

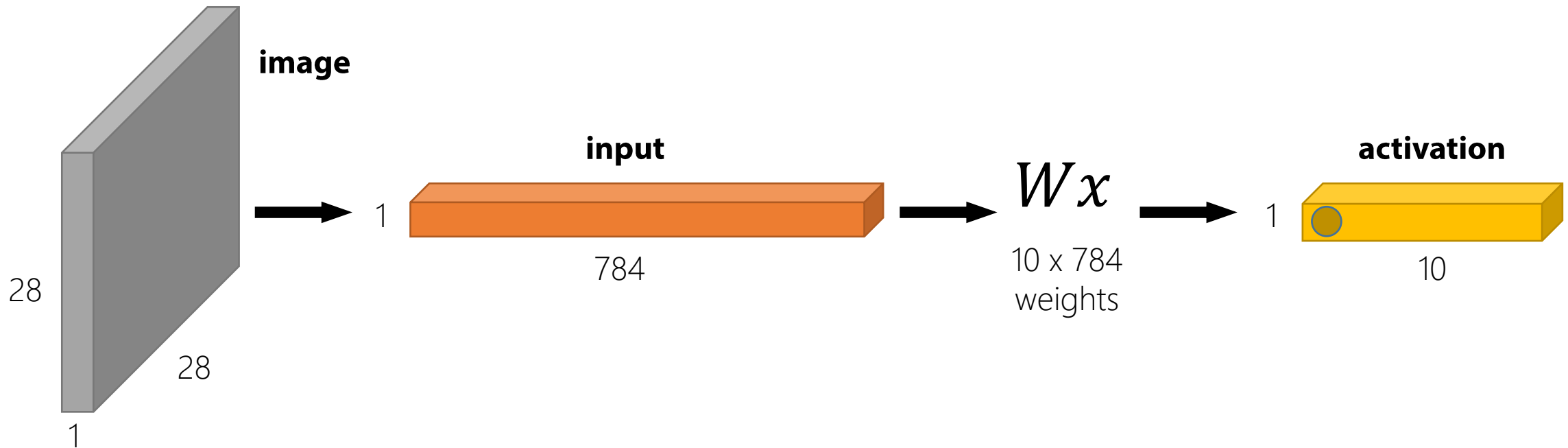


Convolutional Neural Networks

Stephen Baek

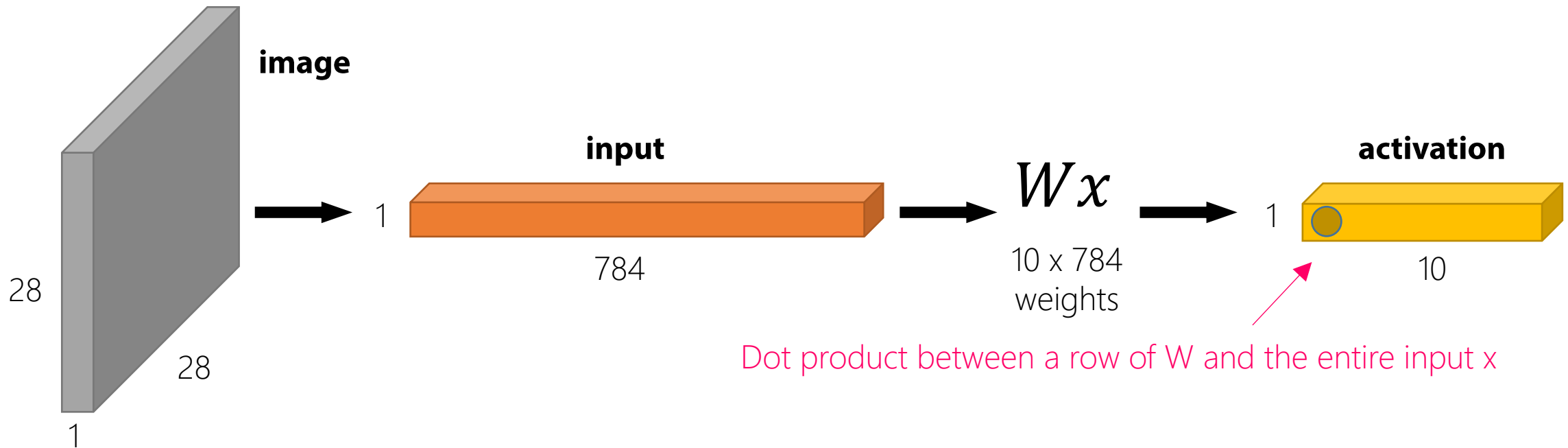
Fully Connected Layer

- 28 x 28 image → stretch to 784 x 1
- 64 x 64 x 3 image → stretch to 12288 x 1
- ...



Fully Connected Layer

- 28 x 28 image → stretch to 784 x 1
- 64 x 64 x 3 image → stretch to 12288 x 1
- ...



Draw your number here

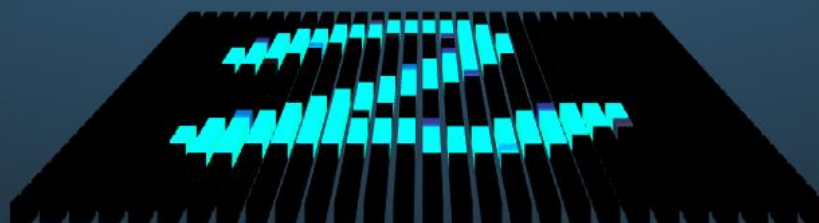


Downsampled drawing:

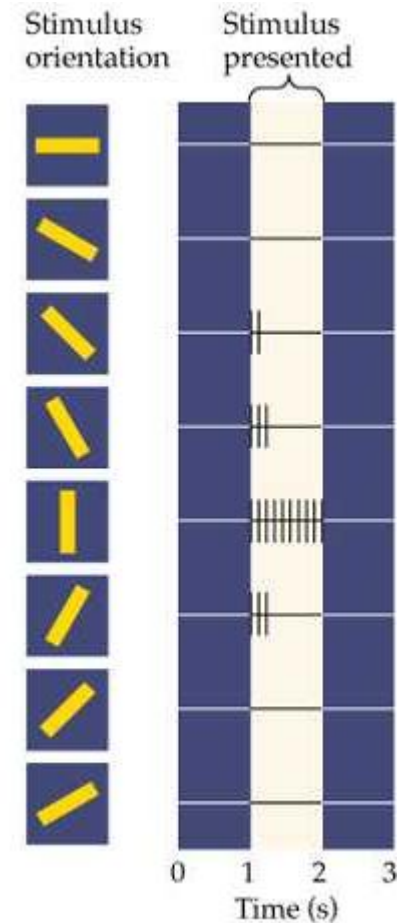
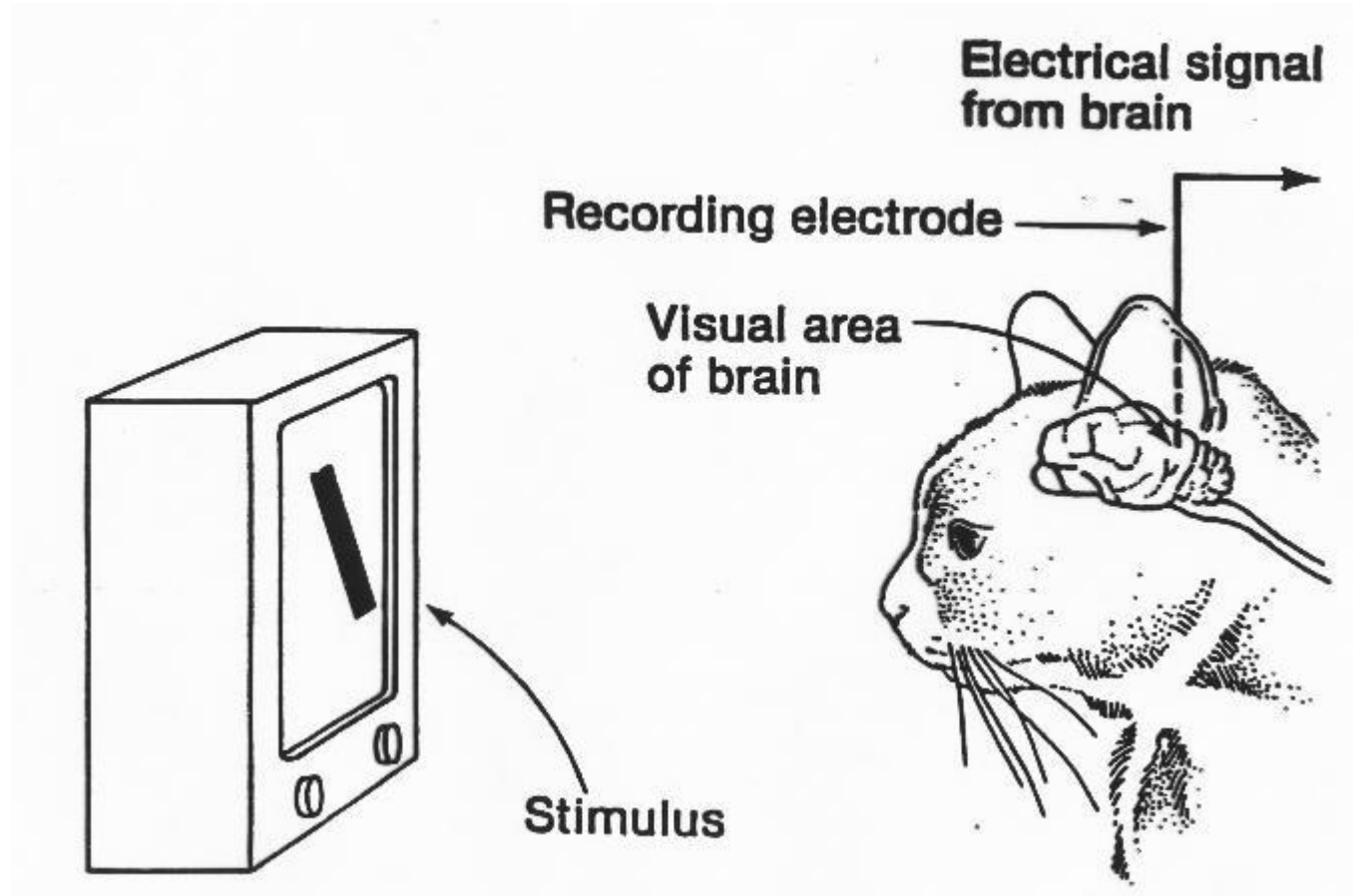
First guess:

Second guess:

0 1 2 3 4 5 6 7 8 9
■ ■ ■ ■ ■ ■ ■ ■ ■ ■



Hubel & Wiesel (1959 ~)



Neurons in the visual cortex respond selectively to oriented edges. Neurons in visual cortex typically respond vigorously to a bar of light oriented at a particular angle and weakly or not at all to other orientations.

Hubel & Wiesel (1959 ~)

Warning!! Visually Disturbing

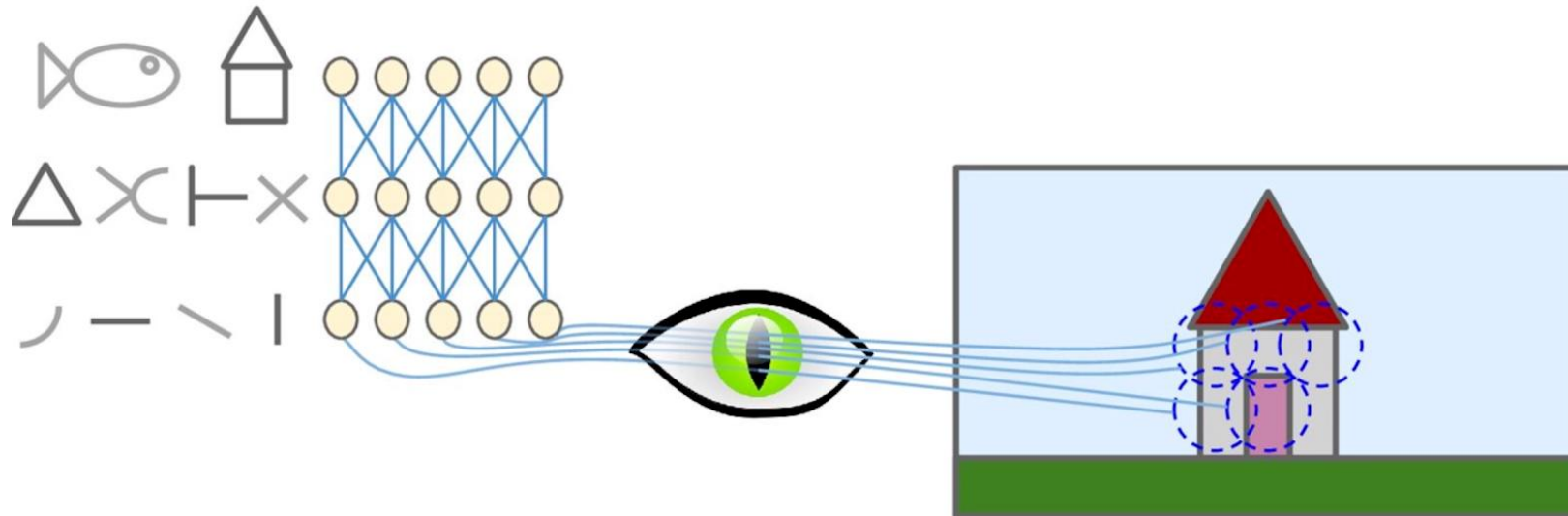
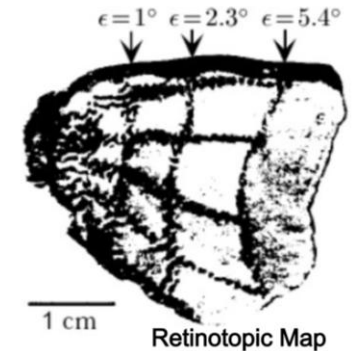


Topographic Maps in Cortex

- Each visual sensitive cell only responds to stimuli of a limited region (receptive field)
- Neighboring cells have partially overlapping receptive fields
- Neighboring points in a visual image evoke activity in neighboring regions of visual cortex
- In this manner, the visual system easily maintain the information of the spatial location of stimulus



(Dayan and Abbott 2001)



The human visual system

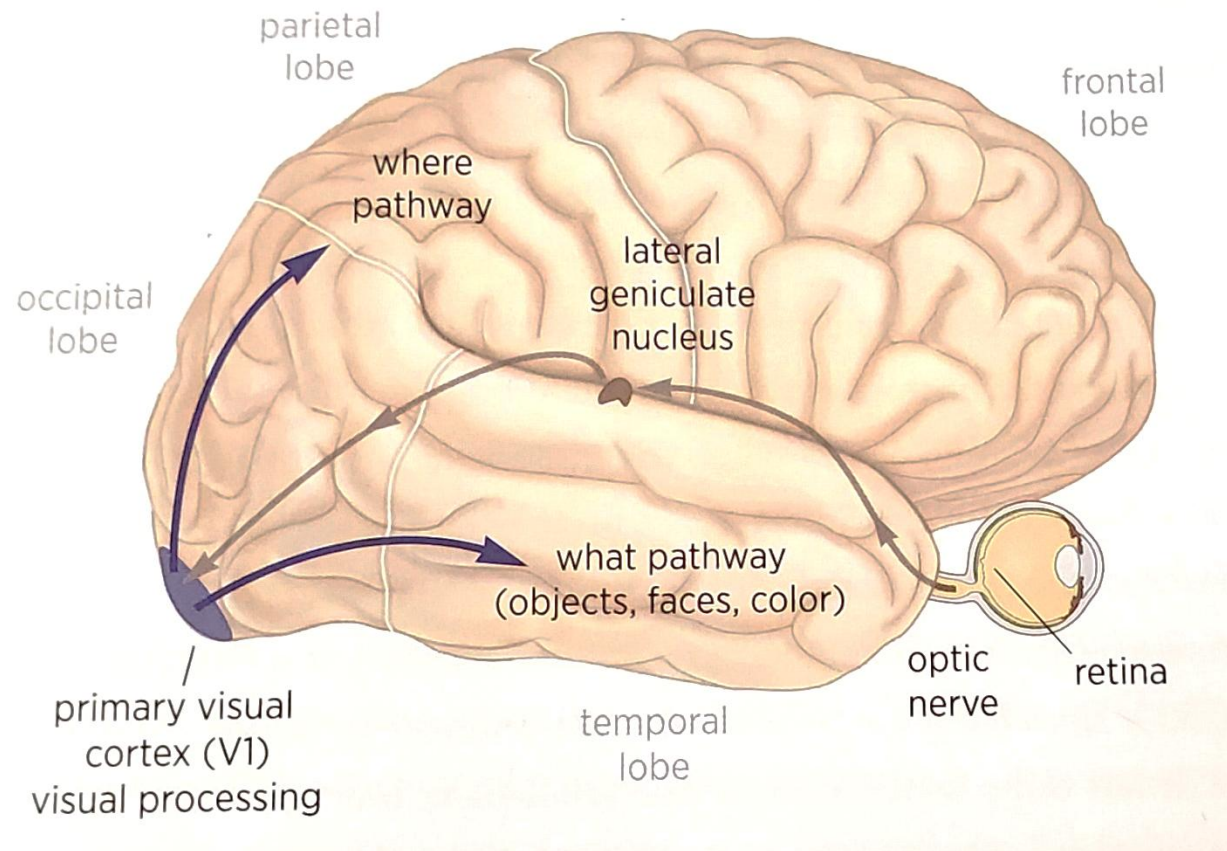
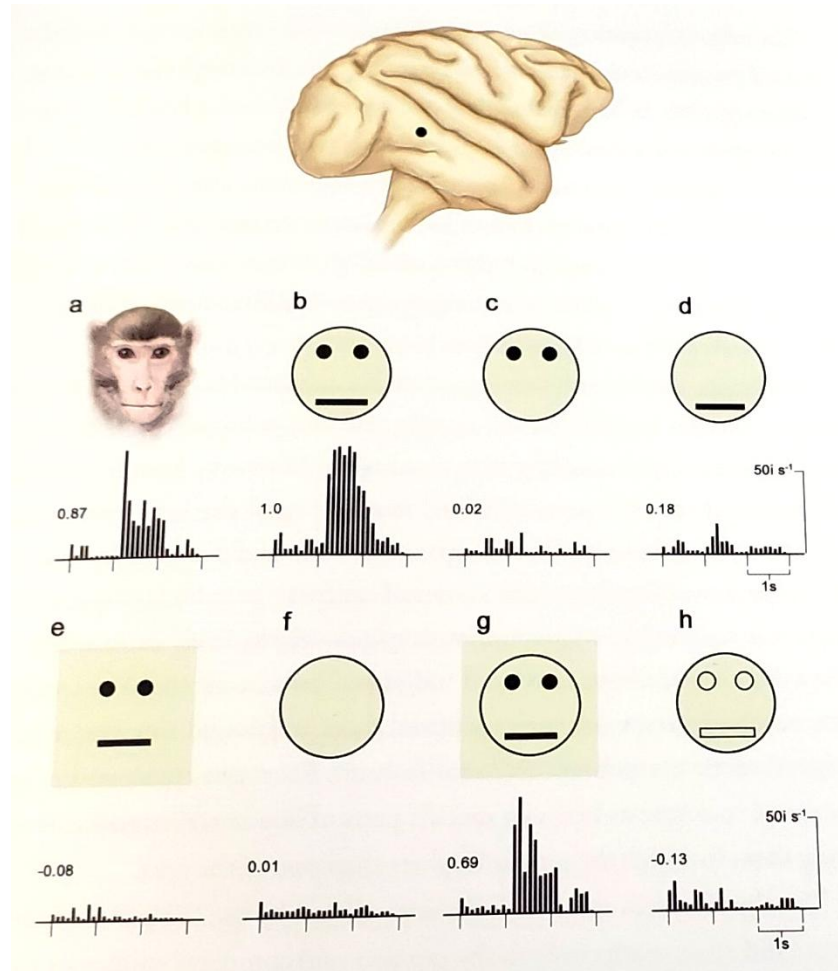


Image Courtesy: Kandel, "Reductionism in Art and Brain Science," 2016

- Retina: visual input
- Retina → Lateral Geniculate Nucleus
 - Visual information flows through the optic nerve
- Lateral Geniculate Nucleus (LGN):
 - A small, ovoid object at the end of the optic tract
 - One on each side of the brain
 - In humans, each LGN has six layers of neurons
 - Sends information to the primary visual cortex (V1)

e.g. Facial Recognition

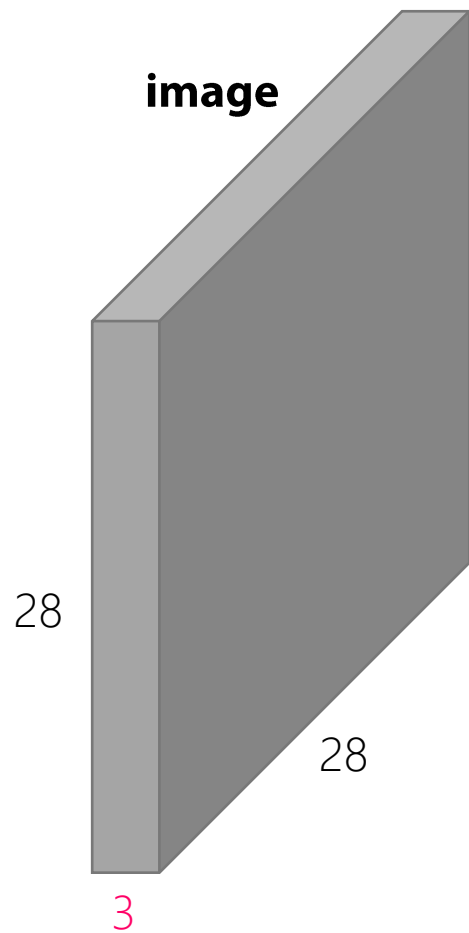


- Holistic face detection
 - “Face cell” in the inferior temporal cortex
 - Fires when there is a face-like object

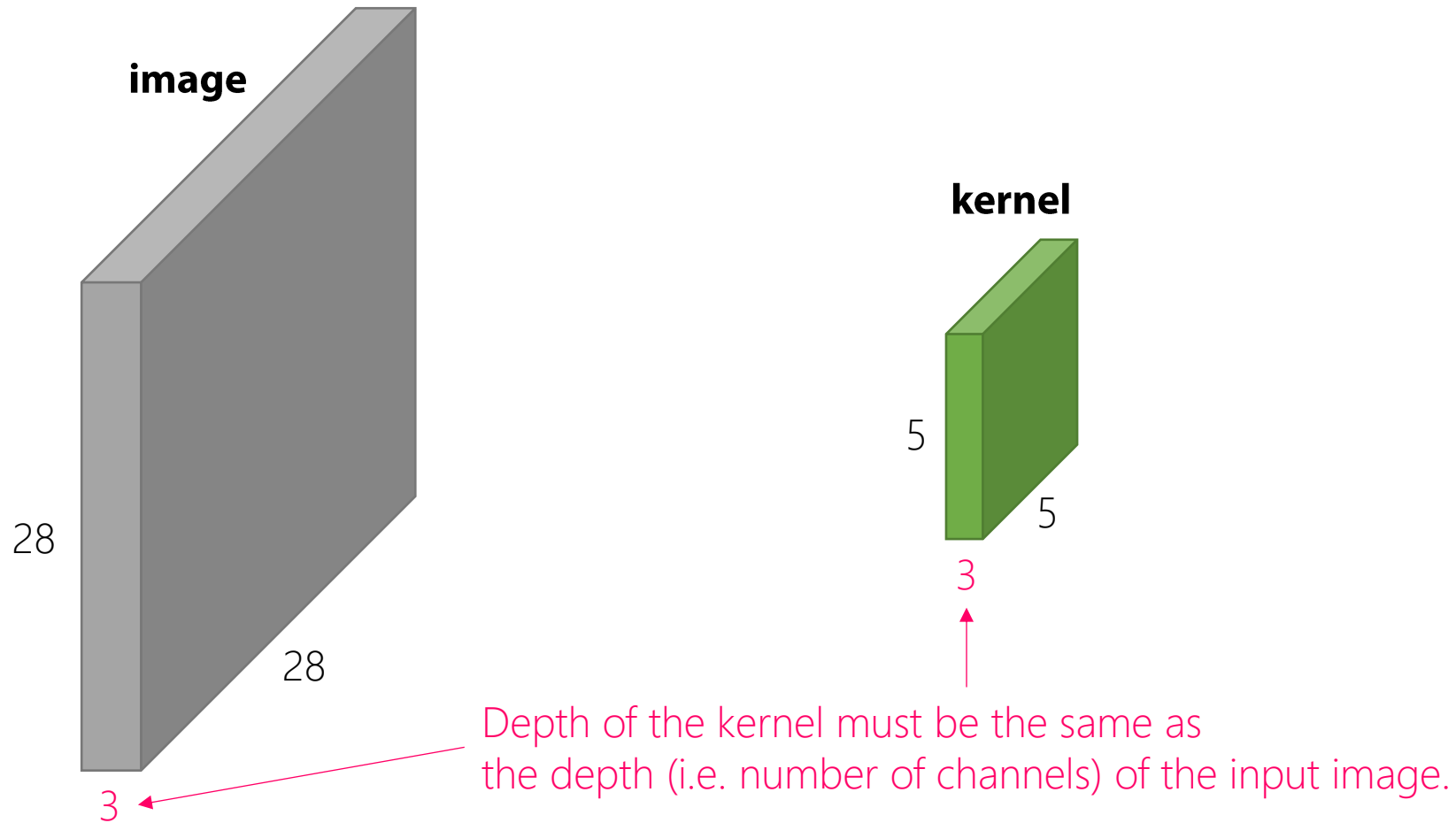
Convolutional Neural Networks

- Key idea:
 - As like how we (humans) understand a visual scene, if neural nets could see small pieces, understand patterns and textures, combine the pieces to see a bigger picture, computers should be able to recognize images.
- A bonus:
 - Typical neural networks are “fully connected”.
 - In an image domain, this means all the pixels are interconnected.
 - However, pixels far apart have no significant meaning...
 - By connecting only the nearing neighbors, computational load could be much lower.

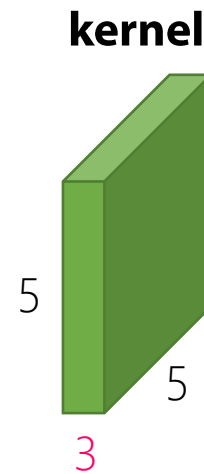
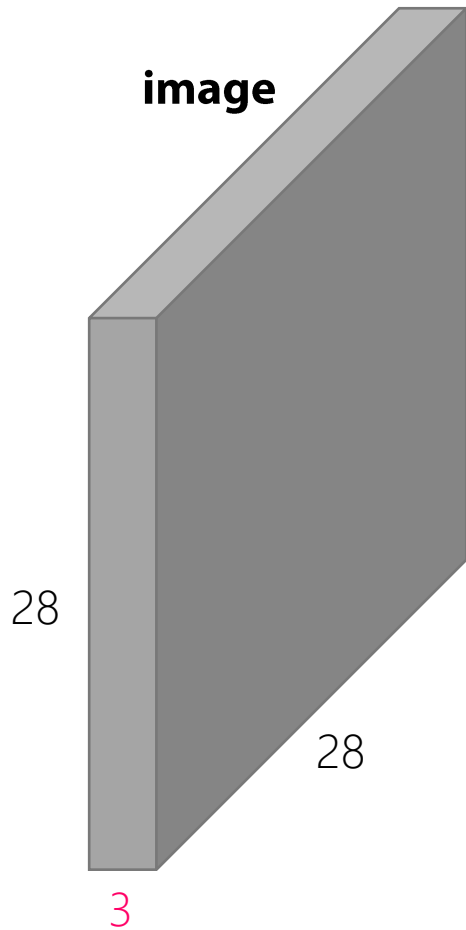
Convolution Layer



Convolution Layer

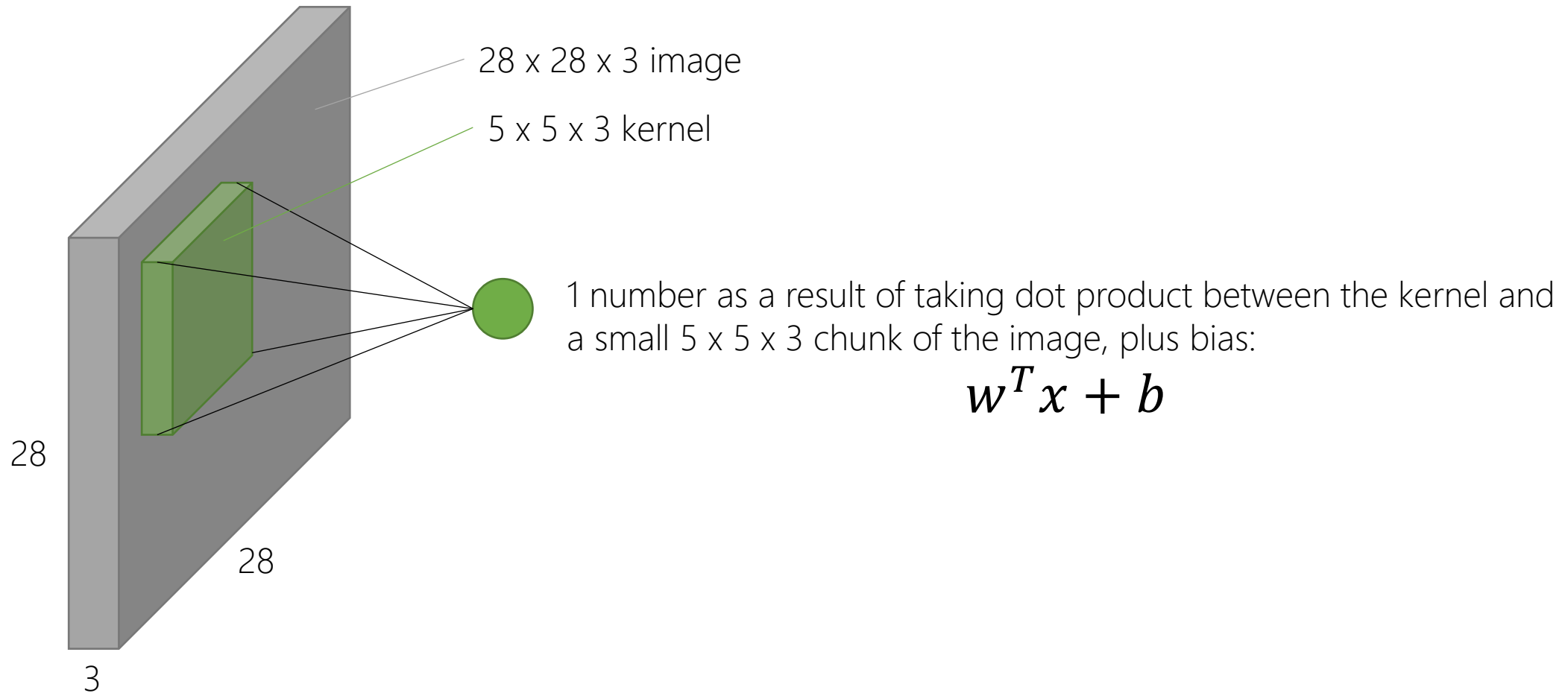


Convolution Layer

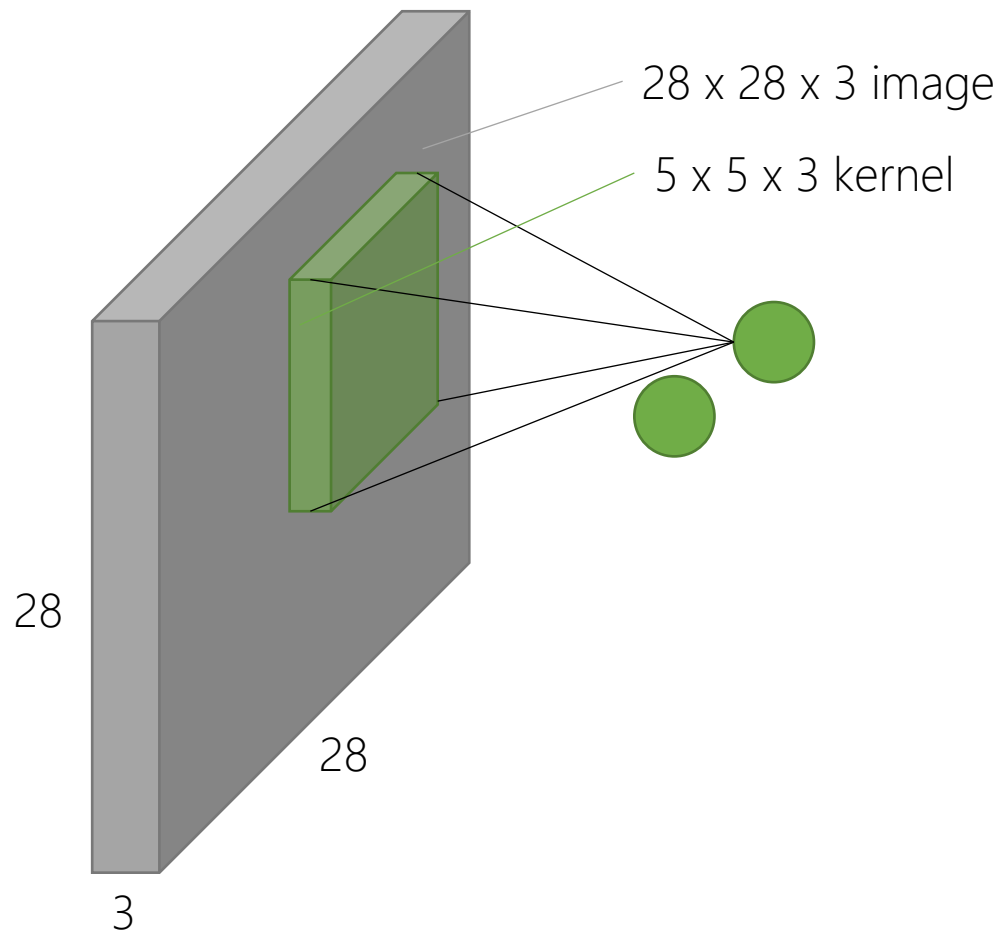


Convolve the kernel with the image!
In other words...
"slide the kernel over the image,
compute dot products each time."

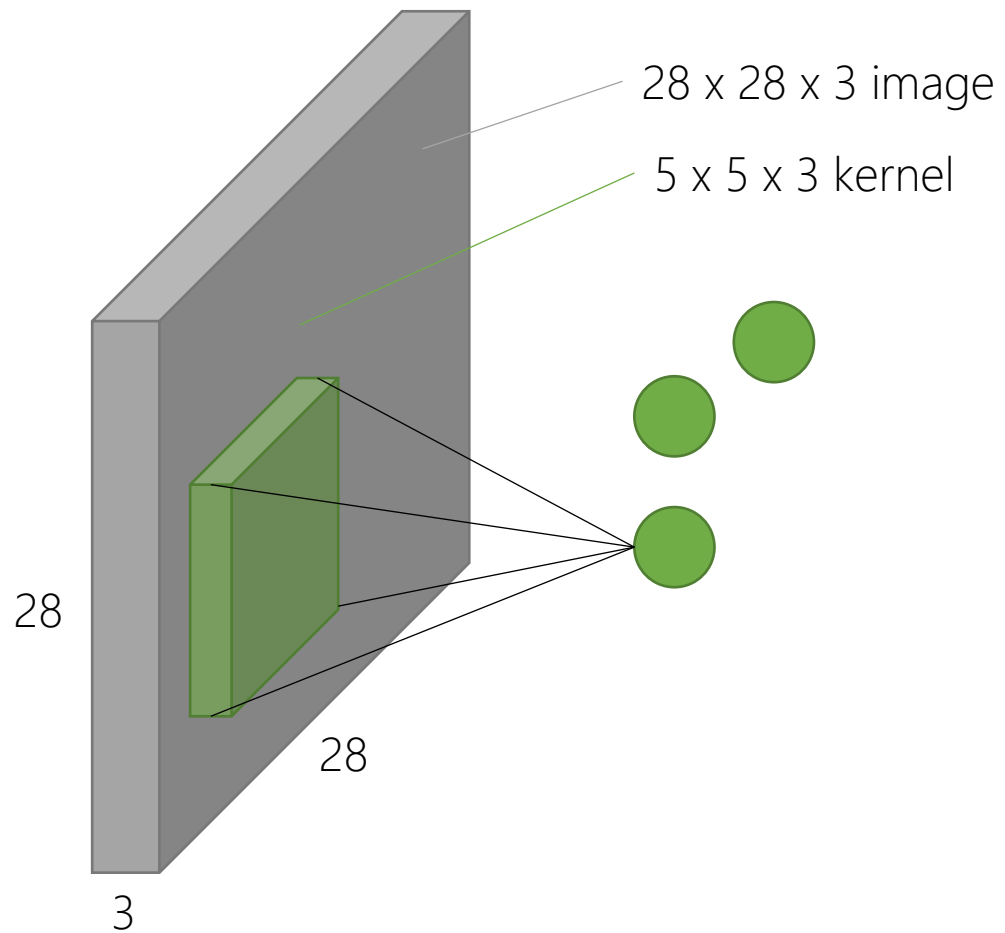
Convolution Layer



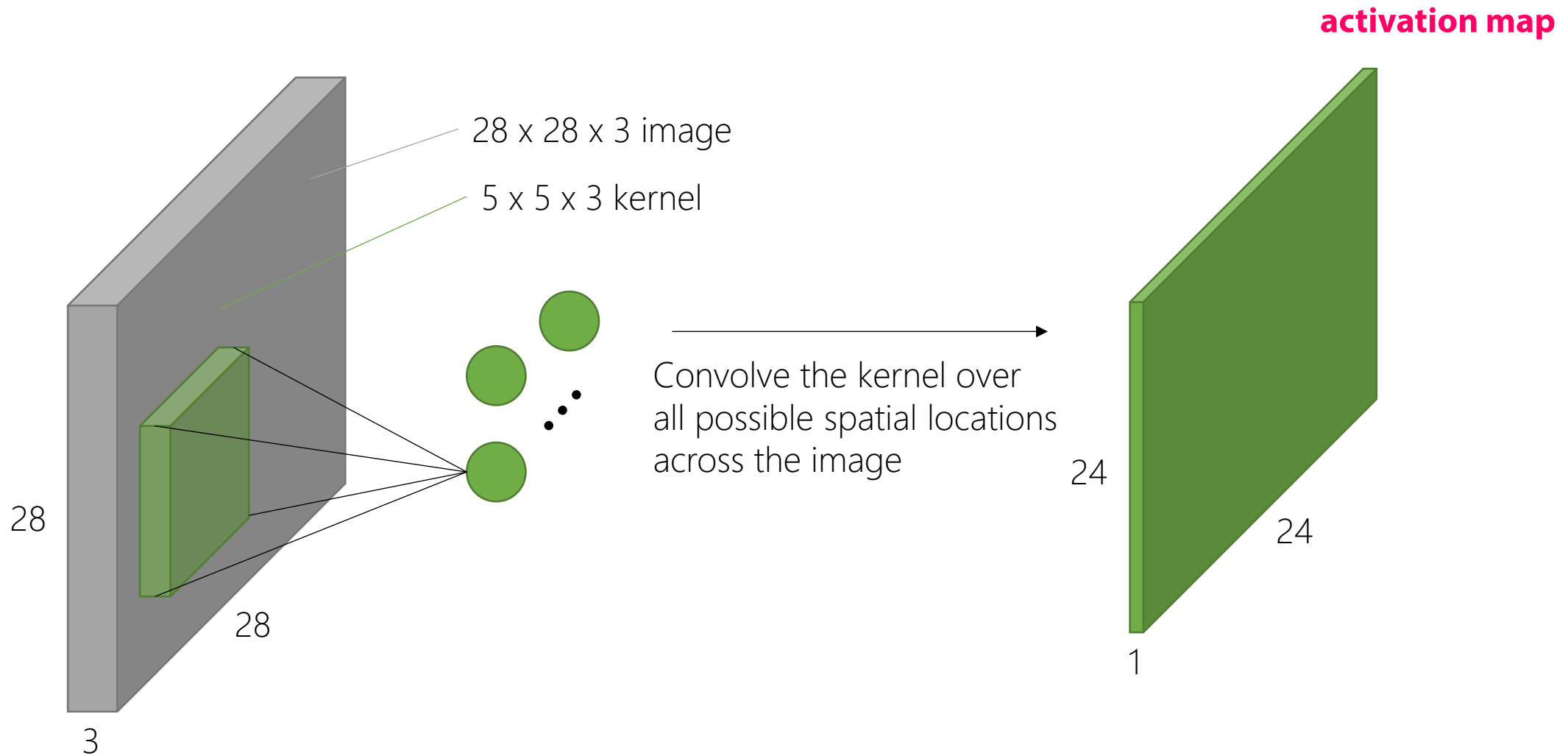
Convolution Layer



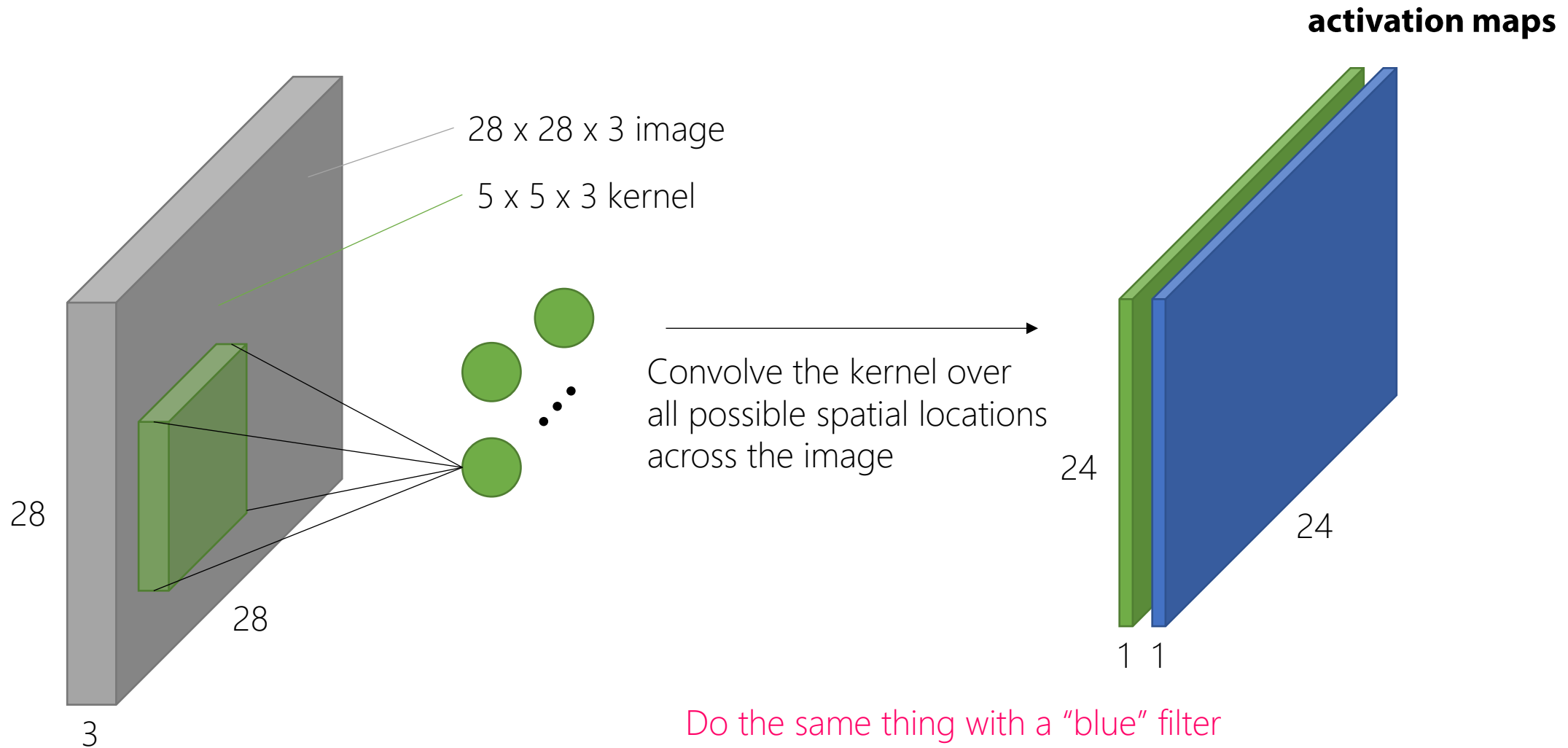
Convolution Layer



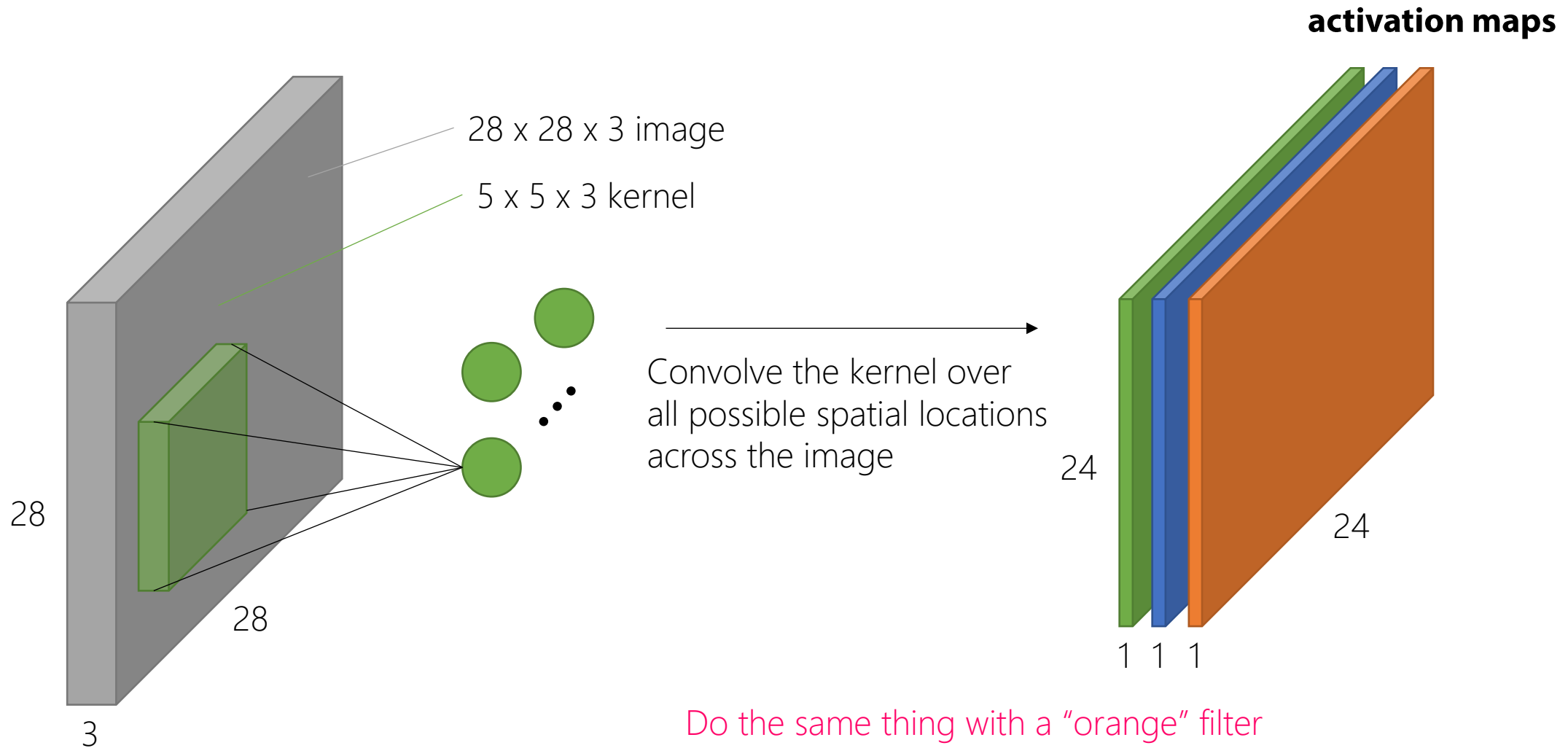
Convolution Layer



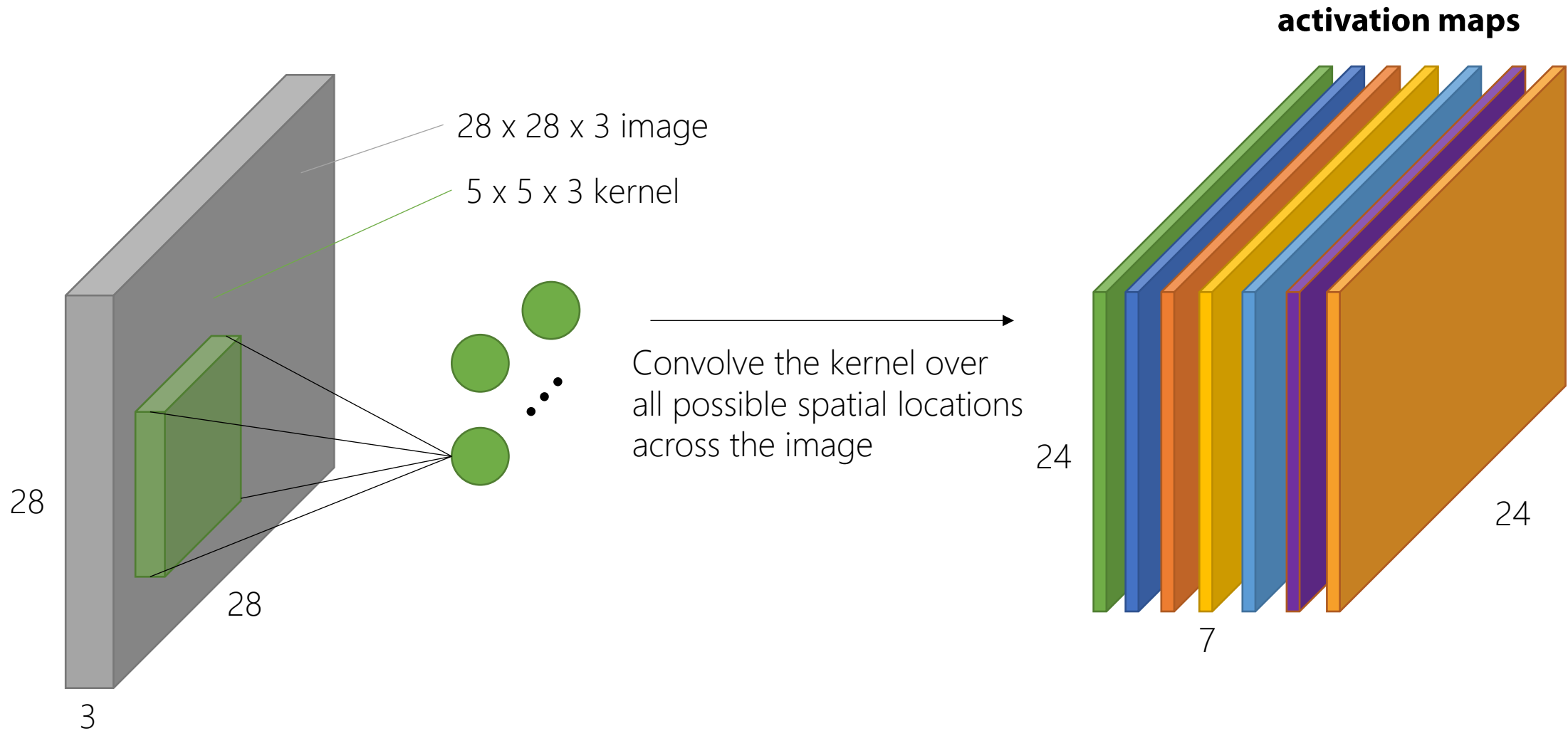
Convolution Layer



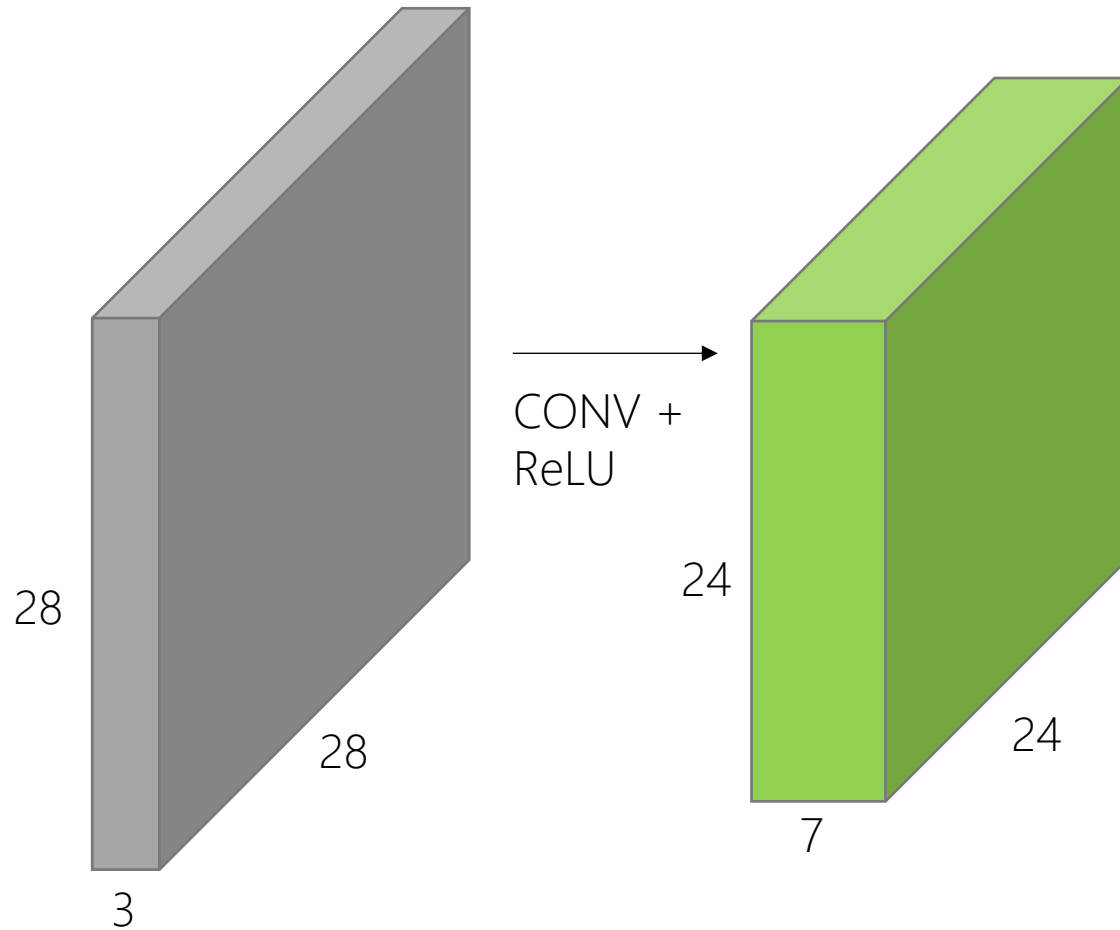
Convolution Layer



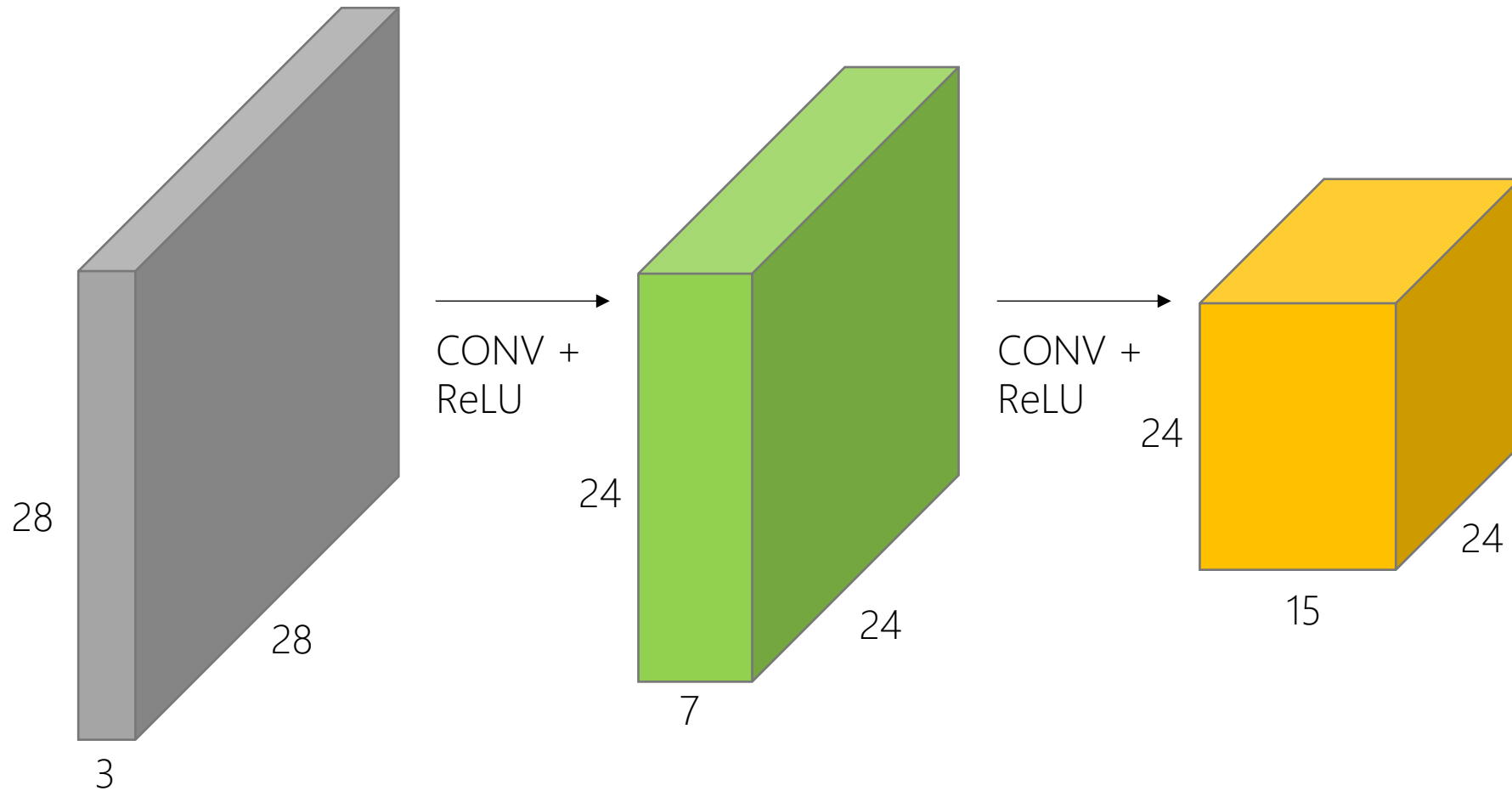
Convolution Layer



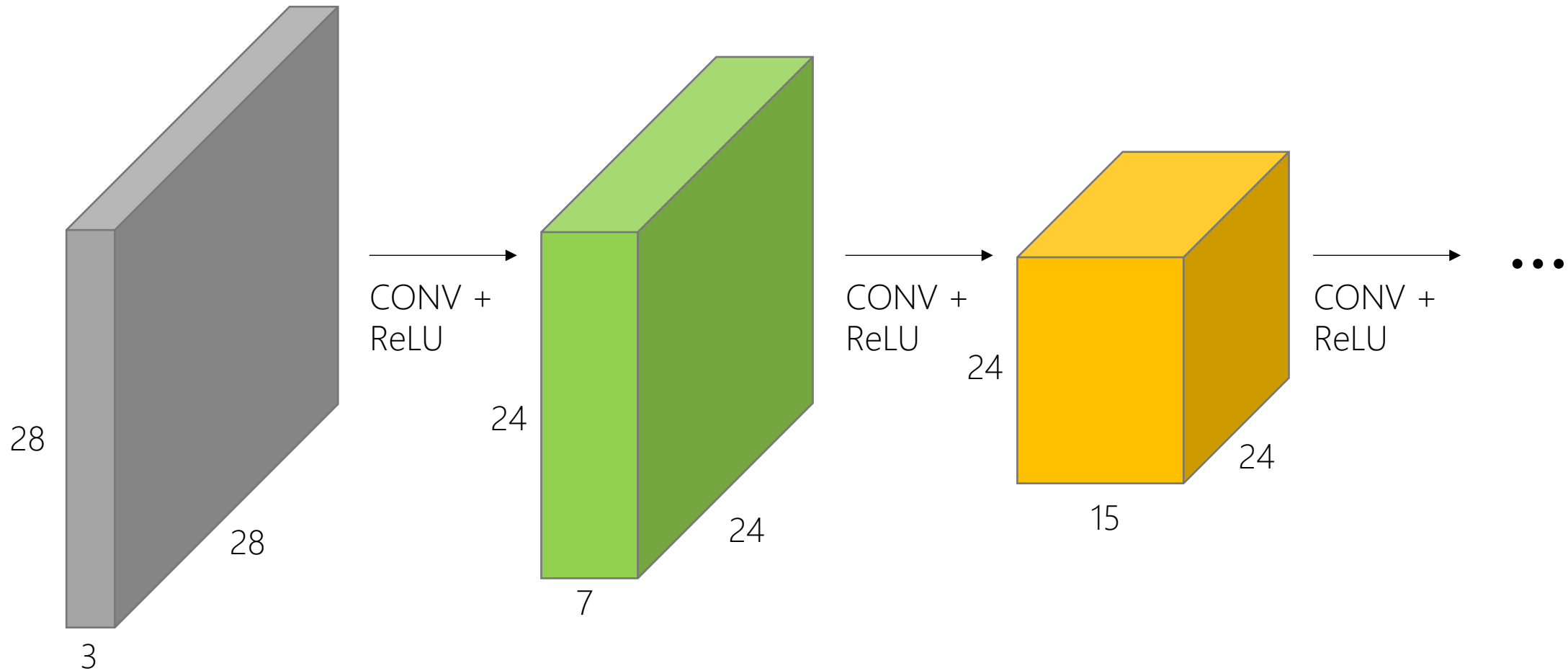
ConvNet is a sequence of convolution layers



ConvNet is a sequence of convolution layers



ConvNet is a sequence of convolution layers



Draw your number here

2

0 1 2 3 4 5 6 7 8 9



Downsampled drawing: 2

First guess: 2

Second guess: 1

Layer visibility

Input layer ☐ Show

Convolution layer 1 ☐ Show

Downsampling layer 1 ☐ Show

Convolution layer 2 ☐ Show

Downsampling layer 2 ☐ Show

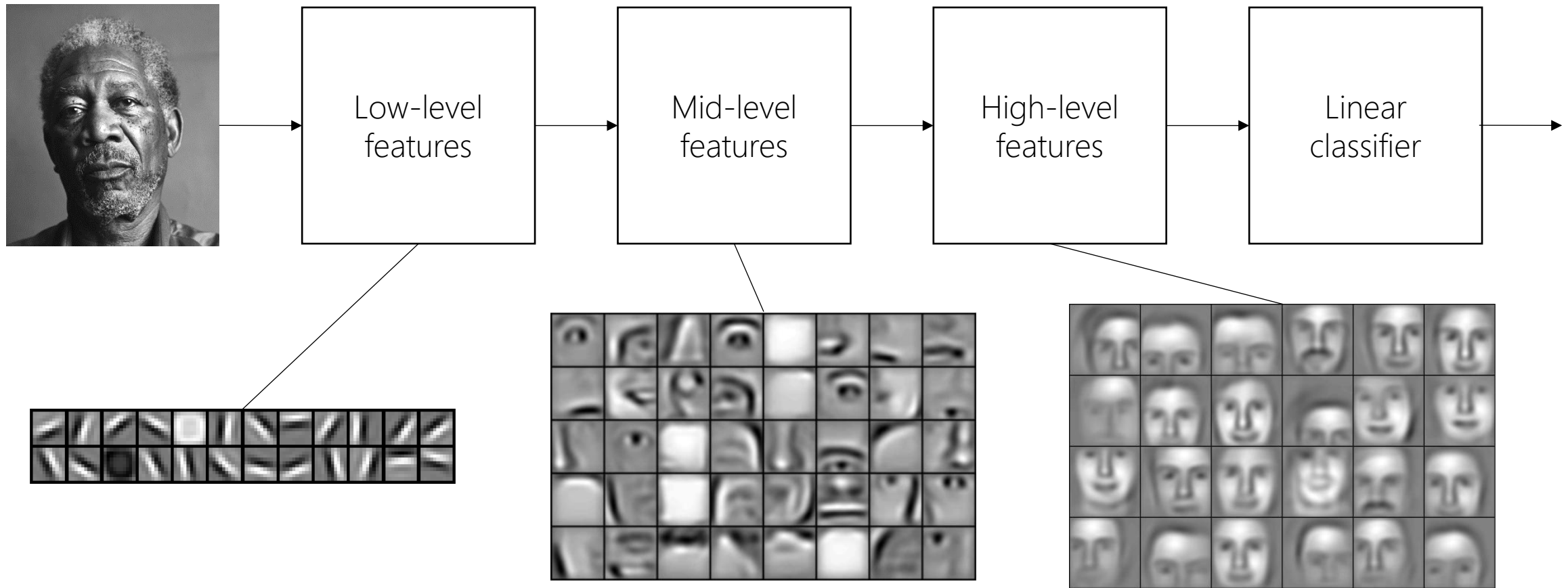
Fully-connected layer 1 ☐ Show

Fully-connected layer 2 ☐ Show

Output layer ☐ Show

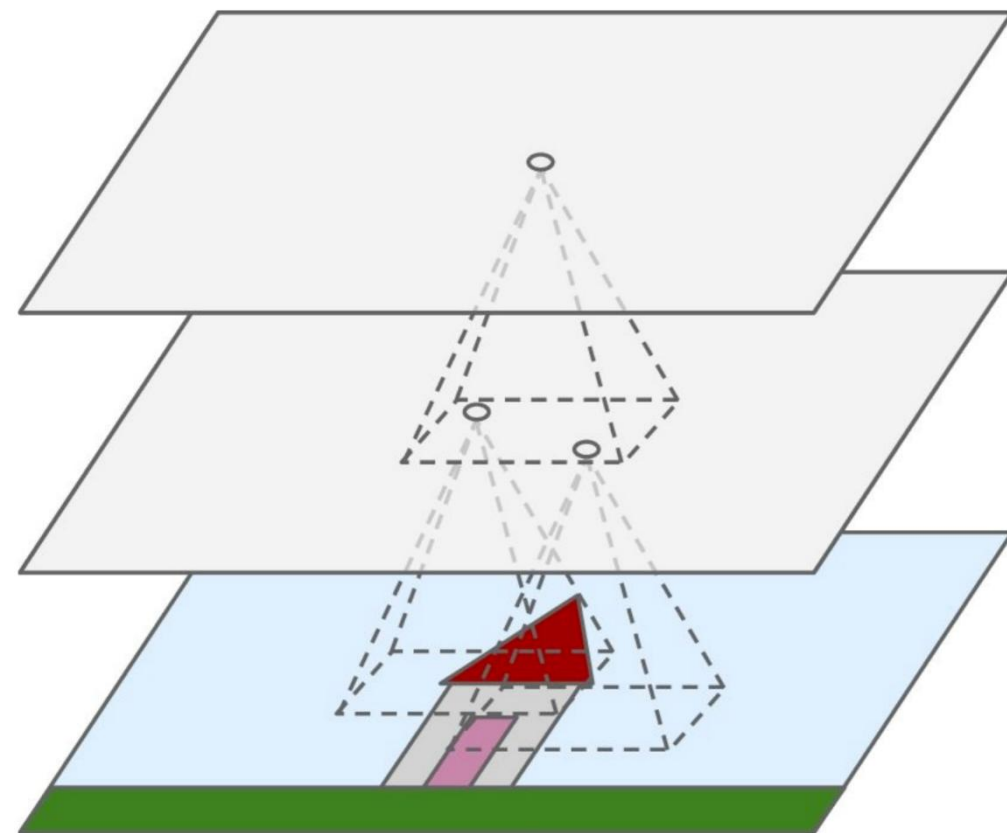
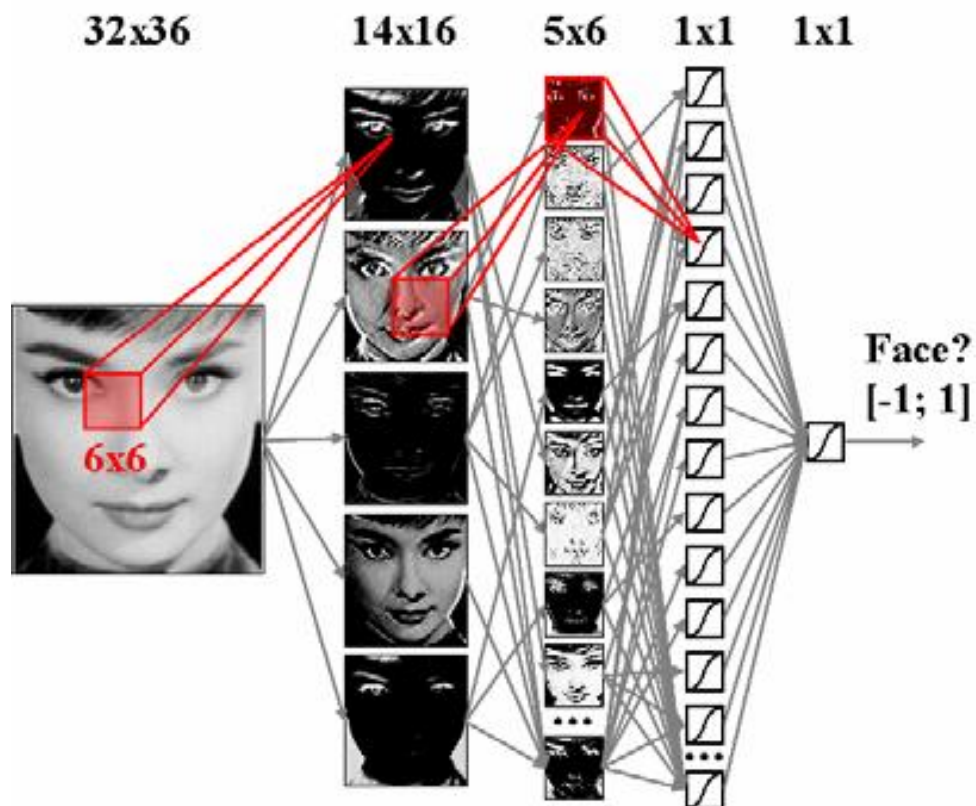
Made by [Adam Harley](#). [Project details](#)

ConvNet is a sequence of convolution layers

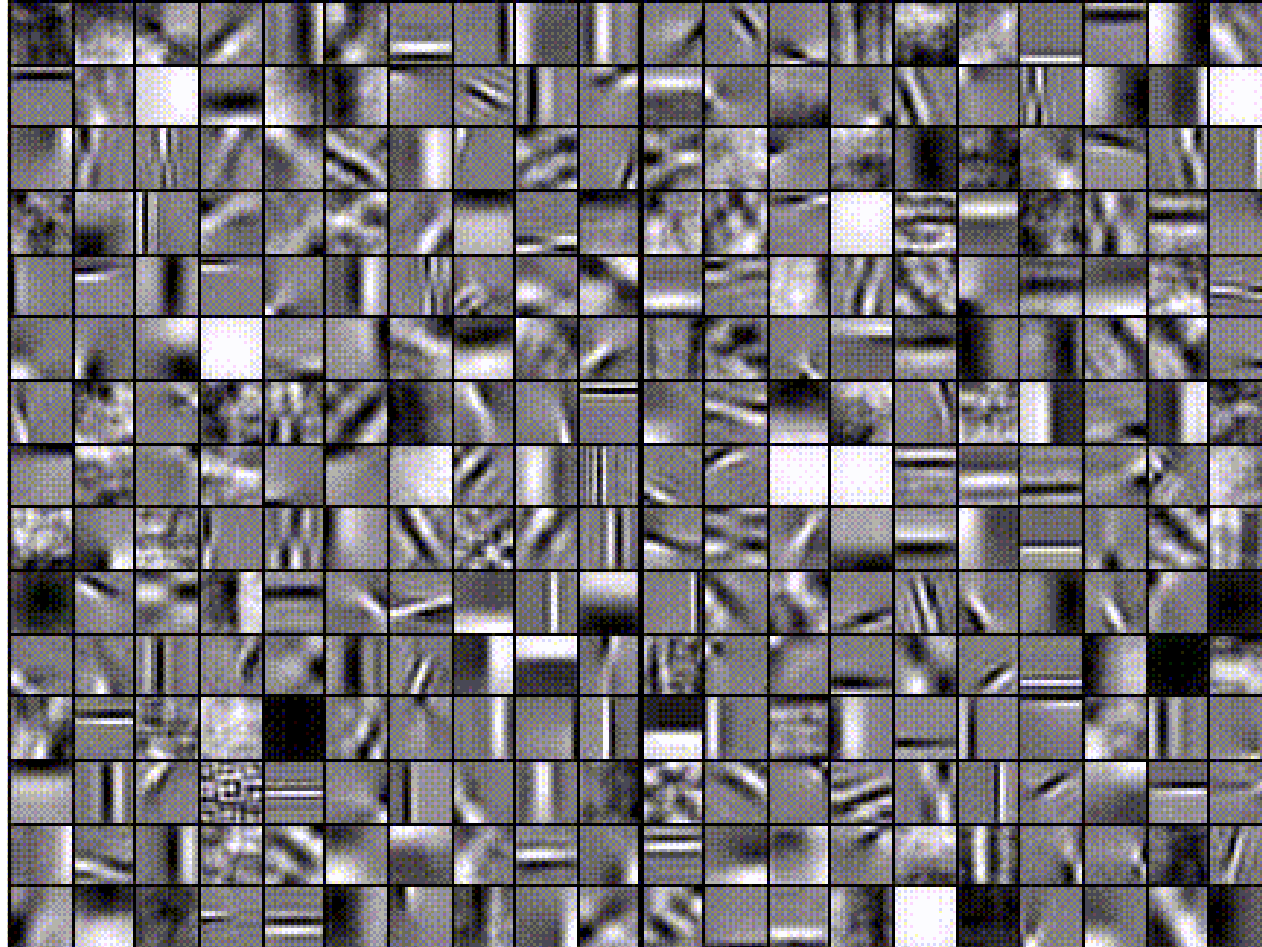


Lee et al. (2009)

ConvNet is a sequence of convolution layers



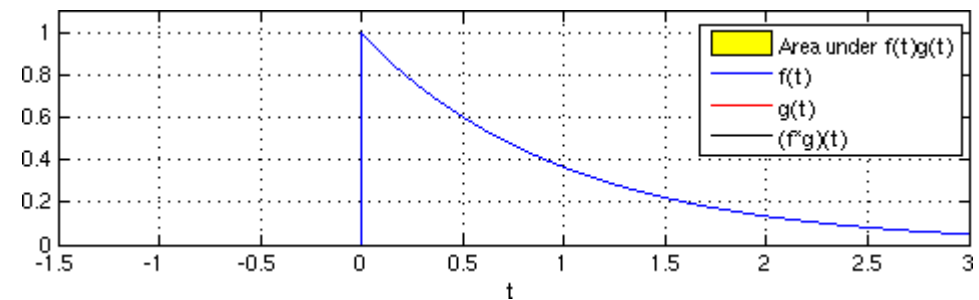
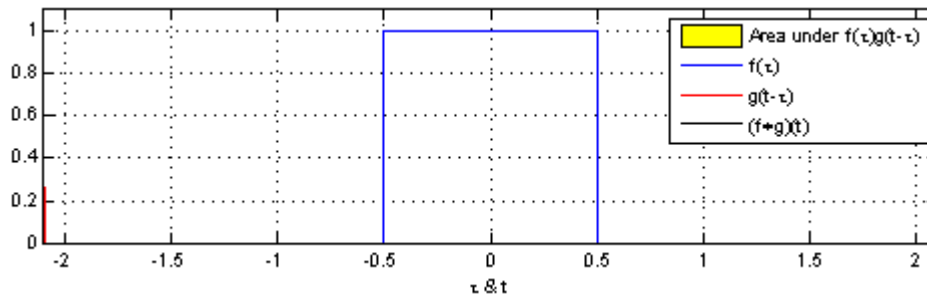
Conv kernels are trainable



A closer look...

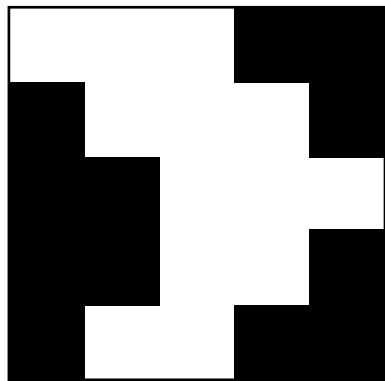
- Convolution

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



A closer look...

- 2D Discrete Convolution



*

$$\begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix}$$



1 _{x1}	1 _{x0}	1 _{x1}	0	0
0 _{x0}	1 _{x1}	1 _{x0}	1	0
0 _{x1}	0 _{x0}	1 _{x1}	1	1
0	0	1	1	0
0	1	1	0	0

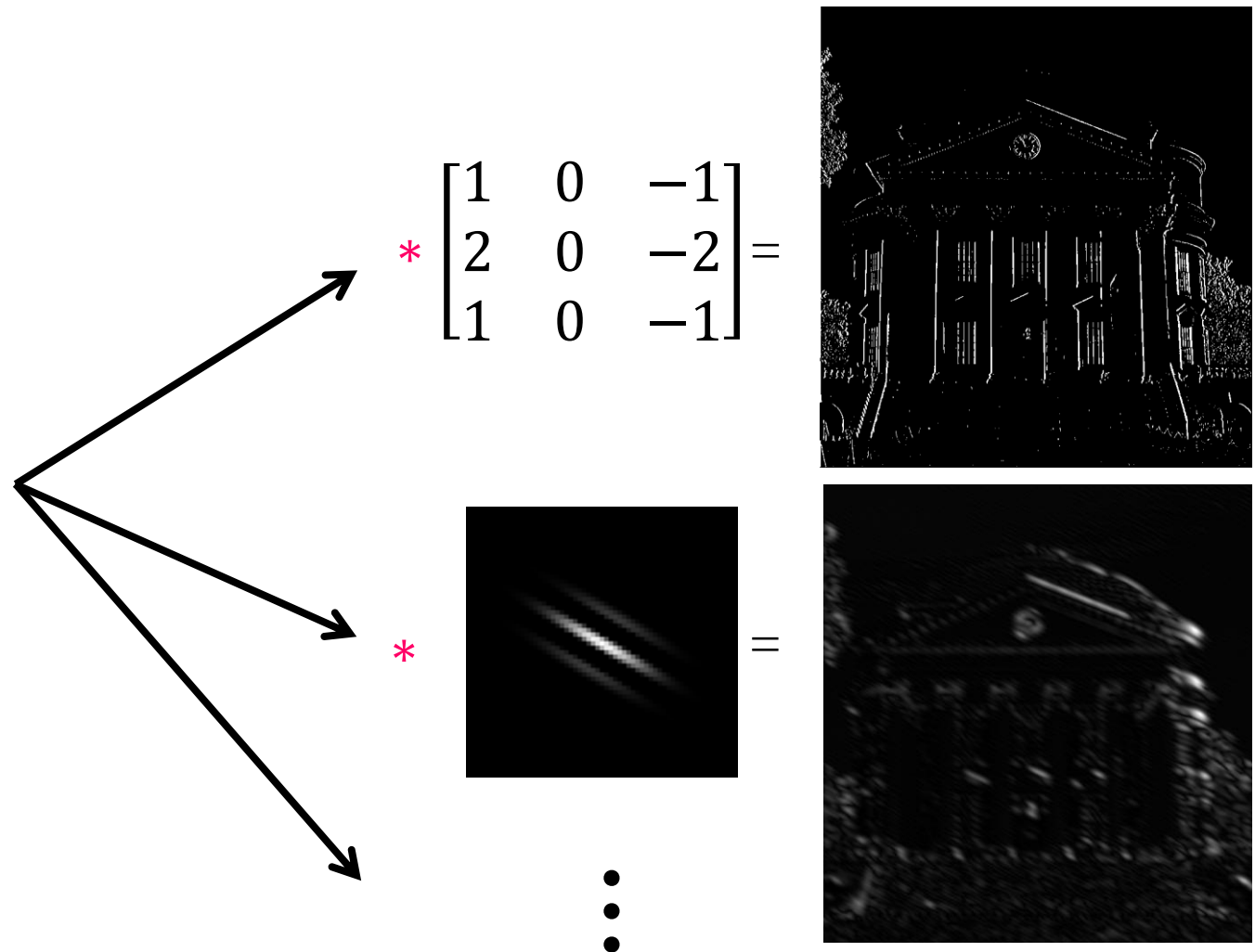
Image

4		

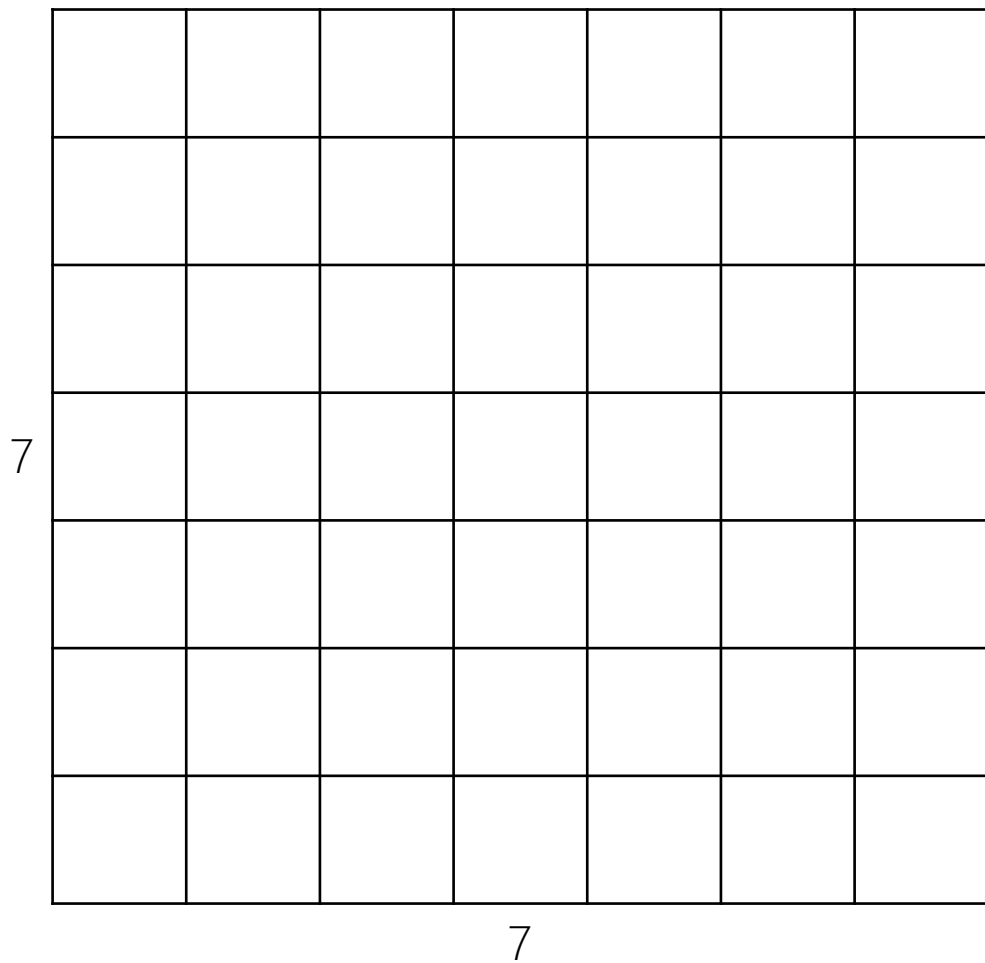
Convolved
Feature

A closer look...

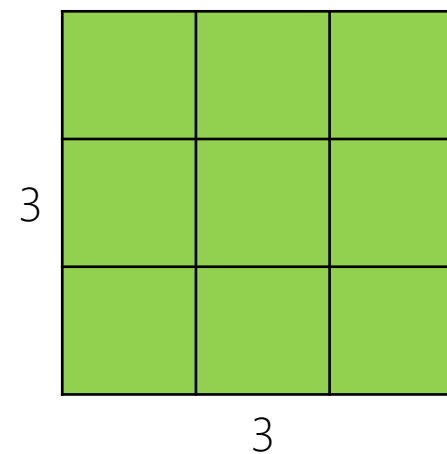
- 2D Discrete Convolution



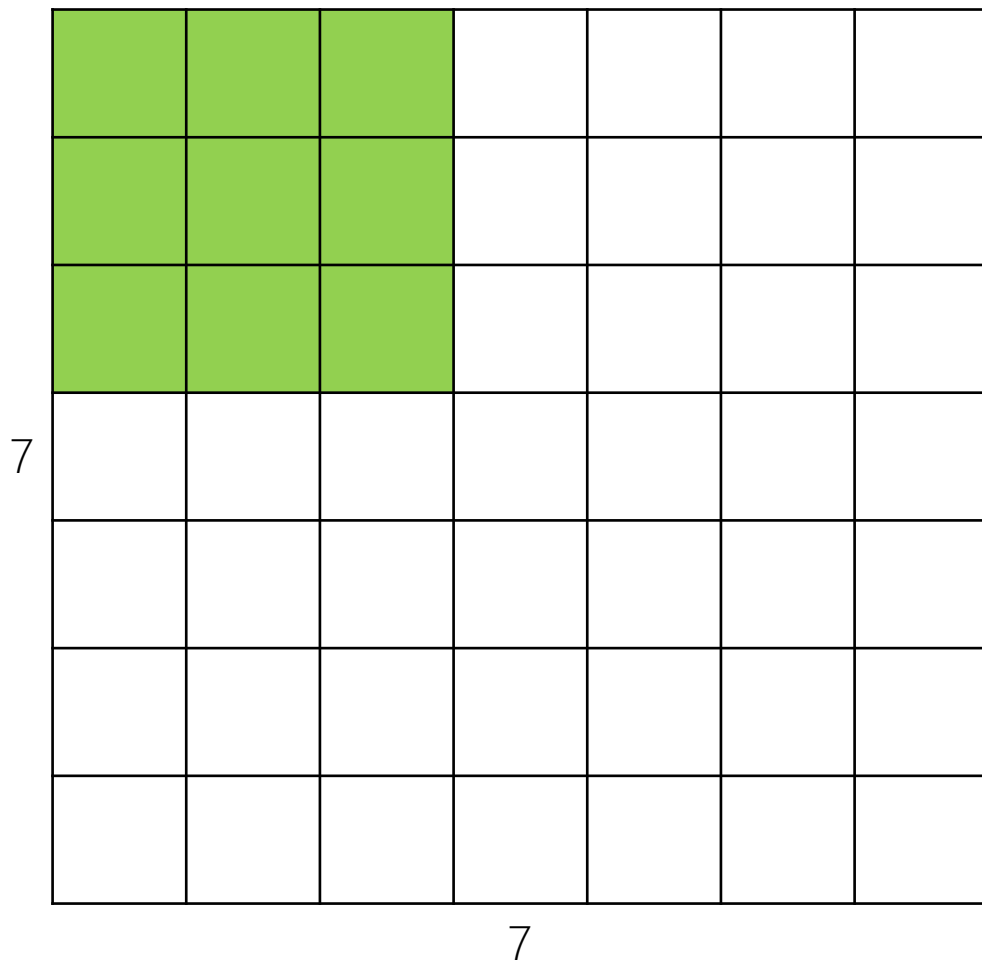
Stride & Padding



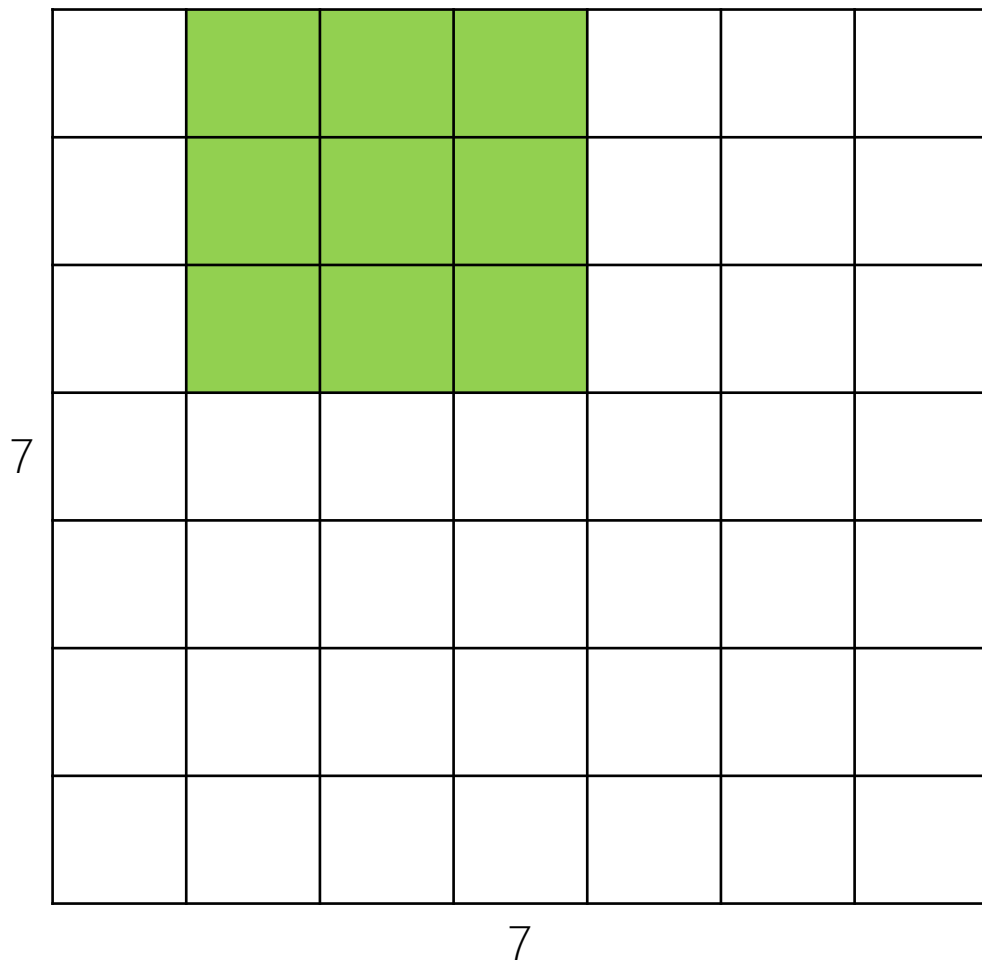
*



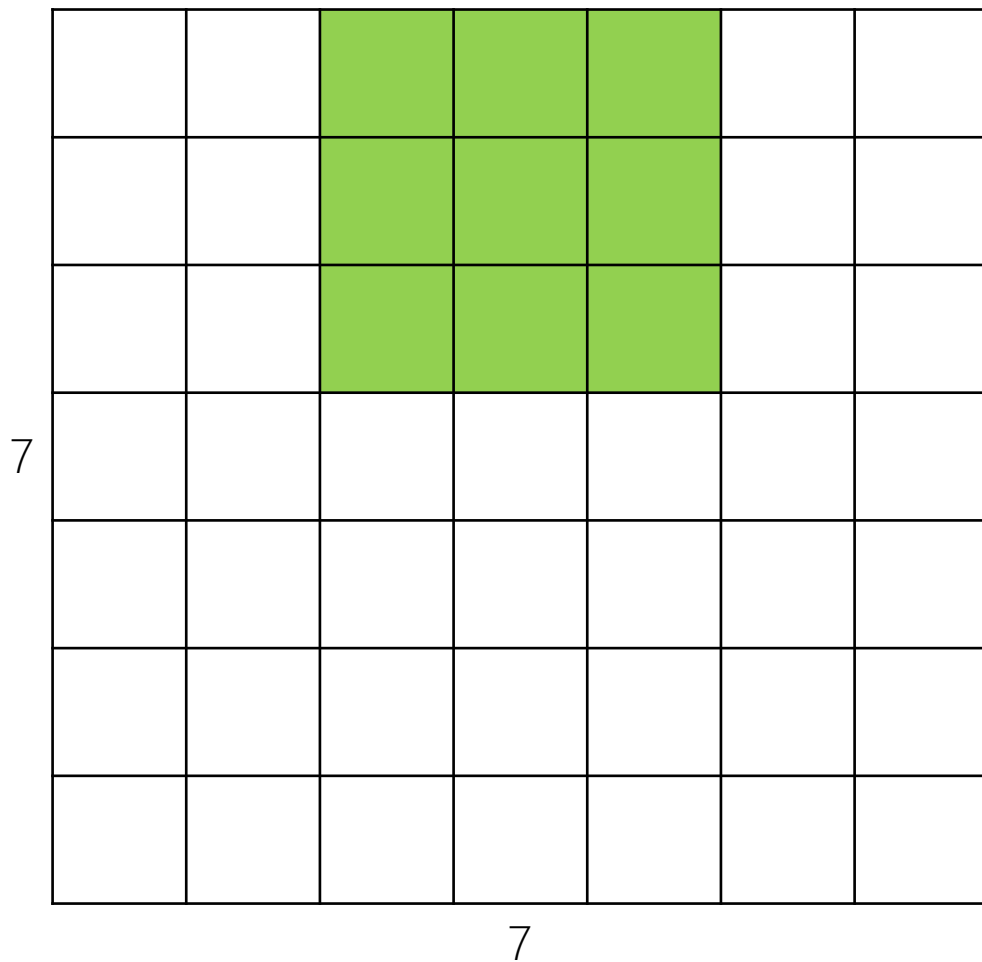
Stride & Padding



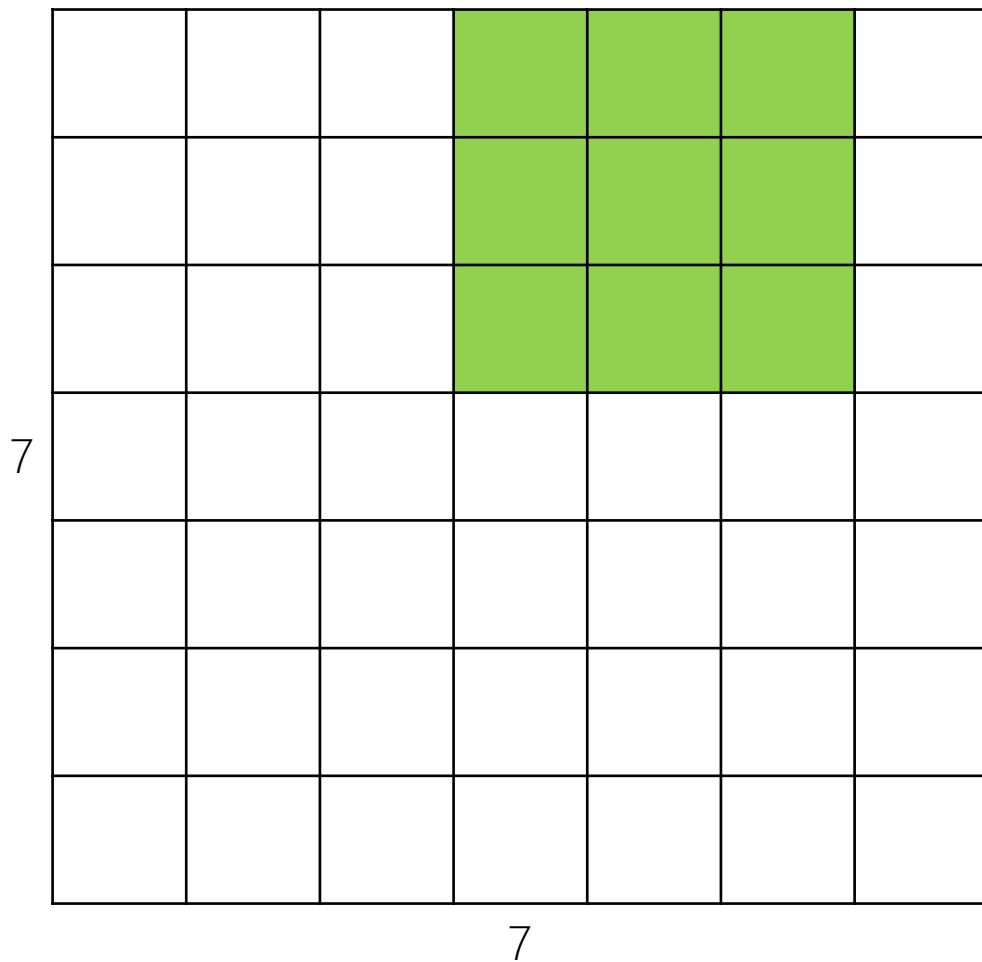
Stride & Padding



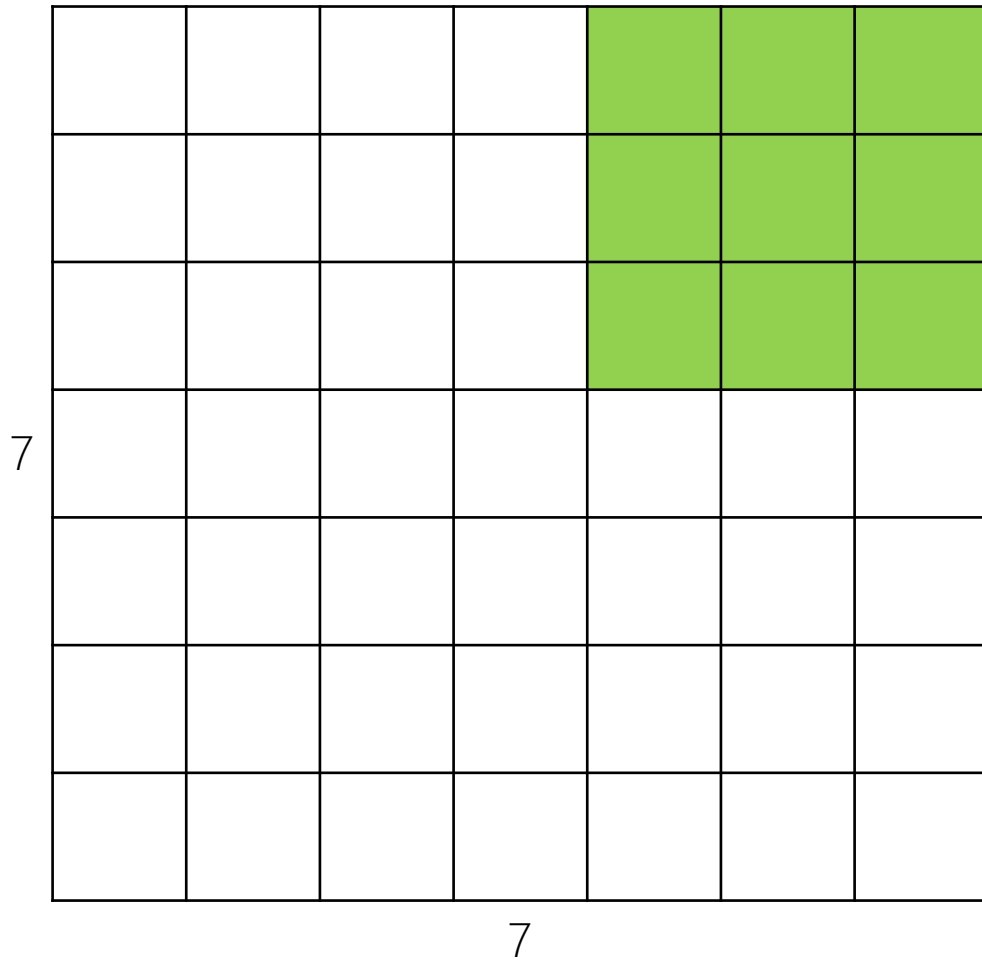
Stride & Padding



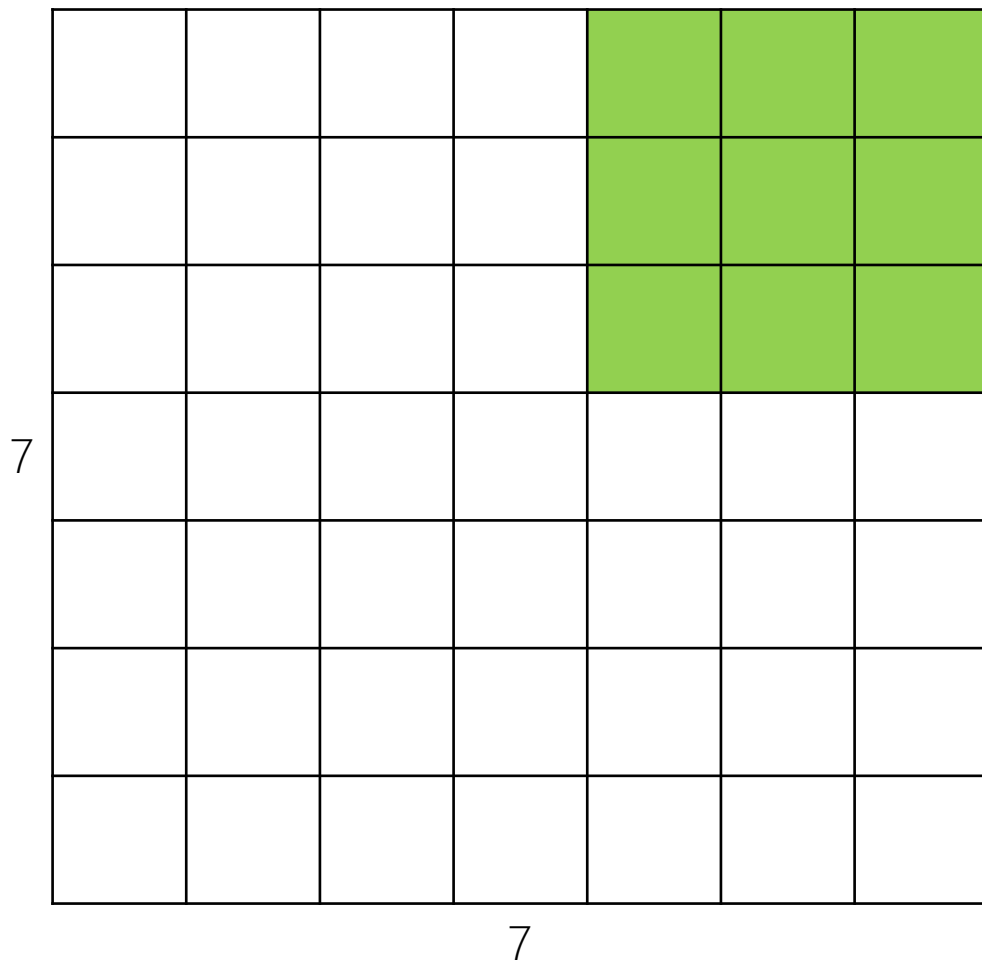
Stride & Padding



Stride & Padding

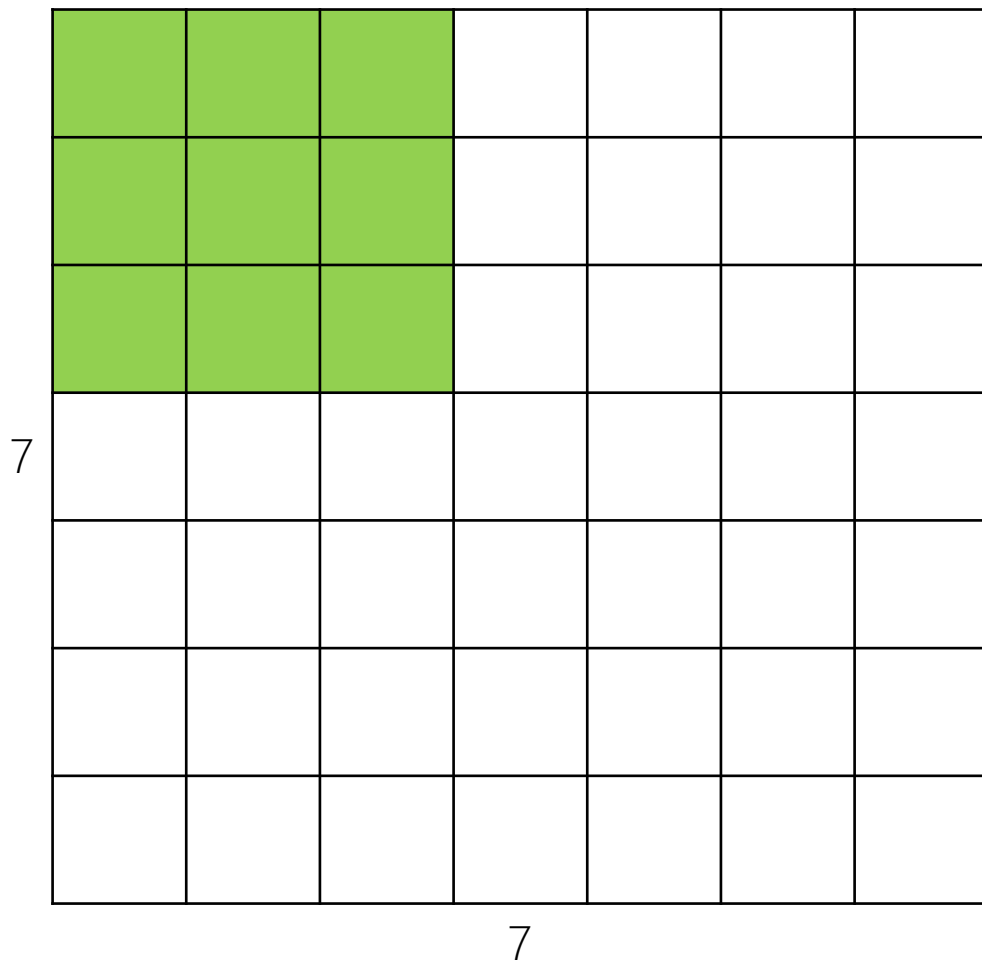


Stride & Padding



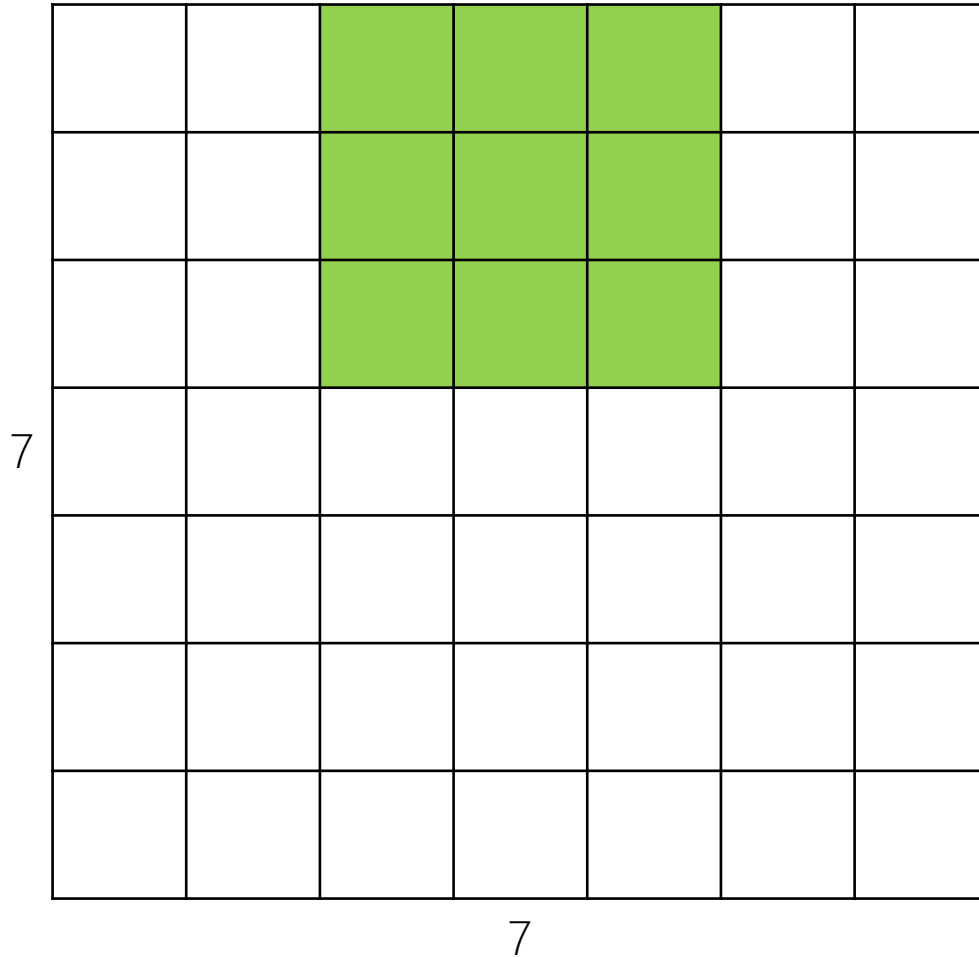
→ 5x5 output

Stride & Padding



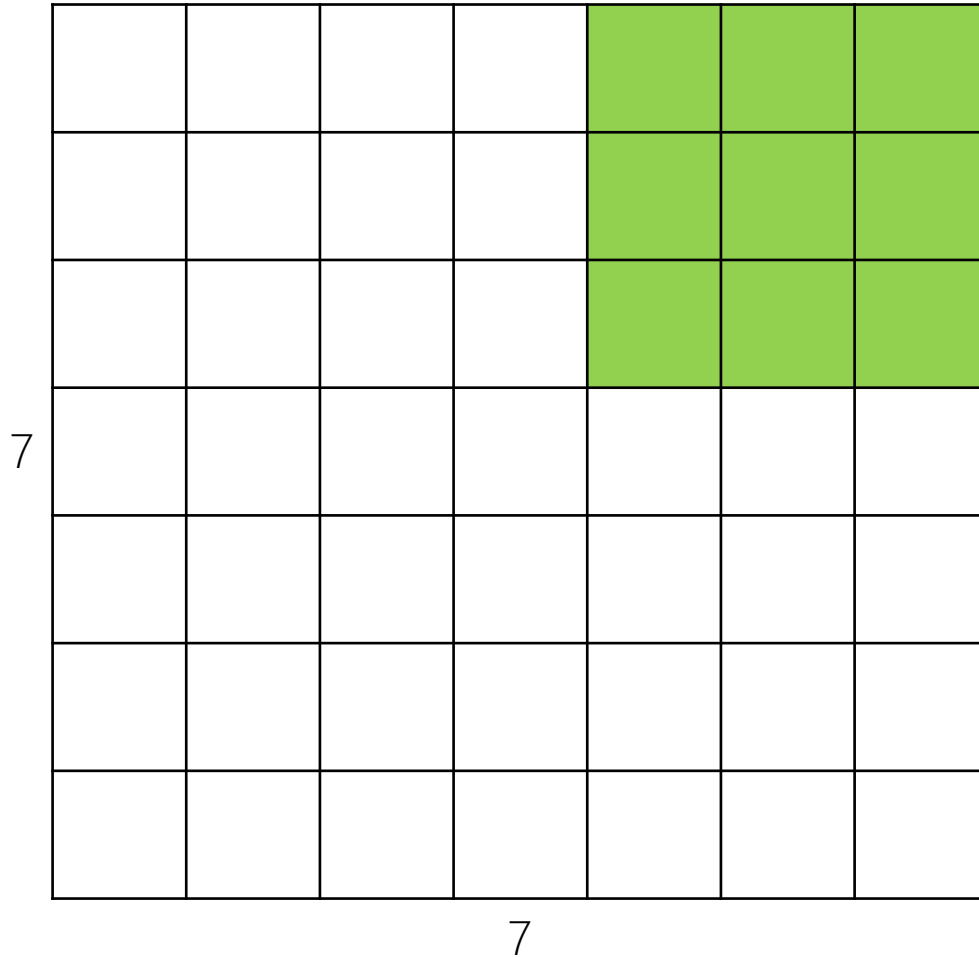
Now with **stride 2**

Stride & Padding



Now with stride 2

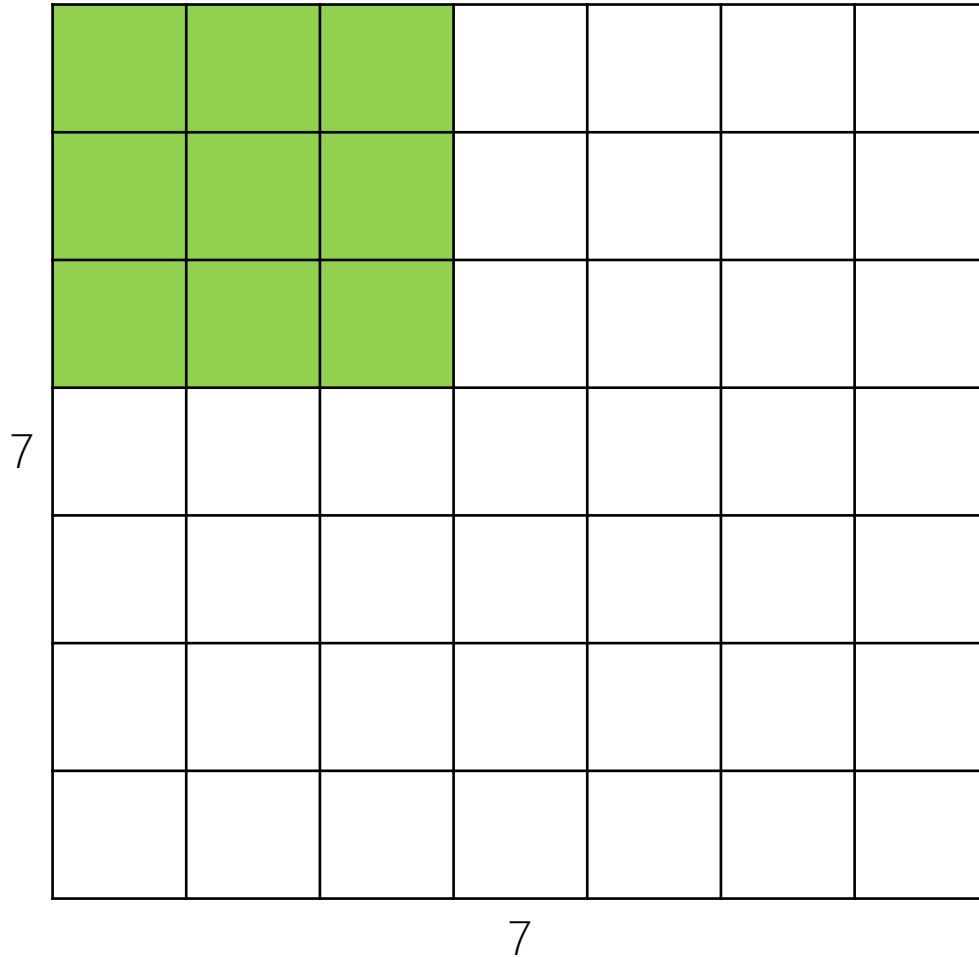
Stride & Padding



Now with **stride 2**

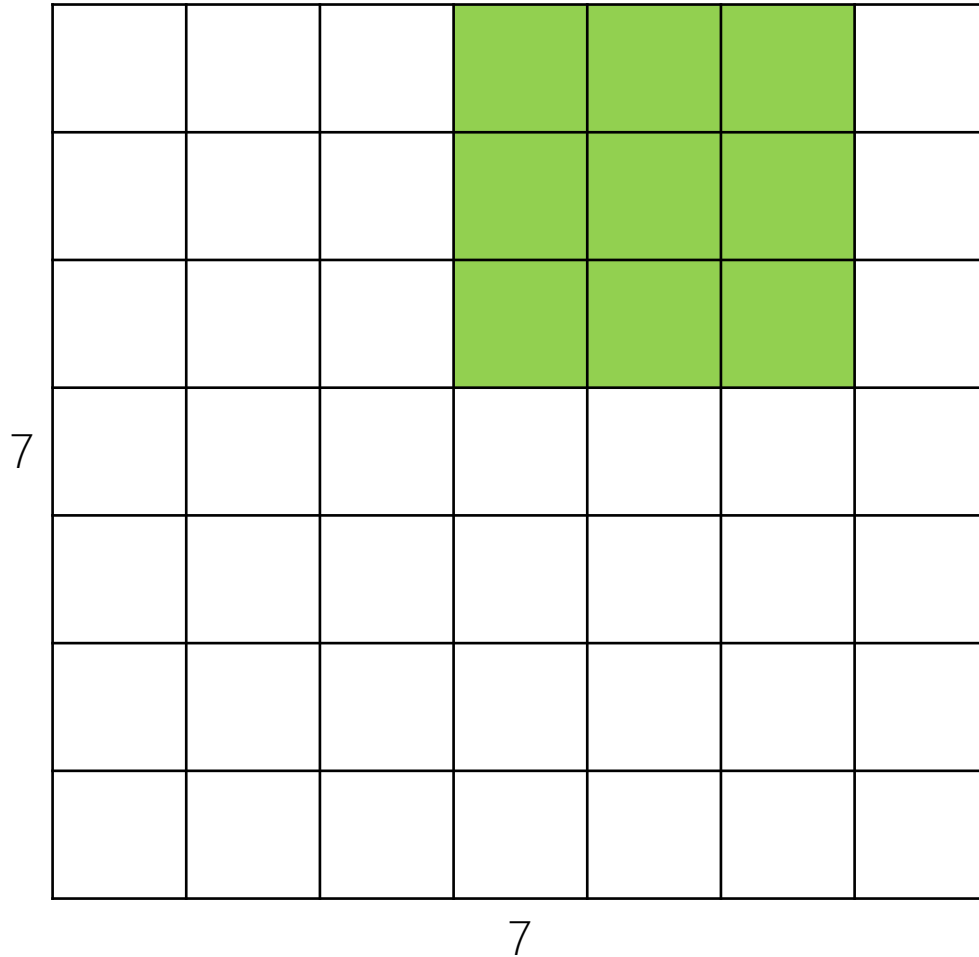
→ 3x3 output

Stride & Padding



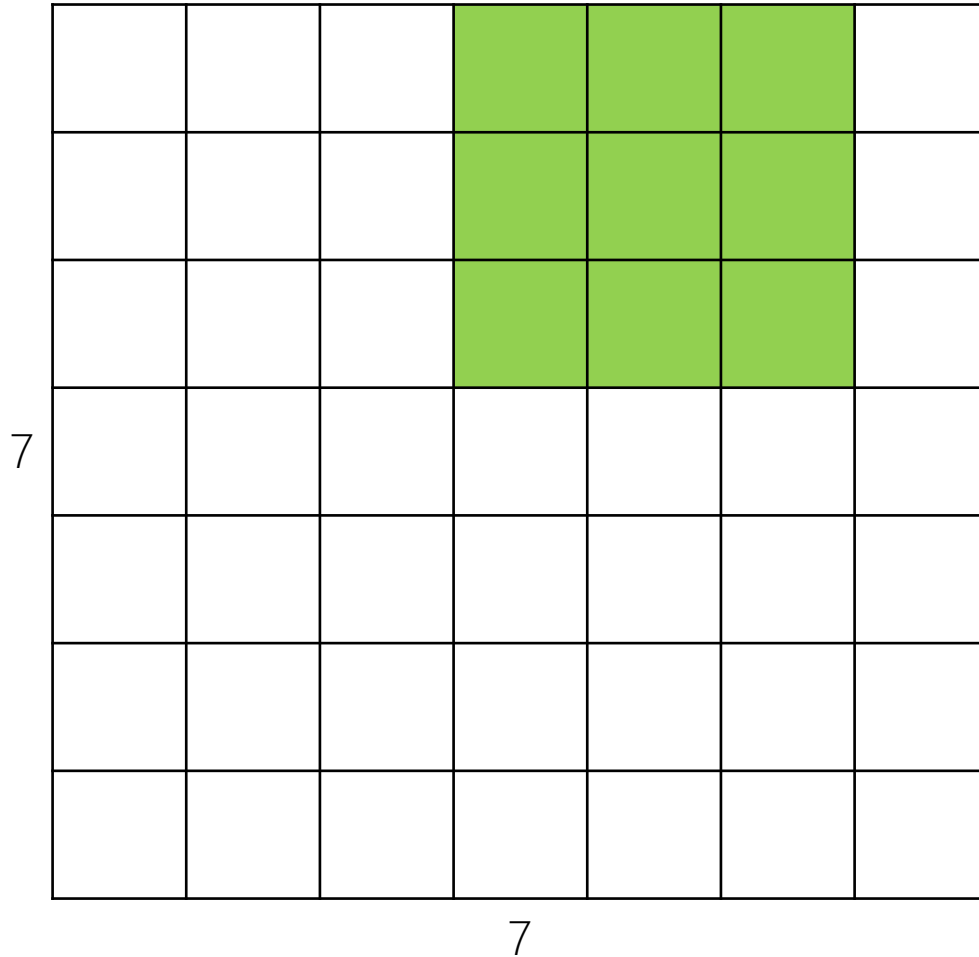
How about with **stride 3**?

Stride & Padding



How about with **stride 3**?

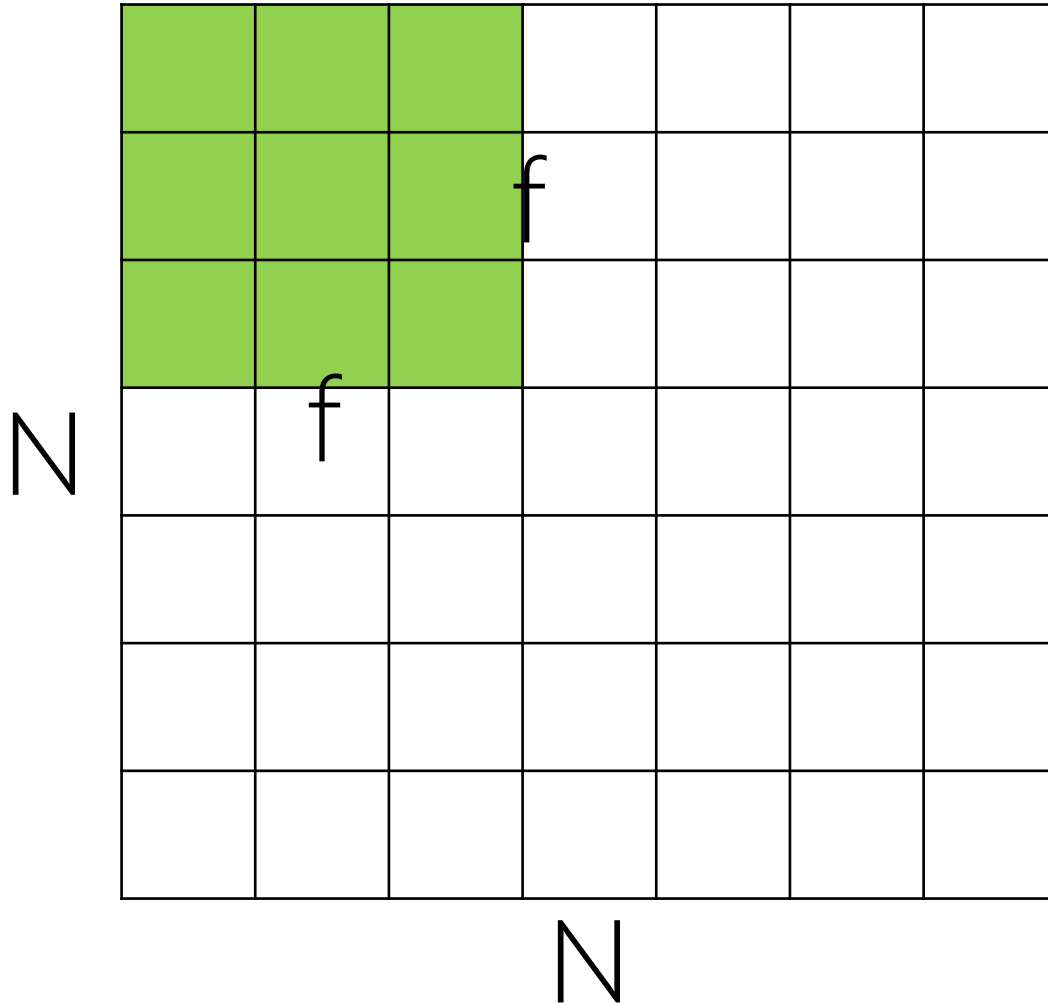
Stride & Padding



How about with **stride 3**?

→ ...?

Stride & Padding



Output size:
 $(N-f)/\text{stride} + 1$

Example, $N=7$, $f=3$:

Stride 1: $(7-3)/1 + 1 = 5$

Stride 2: $(7-3)/2 + 1 = 3$

Stride 3: $(7-3)/3 + 1 = 2.3333$

Stride & Padding

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

Output size:

$$(N-f)/\text{stride} + 1$$

Q. 7x7 input, 3x3 kernel, with stride 1,
padded with 1 pixel. Output size?

Stride & Padding

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

Output size:
 $(N-f)/\text{stride} + 1$

Q. 7x7 input, 3x3 kernel, with stride 1,
padded with 1 pixel. Output size?

→ 7x7

Stride & Padding

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

Output size:

$$(N-f)/\text{stride} + 1$$

Q. 7x7 input, 3x3 kernel, with stride 1, padded with 1 pixel. Output size?

→ 7x7

Q. 7x7 input, 3x3 kernel, with stride 3. You want to make 3x3 output, padding?

Stride & Padding

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

Output size:

$$(N-f)/\text{stride} + 1$$

Q. 7x7 input, 3x3 kernel, with stride 1,
padded with 1 pixel. Output size?

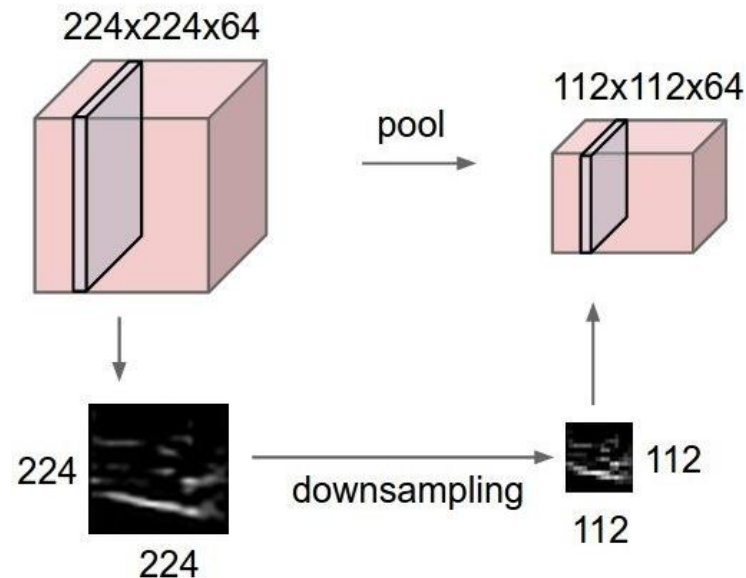
→ 7x7

Q. 7x7 input, 3x3 kernel, with stride 3.
You want to make 3x3 output, padding?

→ 1

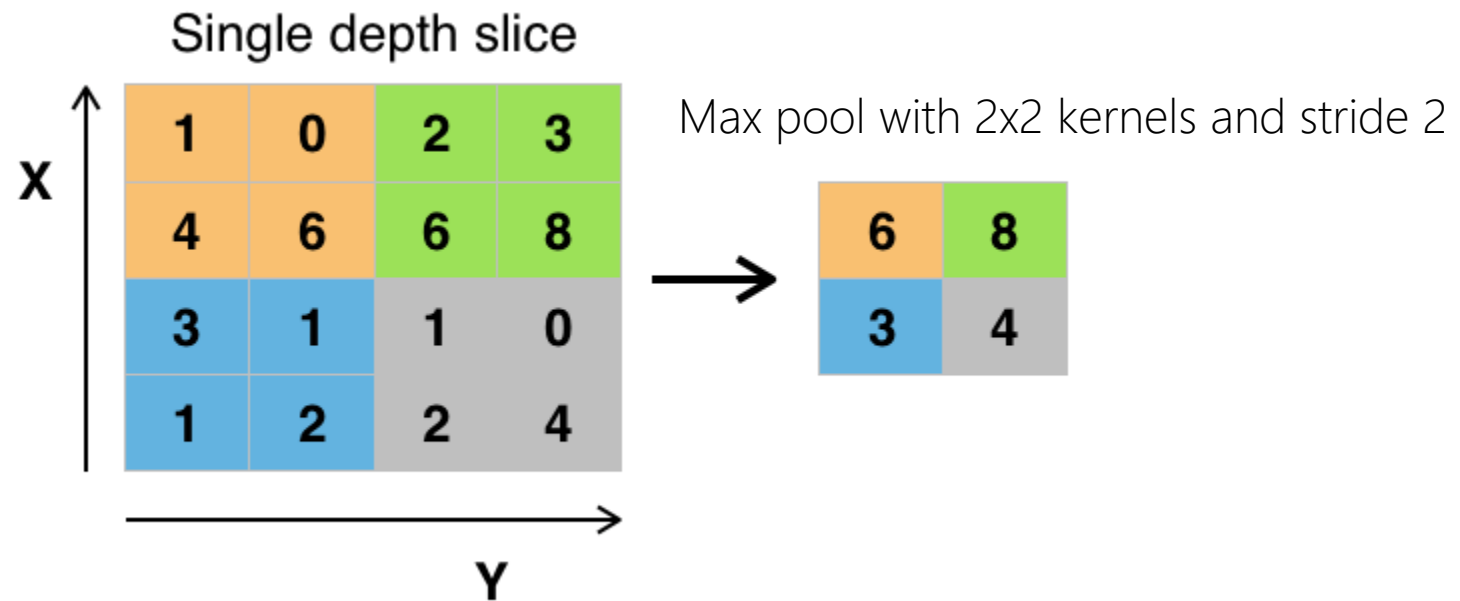
Pooling Layers

- Pooling == ConvNet way of downsampling
- “Reduce the image size” to obtain smaller but more manageable features
- Receptive fields becomes relatively larger as the image size gets smaller
- Operates over each activation map independently

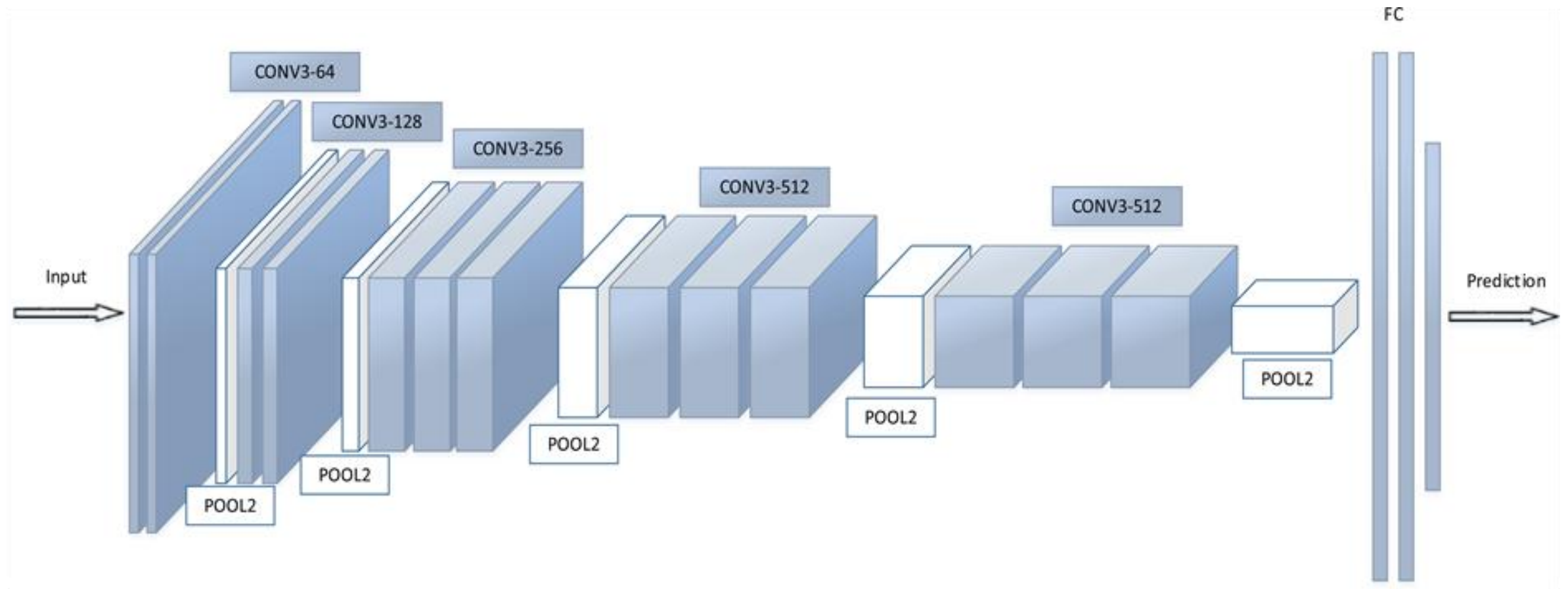


Max Pooling

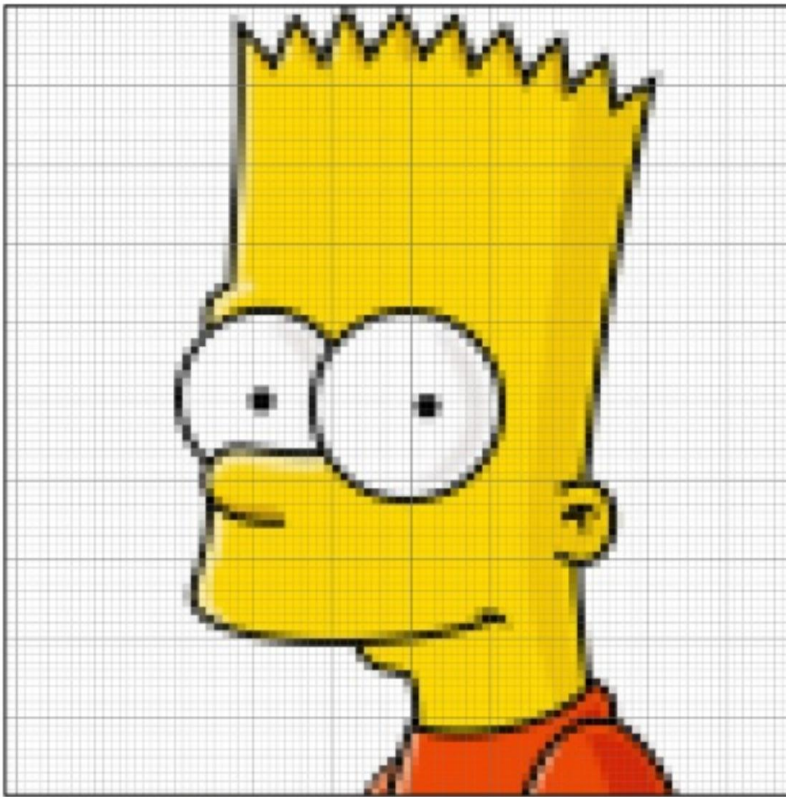
- A non-linear down-sampling method
- An image is partitioned into a set of (non-overlapping) rectangles.
- The maximum of each such sub-region is sampled.



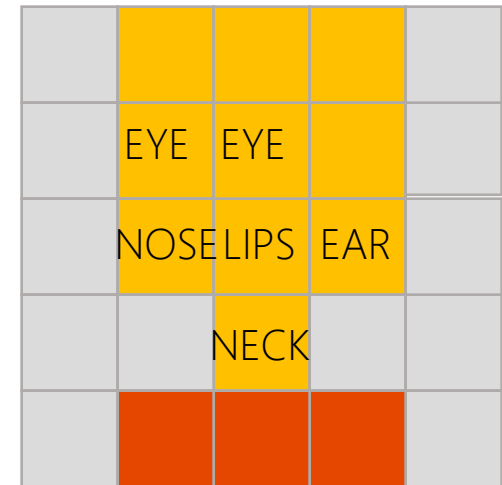
VGG Networks



Abstraction of an Image

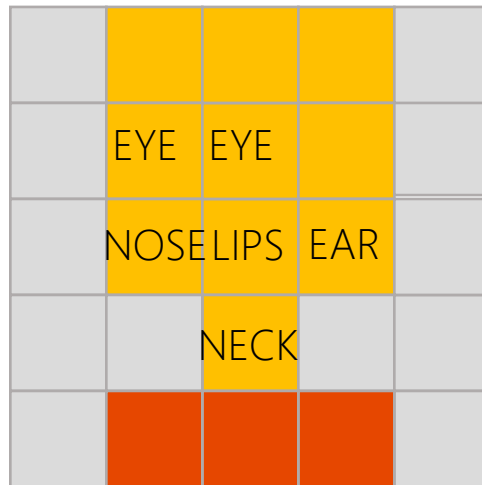


Conv2D + MaxPooling
+ Conv2D + MaxPooling
+ Conv2D + MaxPooling
+ ...



Abstraction of an Image

- Final decision is made by the fully-connected layers:



"Has two eyes, a nose, lips, an ear, and a neck
and wears t-shirts"



Cat	X
Dog	X
Person	○
Orange	X