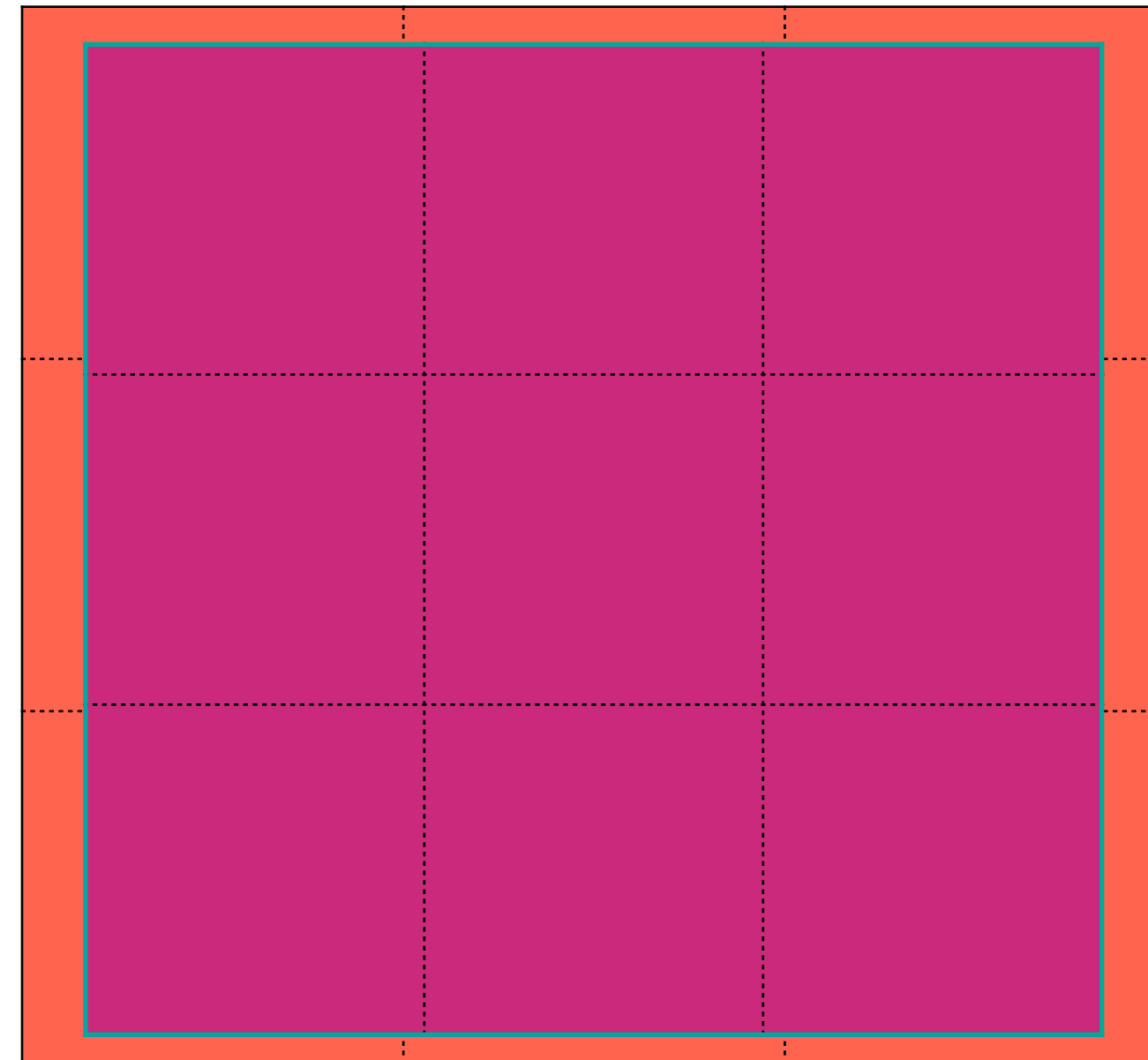
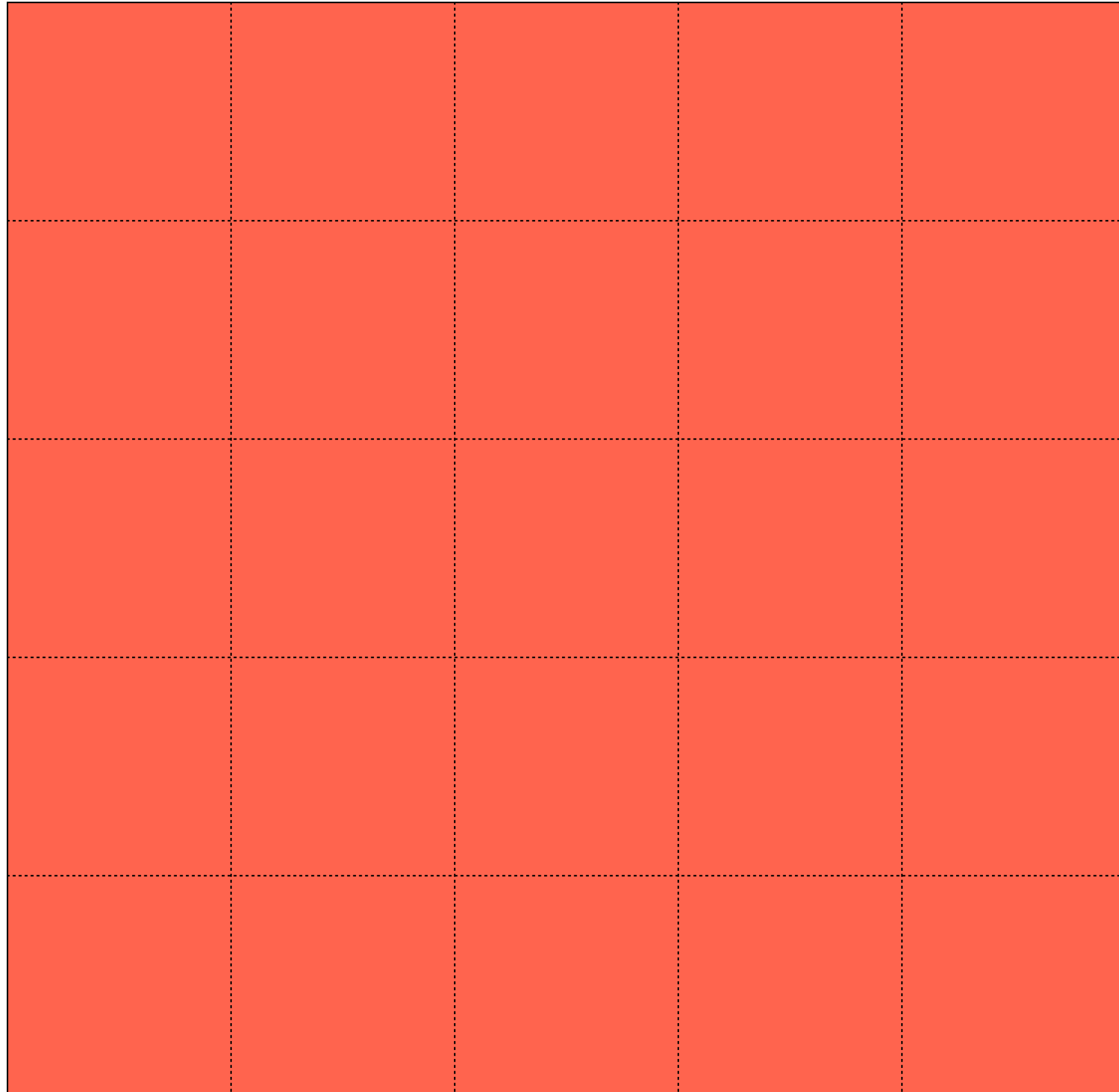


Conv2D

Filter = (3, 3)

Stride = (1, 1)

# Receptive fields



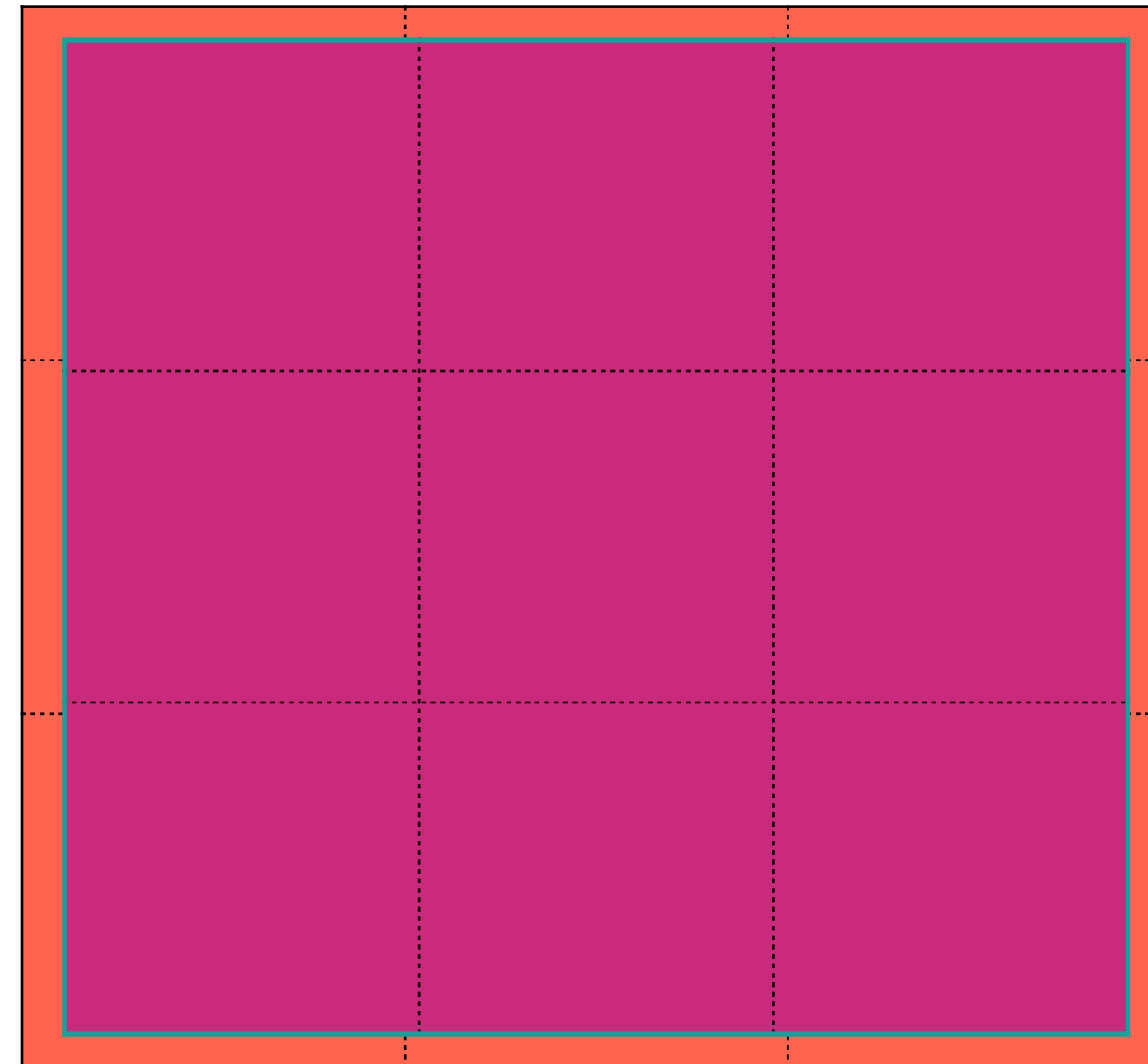
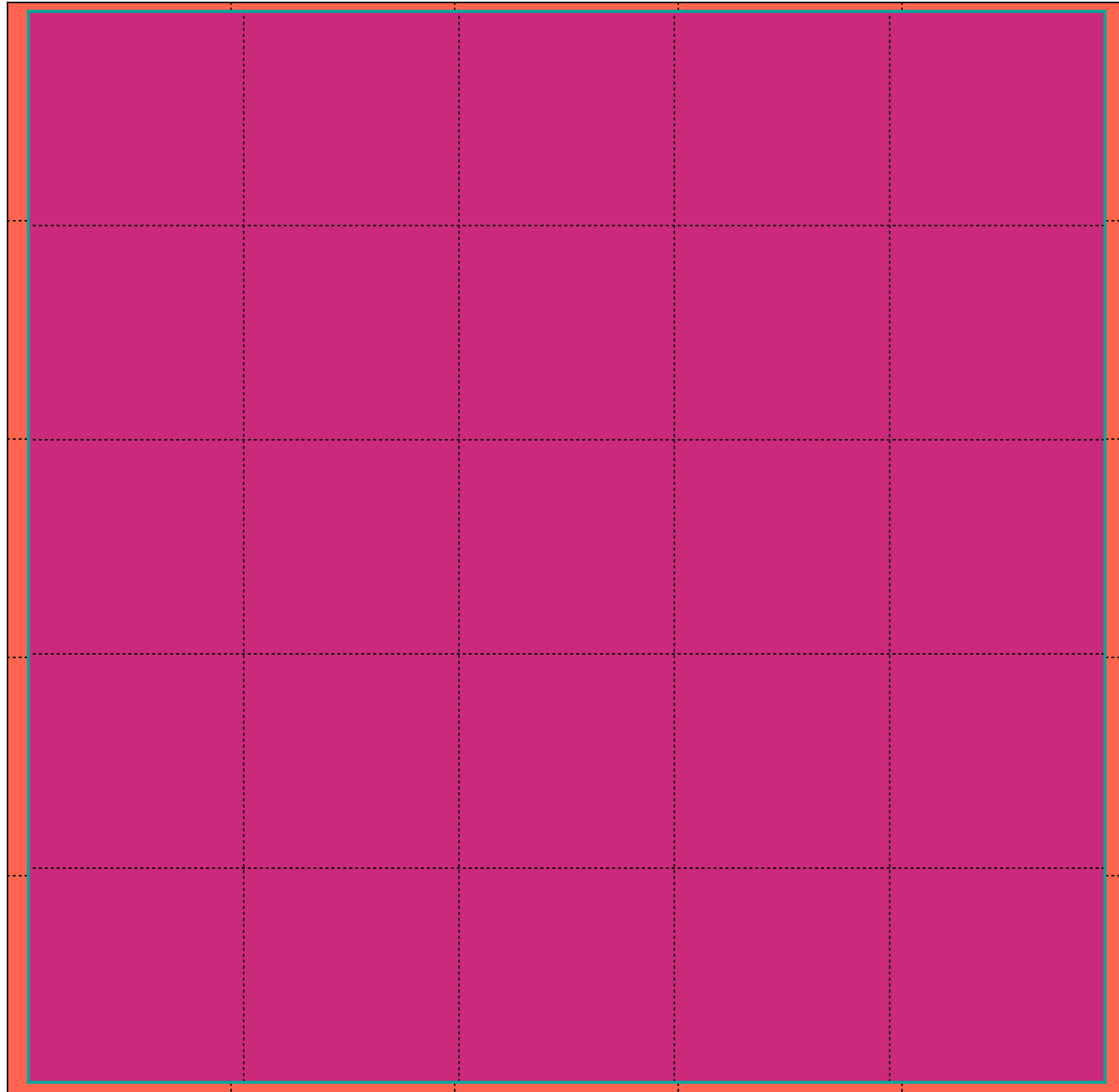
Conv2D  
Filter = (3, 3)  
Stride = (1, 1)



Conv2D  
Filter = (3, 3)  
Stride = (1, 1)

# Receptive fields

Original receptive field: 5x5



Conv2D  
Filter = (3, 3)  
Stride = (1, 1)

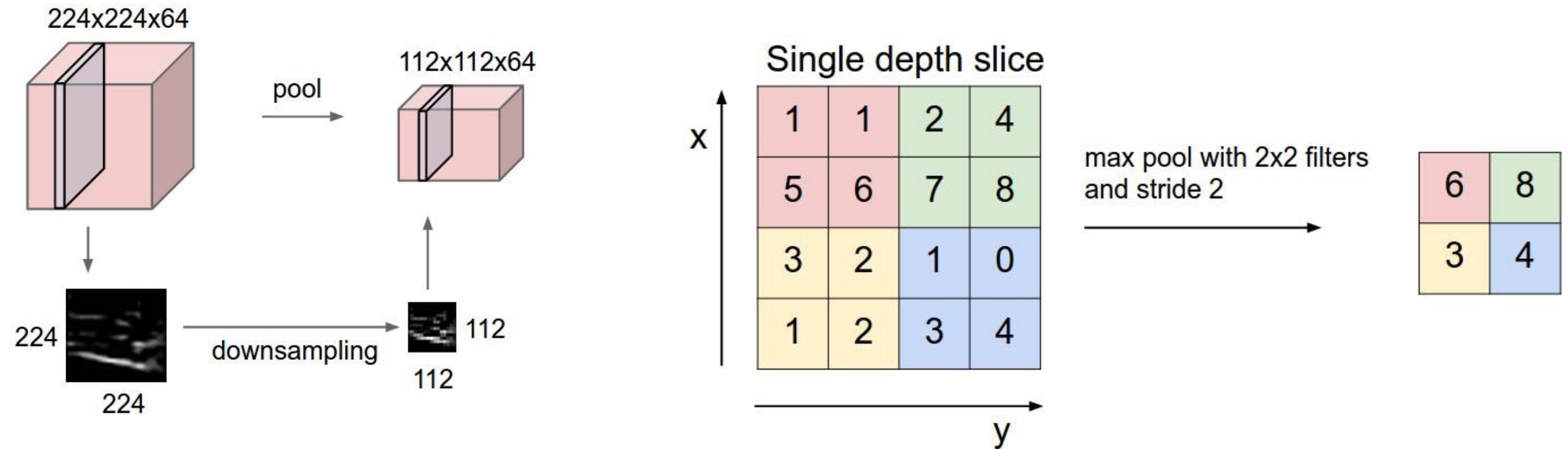


Conv2D  
Filter = (3, 3)  
Stride = (1, 1)

# Receptive fields

- Stacking convolutions one after the other increases the original receptive field: two (3, 3) convs get to a (5, 5) receptive field
- (and tend to perform better than a single (5, 5) conv)

# Pooling



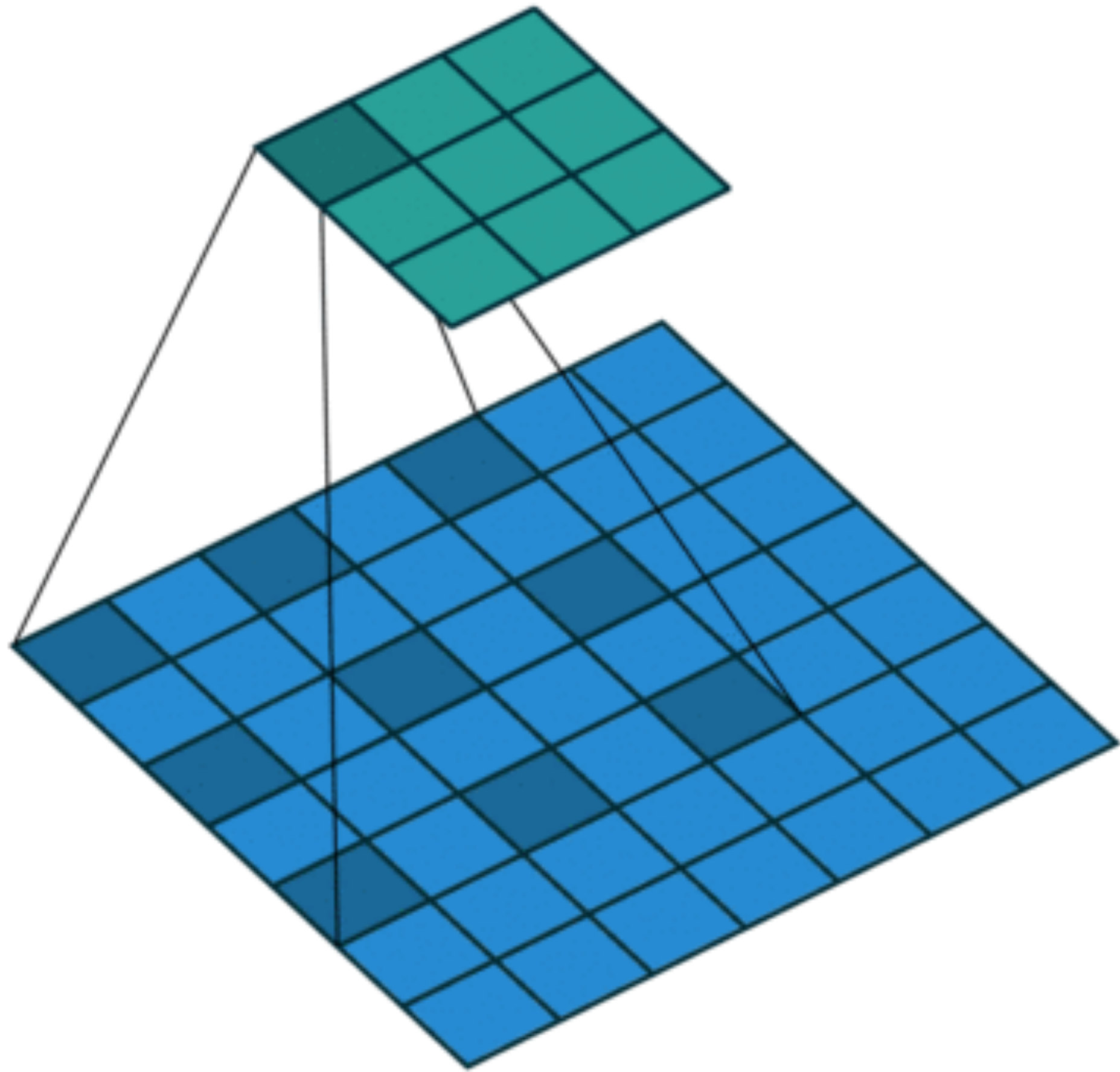
- Pooling subsamples the image by taking average or max value of a pooling region
- $2 \times 2$  max pooling is illustrated, and probably most common
- Very important in early convnet applications, but has recently fallen out of favor (strided convolutions can subsample, and tend to work better in GAN applications).

# Agenda

1. Basics of convolutions, and filter/stride/pooling math
- 2. More advanced convolution types, and computational considerations**
3. Classic convnet architectures

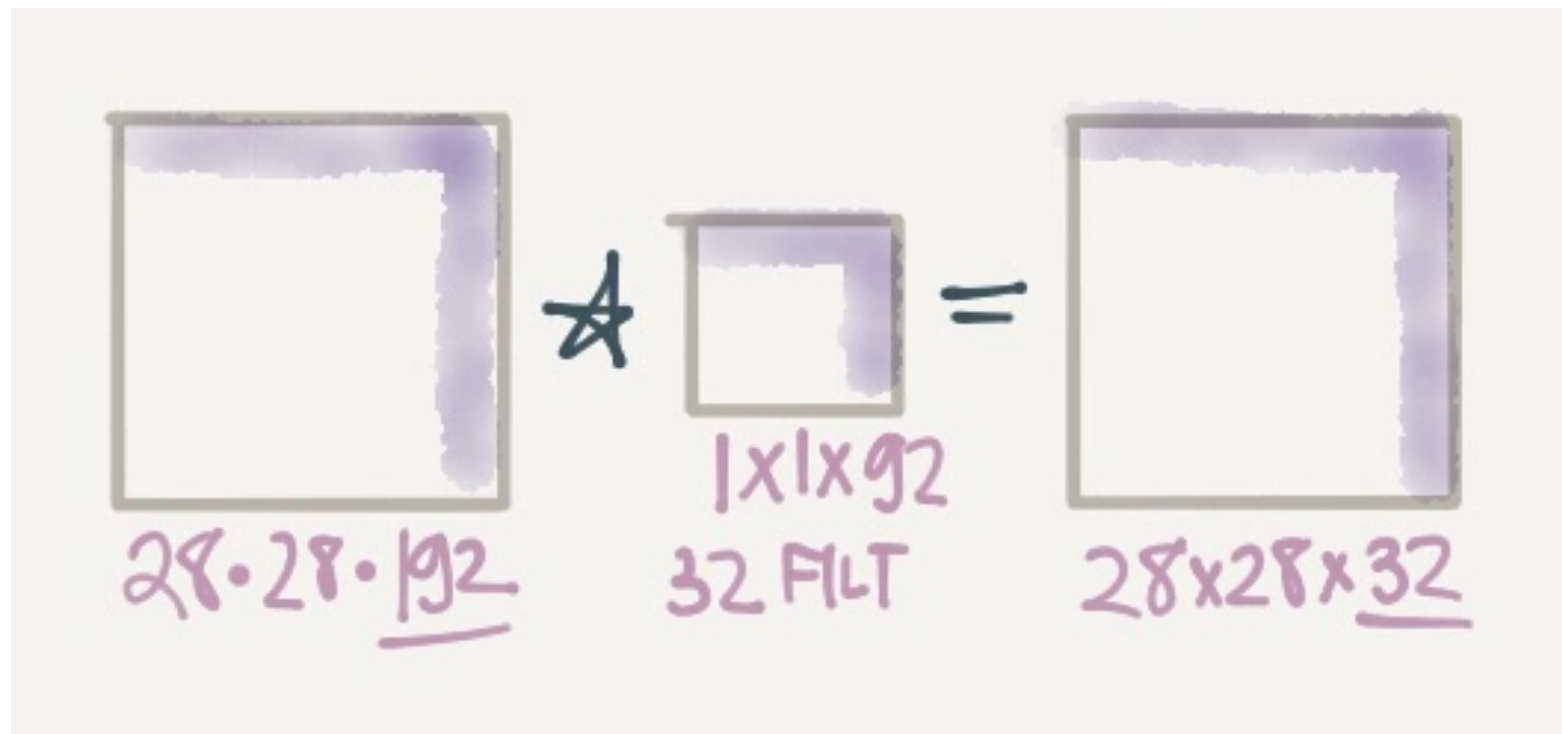


# Dilated Convolution



- Dilated convolutions can “see” a greater portion of the image by skipping pixels
- The (3, 3) 1-dilated convolution illustrated here has a (5, 5) receptive field
- Stacking dilated convolutions up quickly gets to large receptive fields

# 1x1 Convolution



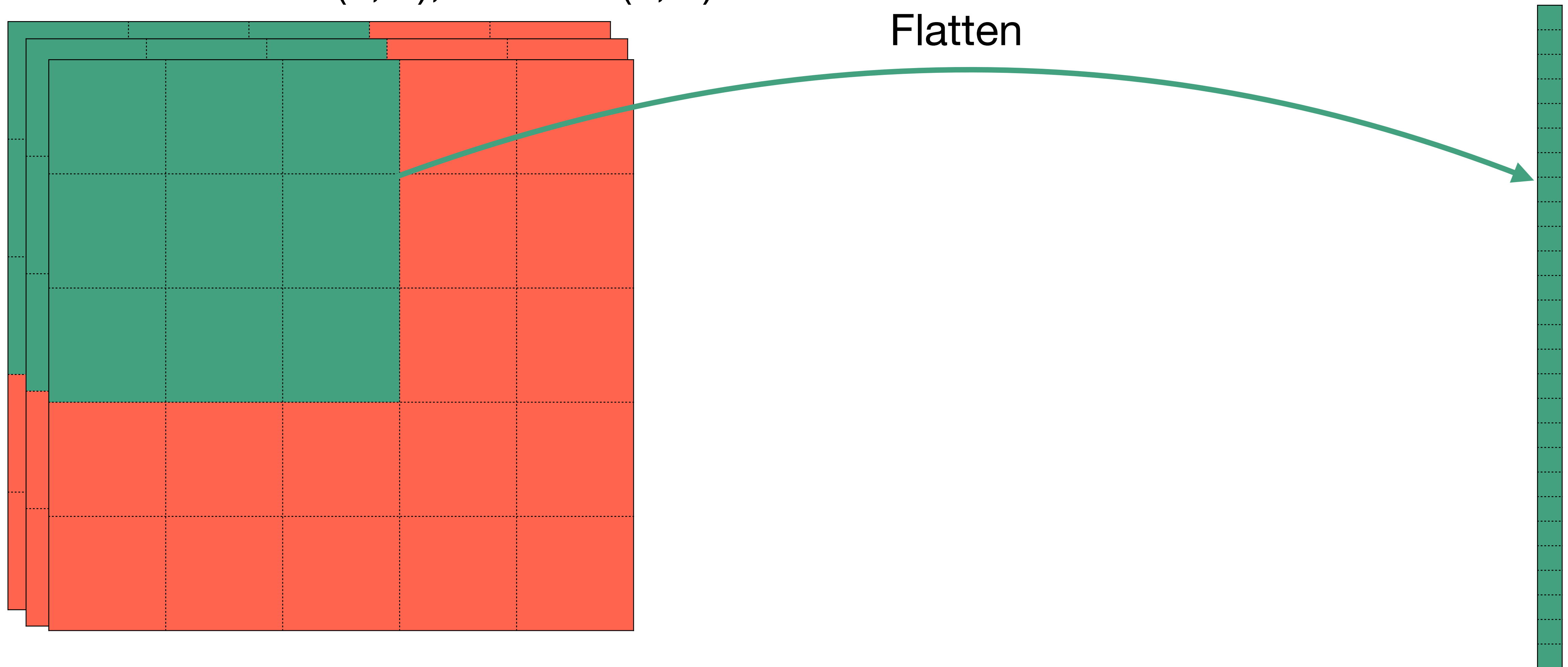
- A way to reduce the “depth” dimension of convolutional outputs
- Corresponds to applying an MLP to every pixel in the convolutional output
- Crucial to popular convnet architectures like Inception (GoogLeNet)

# Implementing the math

Conv2D. Input = (5, 5, 3)

Filters = 32 of size (3, 3), Stride = (1, 1)

$3 * 3 * 3 = 27$ -dimensional

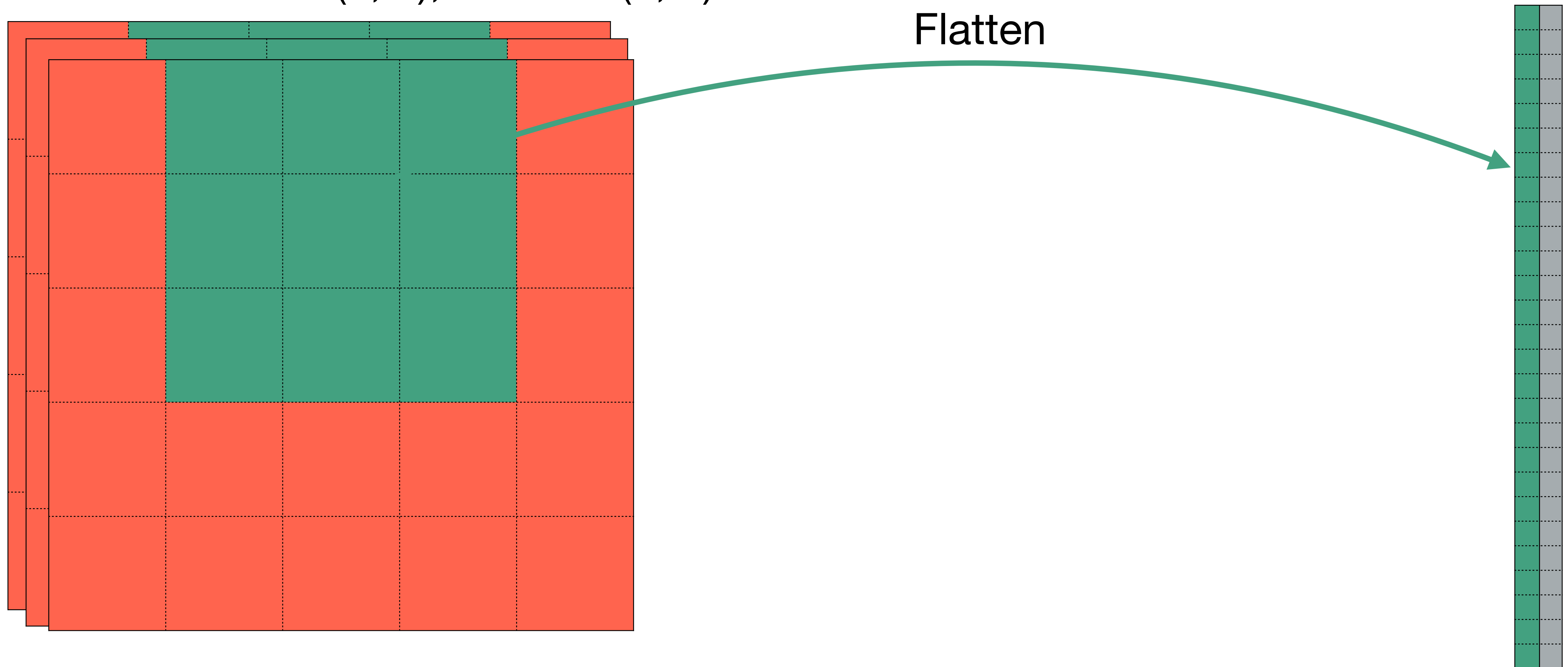


# Implementing the math

Conv2D. Input = (5, 5, 3)

Filters = 32 of size (3, 3), Stride = (1, 1)

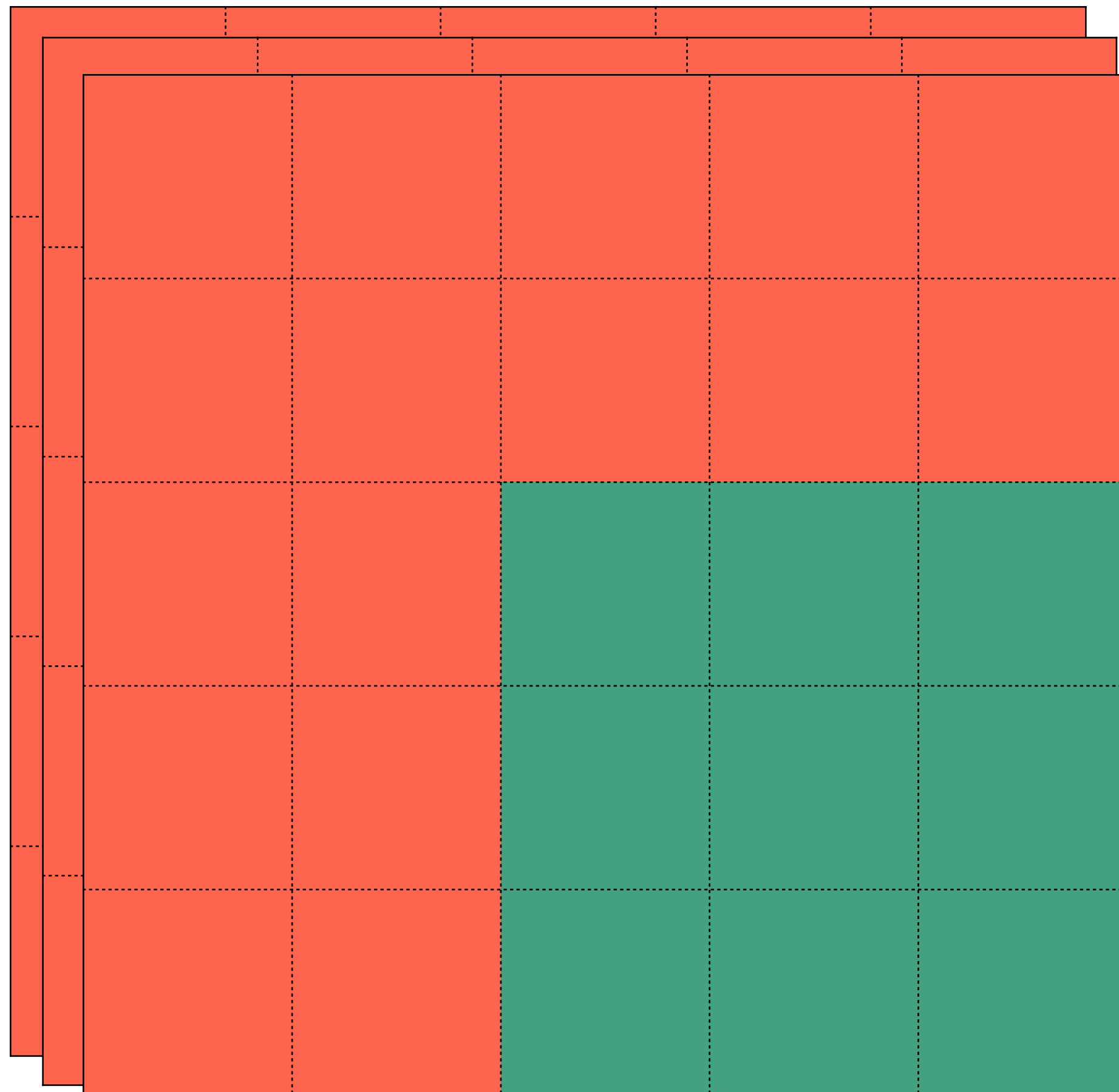
$3 * 3 * 3 = 27$ -dimensional



# Implementing the math

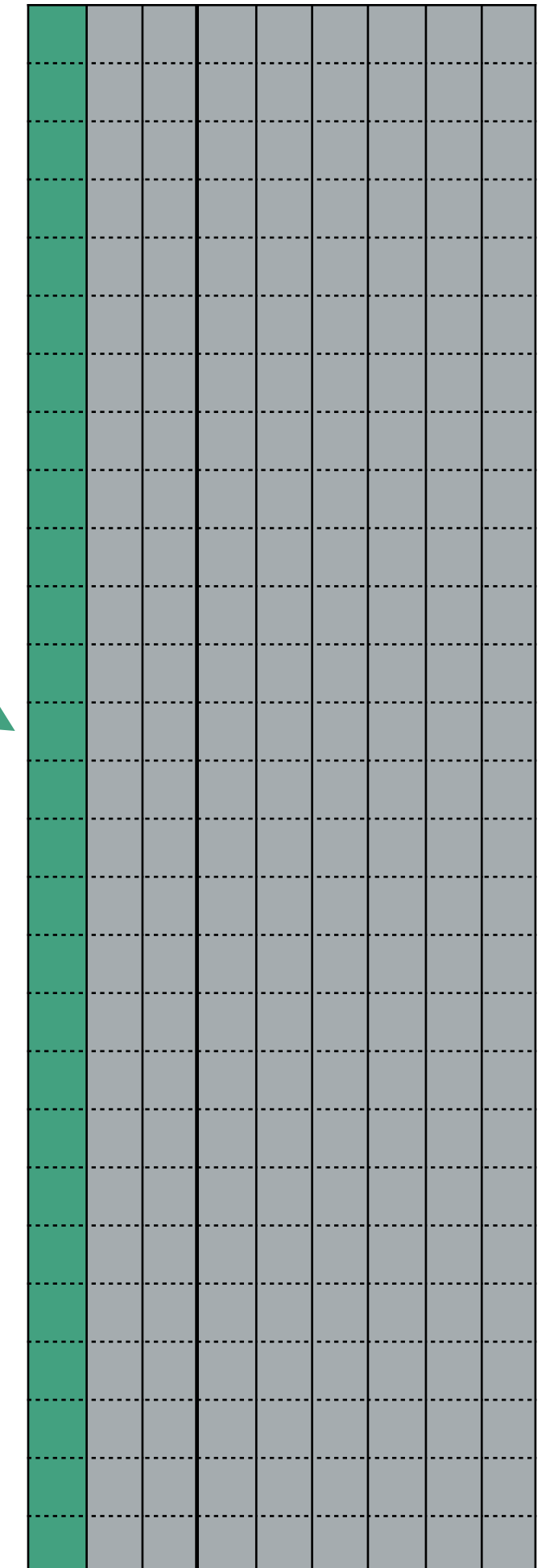
Conv2D. Input = (5, 5, 3)

Filters = 32 of size (3, 3), Stride = (1, 1)



Flatten

X\_col  
(27 x 9)



# Implementing the math

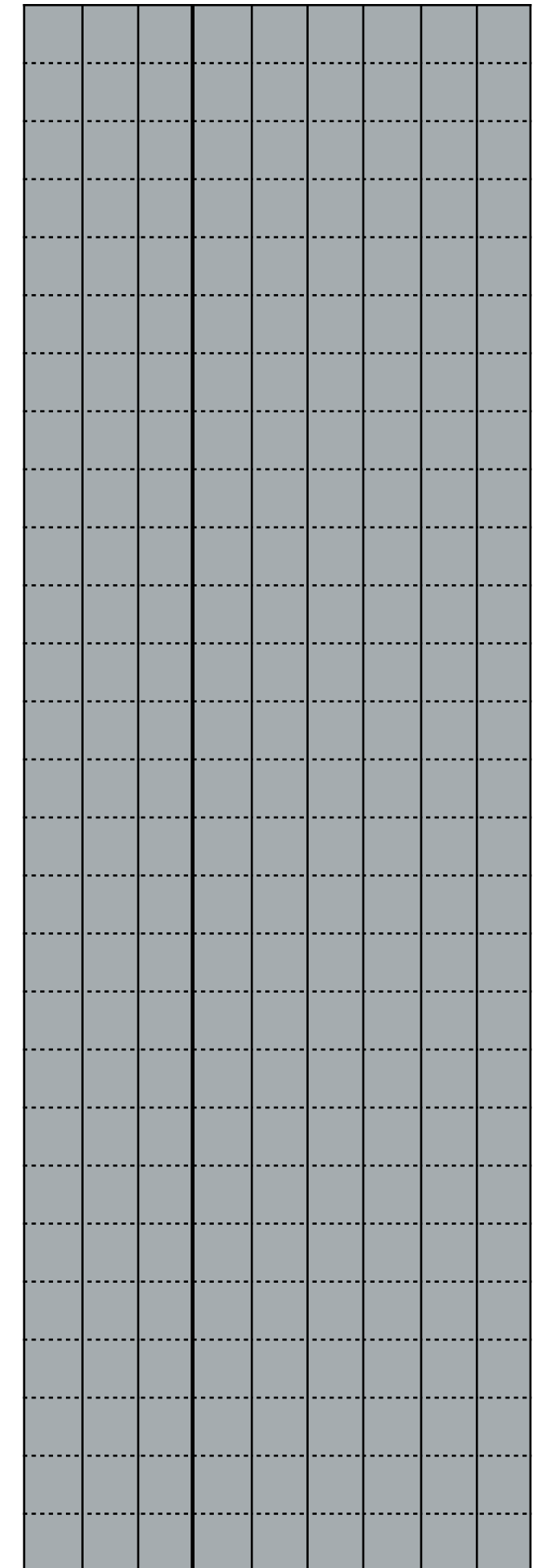
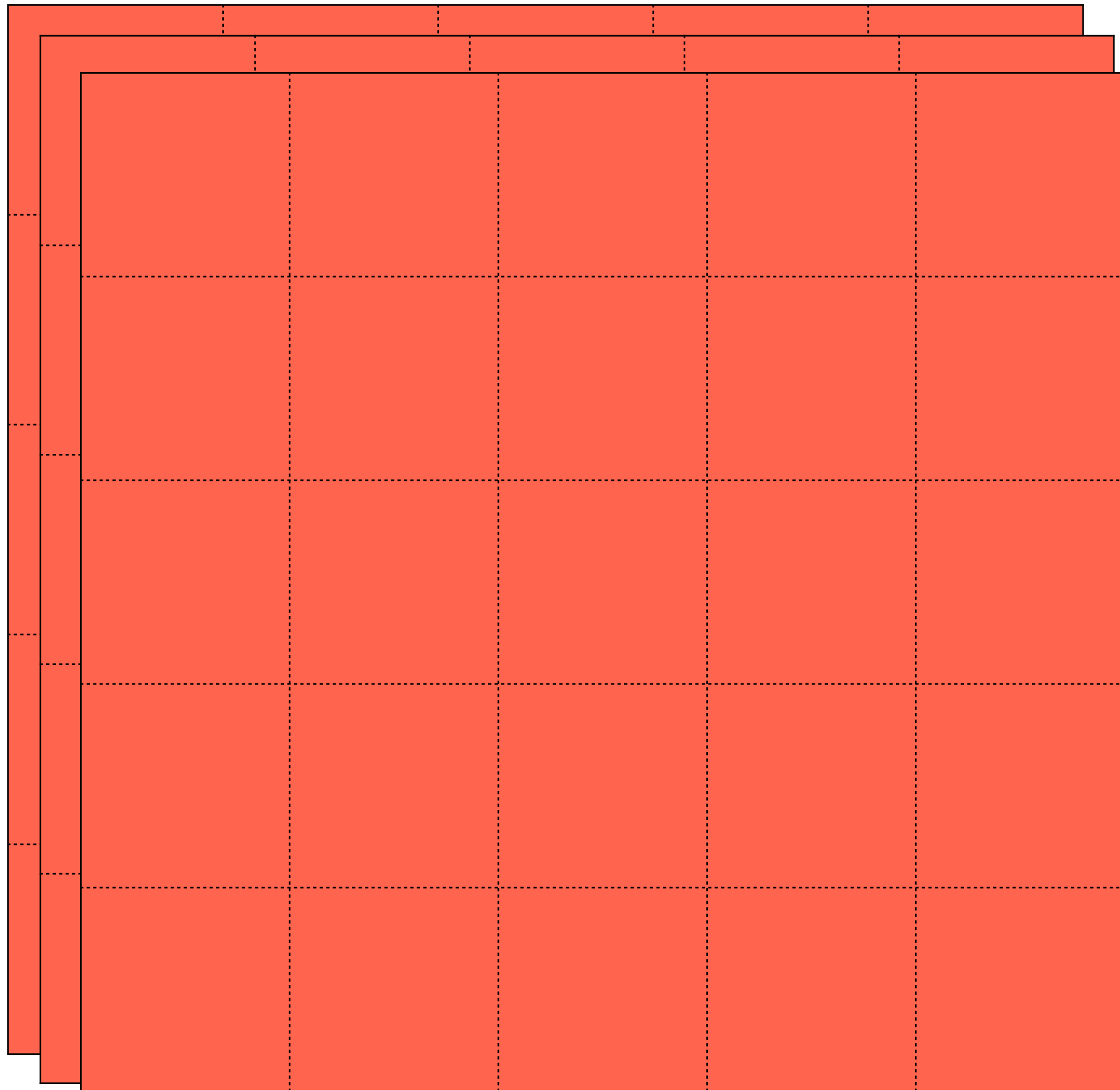
Conv2D. Input = (5, 5, 3)

**Filters = 32 of size (3, 3), Stride = (1, 1)**



W\_row  
(32 x 27)

X\_col  
(27 x 9)

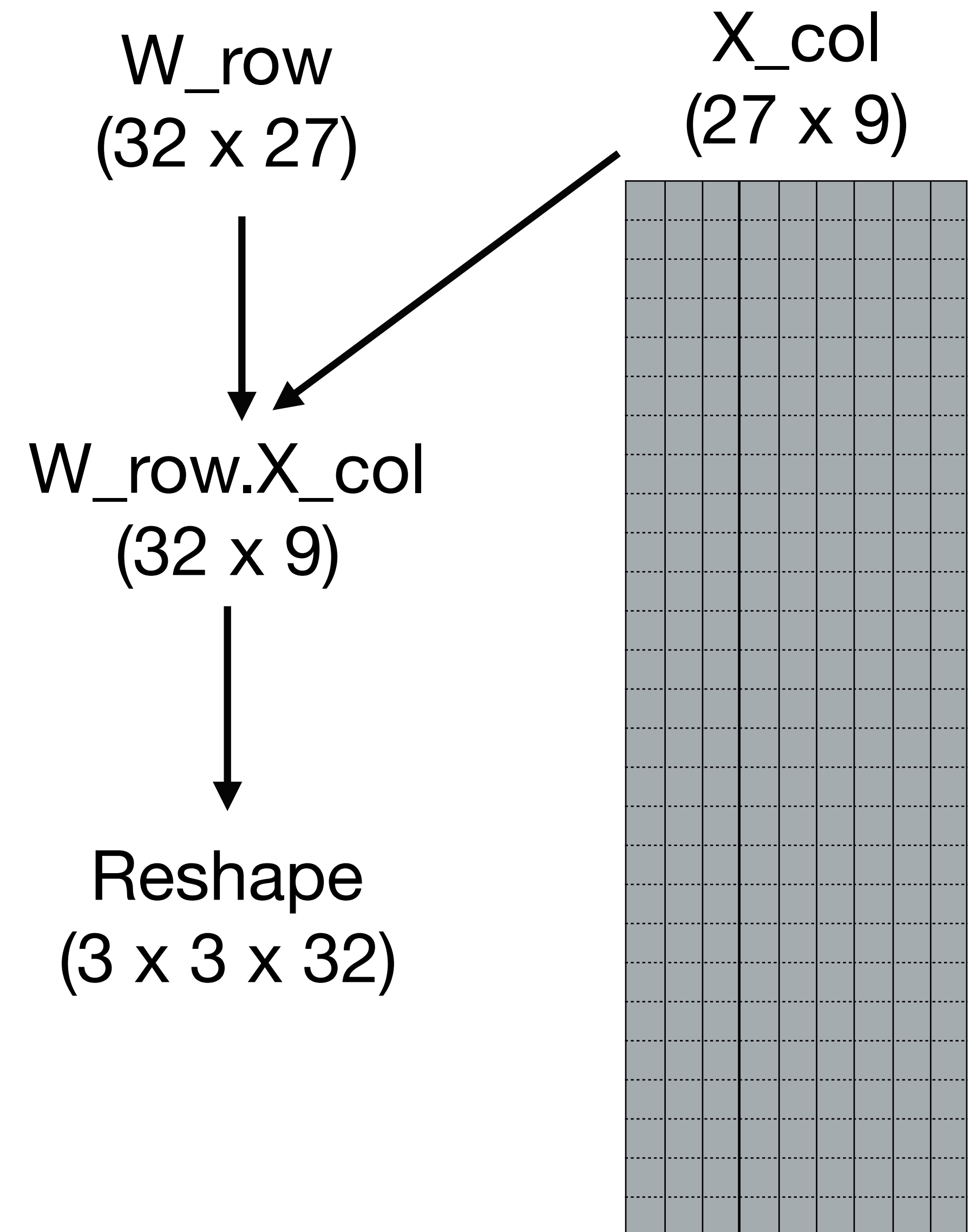
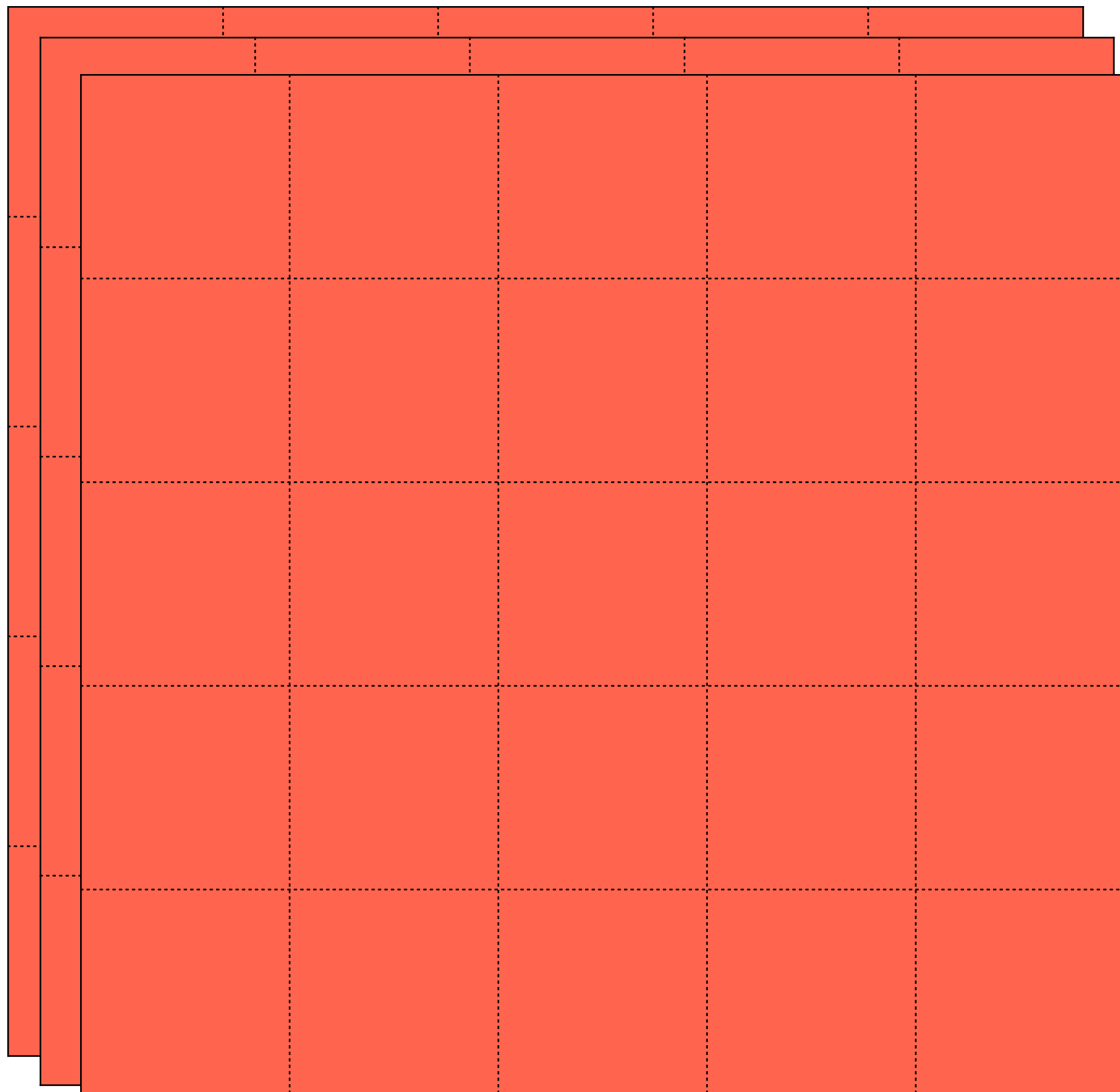




# Implementing the math

Conv2D. Input = (5, 5, 3)

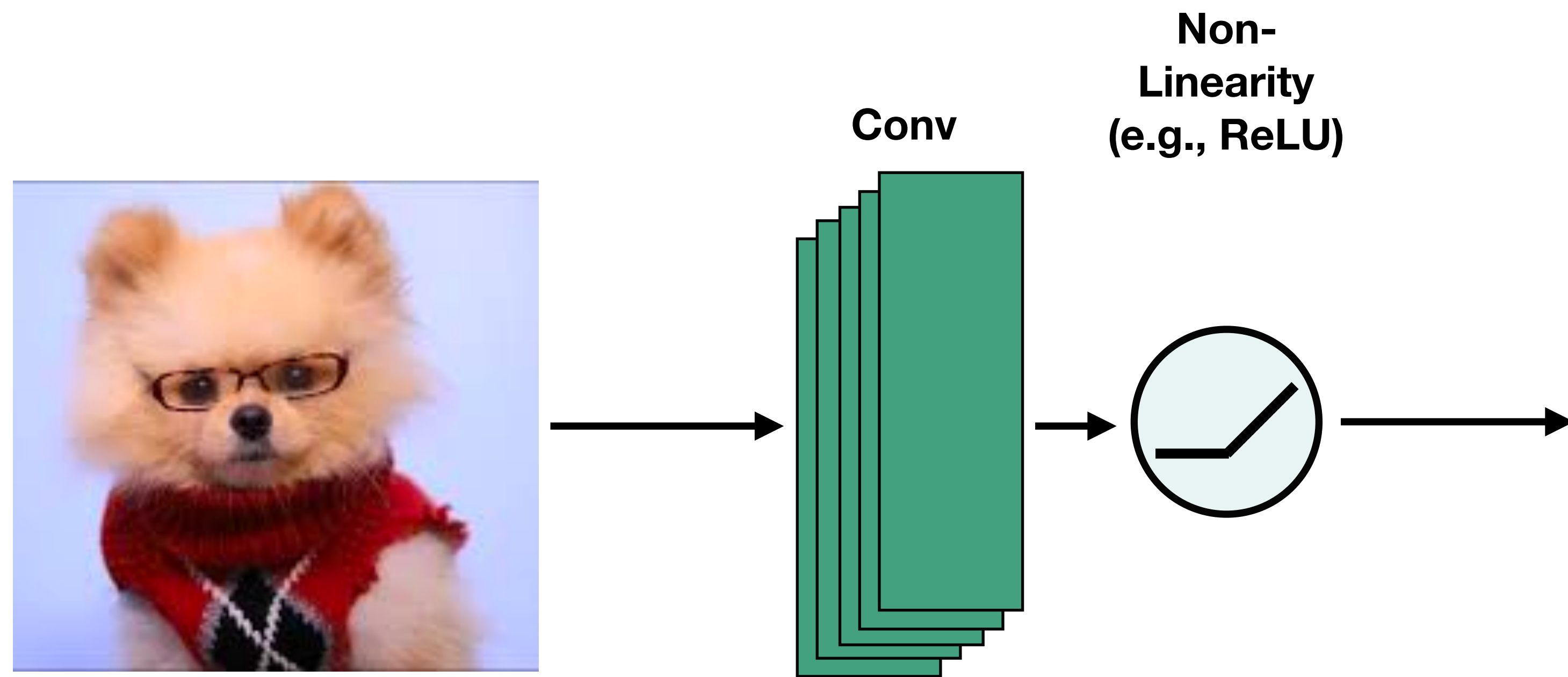
**Filters = 32 of size (3, 3), Stride = (1, 1)**



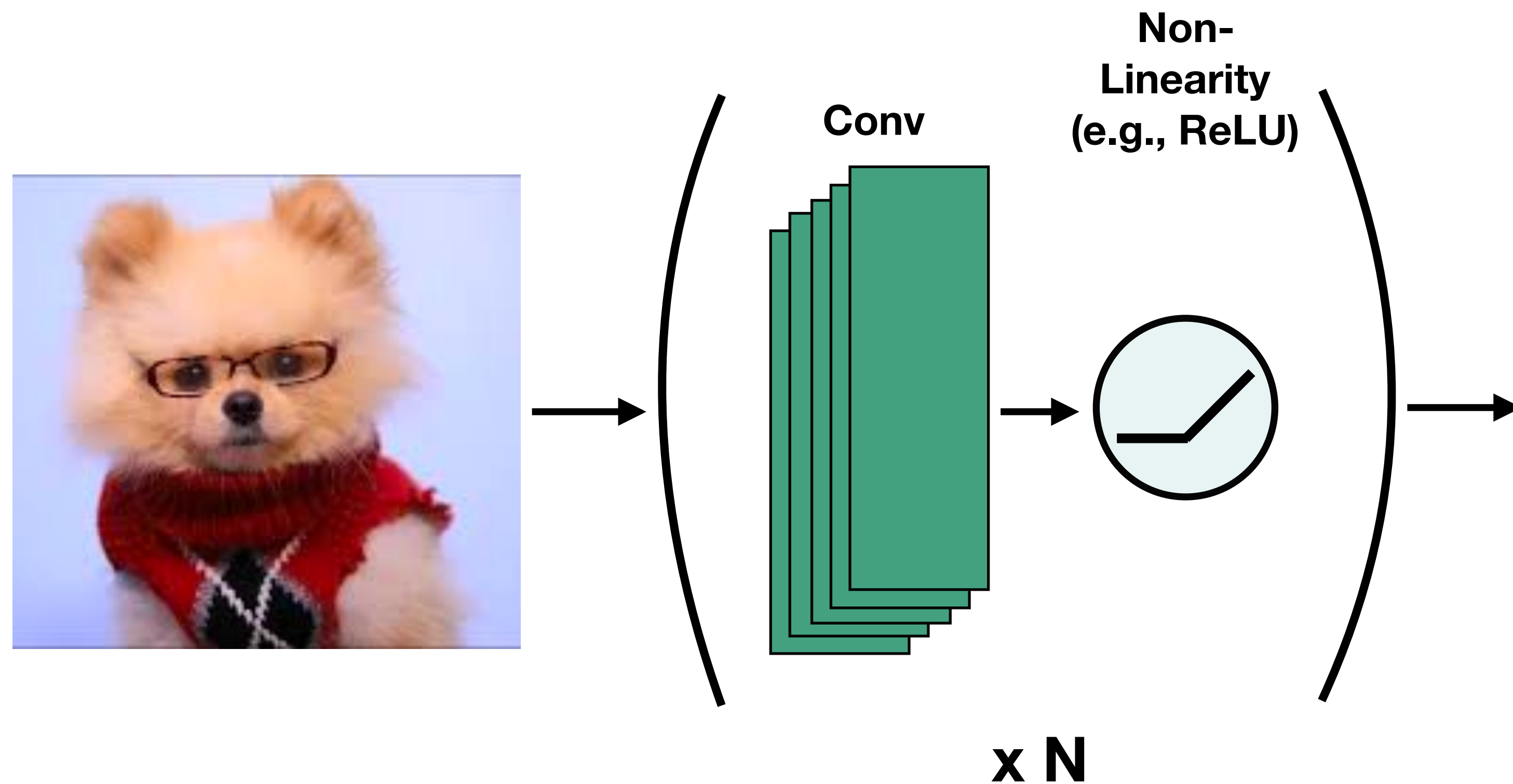
# Agenda

1. Basics of convolutions, and filter/stride/pooling math
2. More advanced convolution types, and computational considerations
3. **Classic ConvNet architectures**

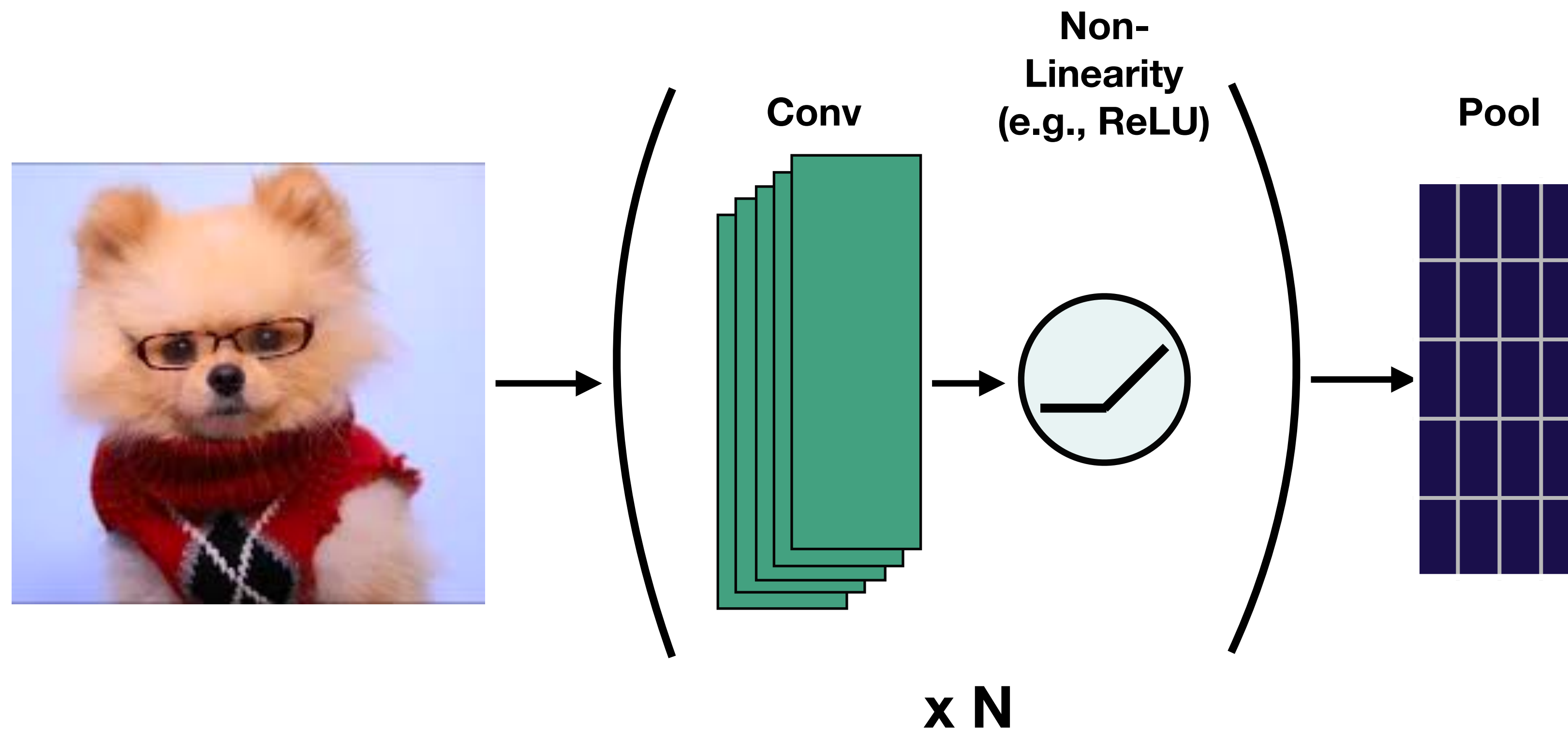
# Most standard: LeNet(-like) architectures



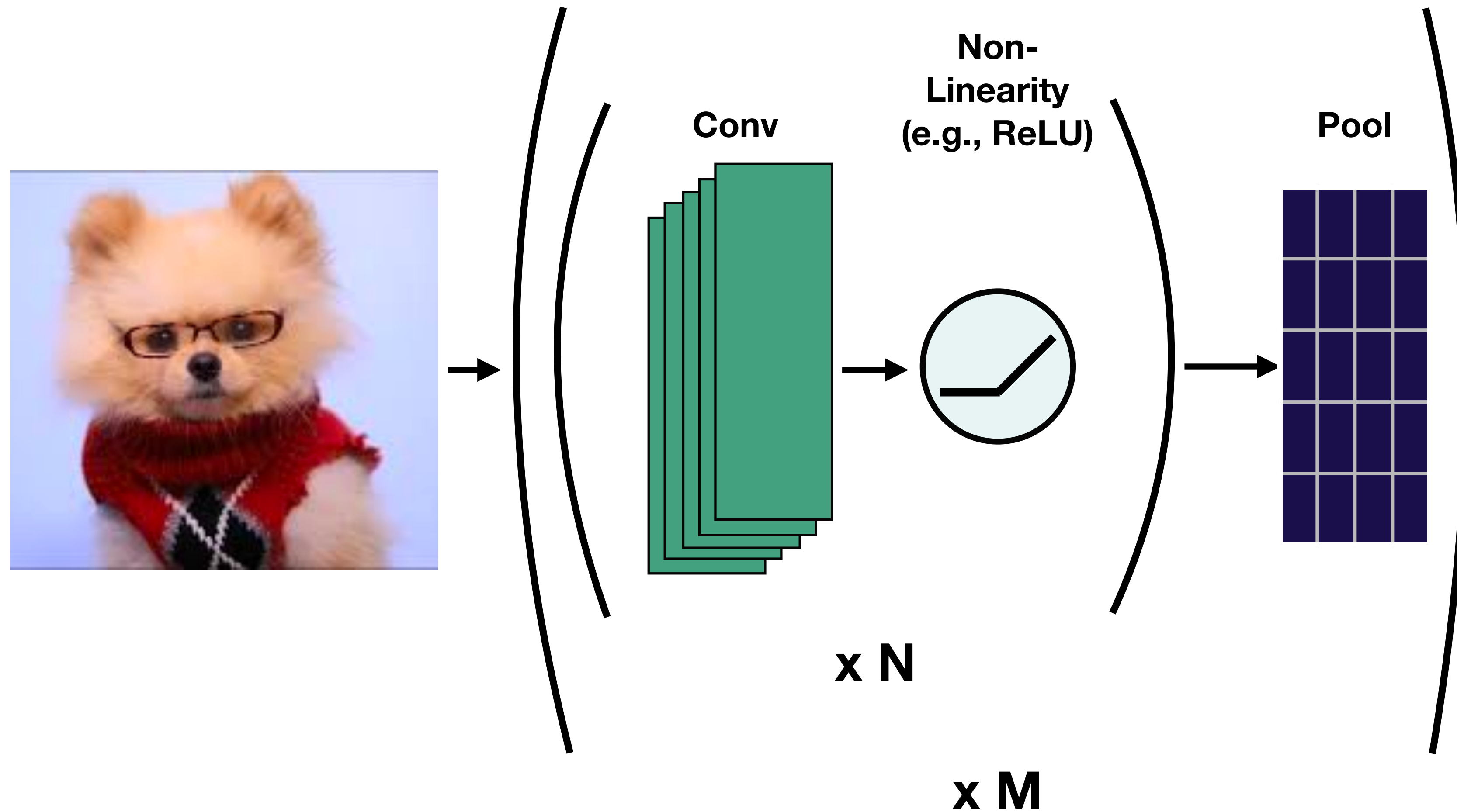
# Most standard: LeNet(-like) architectures



# Most standard: LeNet(-like) architectures

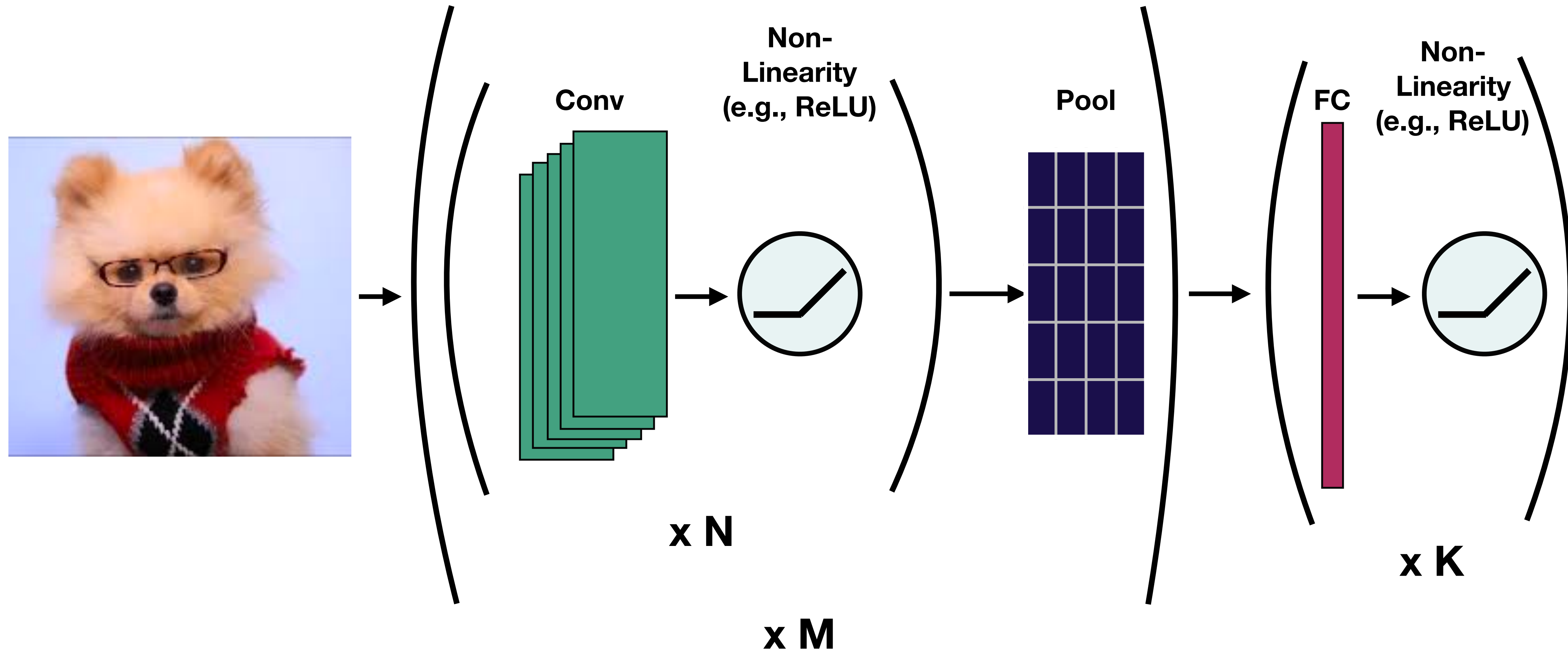


# Most standard: LeNet(-like) architectures

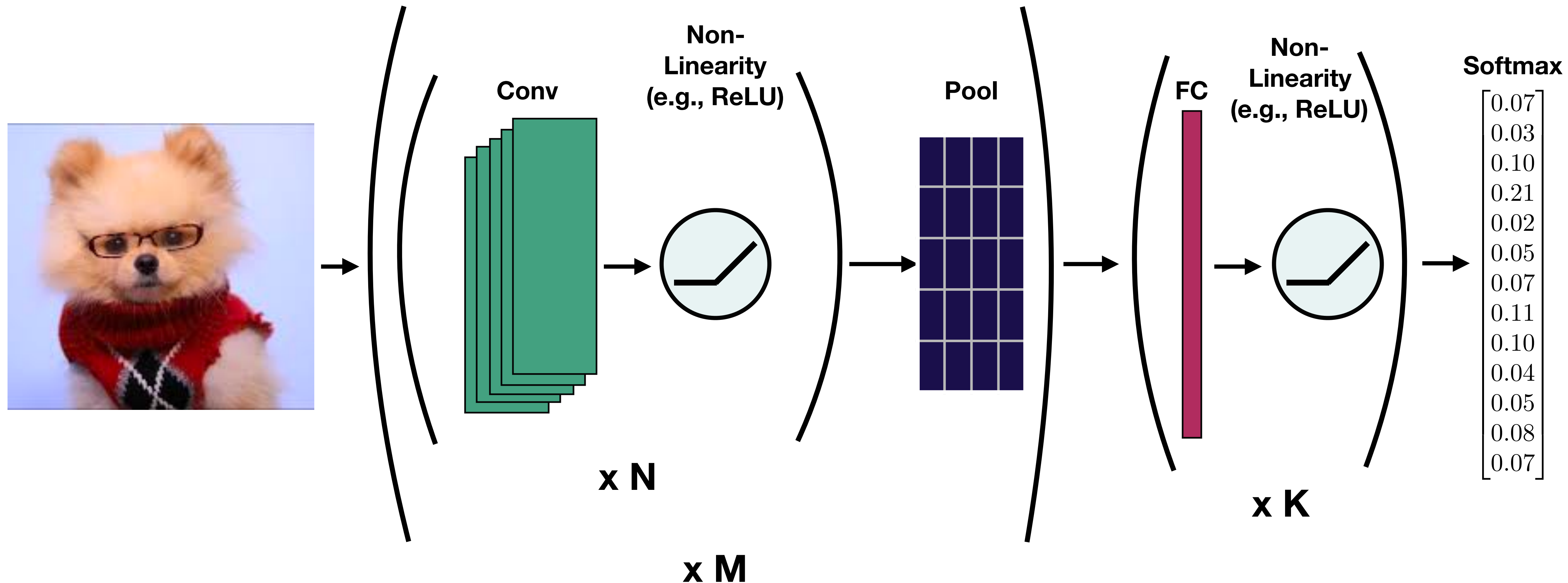




# Most standard: LeNet(-like) architectures

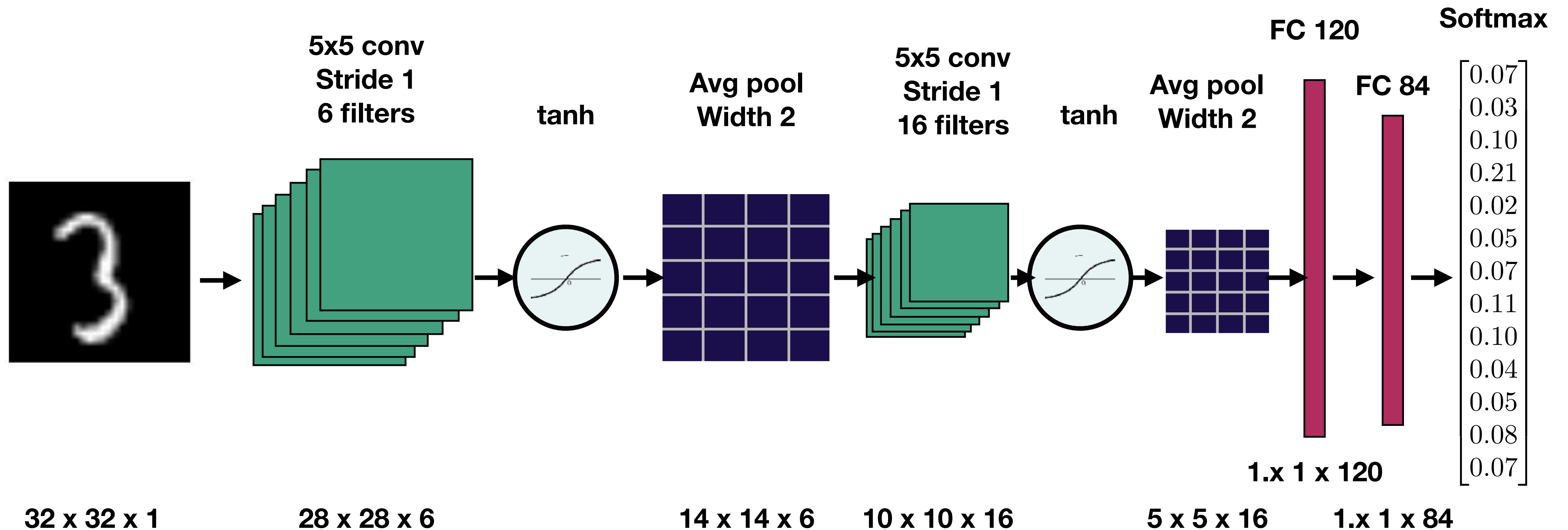


# Most standard: LeNet(-like) architectures

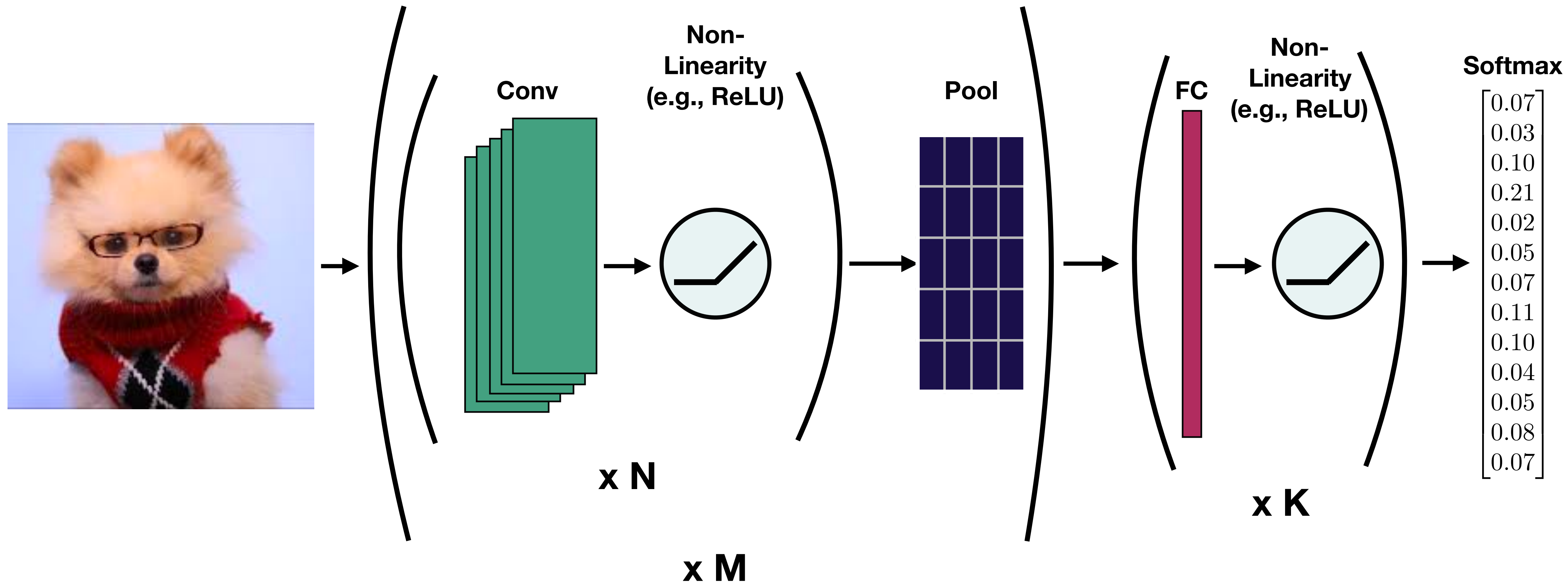




# Classic Convnet Architecture: LeNet



# Most standard: LeNet(-like) architectures



# More modern LeNet-like architectures

- $N = \text{up to } 5$
- $M$  large
- $0 \leq K \leq 2$
- ReLU instead of Tanh

# Where to go to learn more?

- Stanford's CS231n