

Lecture 2.A

Convolutional Networks

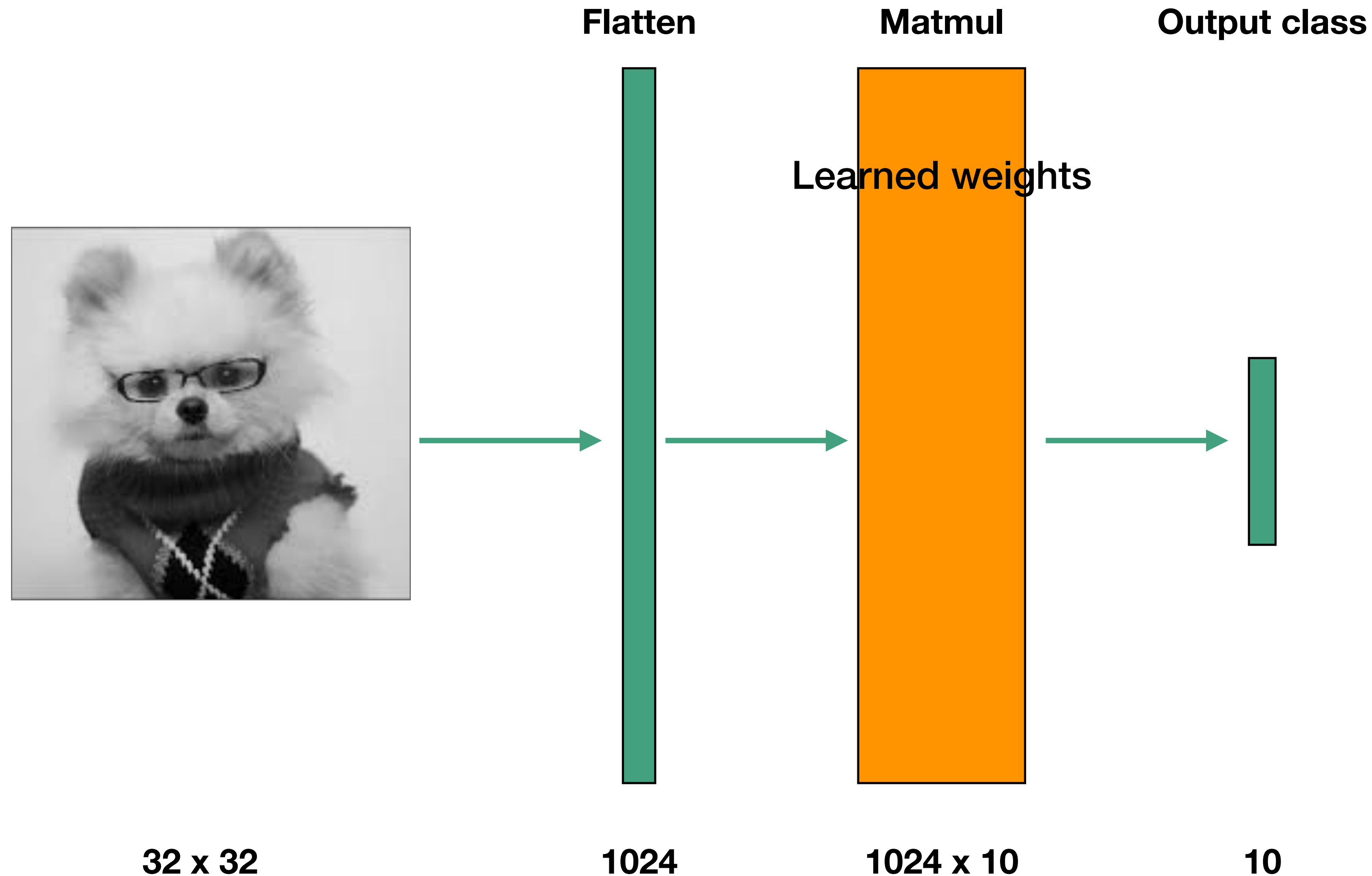
Agenda

- 1. Review of the convolution operation**
2. Other important operations for ConvNets
3. Classic ConvNet architectures

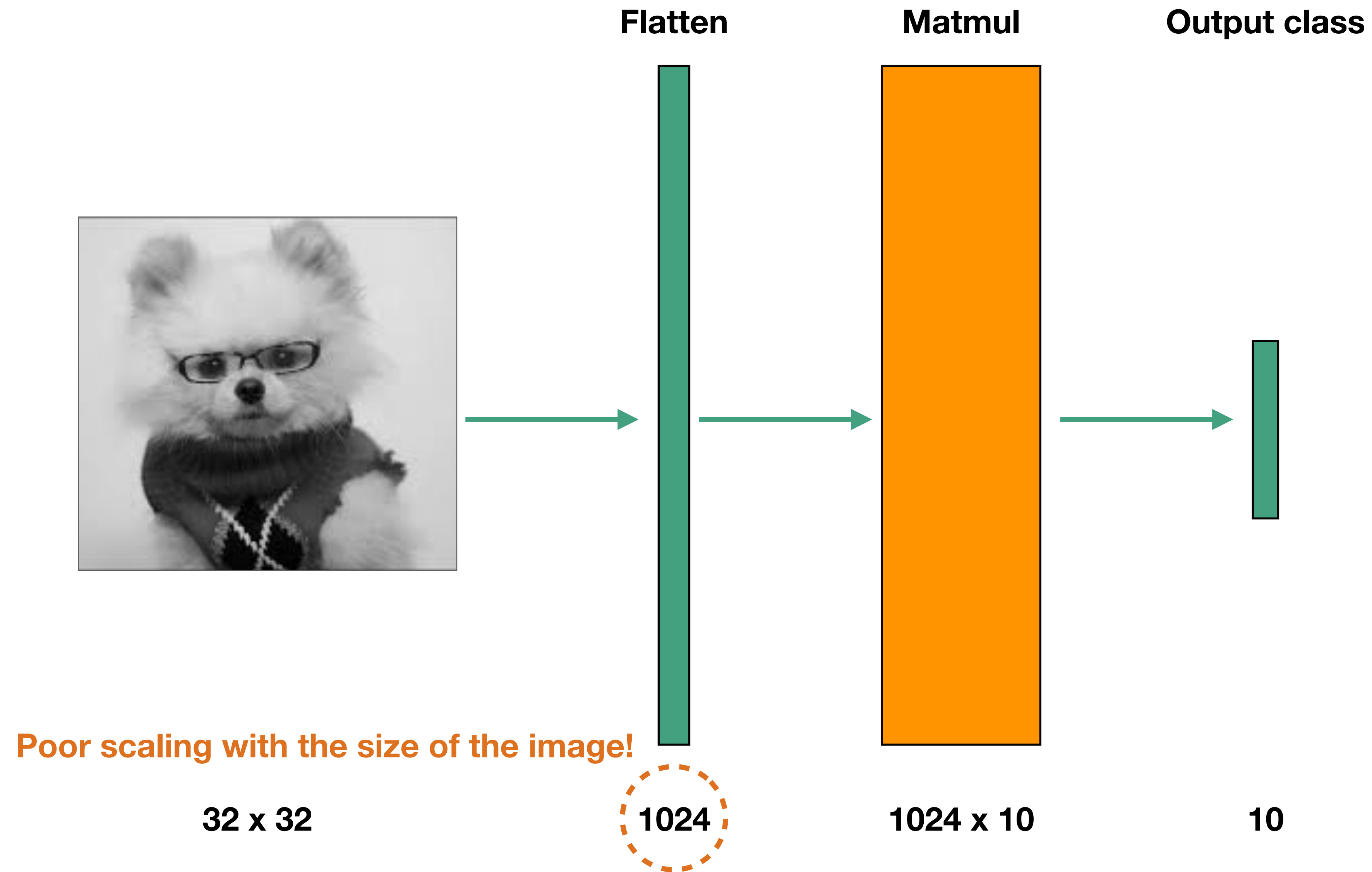
Review of convolutions

- **What's a convolutional filter?**
- Filter stacks and ConvNets
- Strides & padding
- Filter math
- Implementation notes

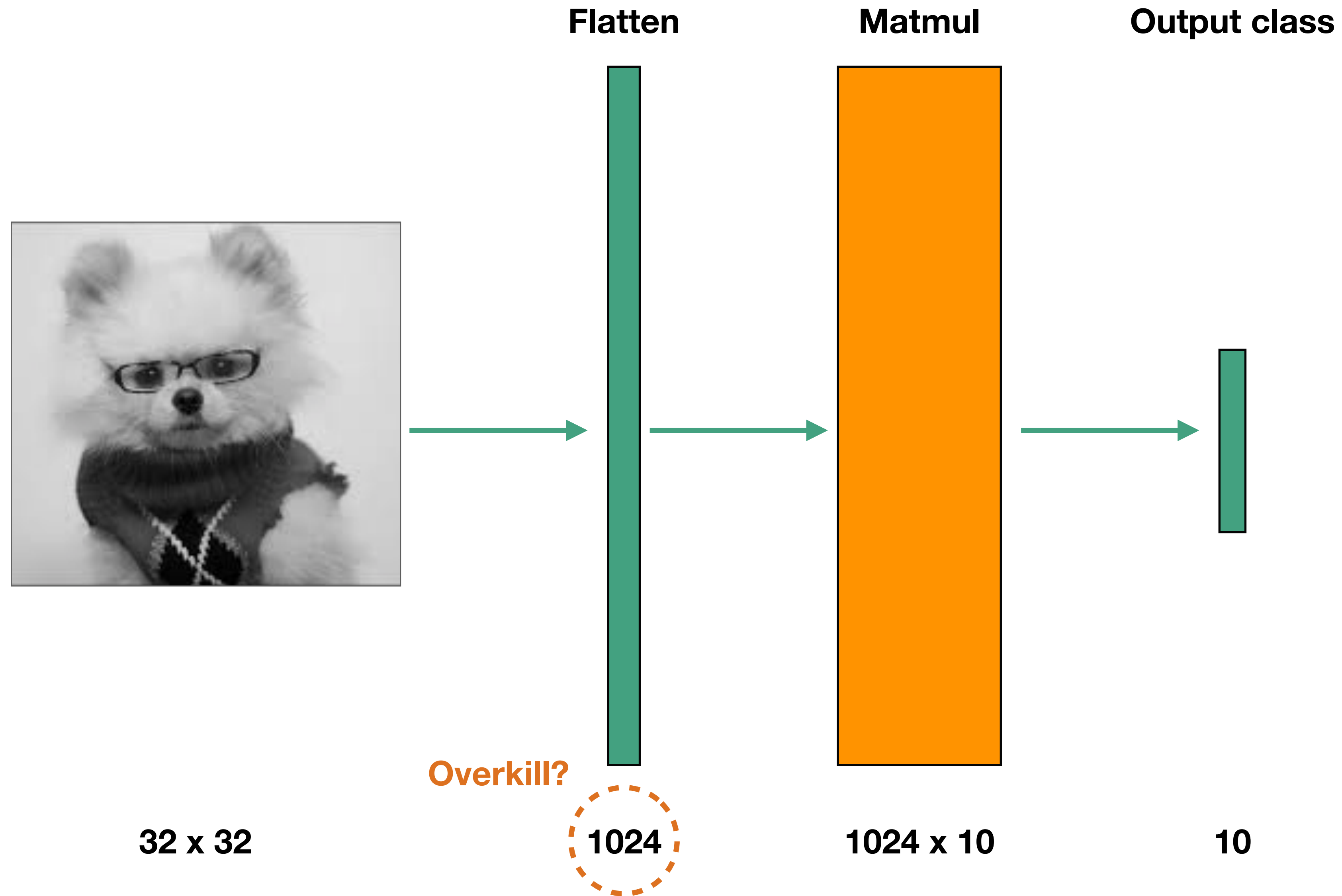
Why not use FC for images?



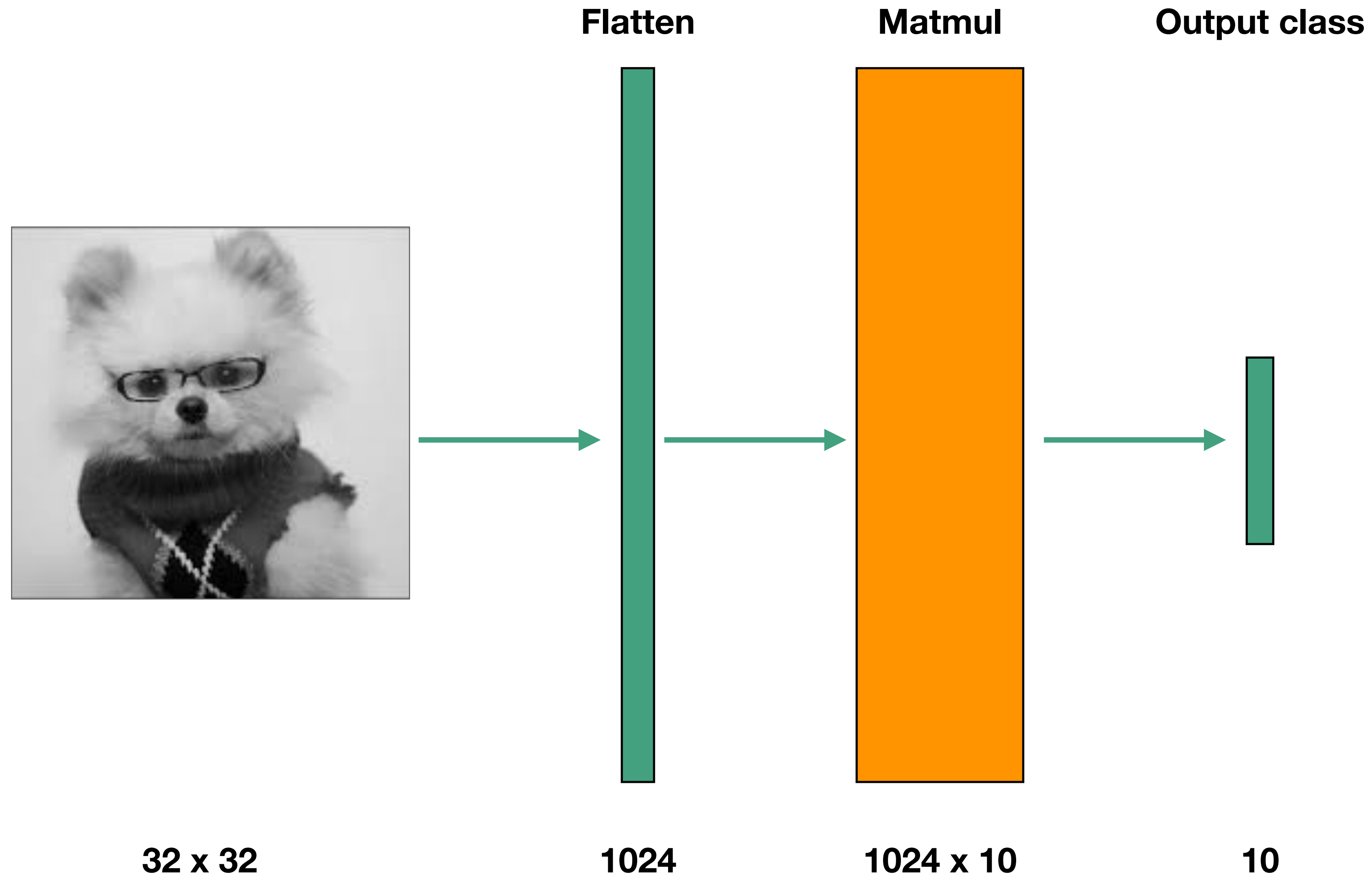
Why not use FC for images?



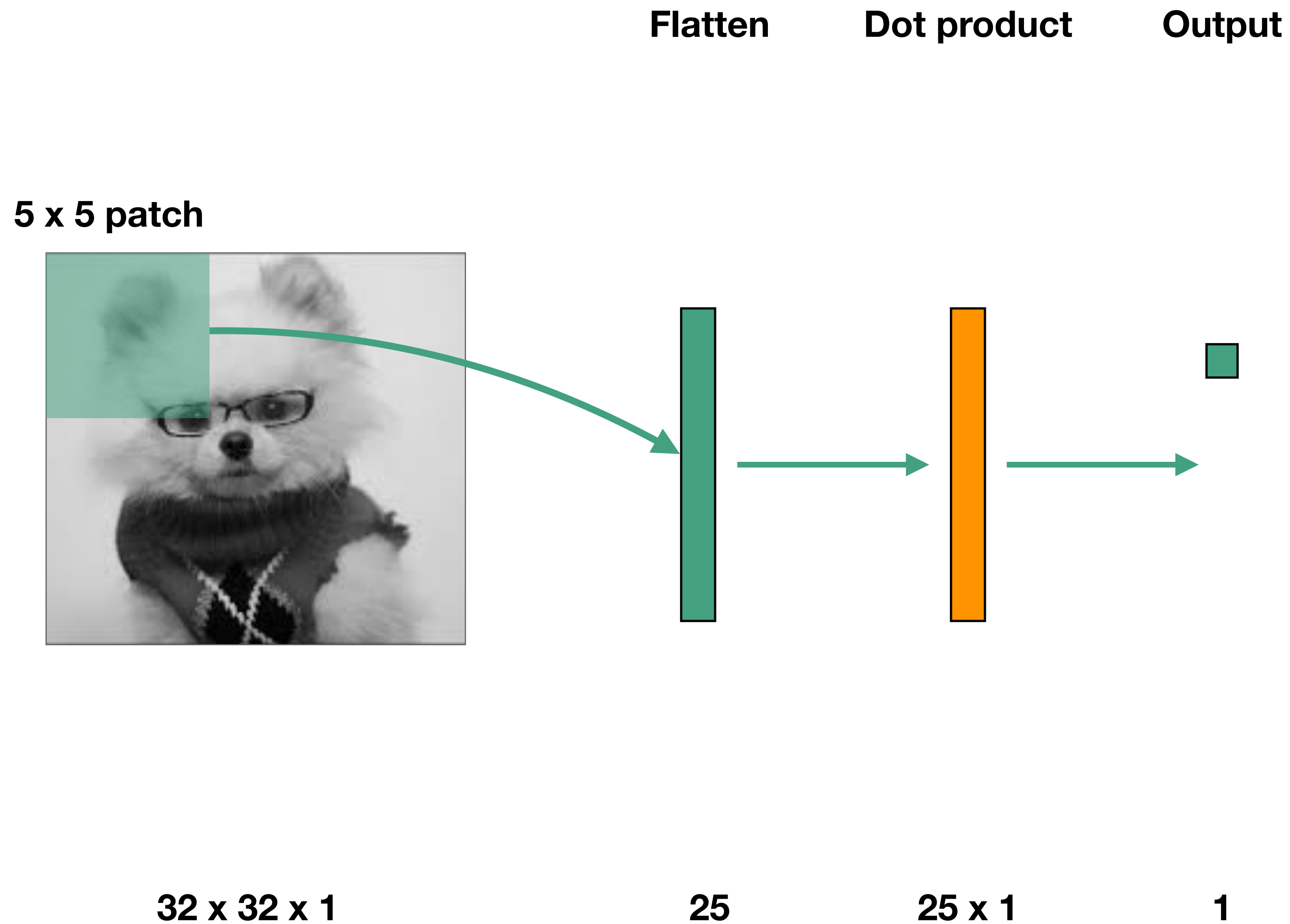
Why not use FC for images?



Why not use FC for images?



Convolutional filters



Convolutional filters

Flatten

Dot product

Output

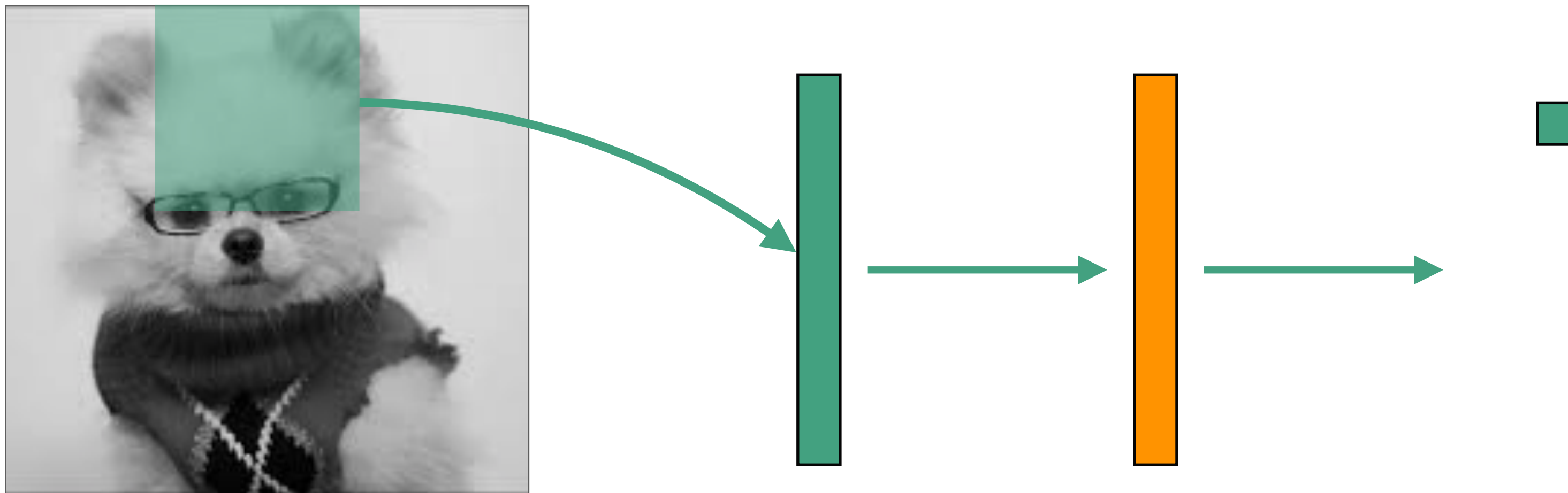


Convolutional filters

Flatten

Dot product

Output

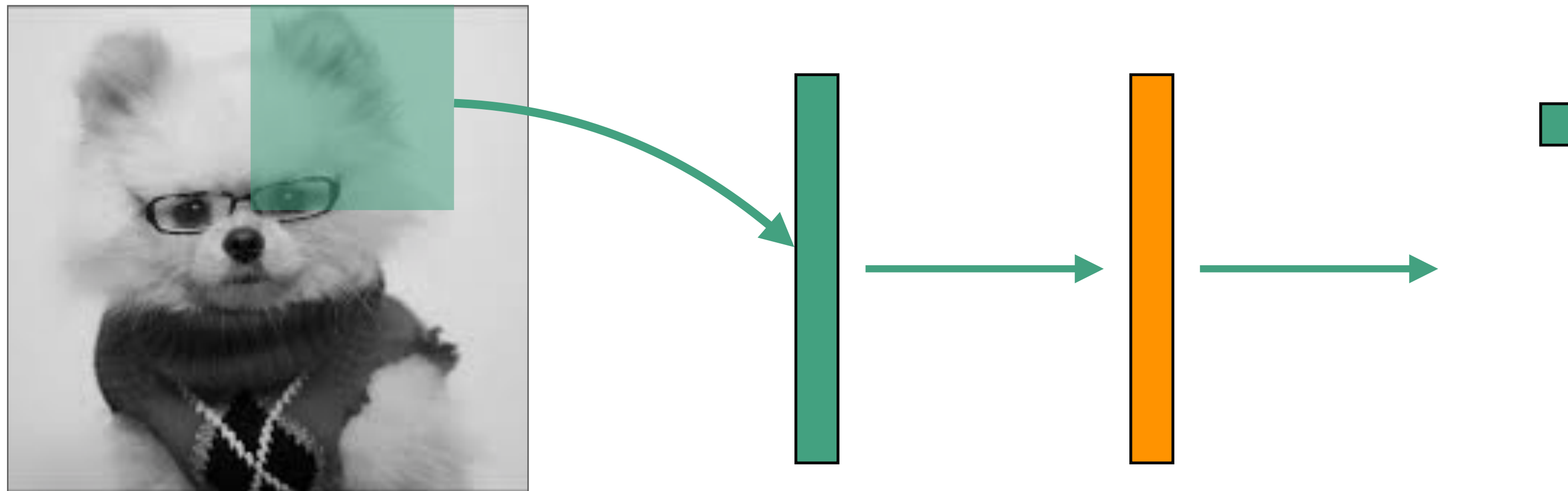


Convolutional filters

Flatten

Dot product

Output

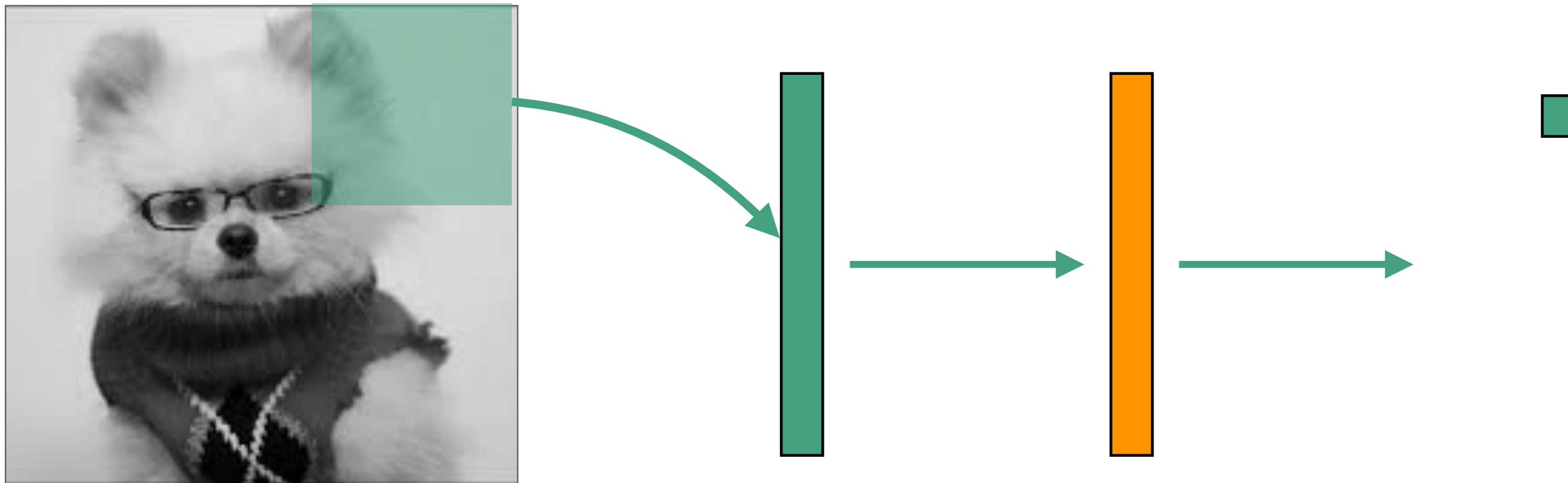


Convolutional filters

Flatten

Dot product

Output



Convolutional filters

Flatten

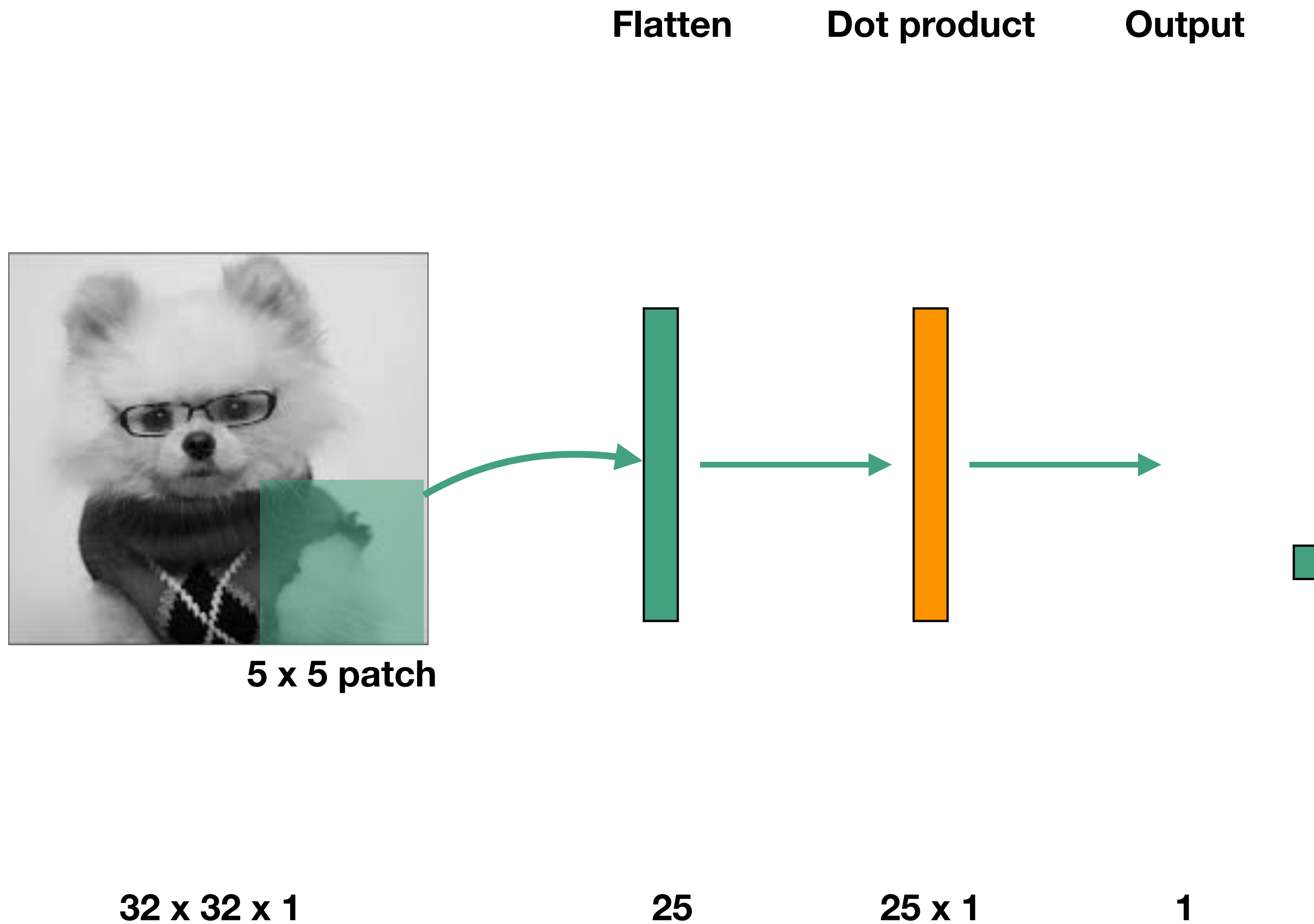
Dot product

Output



...sliding continues...

Convolutional filters

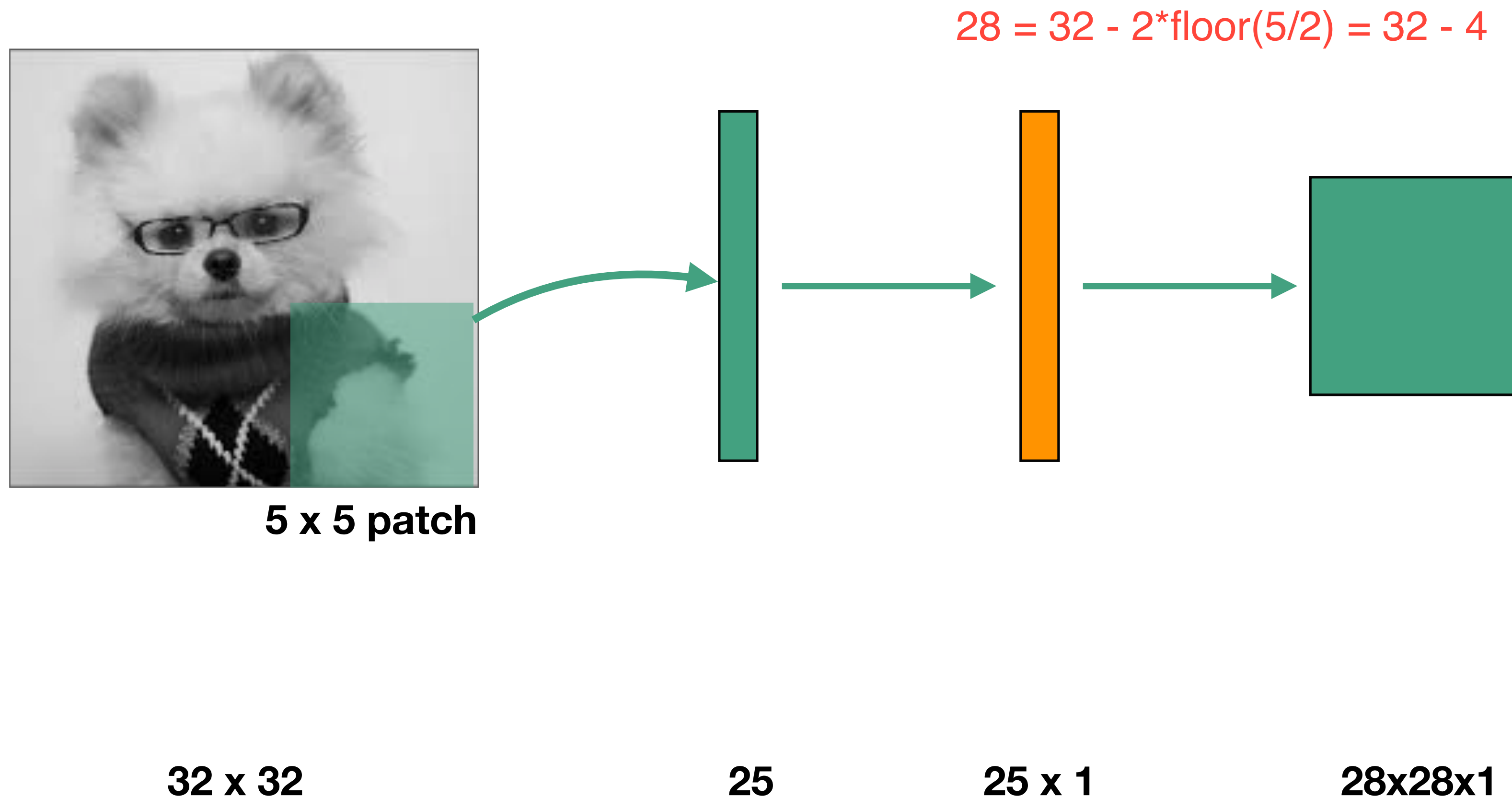


Convolutional filters

Flatten


Dot product

Output



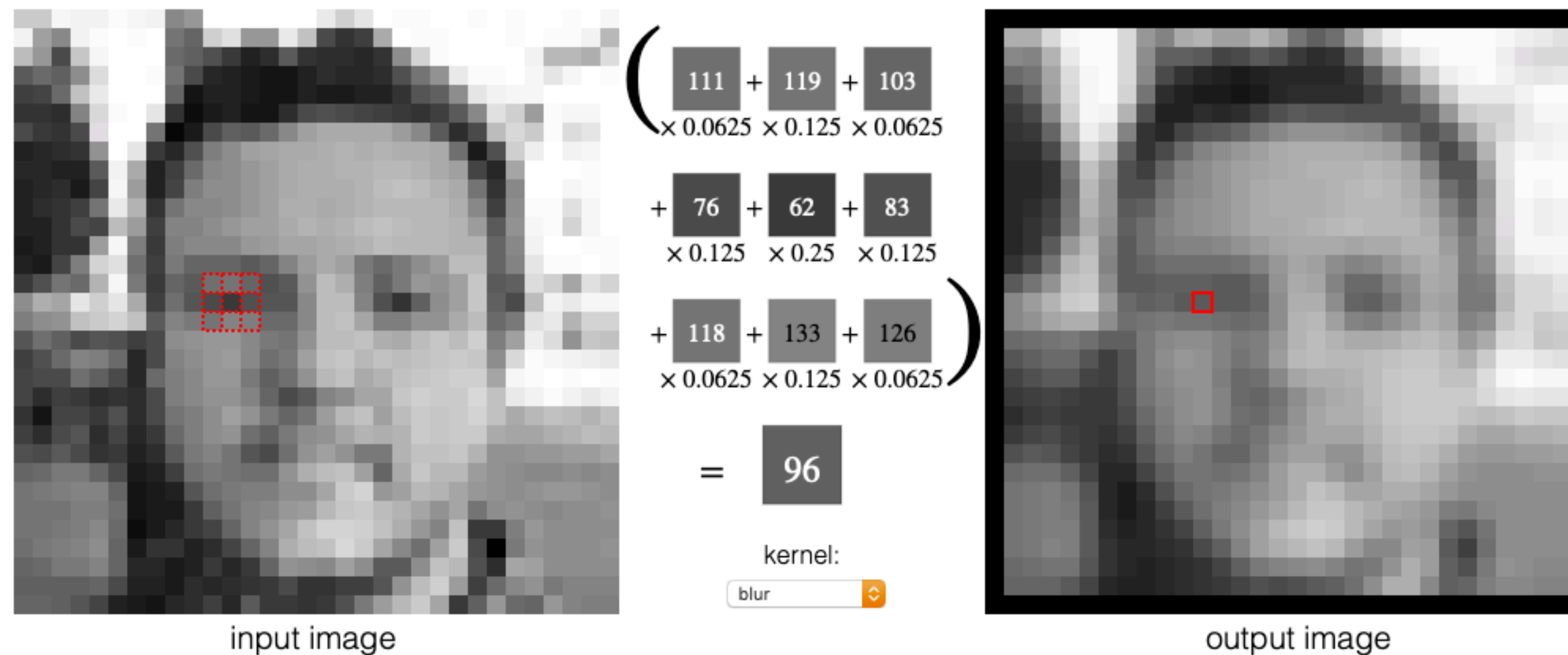
What can a conv filter do?

Let's walk through applying the following 3x3 **blur** kernel to the image of a face from above.

blur 

$$\begin{pmatrix} 0.0625 & 0.125 & 0.0625 \\ 0.125 & 0.25 & 0.125 \\ 0.0625 & 0.125 & 0.0625 \end{pmatrix} \leftarrow \text{Instead of hard-coding the weights, we can learn them!}$$

Below, for each 3x3 block of pixels in the image on the left, we multiply each pixel by the corresponding entry of the kernel and then take the sum. That sum becomes a new pixel in the image on the right. Hover over a pixel on either image to see how its value is computed.

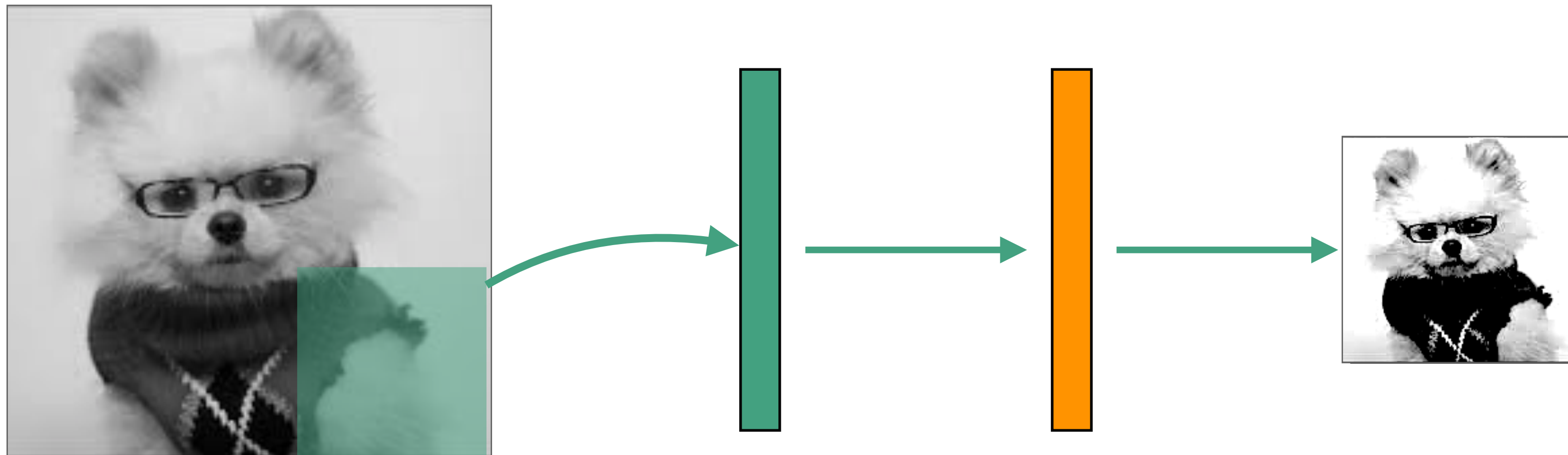


Convolutional filters

Flatten

Dot product

Output



32 x 32 x 3

75

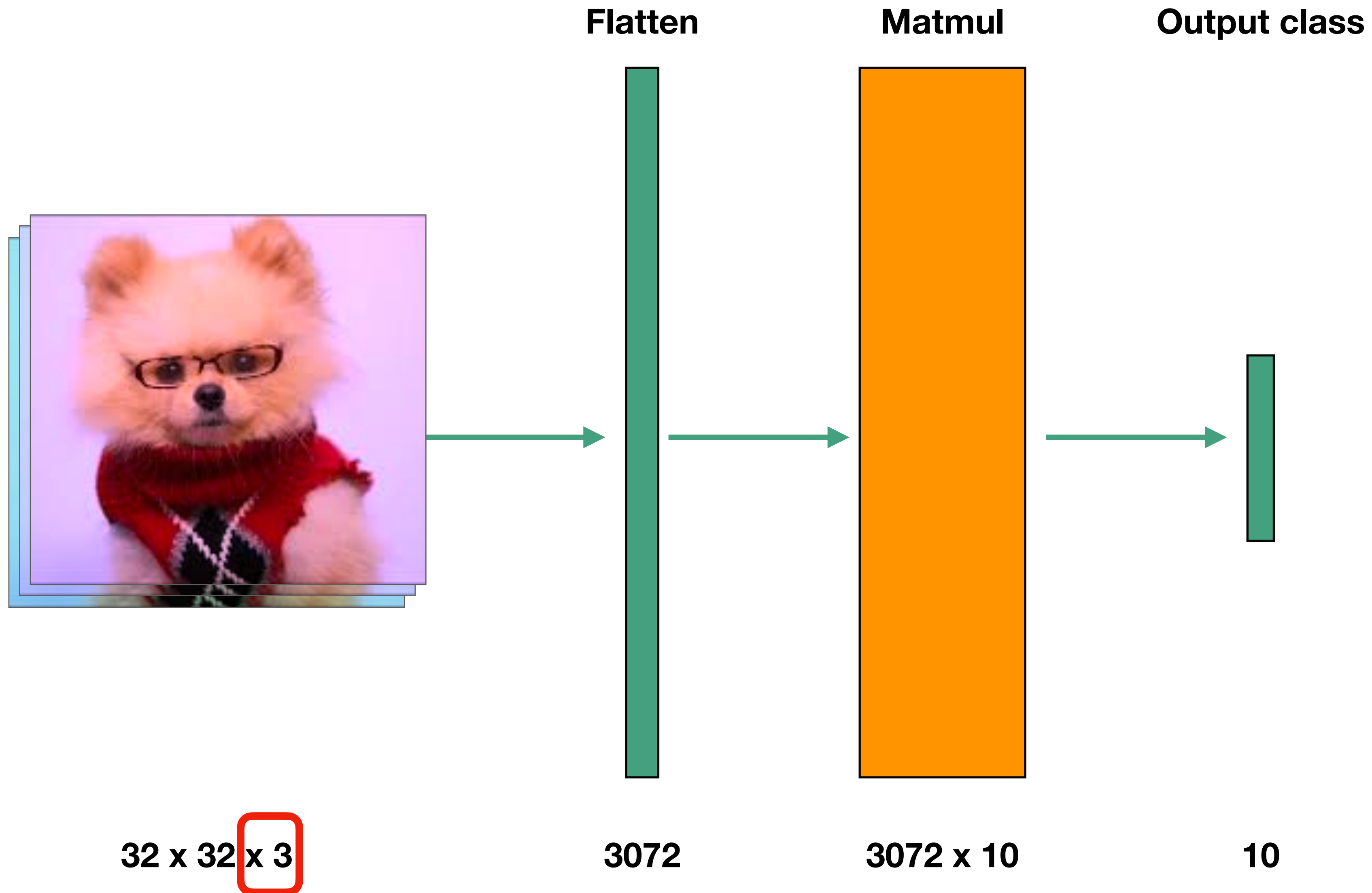
75 x 1

28x28x1

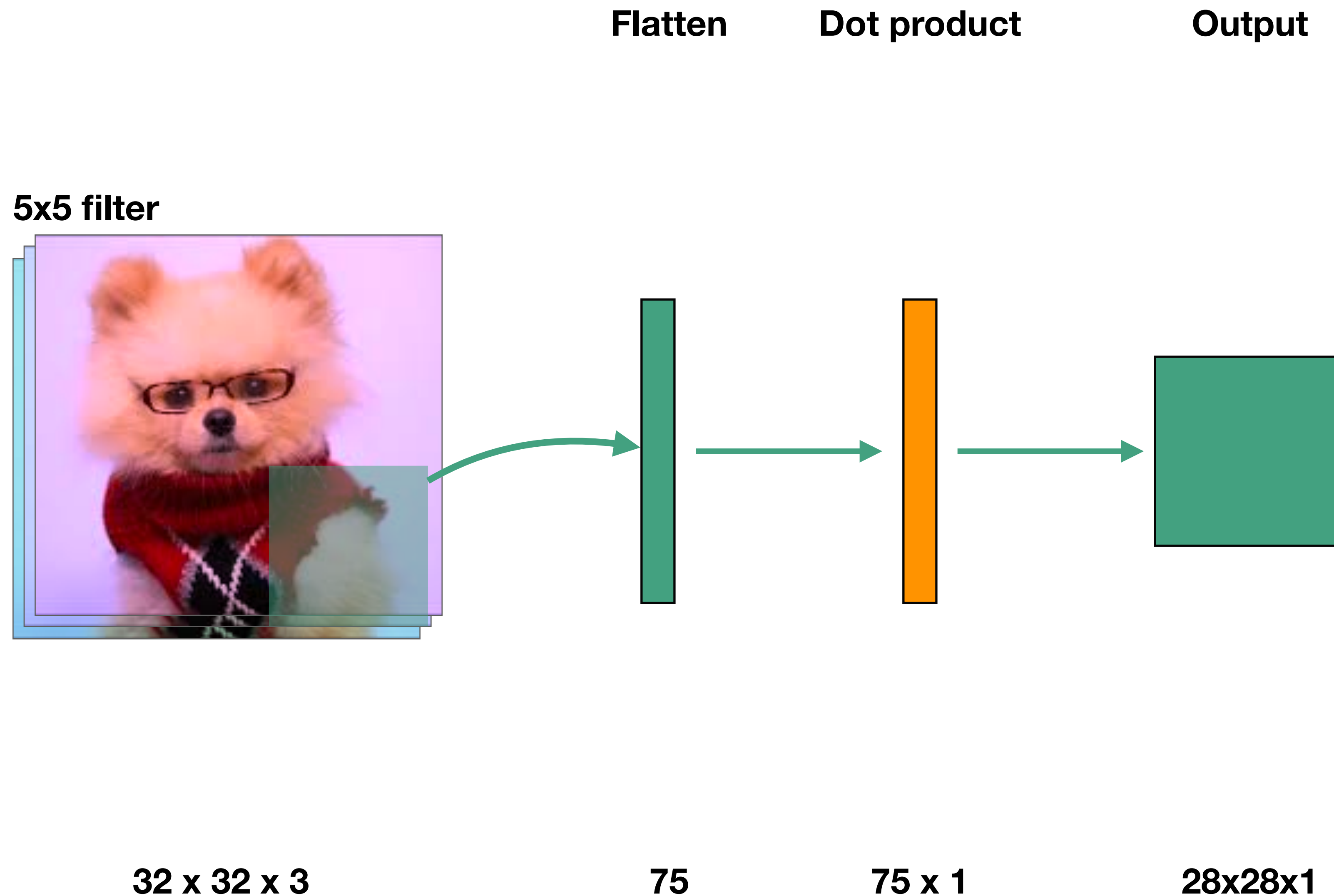
Review of convolutions

- What's a convolutional filter?
- **Filter stacks and ConvNets**
- Strides & padding
- Filter math
- Implementation notes

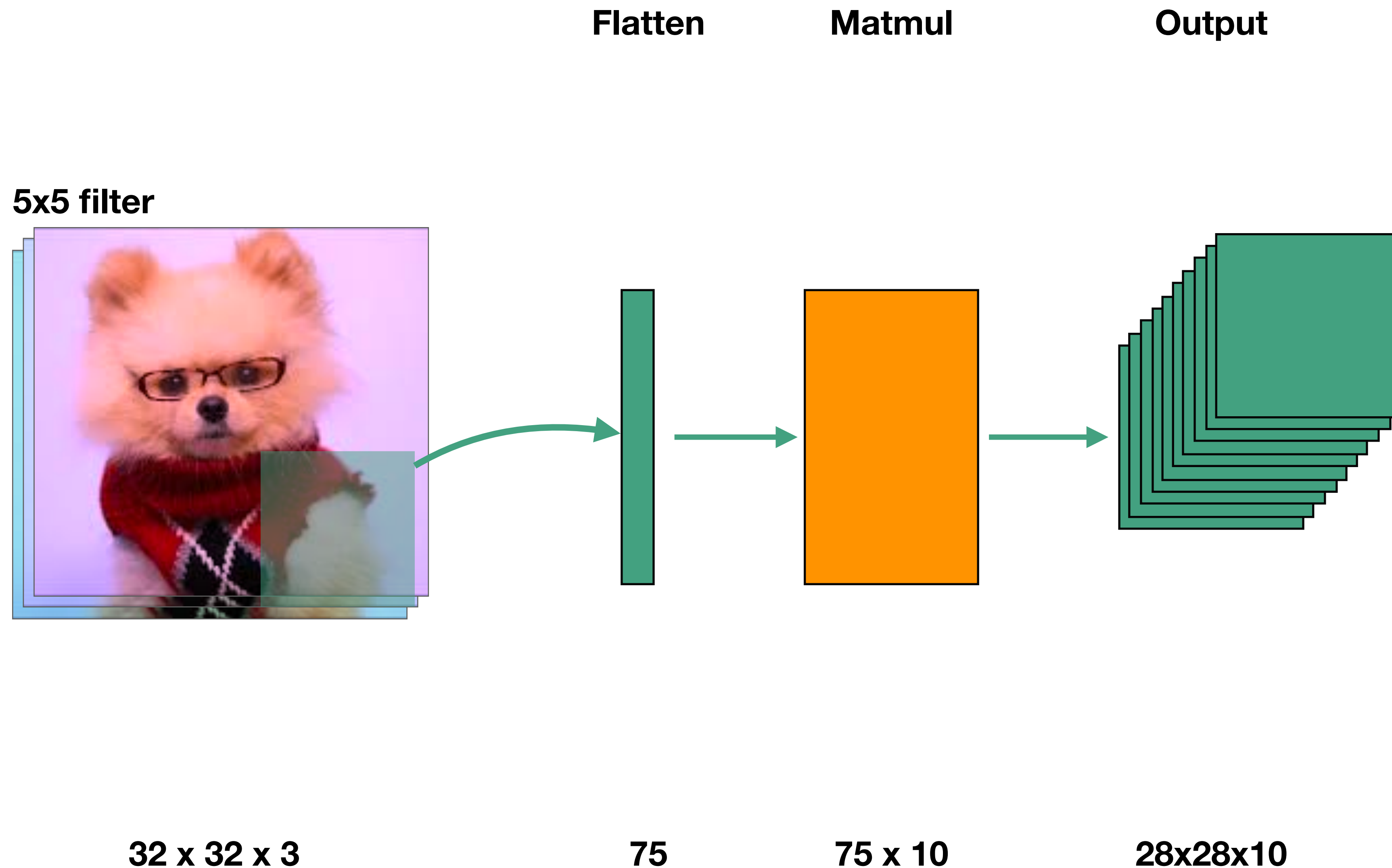
Input can have multiple channels



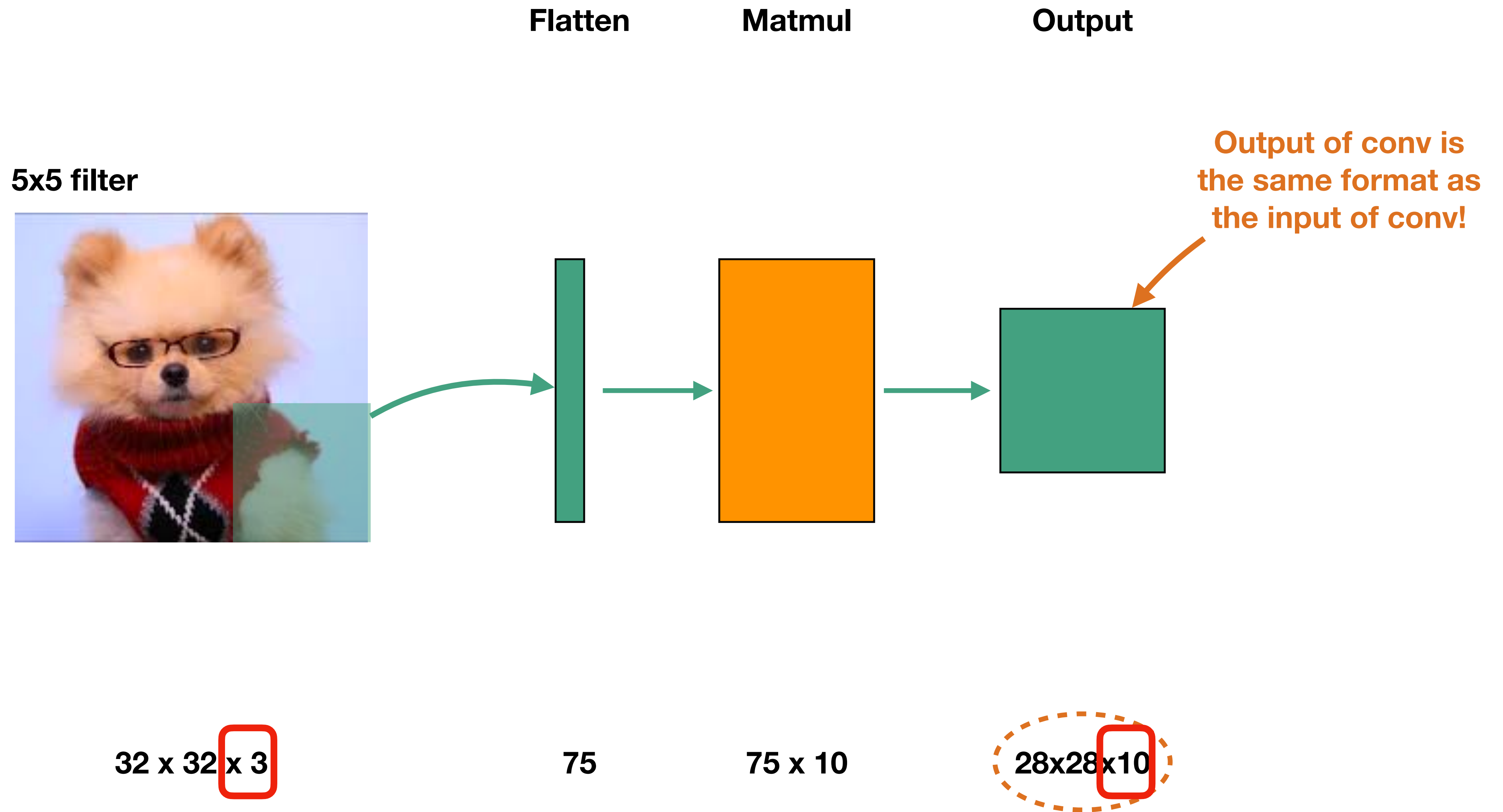
Input can have multiple channels



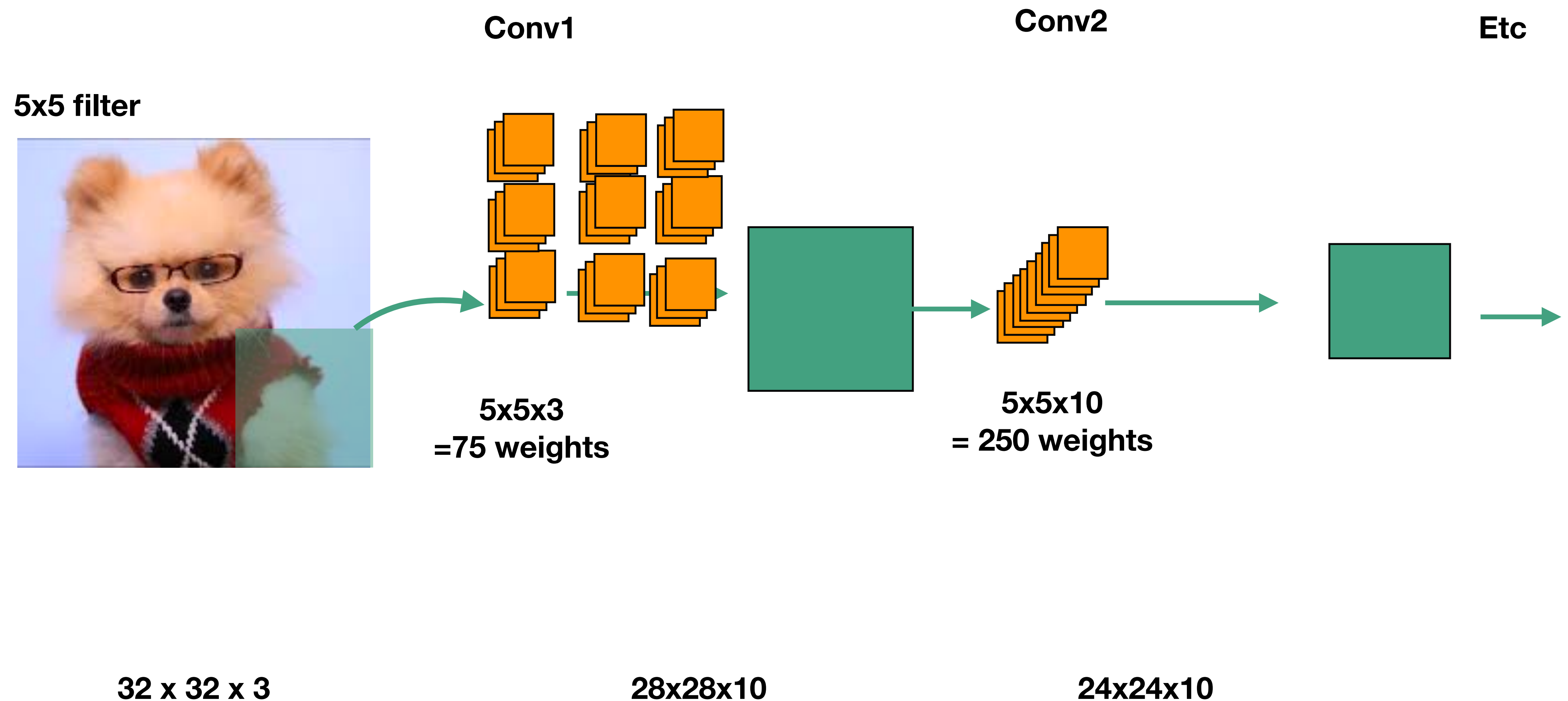
Output can have multiple channels



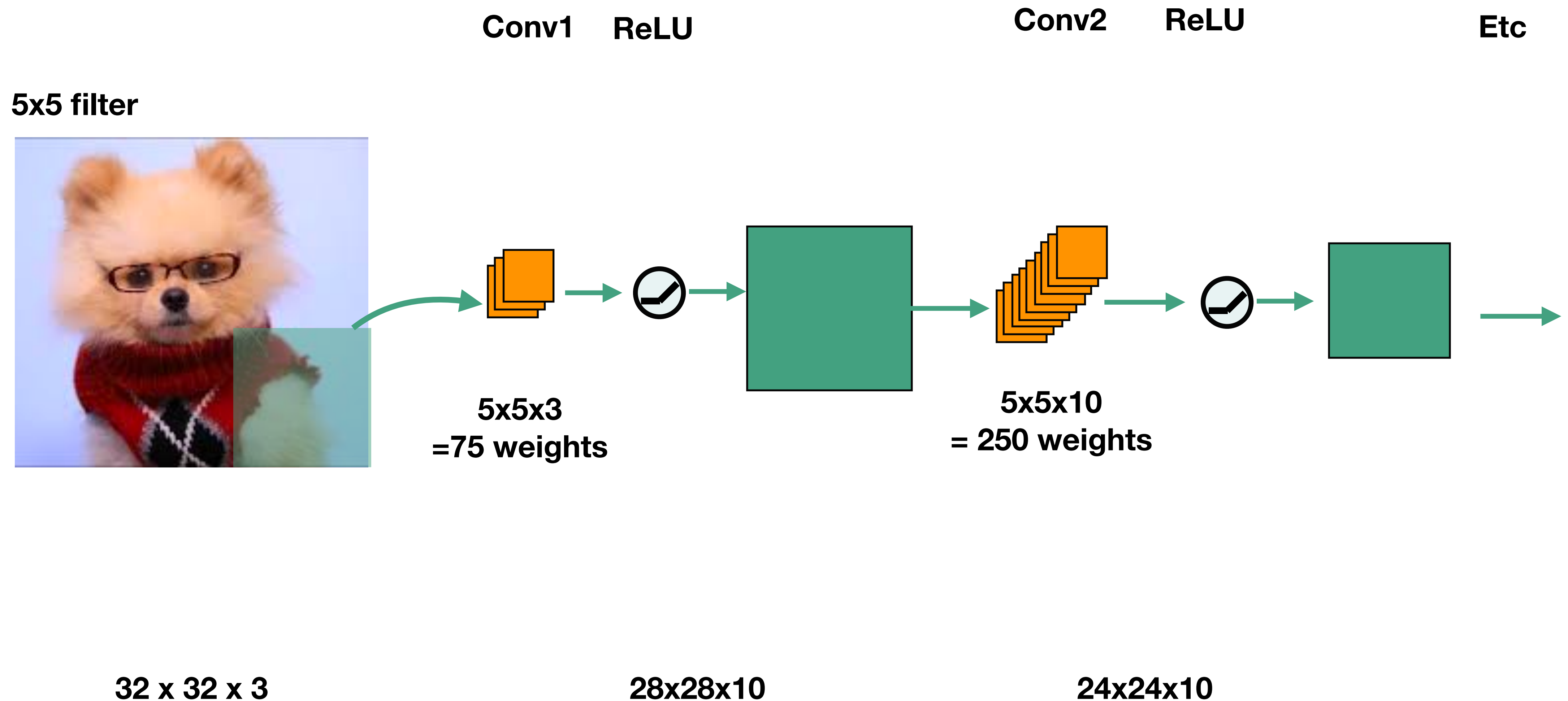
Convolutional filter stacks



Implication —> we can “stack” conv layers



Implication \rightarrow we can “stack” conv layers



Questions?

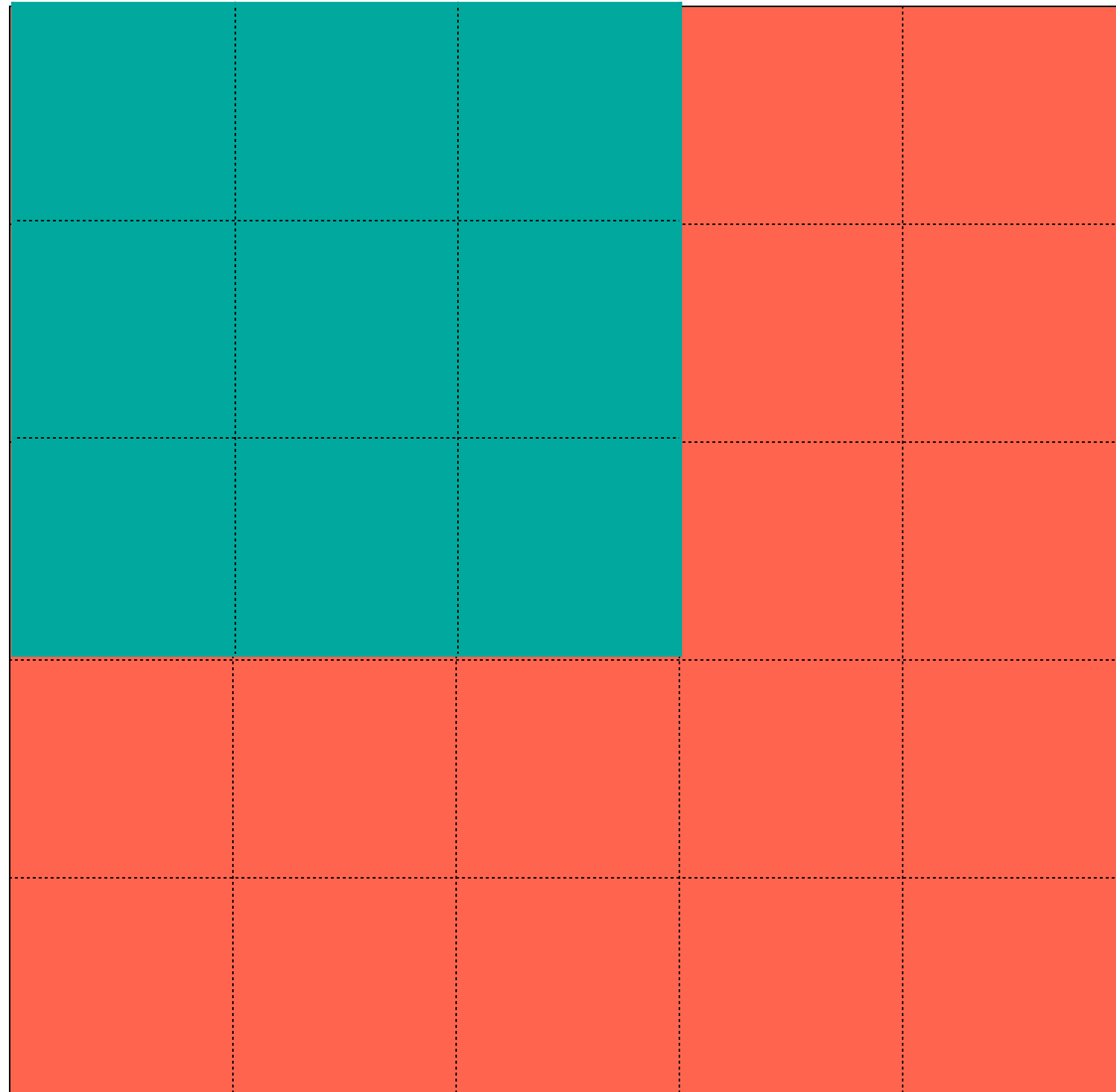
Review of convolutions

- What's a convolutional filter?
- Filter stacks and ConvNets
- **Strides & padding**
- Filter math
- Implementation notes

Strides

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

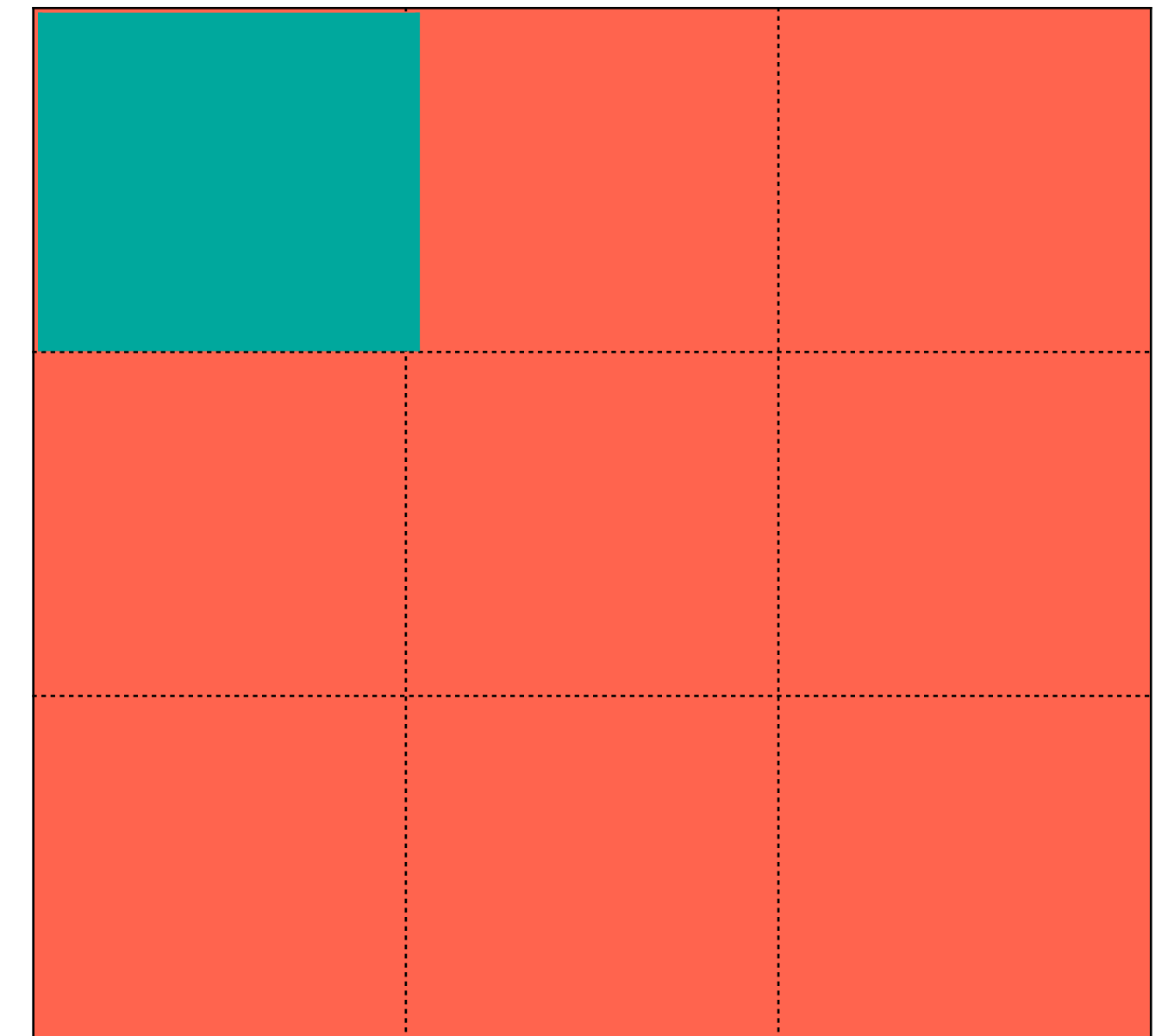
Strides



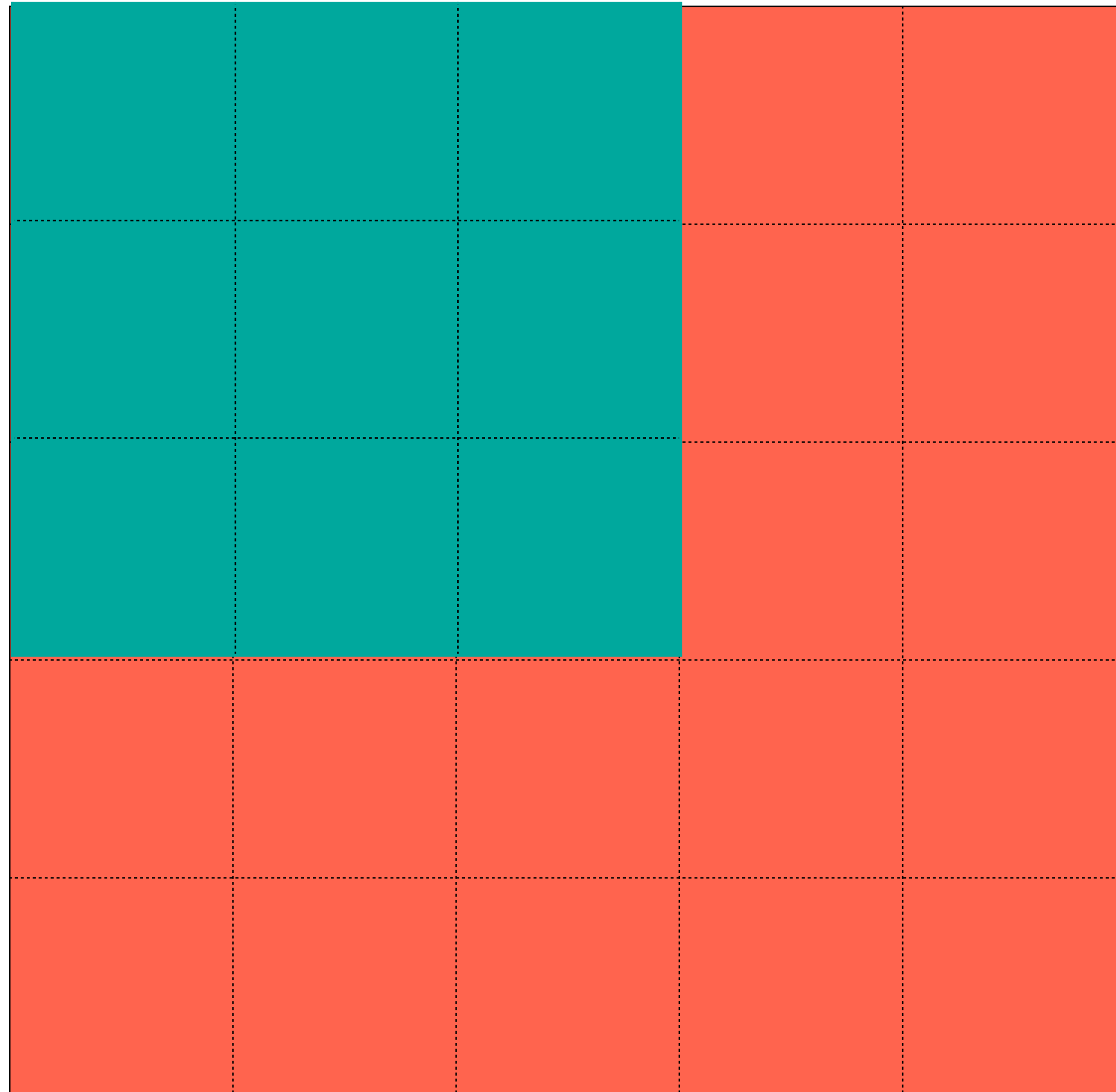
Conv2D

Filter = (3, 3)

Stride = (1, 1)



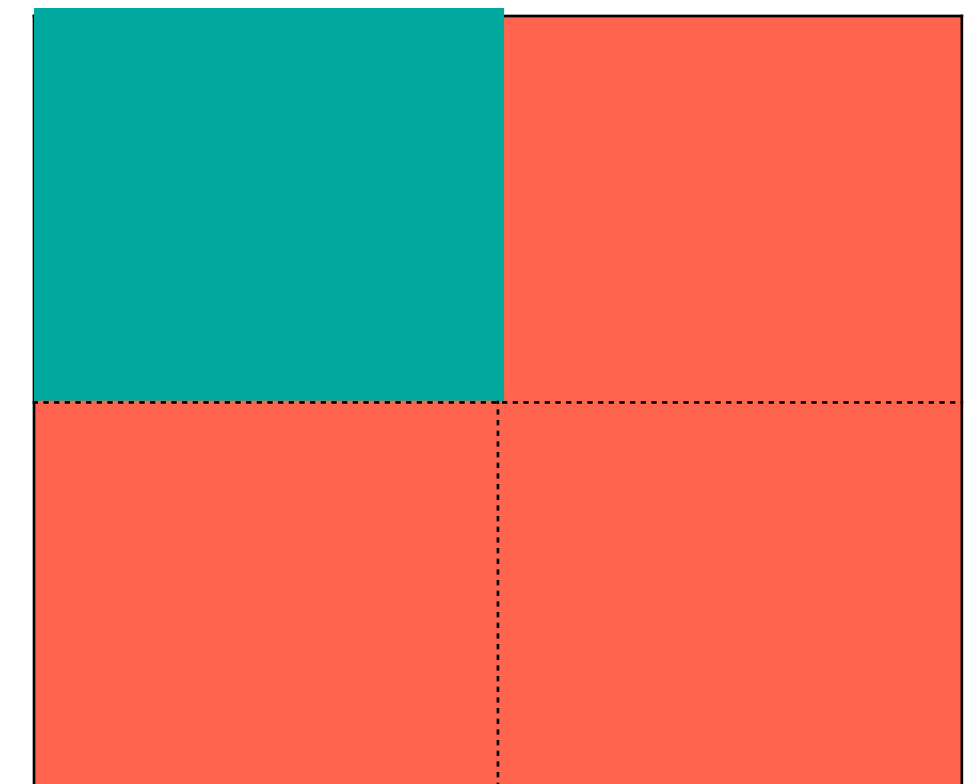
Strides



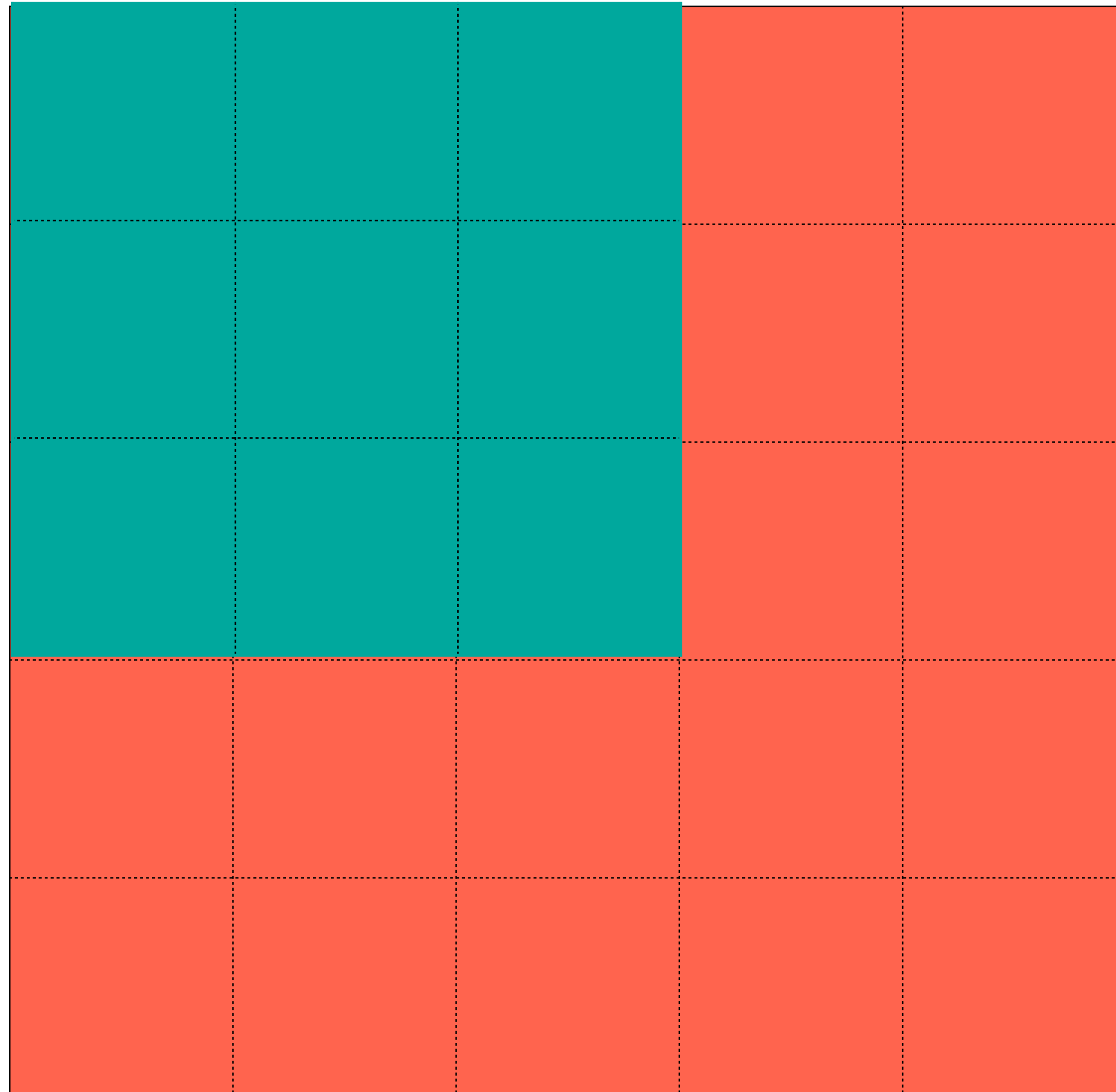
Conv2D

Filter = (3, 3)

Stride = (2, 2)



Strides



Padding can solve this out-of-boundary issue

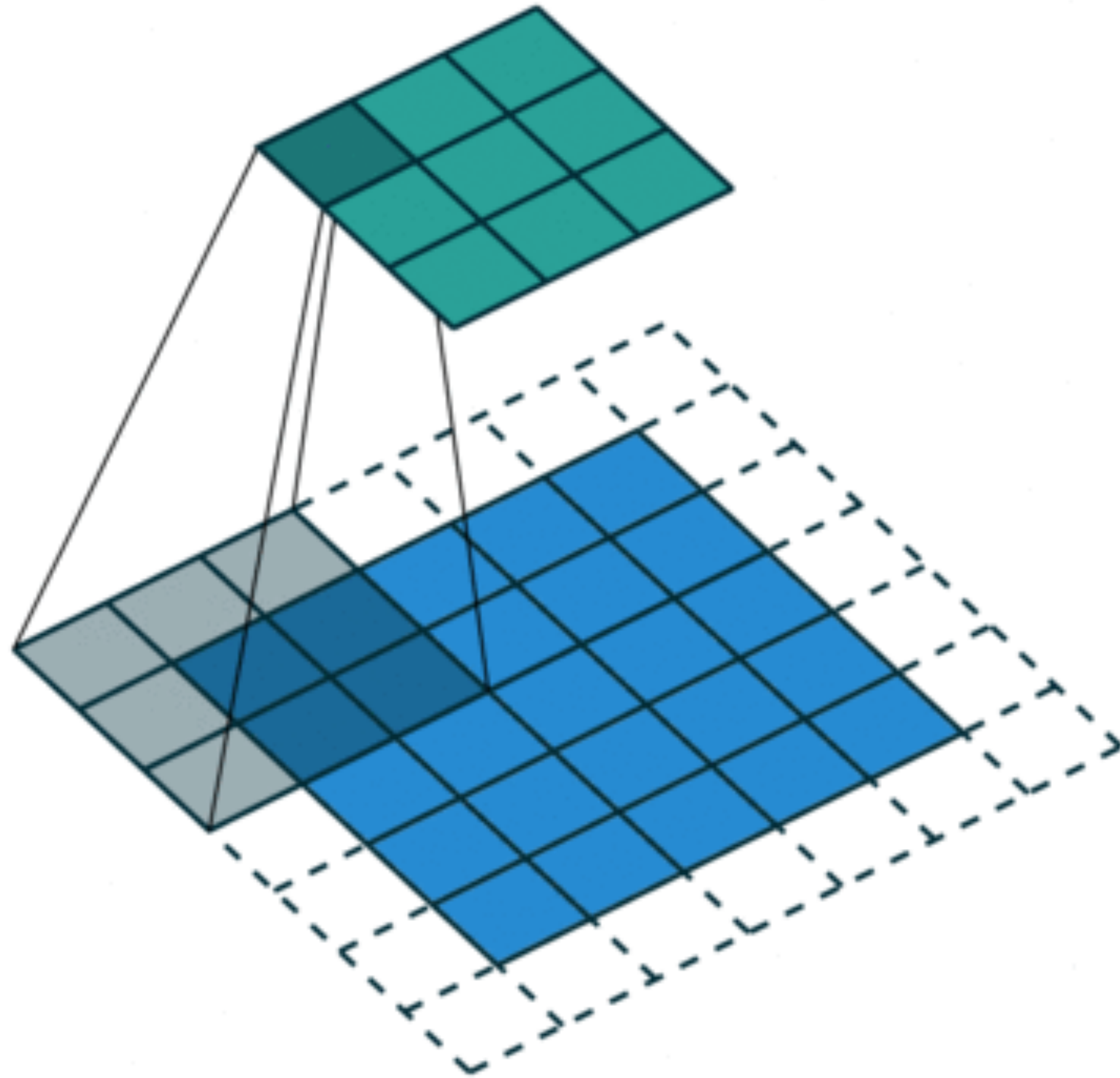
Conv2D

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


Review of convolutions

- What's a convolutional filter?
- Filter stacks and ConvNets
- Strides & padding
- **Filter math**
- Implementation notes


Conv2D Math

- Input: $W \times H \times D$ volume
- Parameters:
 - K filters, each with size (F, F)
 - ...moving at stride (S, S)
 - ...with padding P
- Output: $W' \times H' \times K$ volume
 - $W' = (W - F + 2P) / S + 1$
 - $H' = (H - F + 2P) / S + 1$
- Each filter has $(F * F * D)$ parameters
- $K * (F * F * D)$ total in the layer

Conv2D Math

- Input: $W \times H \times D$ volume
 - Parameters:
 - K filters, each with size (F, F)
 - ...moving at stride (S, S)
 - ...with padding P
 - Output: $W' \times H' \times K$ volume
 - $W' = (W - F + 2P) / S + 1$
 - $H' = (H - F + 2P) / S + 1$
 - Each filter has $(F * F * D)$ parameters
 - $K * (F * F * D)$ total in the layer
- 
- **Commonly set to powers of 2 (e.g. 32, 64, 128)**

Conv2D Math

- Input: $W \times H \times D$ volume
 - Parameters:
 - K filters, each with size (F, F)
 - ...moving at stride (S, S)
 - ...with padding P
 - Output: $W' \times H' \times K$ volume
 - $W' = (W - F + 2P) / S + 1$
 - $H' = (H - F + 2P) / S + 1$
 - Each filter has $(F * F * D)$ parameters
 - $K * (F * F * D)$ total in the layer
- 
- A diagram consisting of a vertical orange line segment on the left, a horizontal orange arrow pointing to the right, and a vertical orange line segment on the right. The horizontal arrow connects the 'Parameters' section to the 'Commonly' section.
- **Commonly $(5, 5)$, $(3, 3)$, $(2, 2)$, $(1, 1)$**

Conv2D Math

- Input: $W \times H \times D$ volume
 - Parameters:
 - K filters, each with size (F, F)
 - ...moving at stride (S, S)
 - ...with padding P
 - Output: $W' \times H' \times K$ volume
 - $W' = (W - F + 2P) / S + 1$
 - $H' = (H - F + 2P) / S + 1$
 - Each filter has $(F * F * D)$ parameters
 - $K * (F * F * D)$ total in the layer
- • **'SAME' sets it automatically**

A guide to convolution arithmetic for deep learning

Vincent Dumoulin^{1★} and Francesco Visin^{2★†}

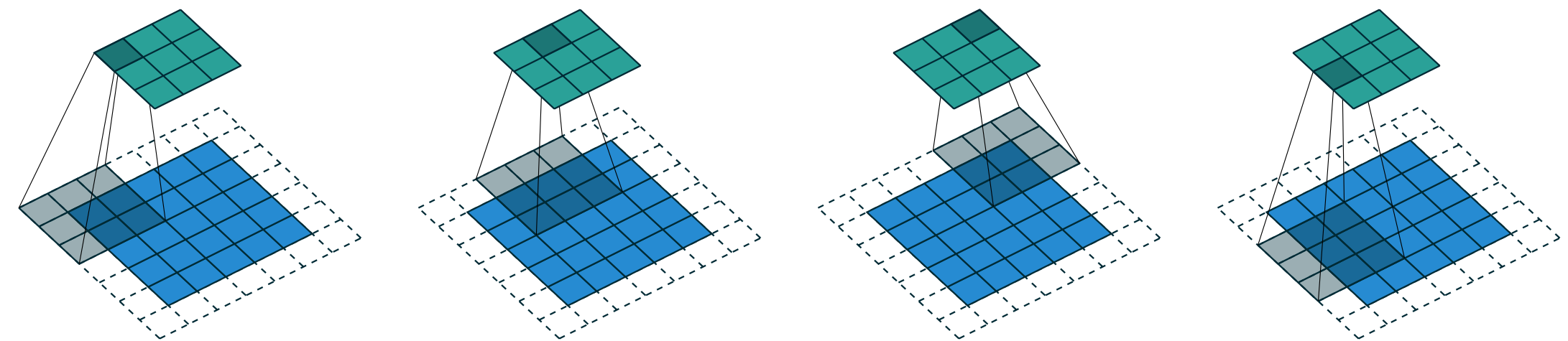


Figure 2.6: (Arbitrary padding and strides) Convolving a 3×3 kernel over a 5×5 input padded with a 1×1 border of zeros using 2×2 strides (i.e., $i = 5$, $k = 3$, $s = 2$ and $p = 1$).

- Lots of cool visualizations and comforting equations

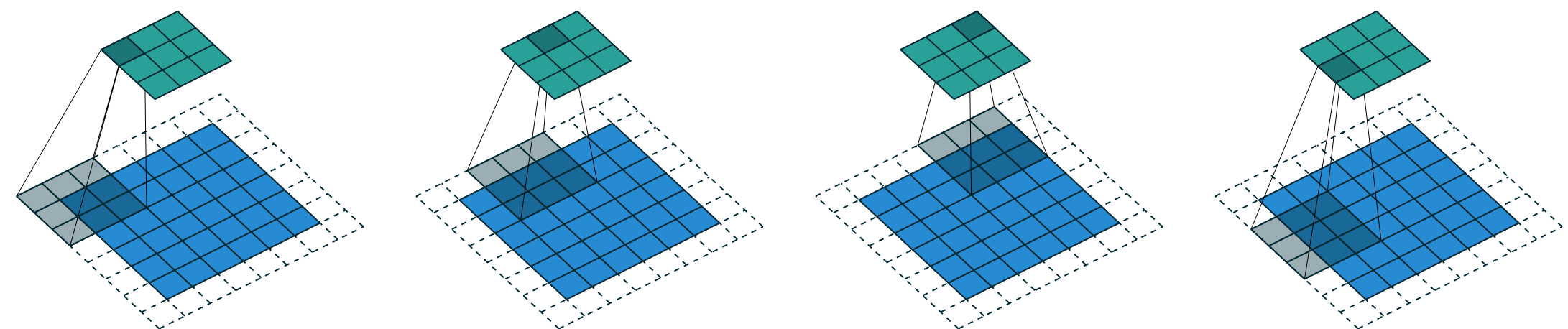


Figure 2.7: (Arbitrary padding and strides) Convolving a 3×3 kernel over a 6×6 input padded with a 1×1 border of zeros using 2×2 strides (i.e., $i = 6$, $k = 3$, $s = 2$ and $p = 1$). In this case, the bottom row and right column of the zero padded input are not covered by the kernel.

Review of convolutions

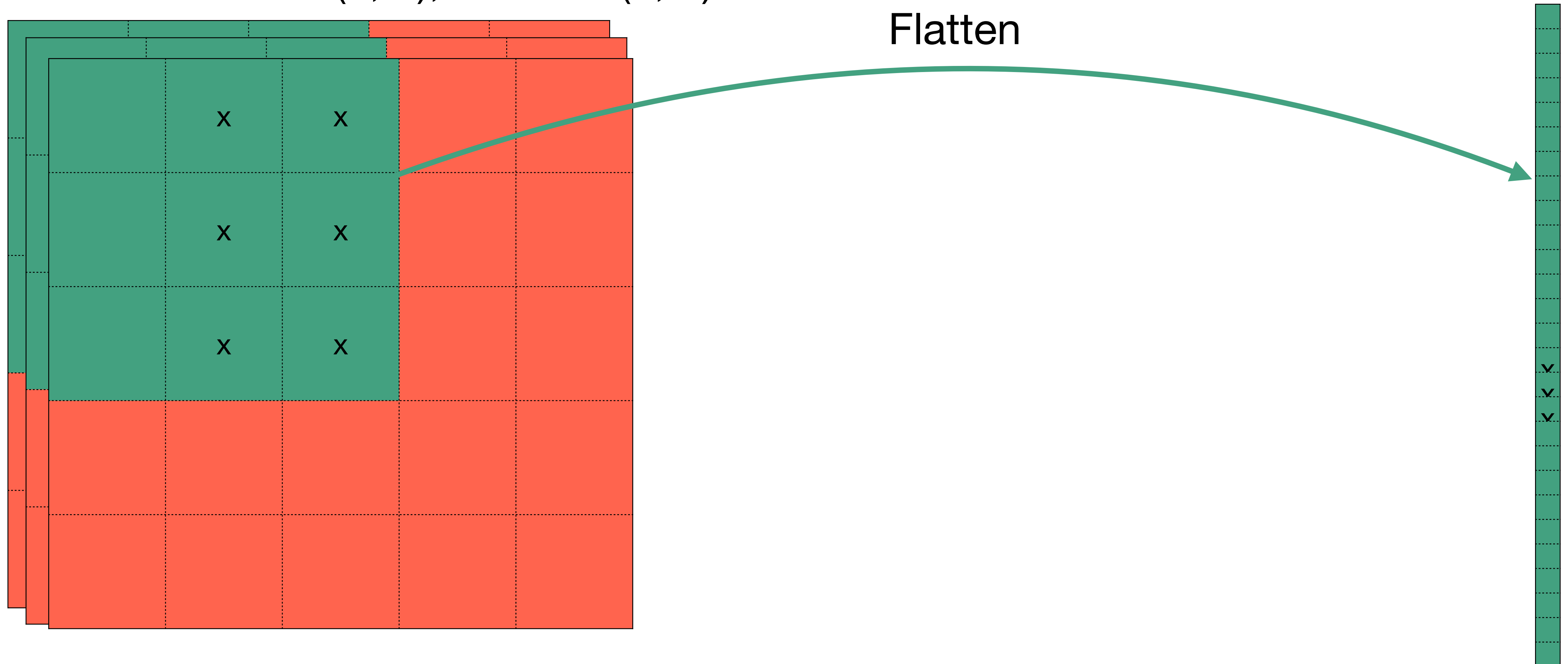
- What's a convolutional filter?
- Filter stacks and ConvNets
- Strides & padding
- Filter math
- **Implementation notes**

Convolution implementation

Conv2D. Input = (5, 5, 3)

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

$3 * 3 * 3 = 27$ -dimensional

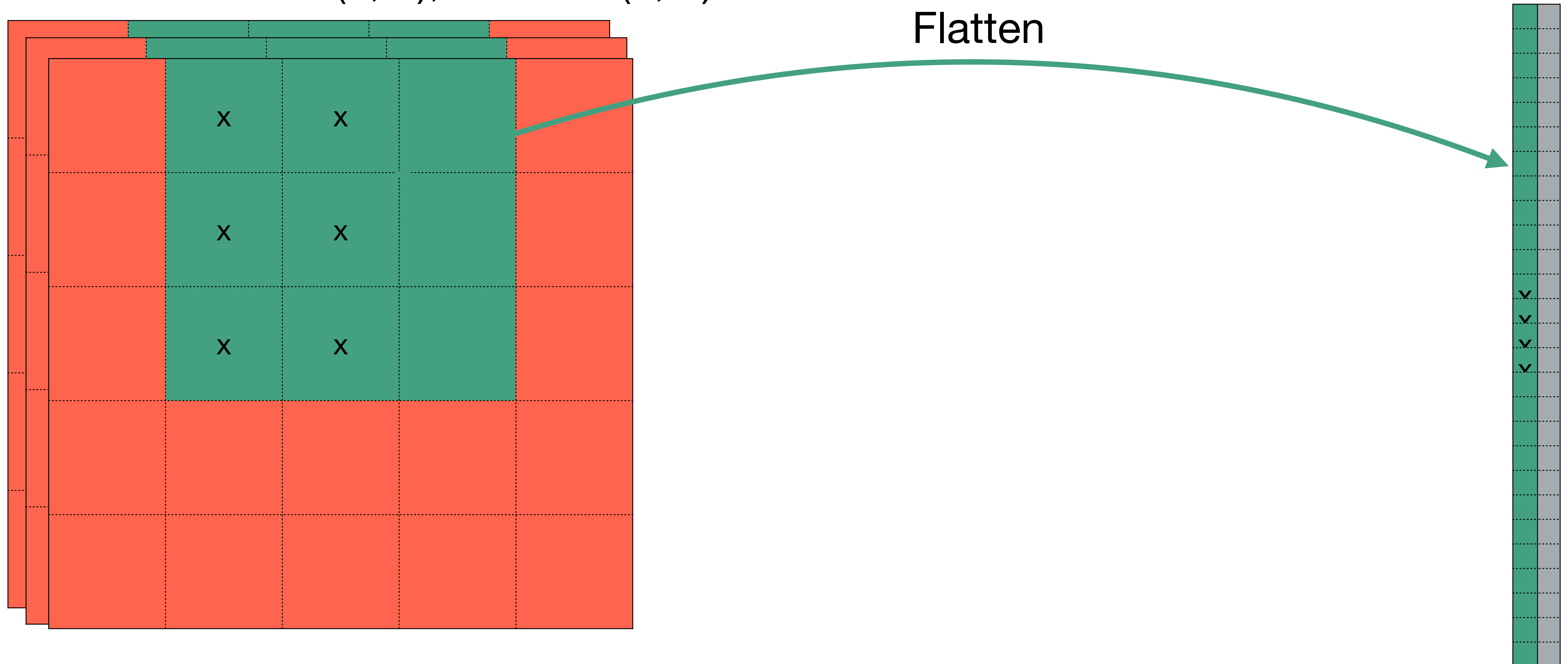


Convolution implementation

Conv2D. Input = (5, 5, 3)

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

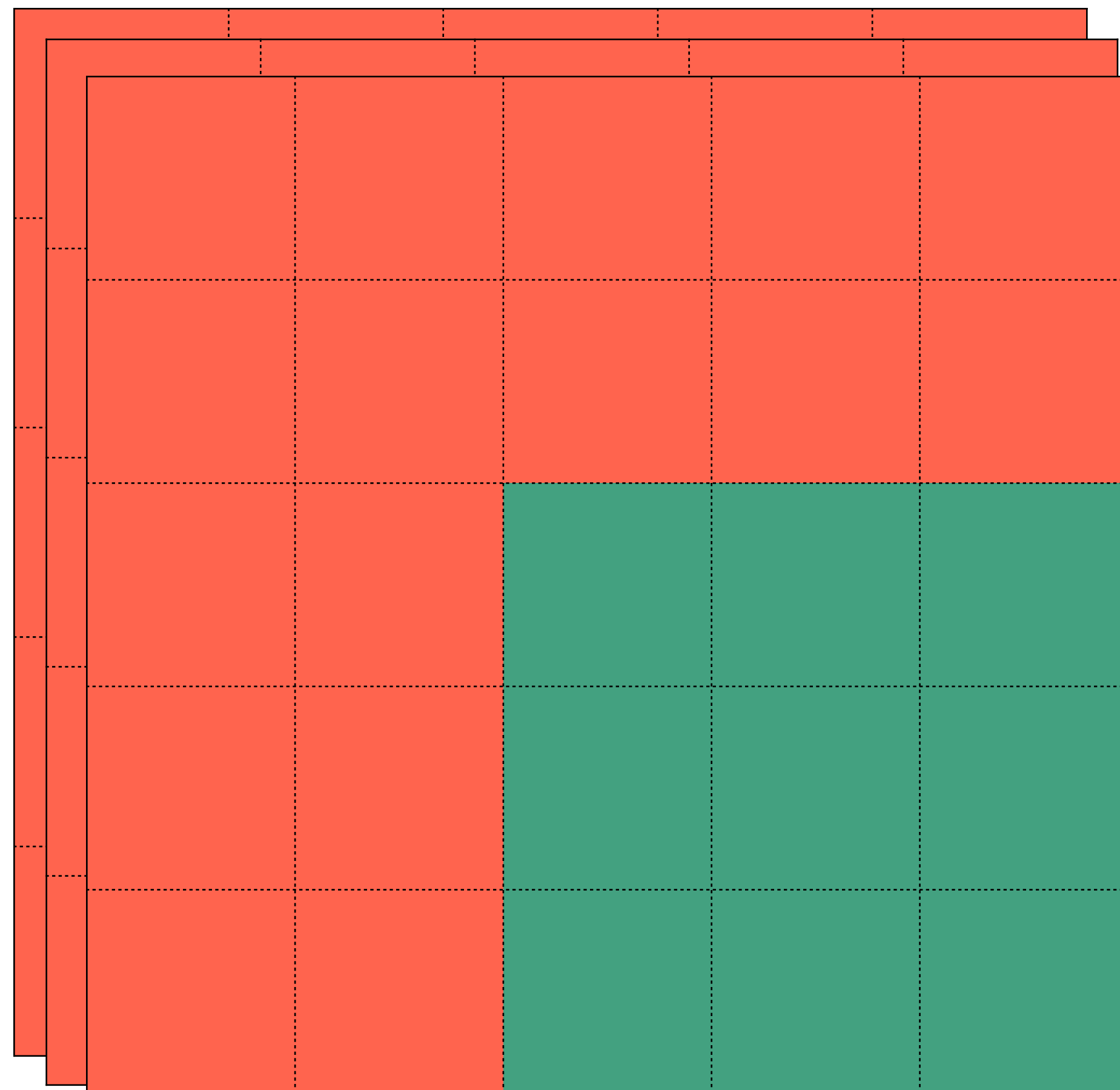
$3 * 3 * 3 = 27$ -dimensional



...sliding continues...

Convolution implementation

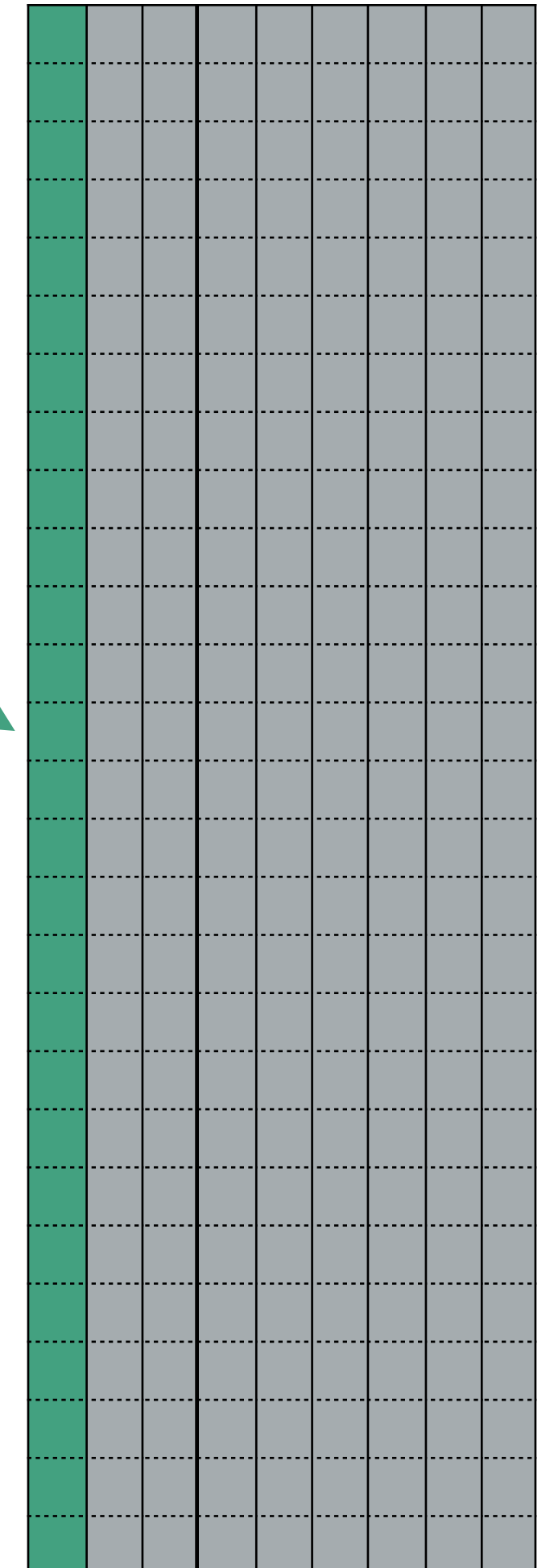
Conv2D. Input = (5, 5, 3)
Filters = 32 of size (3, 3), Stride = (1, 1)



Flatten

Im2Col

X_col
(27 x 9)



3x3 filter, 3-channel input

Conv2D. Input = (5, 5, 3)

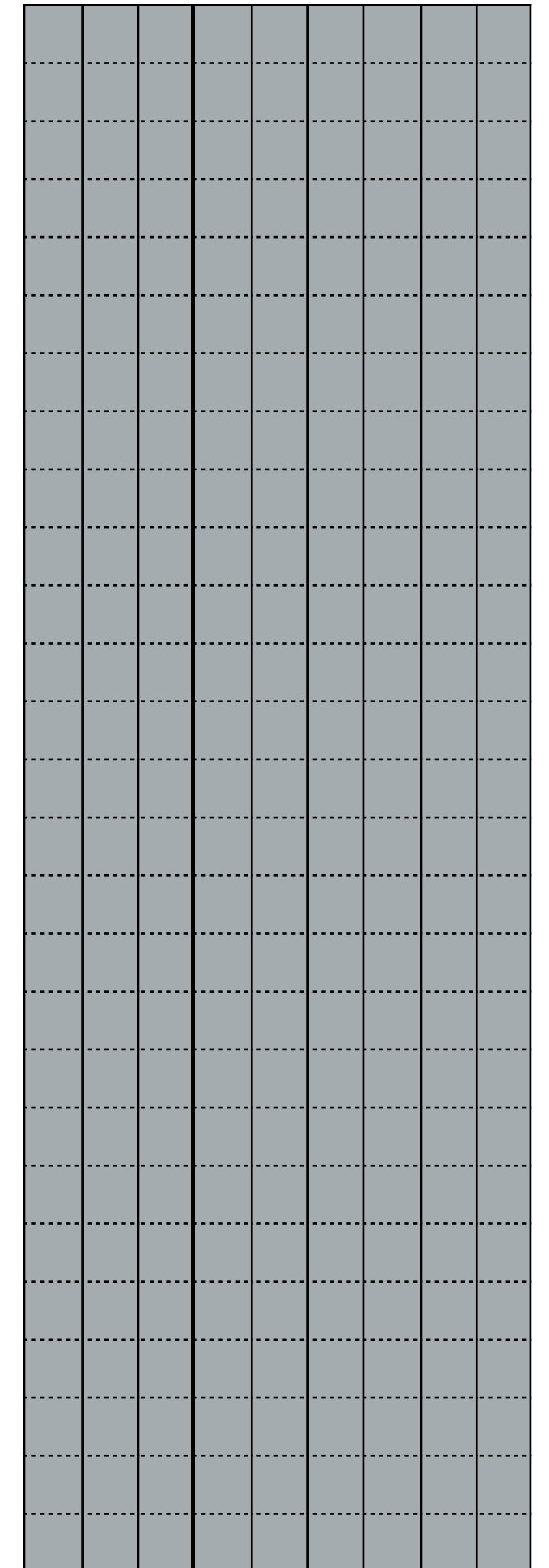
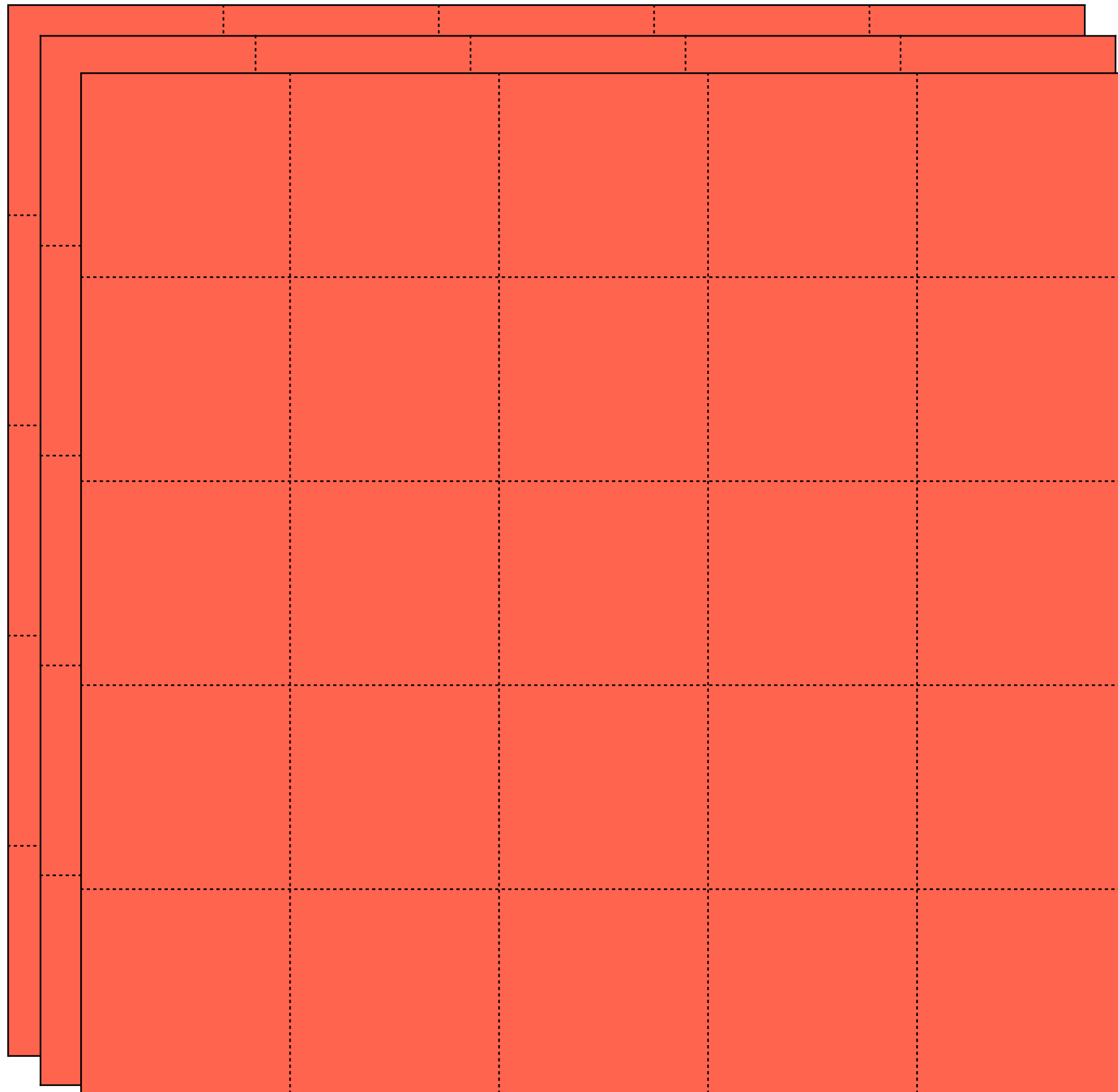
Filters = 32 of size (3, 3), Stride = (1, 1)



W_row
(32 x 27)



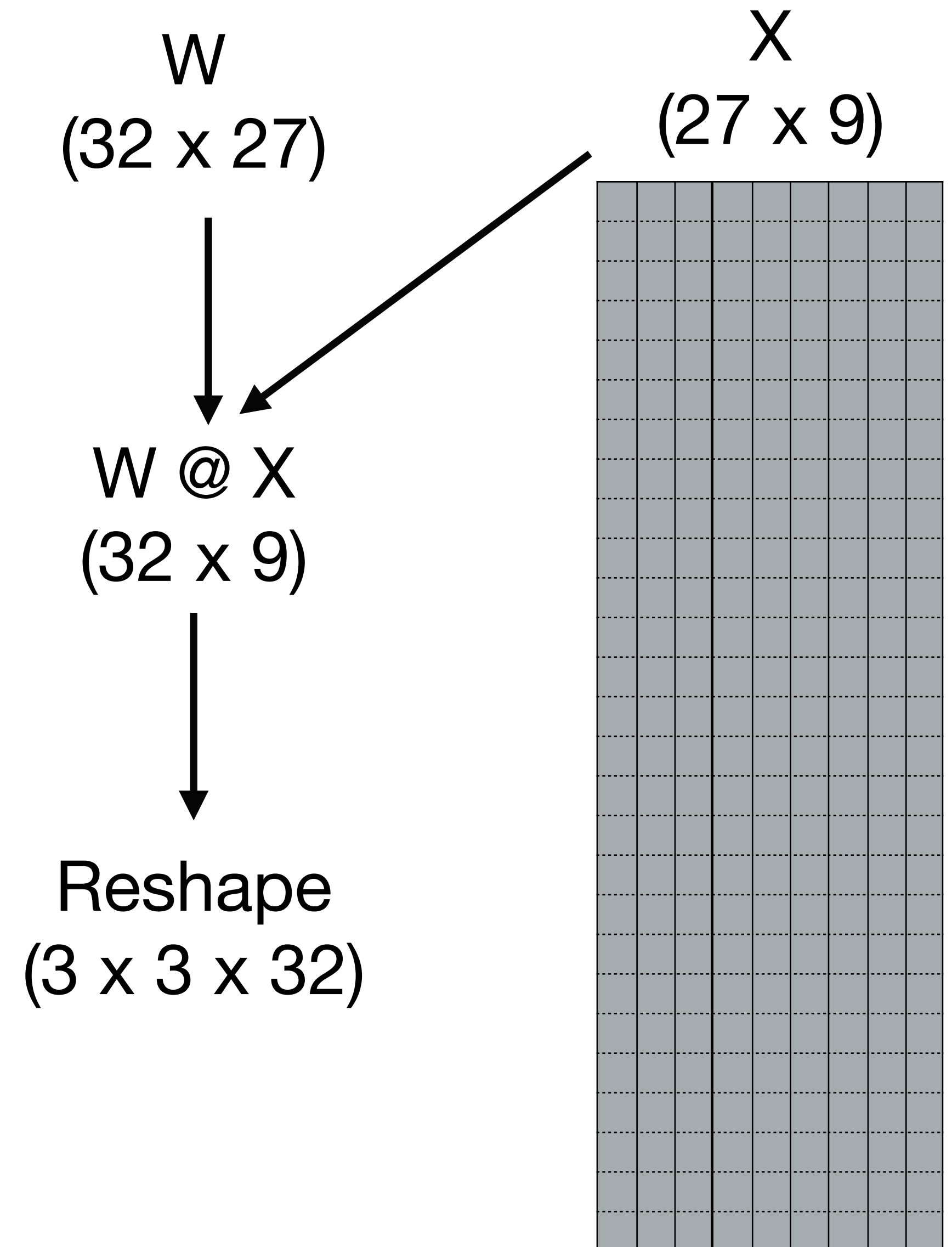
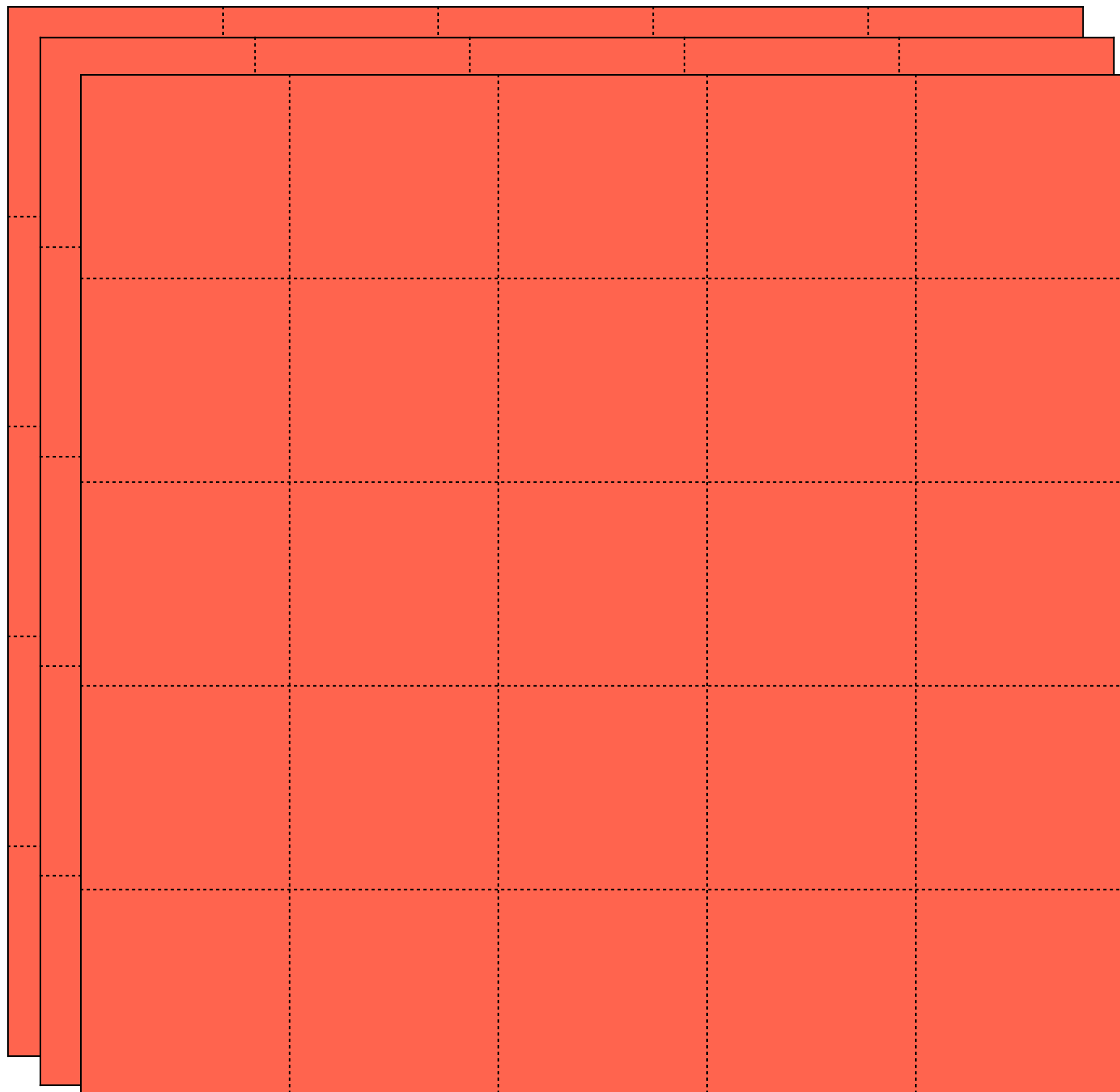
X_col
(27 x 9)



Convolution implementation

Conv2D. Input = (5, 5, 3)

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



Questions?

Agenda

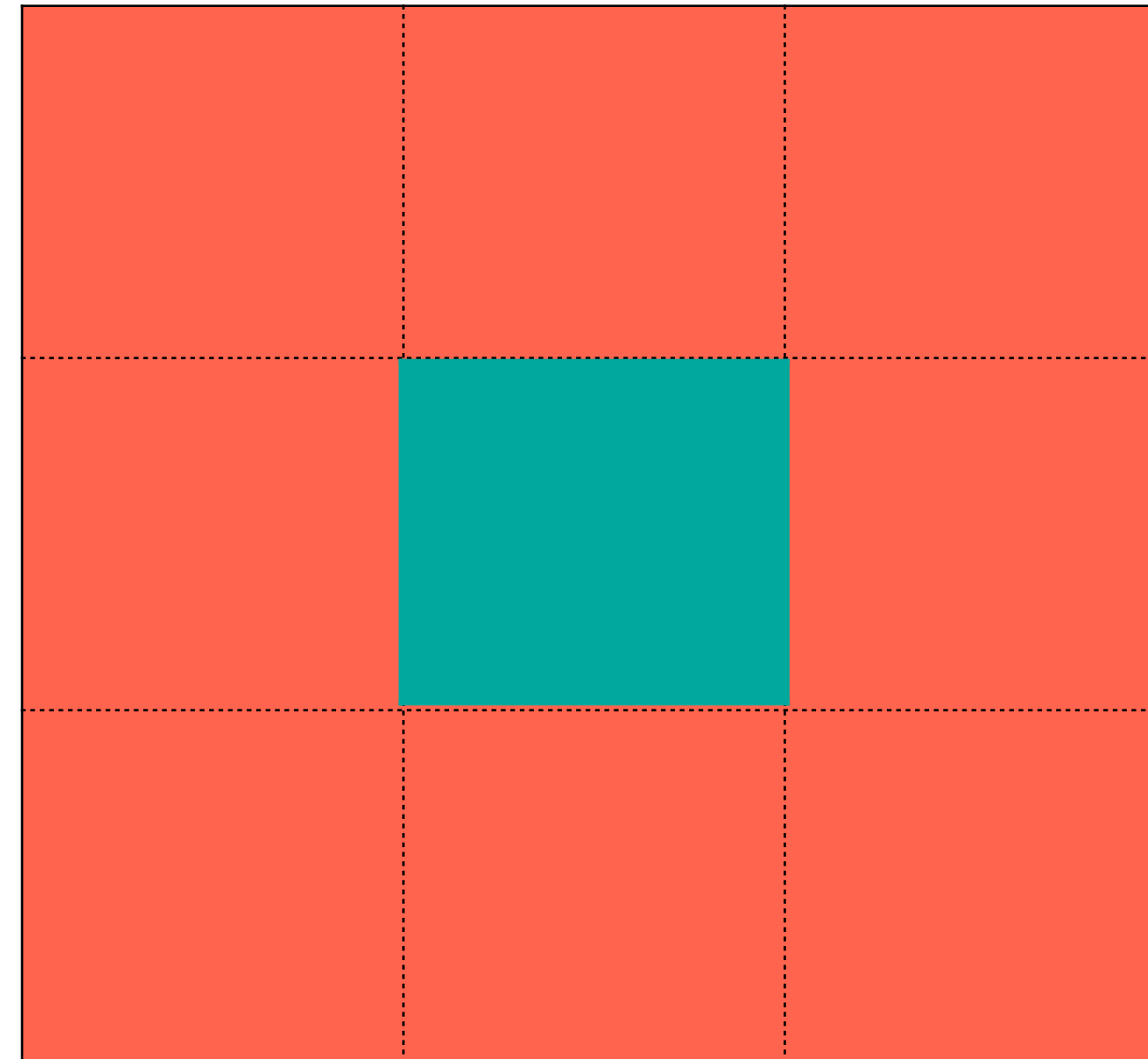
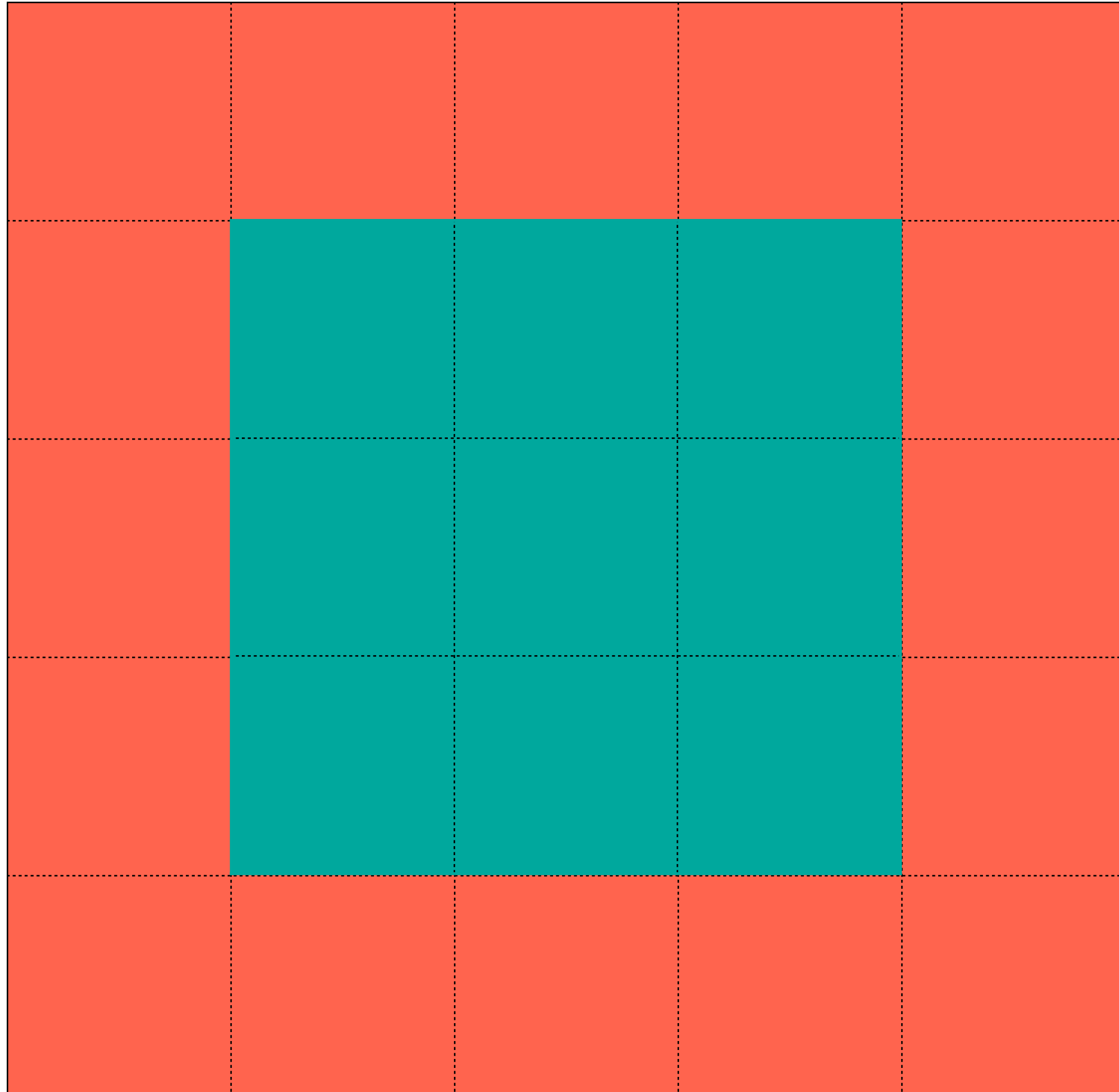
1. Review of the convolution operation
- 2. Other important operations for ConvNets**
3. Classic ConvNet architectures

Other important ConvNet operations

- **Increasing the receptive field (dilated convolutions)**
- Decreasing the size of the tensor
 - Pooling
 - 1x1-convolutions

Receptive fields

Receptive field: 3x3

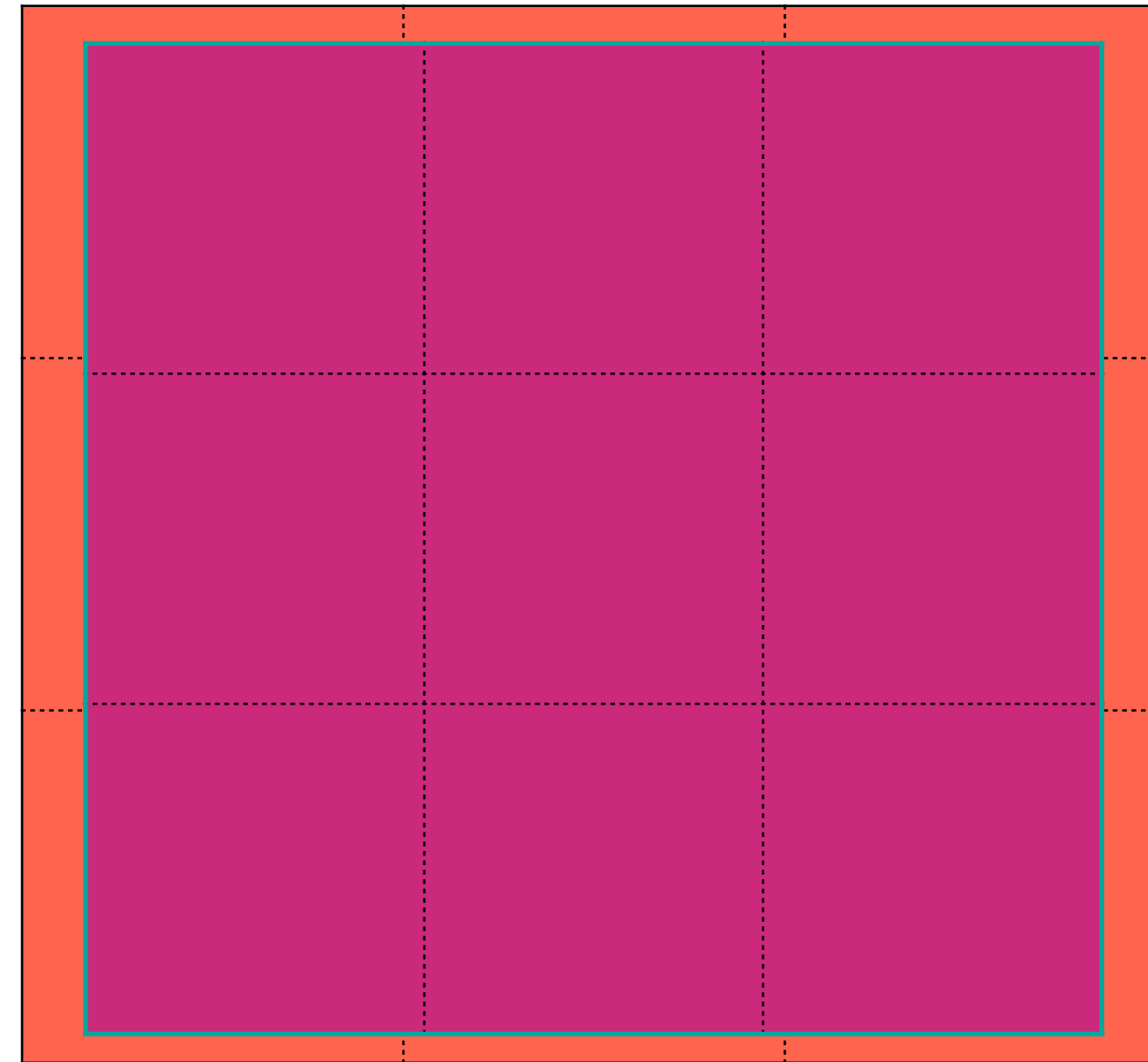
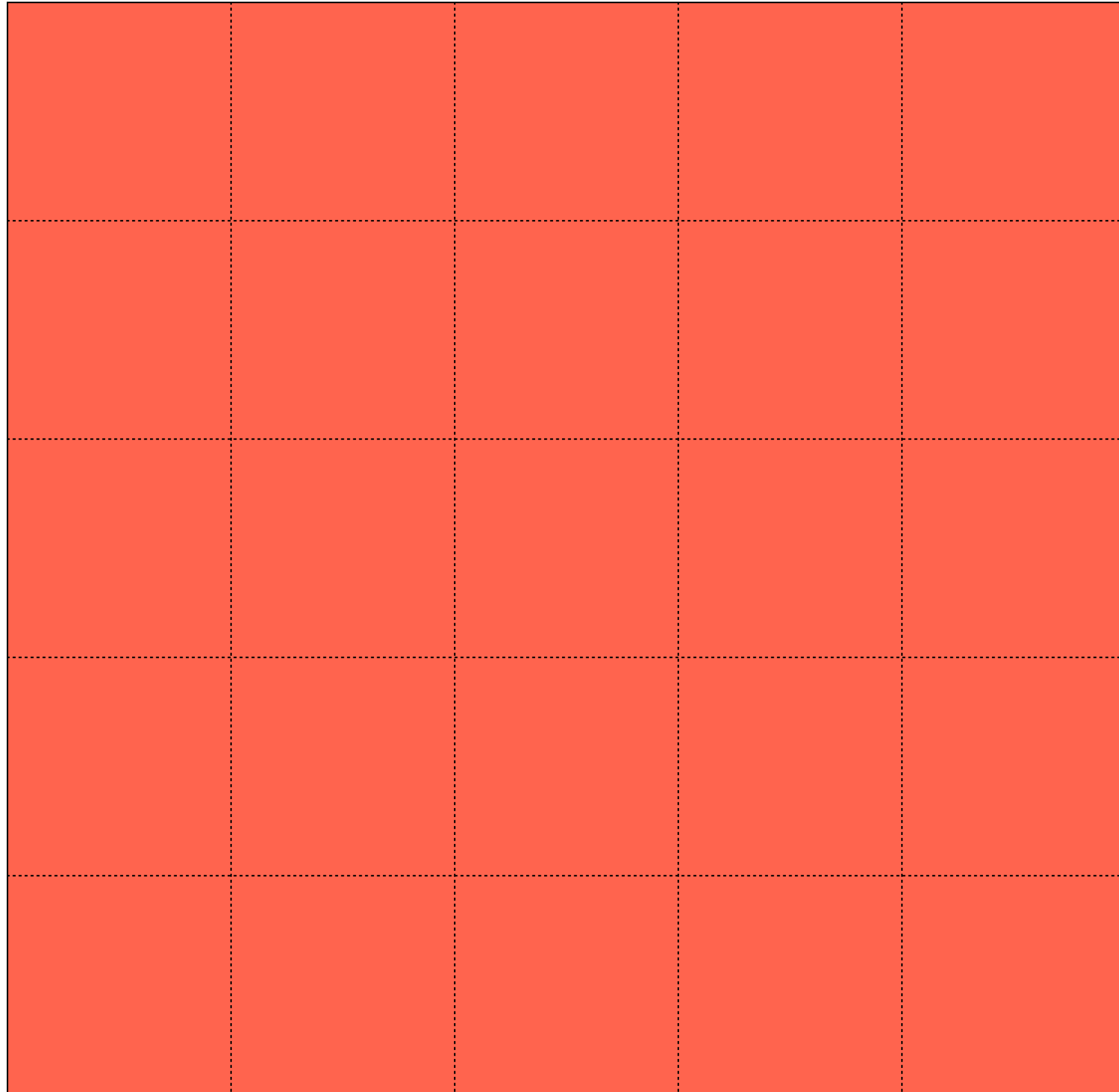


Conv2D

Filter = (3, 3)

Stride = (1, 1)

Receptive fields



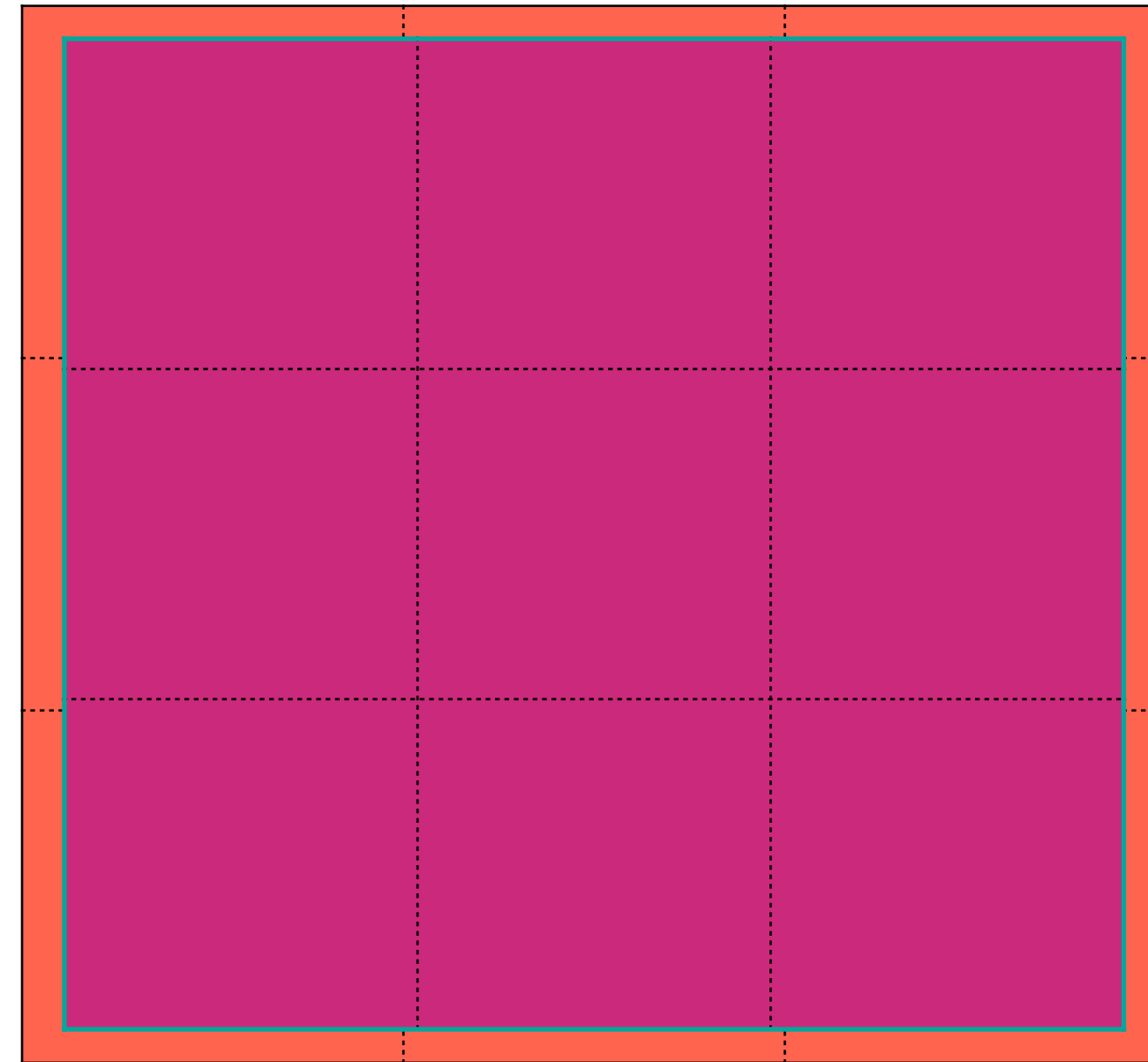
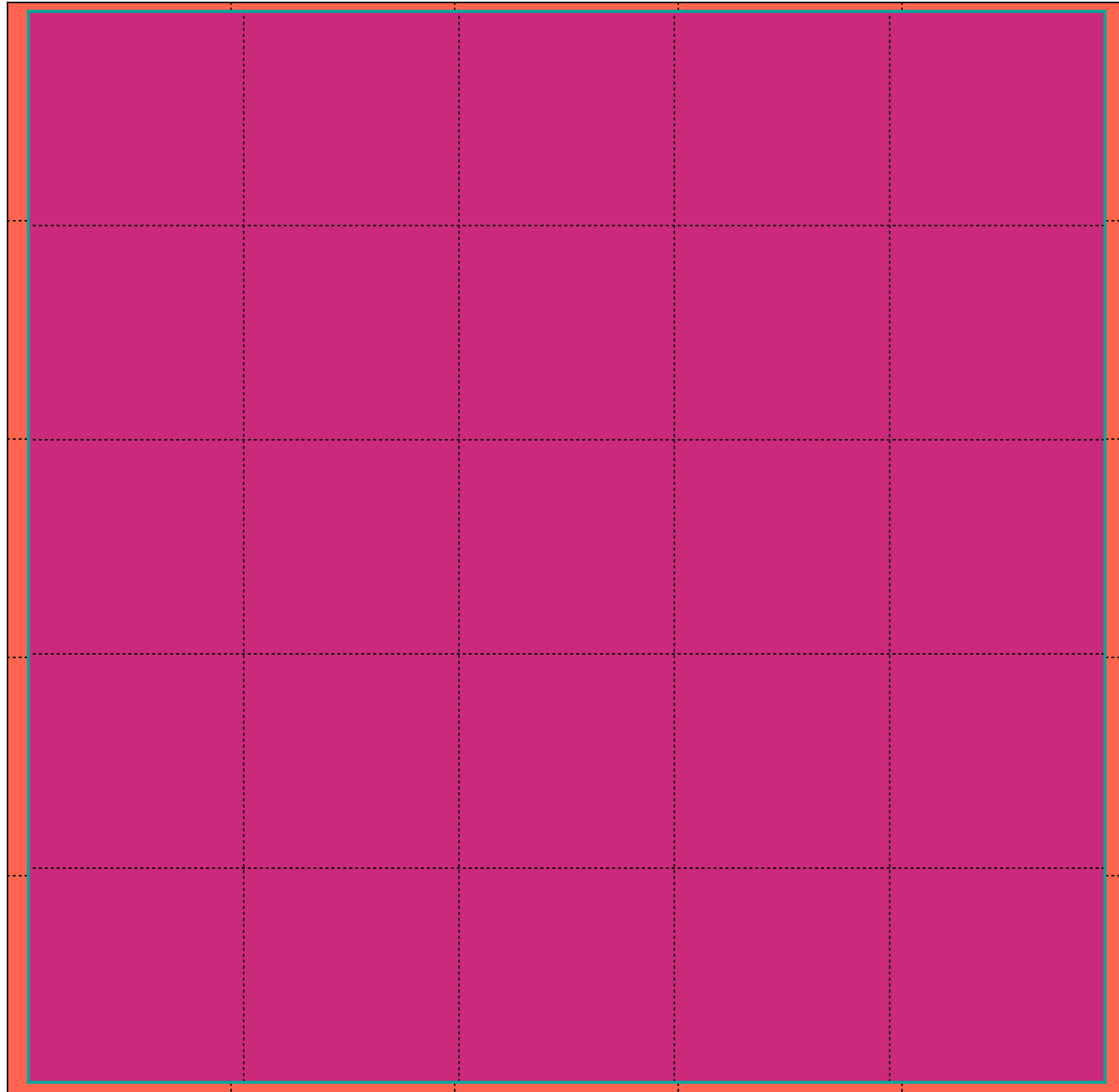
Conv2D
Filter = (3, 3)
Stride = (1, 1)



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

Receptive fields

Original receptive field: 5x5



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

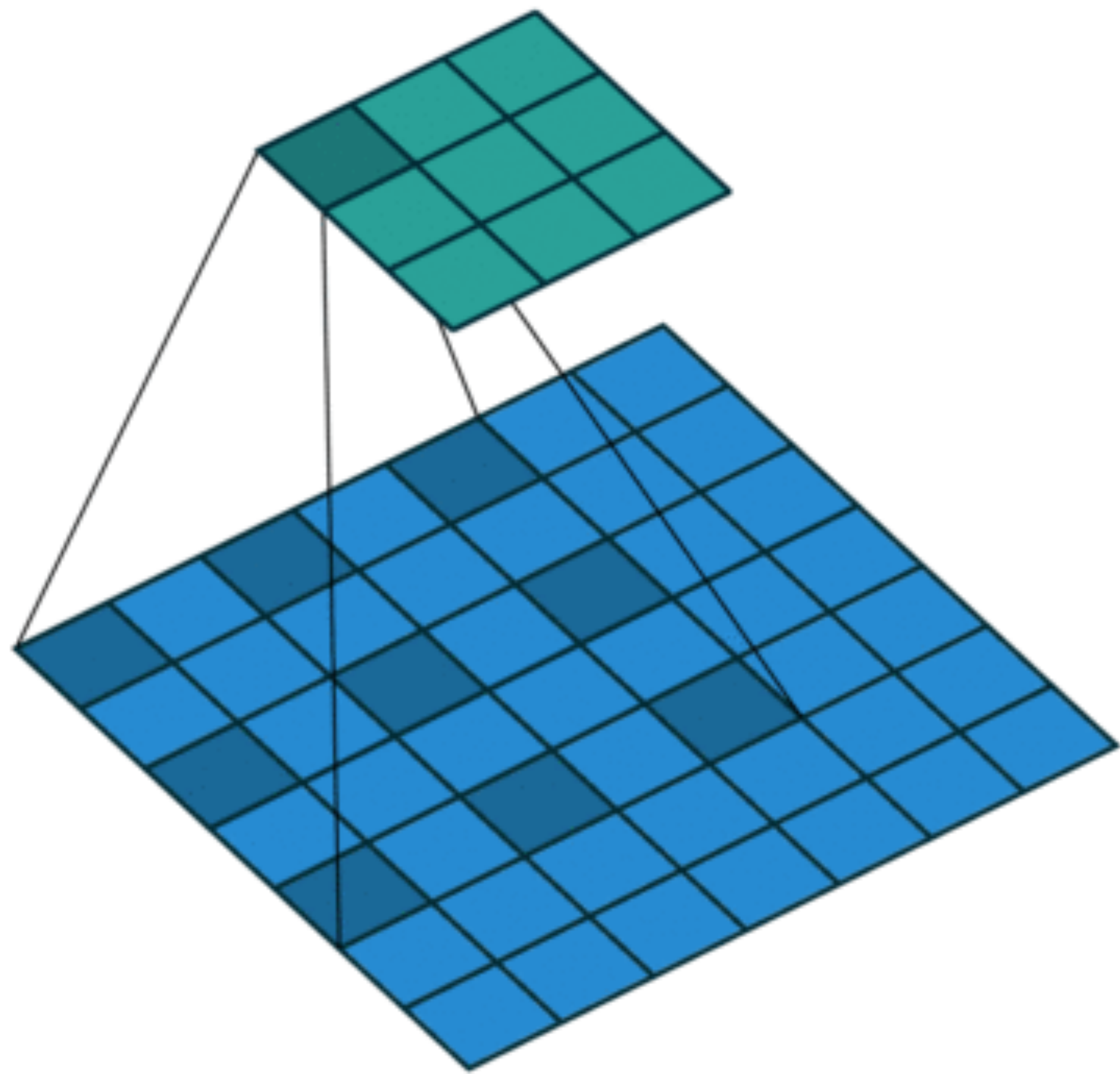


Conv2D
Filter = (3, 3)
Stride = (1, 1)

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)
 - (with fewer parameters!)

Dilated Convolution

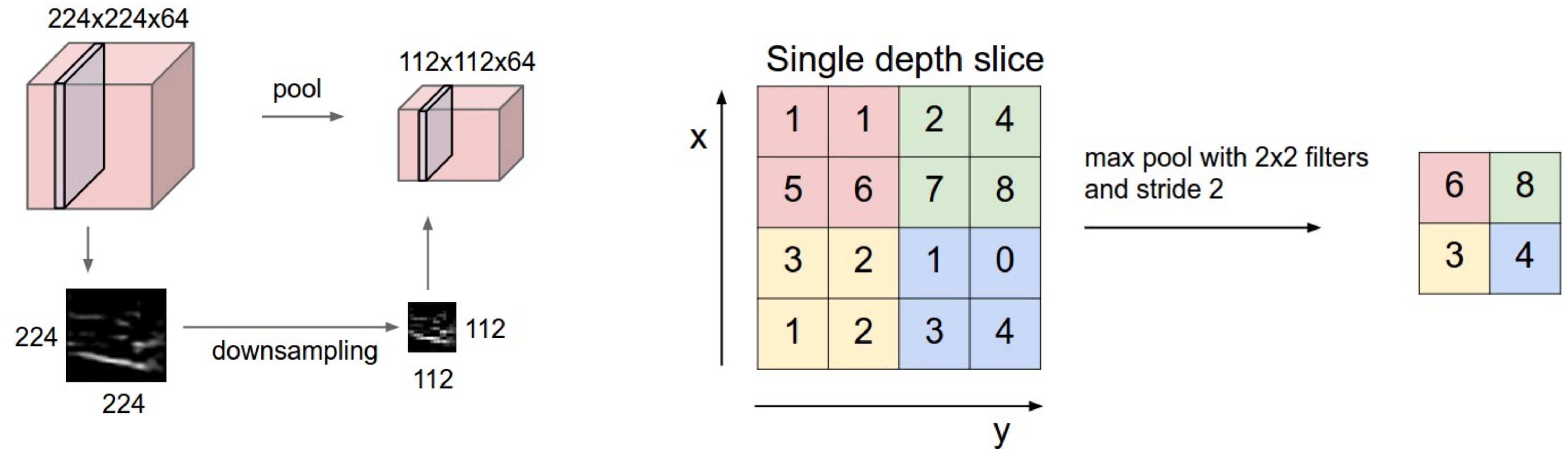


- Dilated convolutions can “see” a greater portion of the image by skipping pixels
- The $(3, 3)$ 1-dilated convolution illustrated here has a $(5, 5)$ receptive field
- Stacking dilated convolutions up quickly gets to large receptive fields

Other important ConvNet operations

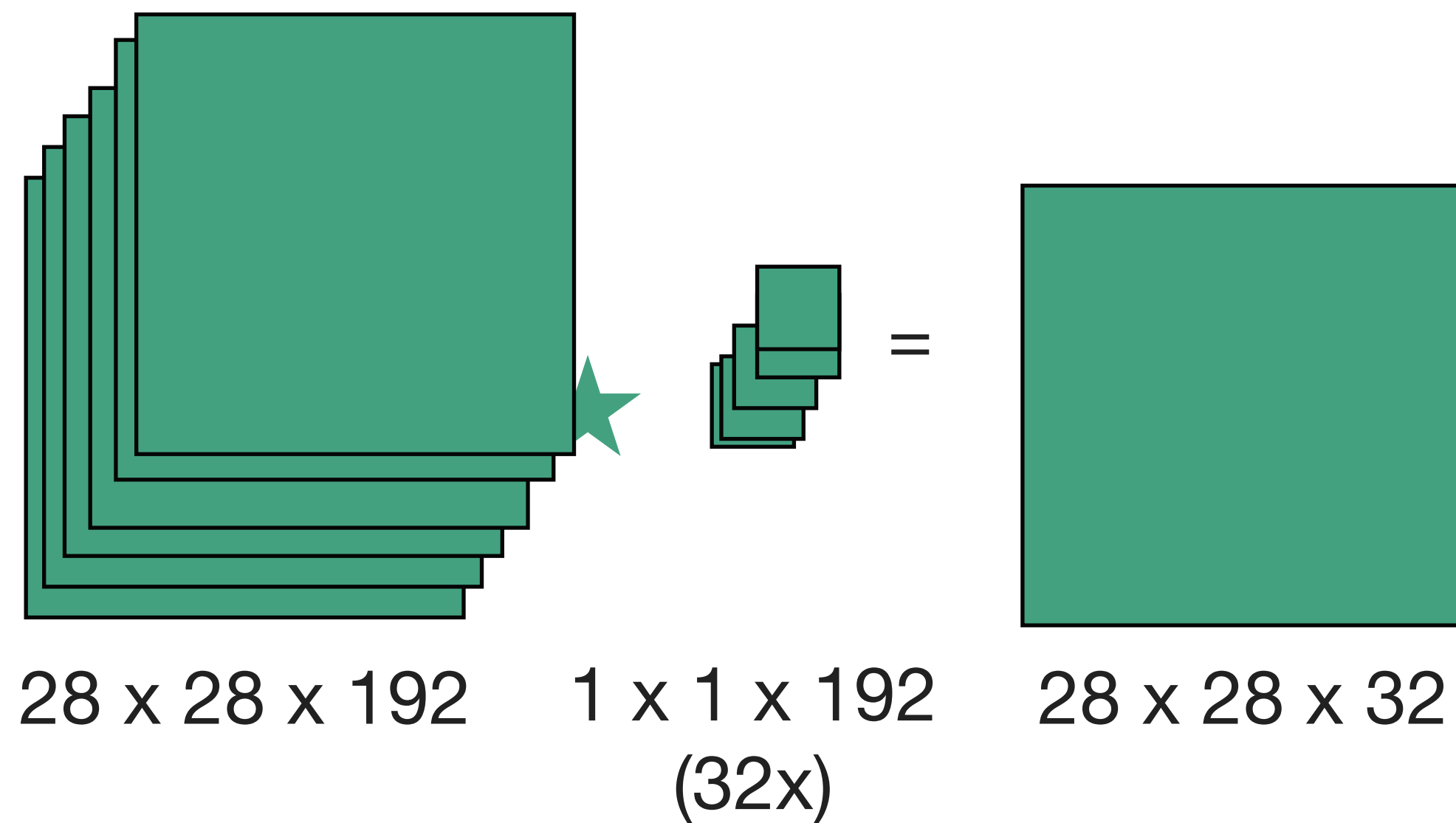
- Increasing the receptive field (dilated convolutions)
- **Decreasing the size** of the tensor
 - Pooling
 - 1x1-convolutions

Pooling



- Subsamples the image through average or max of region
- 2×2 max pooling is most common
- Recently fallen out of favor

1x1 Convolution



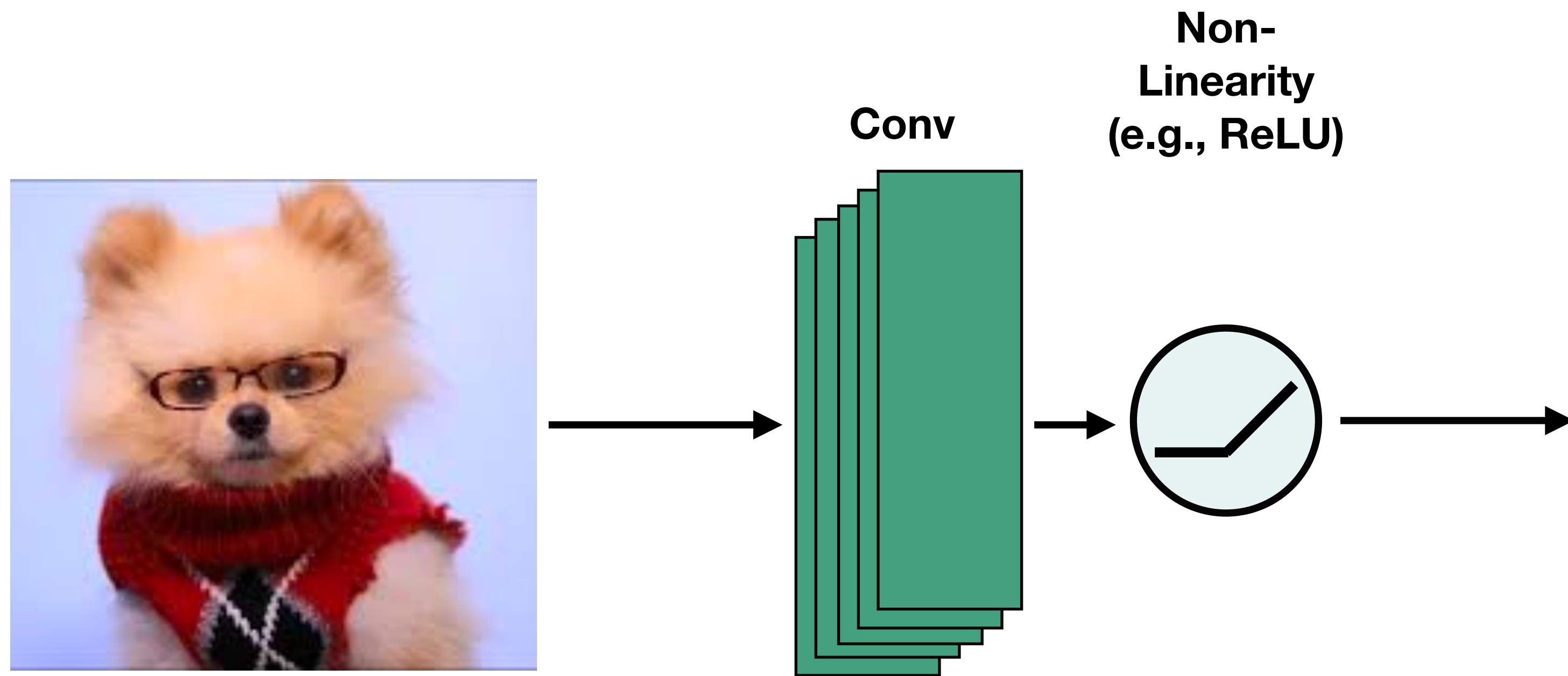
- A way to **reduce the “depth”** dimension of convolutional outputs
- Corresponds to applying an MLP to every pixel in the convolutional output
- Crucial to popular convnet architectures like Inception (GoogLeNet)

Questions?

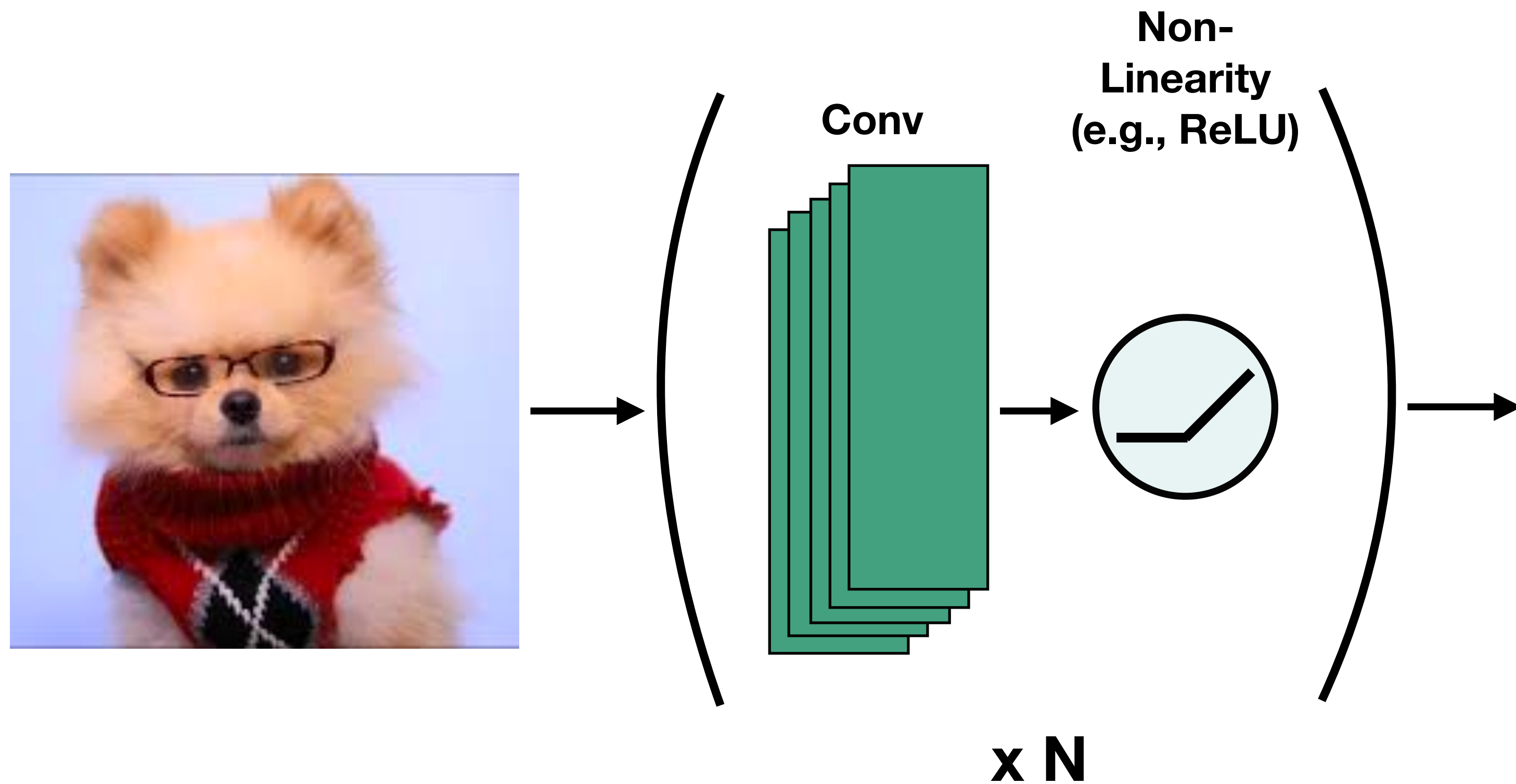
Agenda

1. Review of the convolution operation
2. Other important operations for ConvNets
3. **Classic ConvNet architectures**

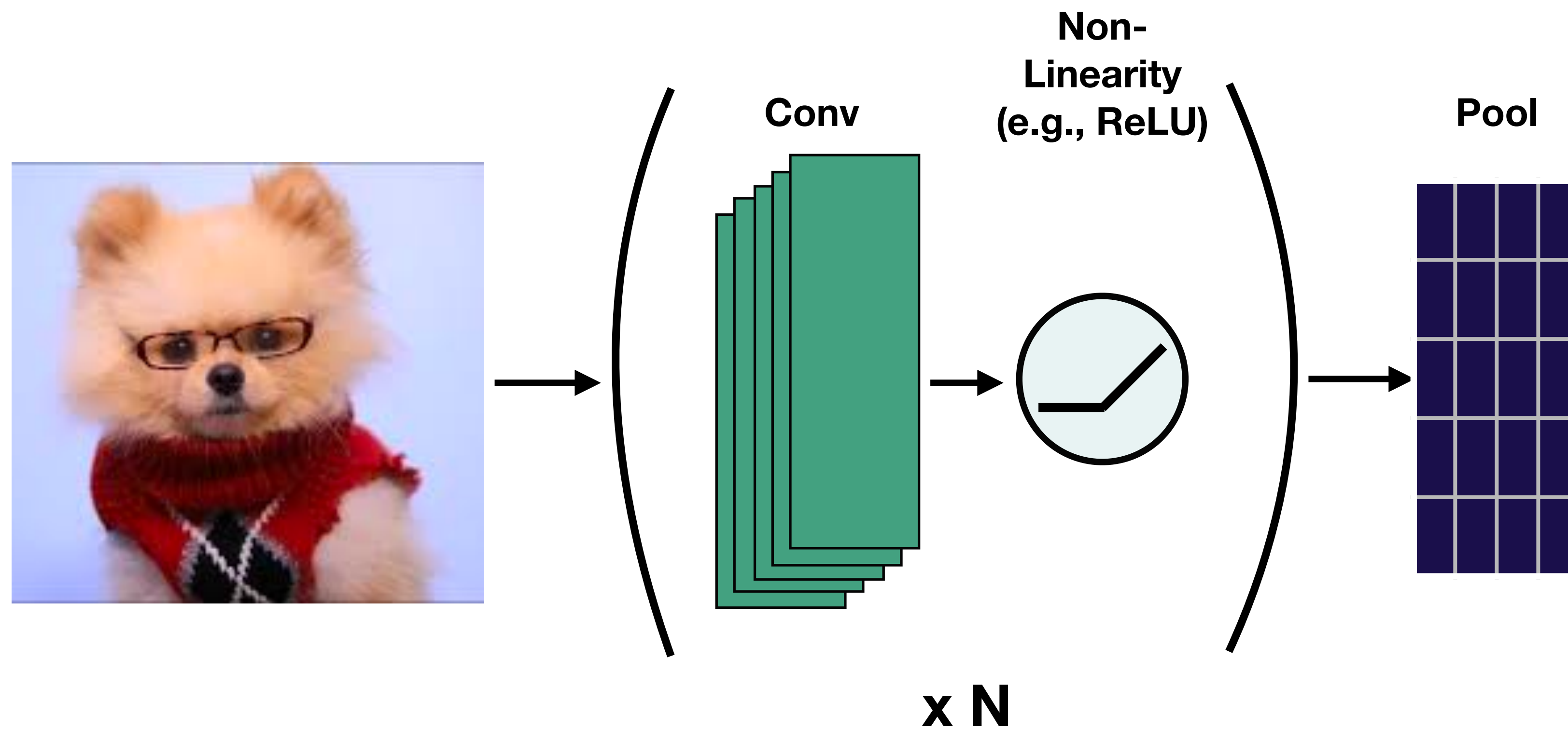
Most standard: LeNet(-like) architectures



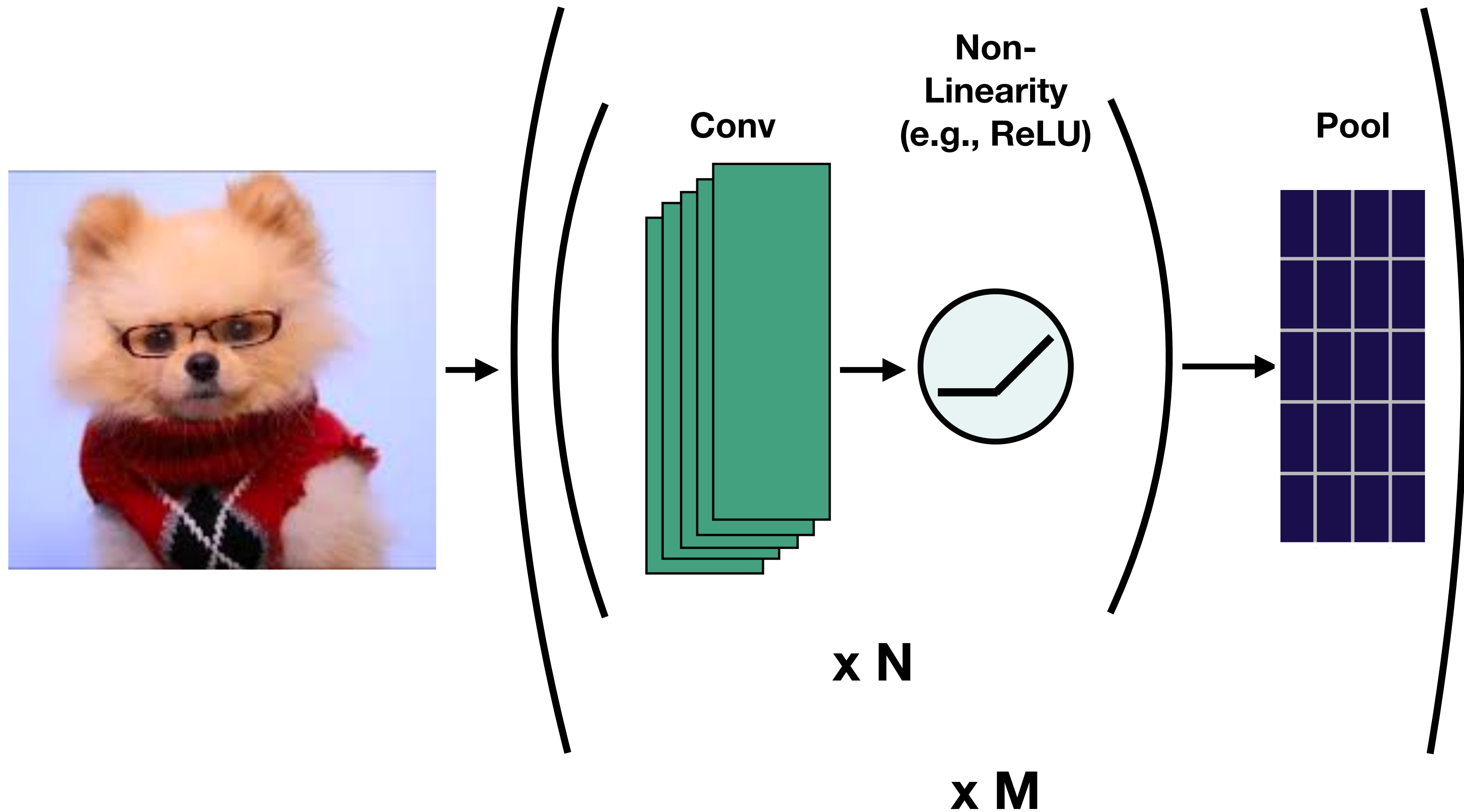
Most standard: LeNet(-like) architectures



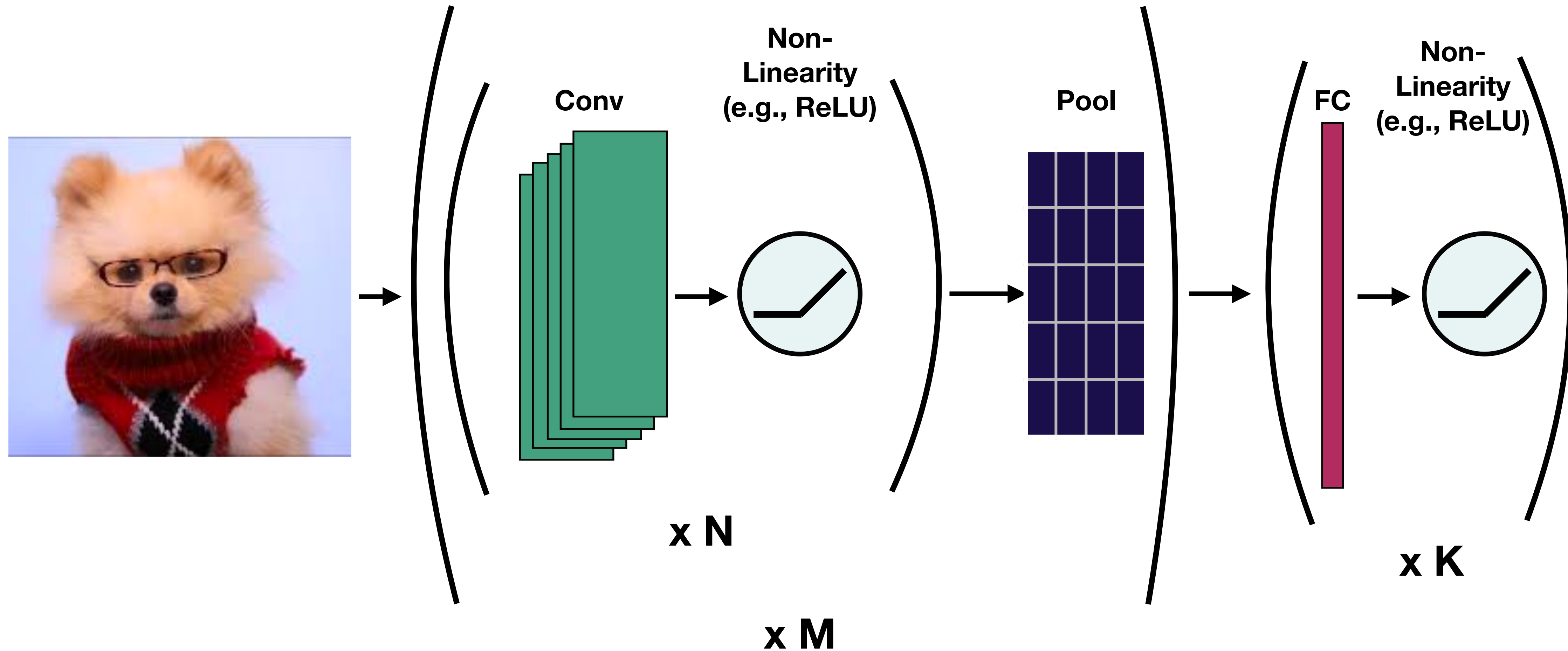
Most standard: LeNet(-like) architectures



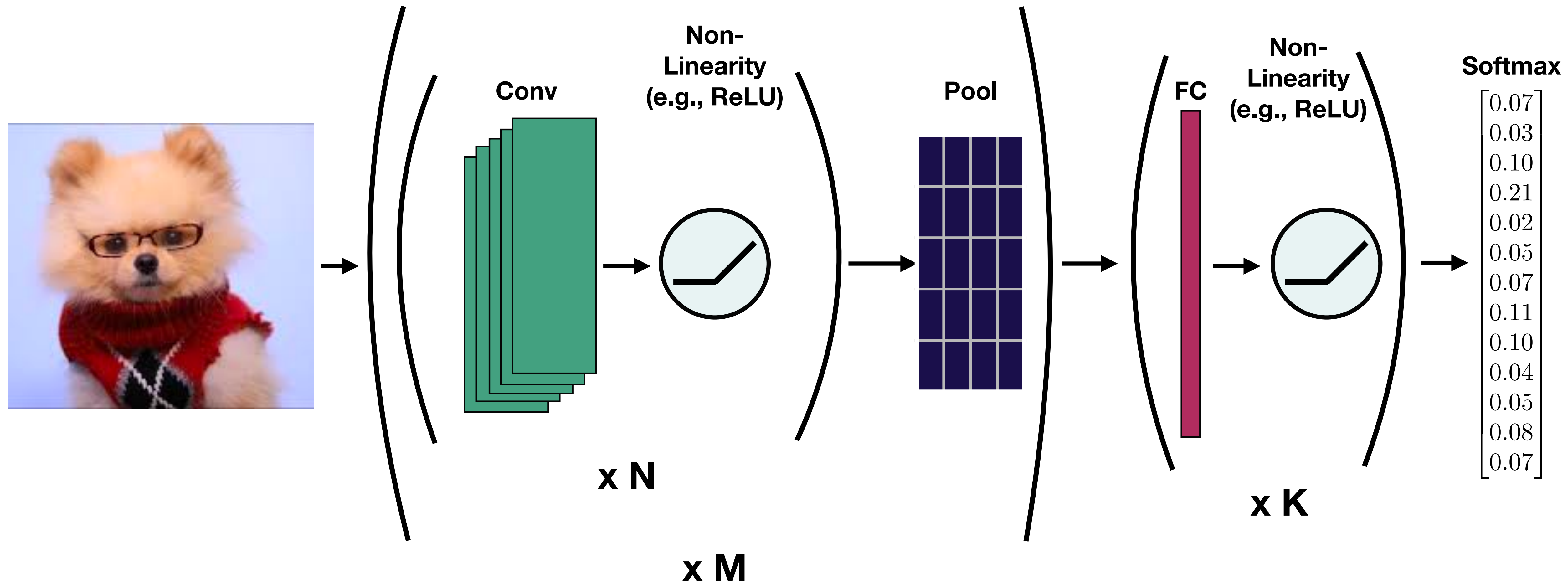
Most standard: LeNet(-like) architectures



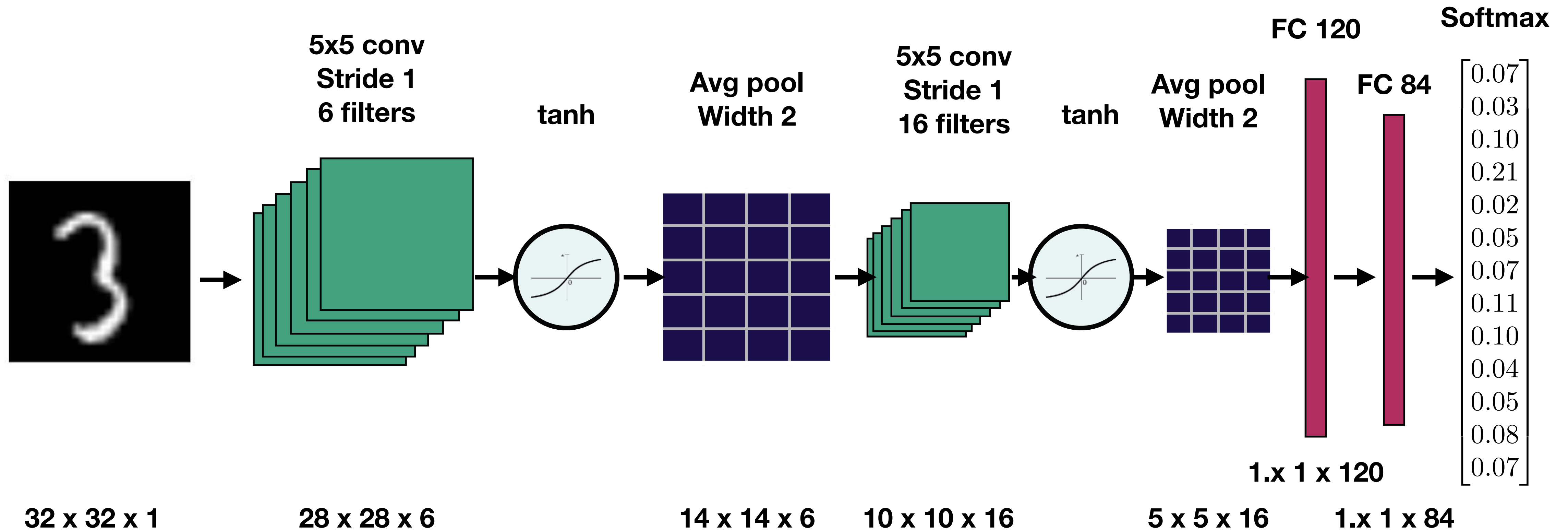
Most standard: LeNet(-like) architectures



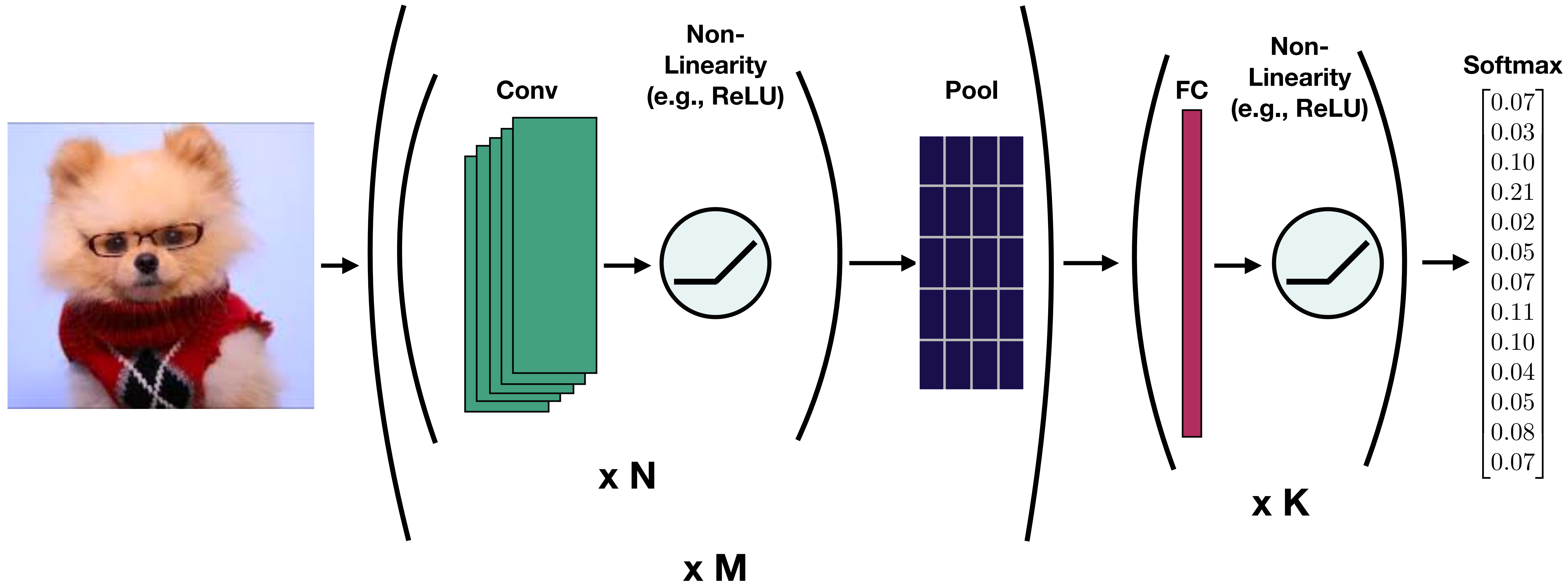
Most standard: LeNet(-like) architectures



Classic Convnet Architecture: LeNet



Most standard: LeNet(-like) architectures



More modern convnet architectures in the Vision Applications lecture!

Questions?