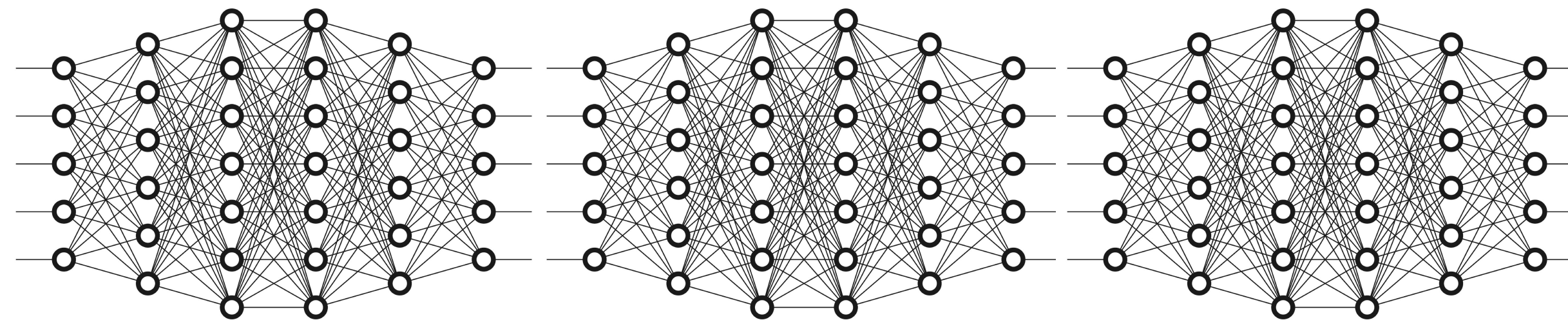




# CPSC 425: Computer Vision



## Lecture 34: Convolutional Neural Networks

# Menu for Today (November 28, 2018)

## Topics:

- Convolutional Layers
- Convolutional Neural Networks
- Pooling Layer
- R-CNN

## Readings:

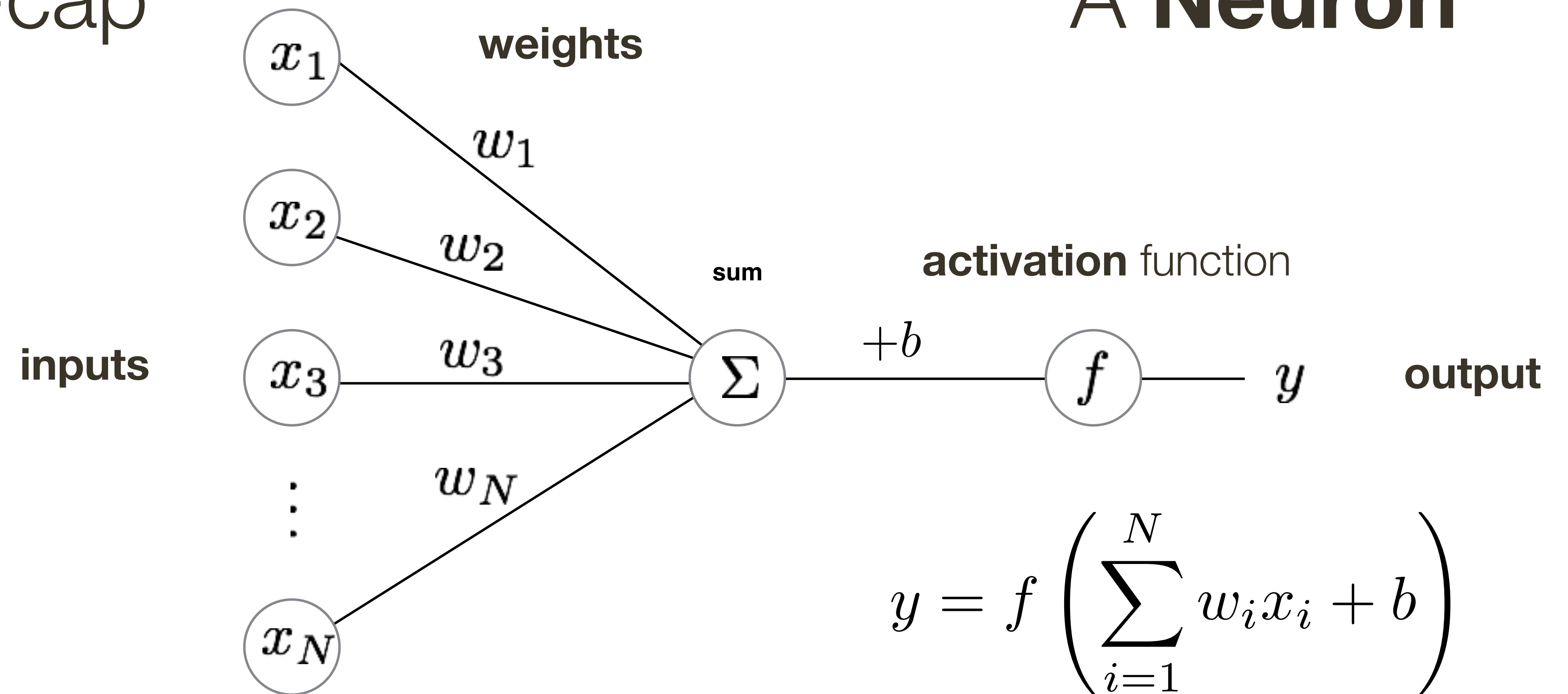
- **Today's** Lecture: N/A
- **Next** Lecture: N/A

## Reminders:

- **Assignment 5:** Scene Recognition with Bag of Words due **last day of classes**

# Lecture 33: Re-cap

## A Neuron



- The basic unit of computation in a neural network is a neuron.
- A neuron accepts some number of input signals, computes their weighted sum, and applies an **activation function** (or **non-linearity**) to the sum.
- Common activation functions include sigmoid and rectified linear unit (ReLU)

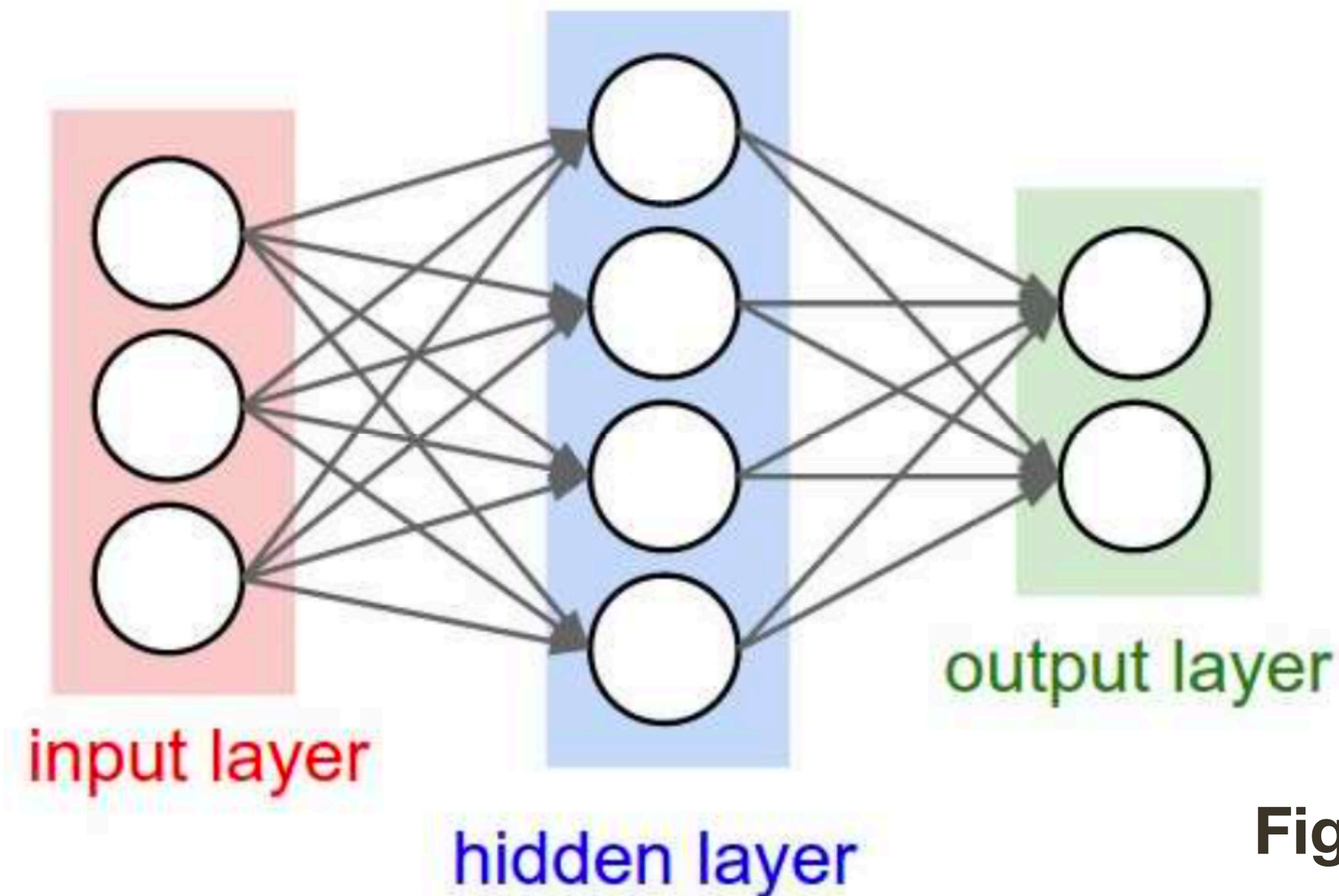
# Lecture 33: Re-cap

# Neural Network

A neural network comprises neurons connected in an acyclic graph

The outputs of neurons can become inputs to other neurons

Neural networks typically contain multiple layers of neurons



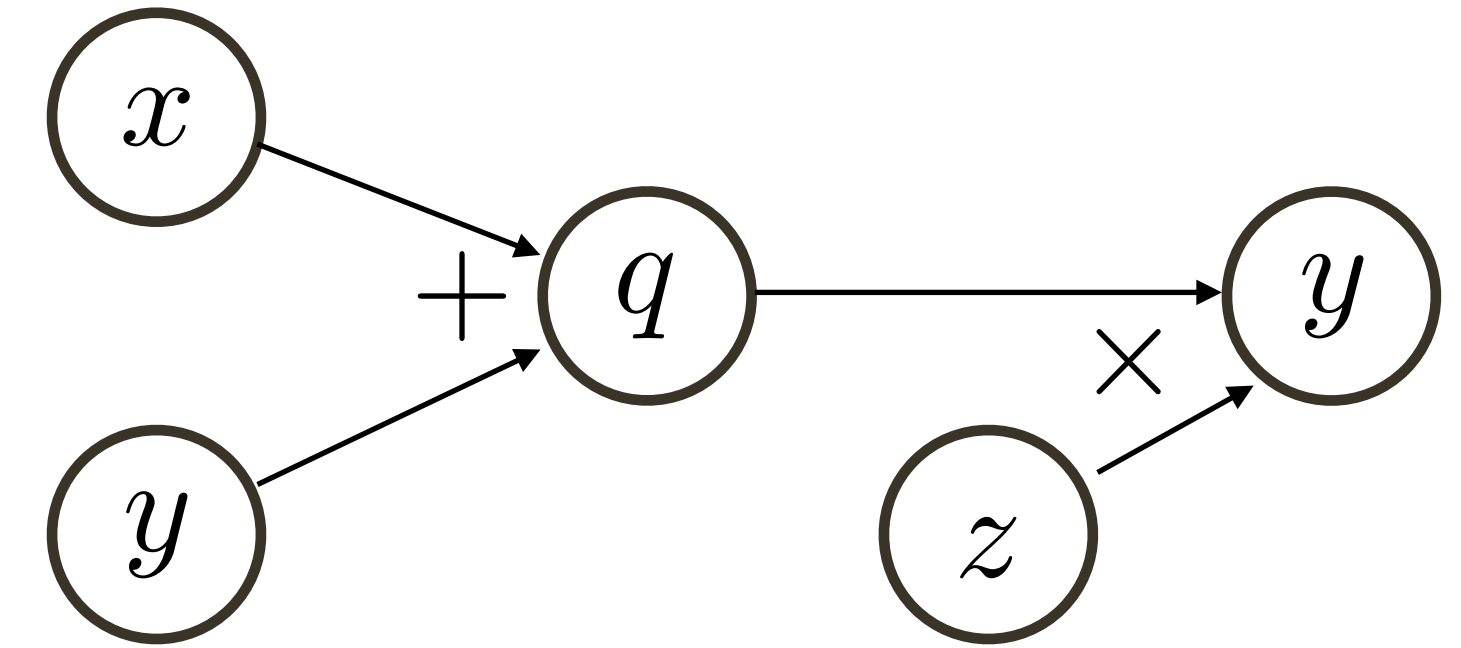
**Figure credit:** Fei-Fei and Karpathy

Example of a neural network with three inputs, a single hidden layer of four neurons, and an output layer of two neurons



# Lecture 33: Re-cap

$$f(x, y, z) = (x + y)z$$



$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x} = \frac{\partial f}{\partial q} \cdot 1$$

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial y} = \frac{\partial f}{\partial q} \cdot 1$$

$$\frac{\partial f}{\partial z} = q$$

Suppose the network input is:  $(x, y, z) = (-2, 5, -4)$

Then:  $q = x + y = 3$        $f = qz = -12$       (**forward** pass)

$$\frac{\partial f}{\partial q} = z = -4$$

$$\frac{\partial f}{\partial x} = -4$$

$$\frac{\partial f}{\partial y} = -4$$

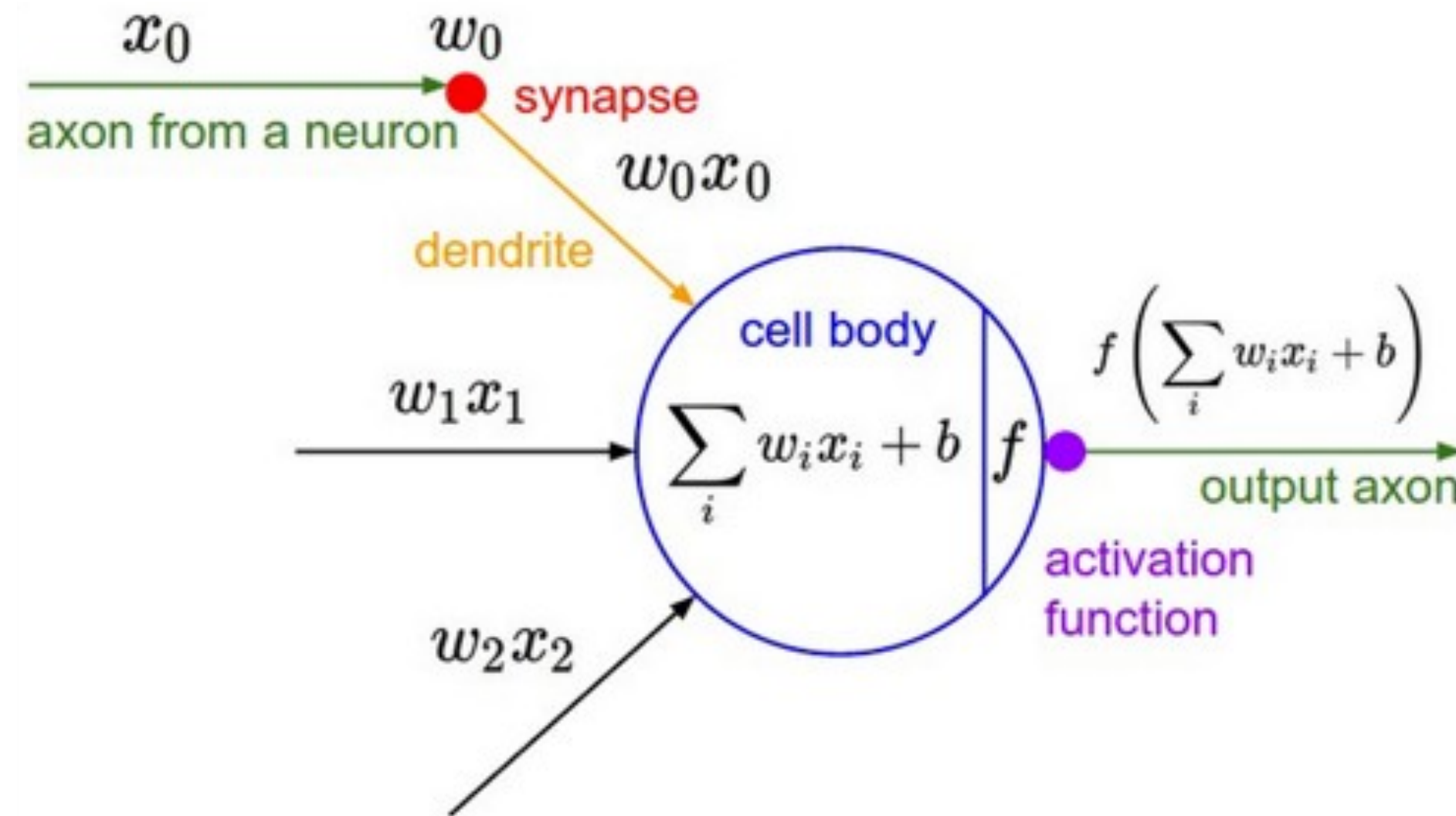
$$\frac{\partial f}{\partial z} = 3$$

(**backward** pass)

# Lecture 33: Re-cap

# Backpropagation

Chaining products and sums may seem like a simple example. But recall the basic unit in a neural network.

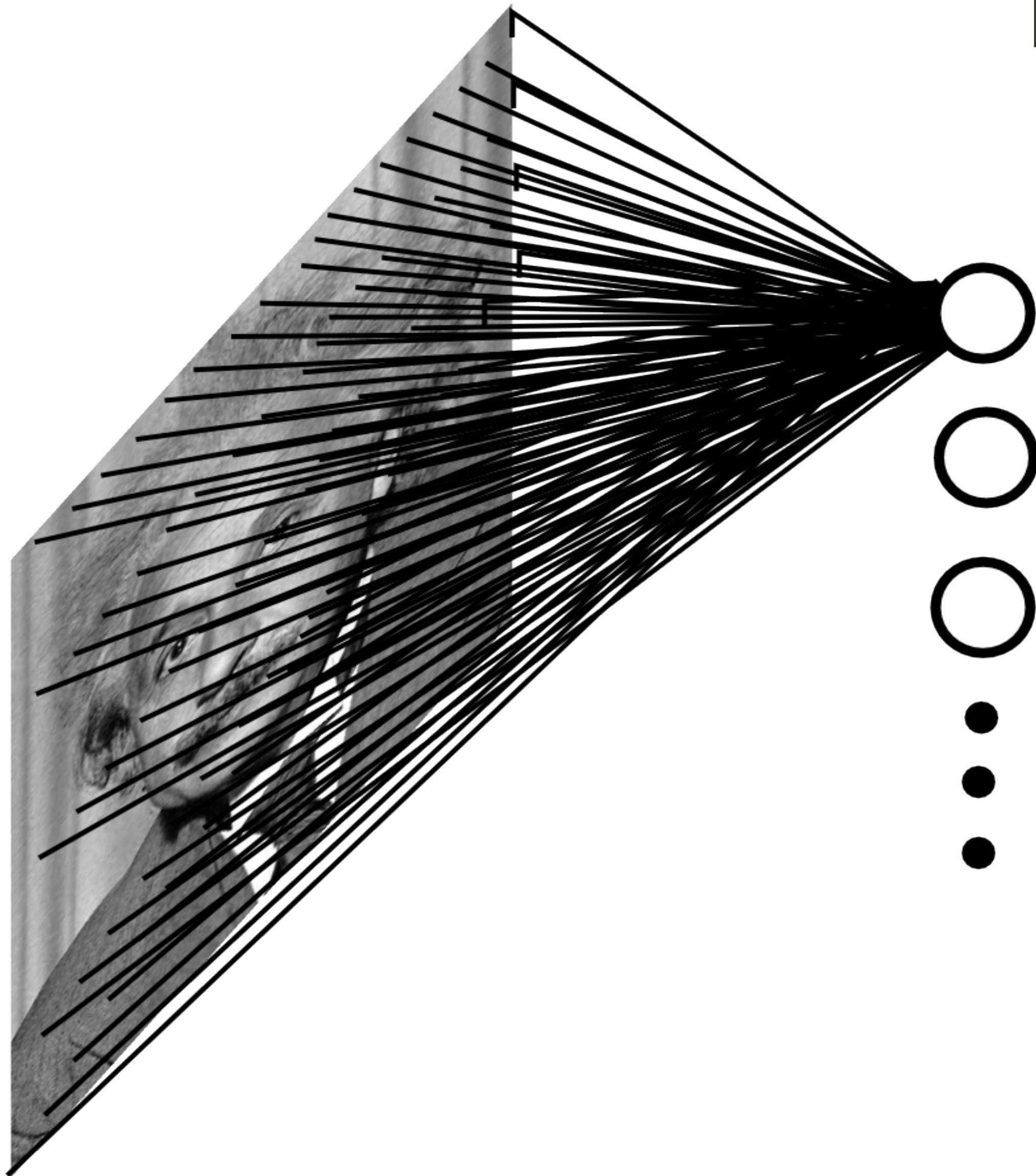


**Figure credit:** Fei-Fei and Karpathy

It consists of products, sums, and activation functions (e.g. ReLU, which is a max), which we can chain together

Common loss functions are also differentiable

# Fully Connected Layer



**Example:** 200 x 200 image (small)  
x 40K hidden units

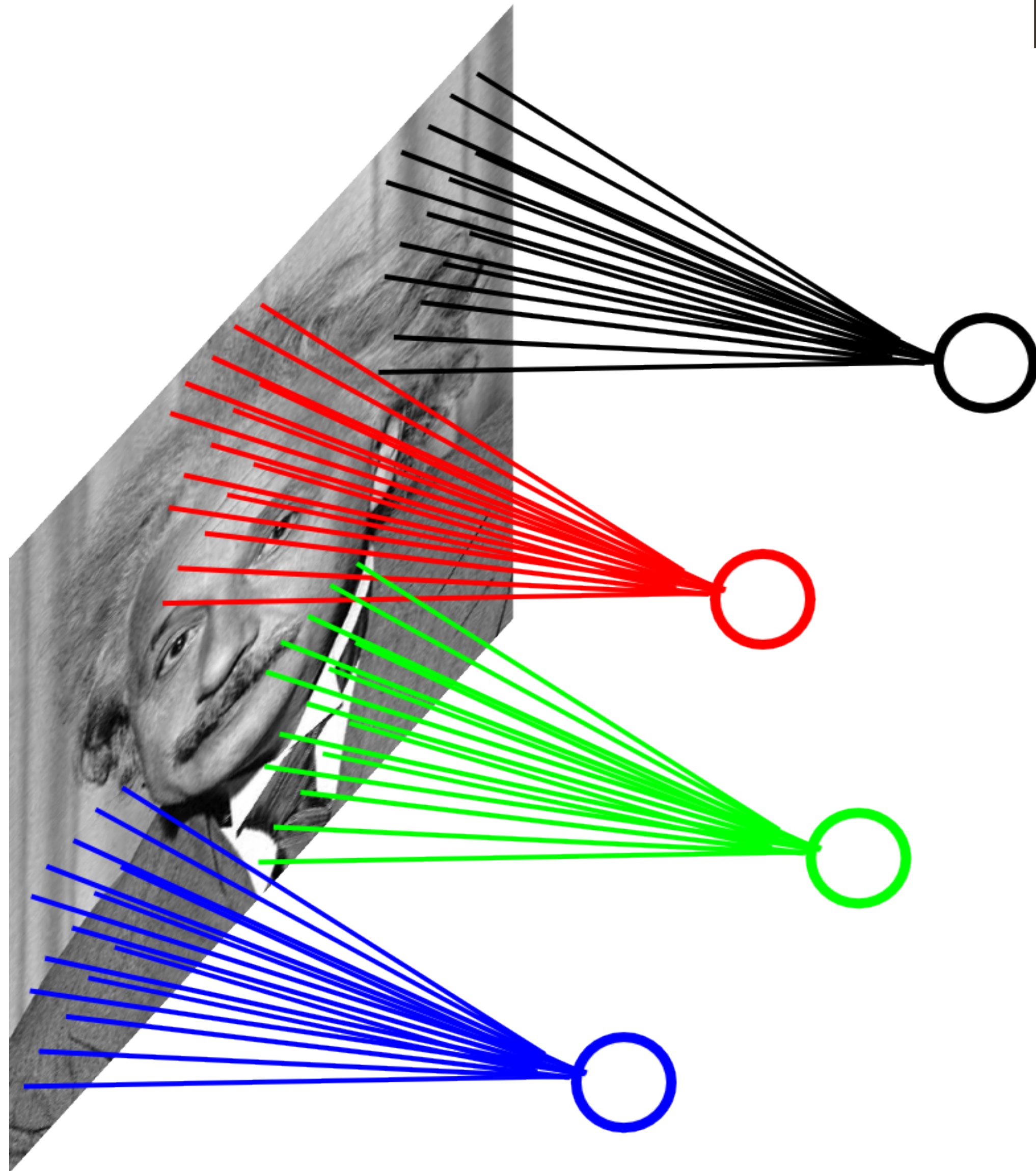
= ~ **2 Billion** parameters (for one layer!)

Spatial correlations are generally local

Waste of resources + we don't have  
enough data to train networks this large



# Locally Connected Layer



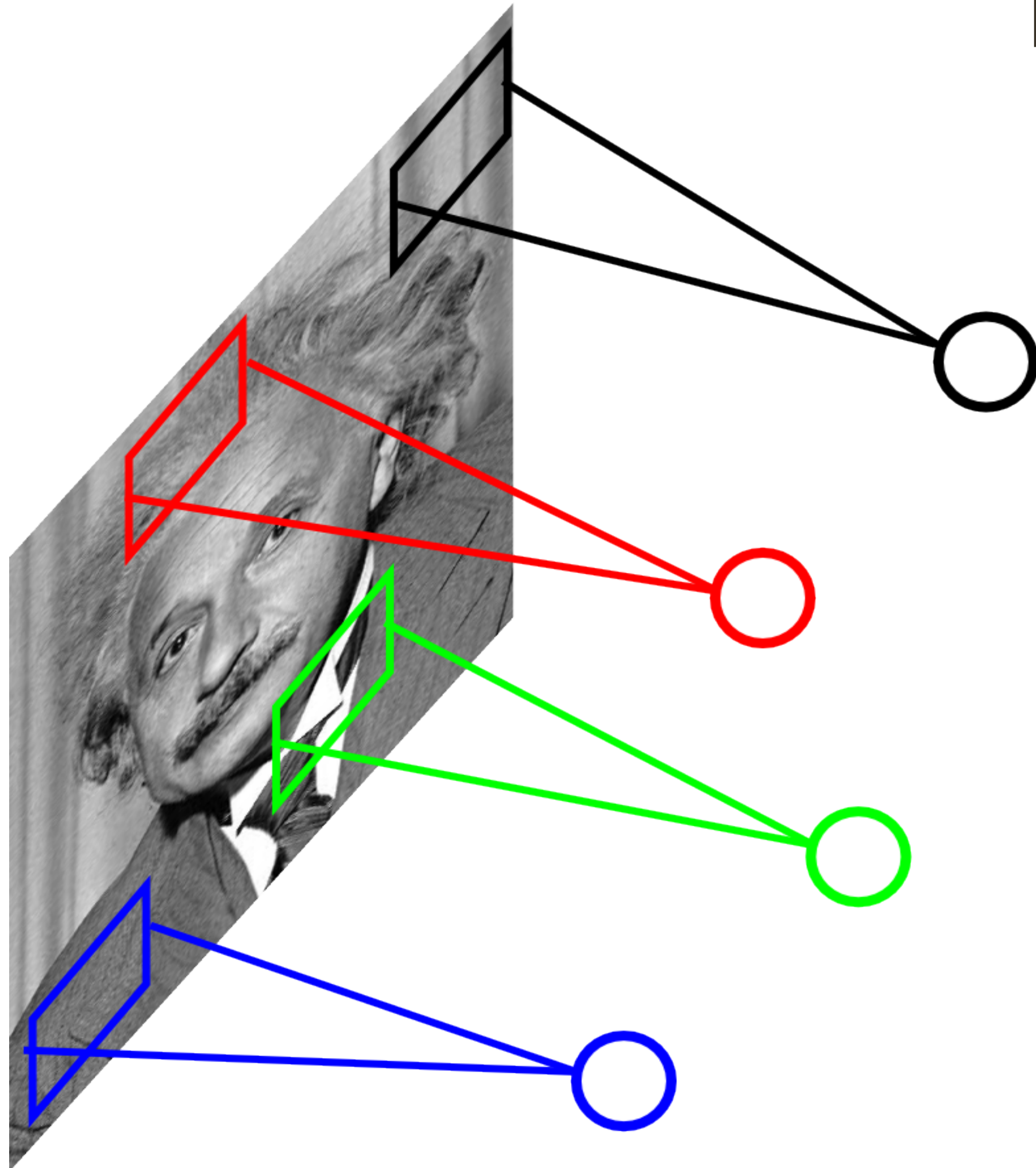
**Example:** 200 x 200 image (small)  
x 40K hidden units

**Filter size:** 10 x 10

= ~ **4 Million** parameters



# Locally Connected Layer



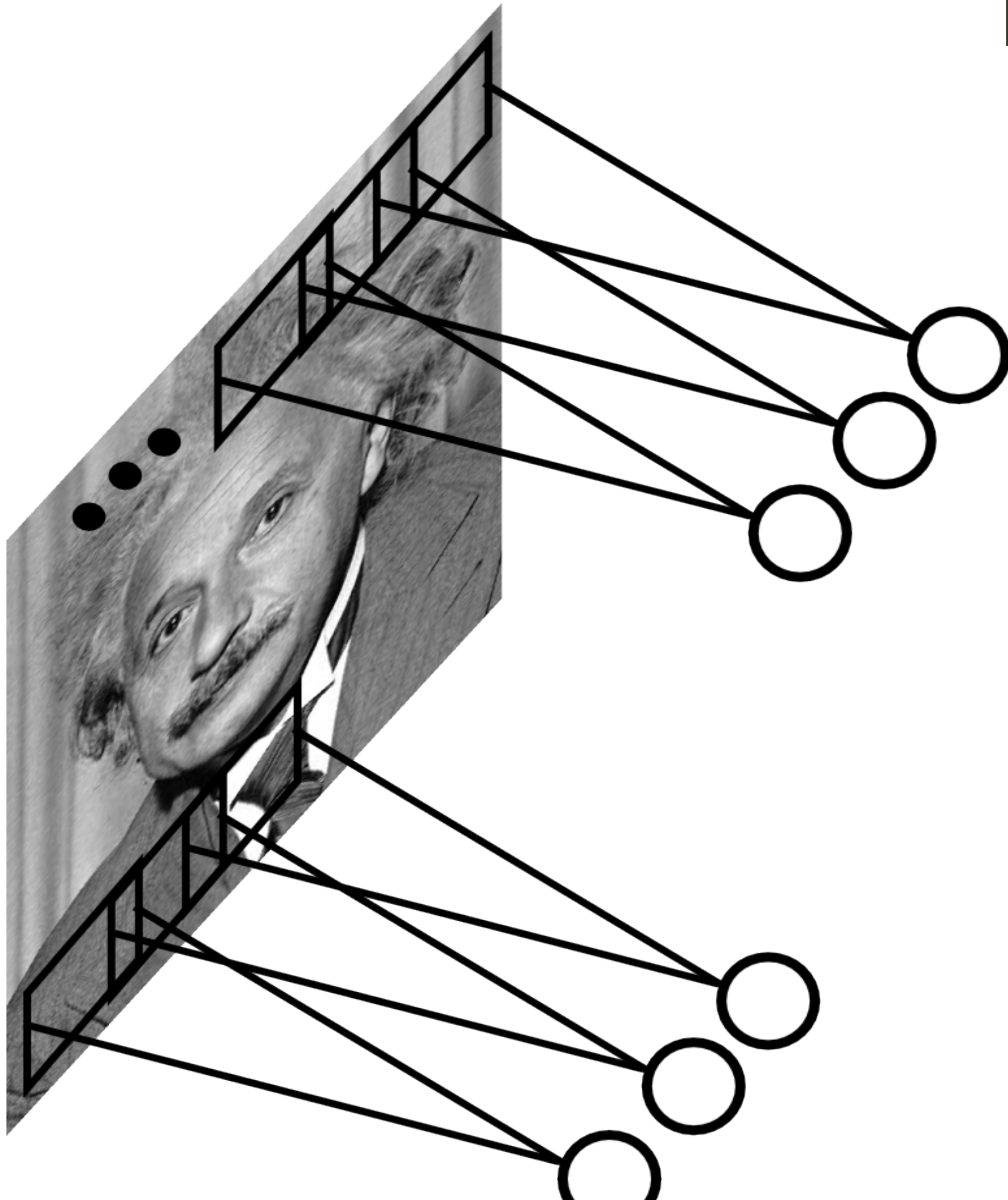
**Example:** 200 x 200 image (small)  
x 40K hidden units

**Filter size:** 10 x 10

= ~ **4 Million** parameters

**Stationarity** — statistics is similar at  
different locations

# Convolutional Layer



**Example:** 200 x 200 image (small)  
x 40K hidden units

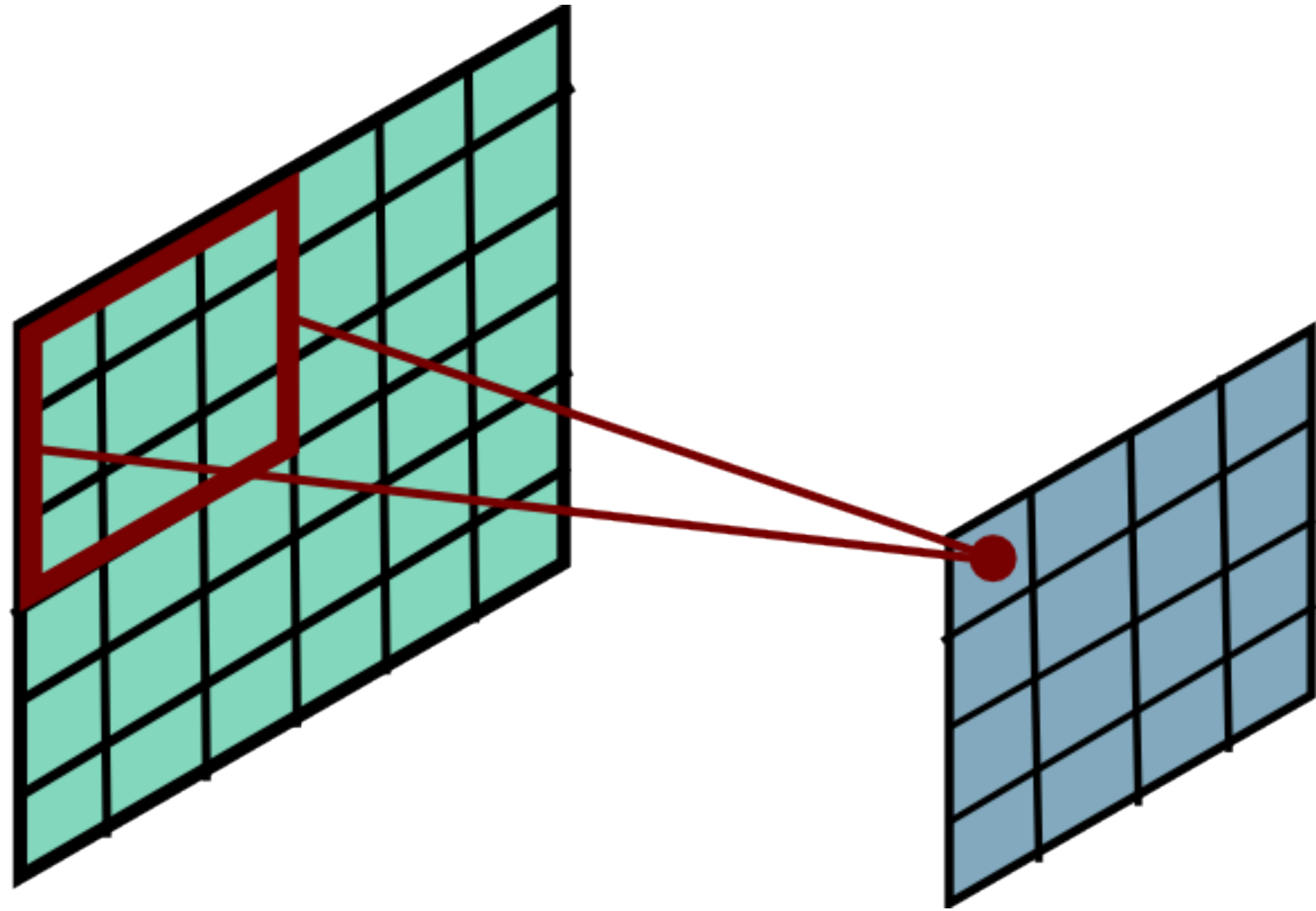
**Filter size:** 10 x 10

= ~ **4 Million** ~~parameters~~

= 100 parameters

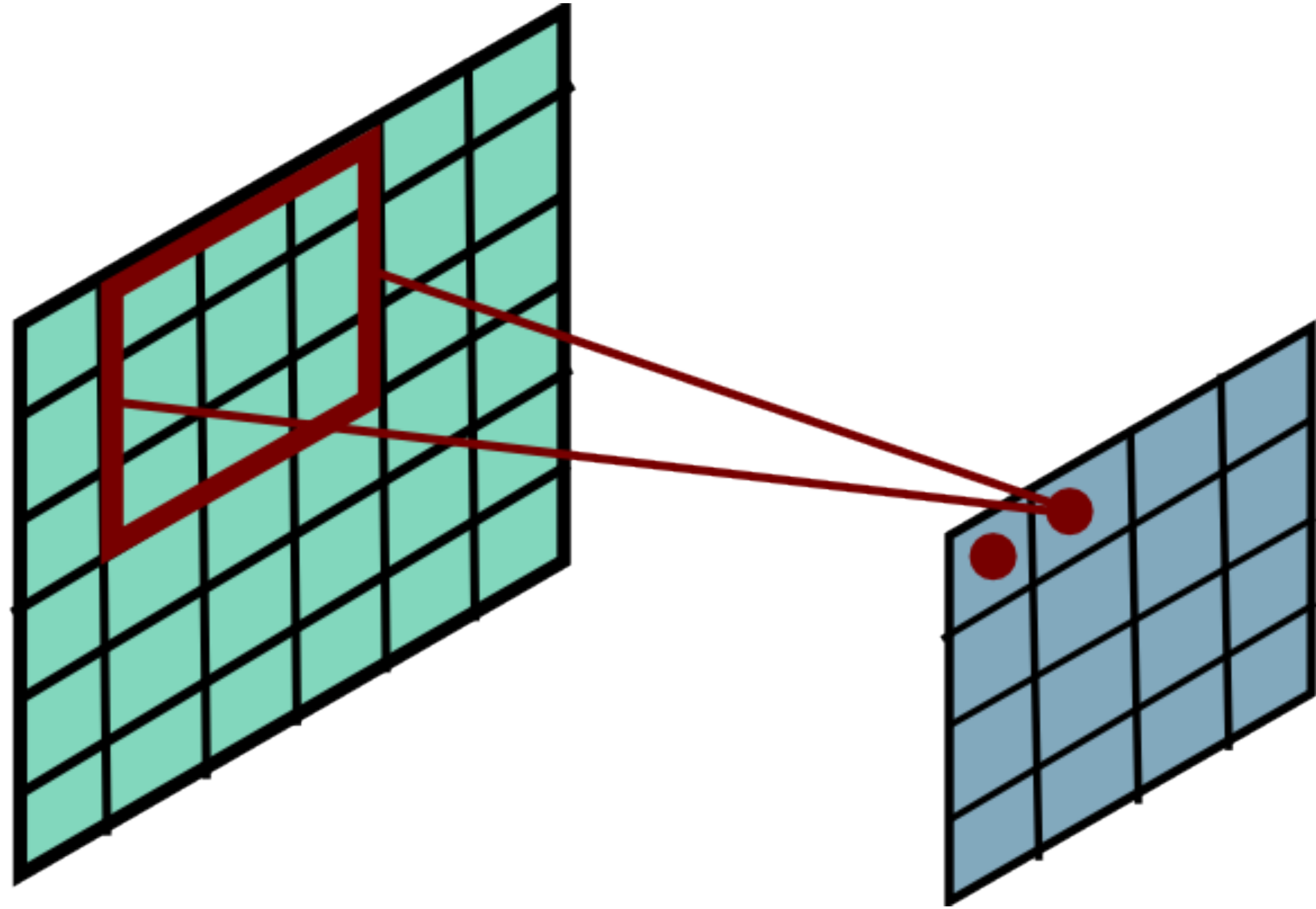
Share the same parameters across the locations (assuming input is stationary)

# Convolutional Layer



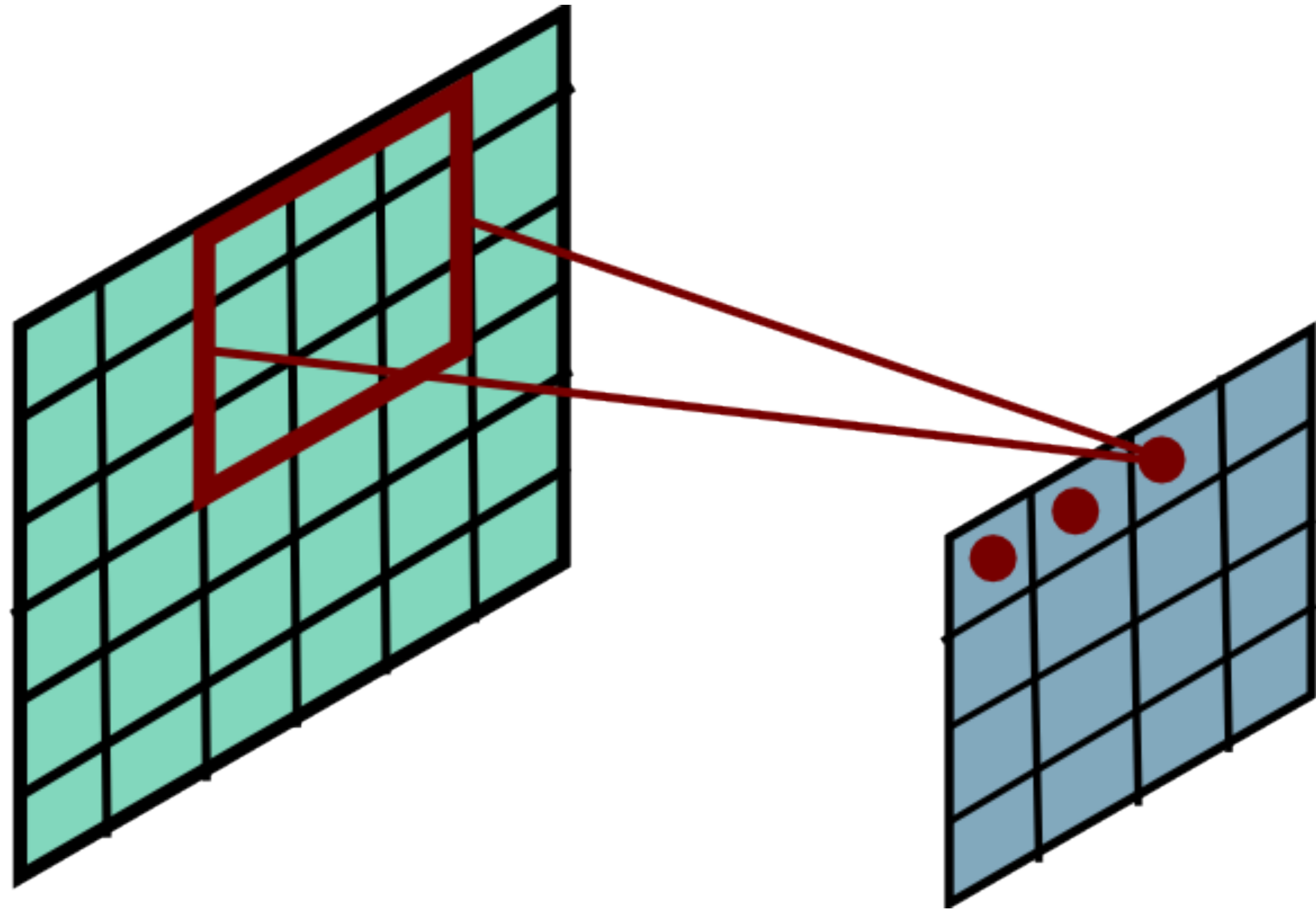


# Convolutional Layer

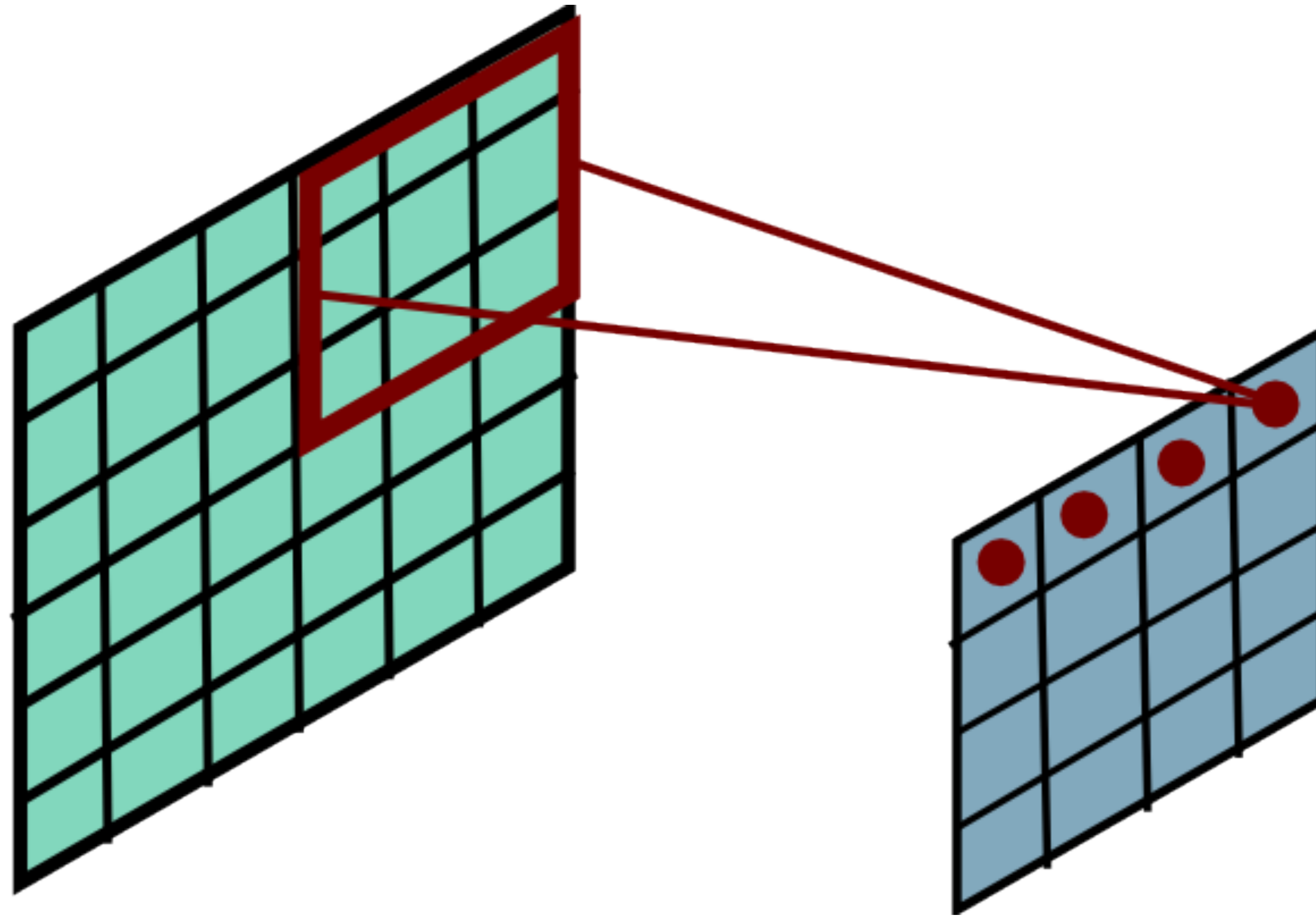




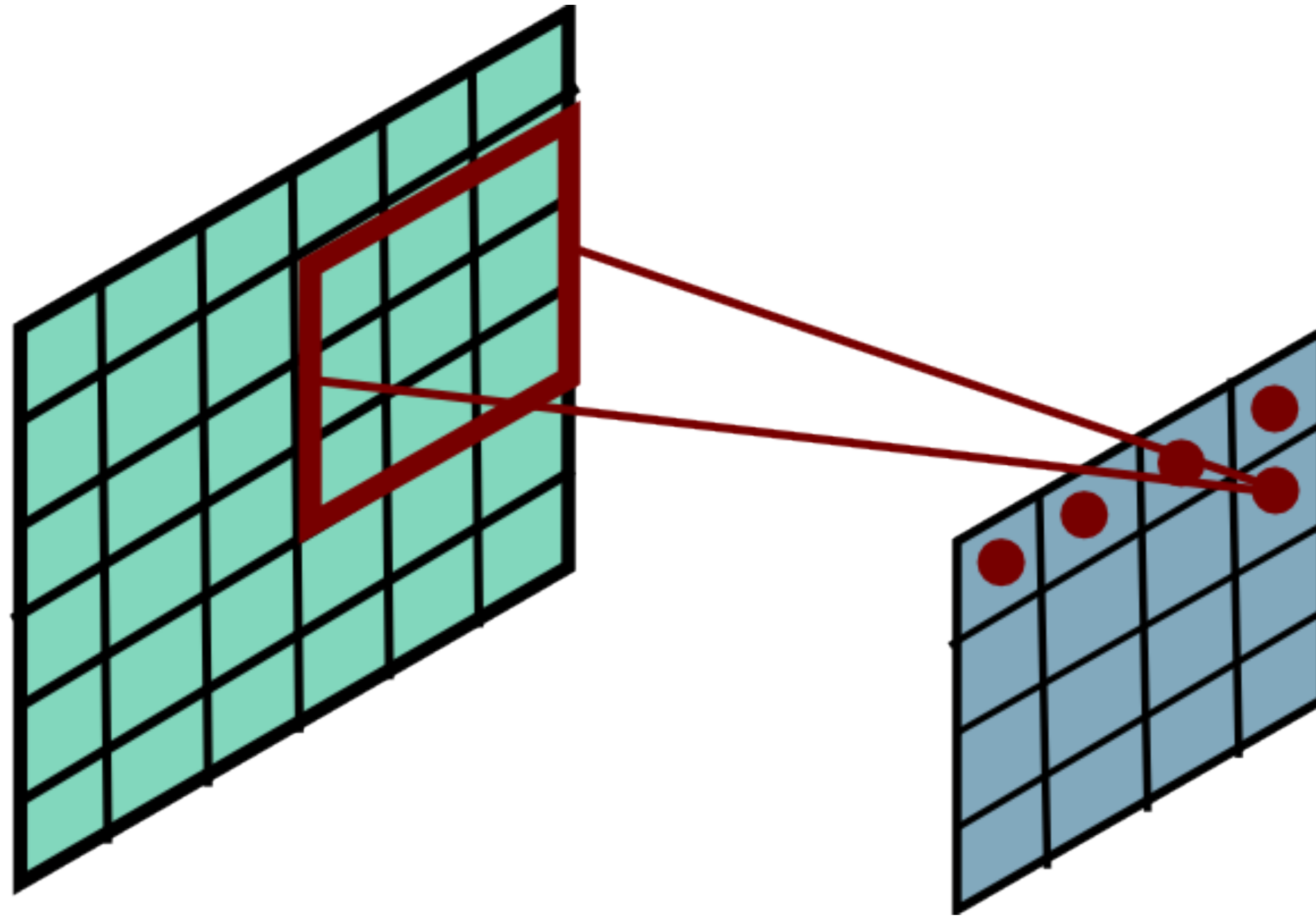
# Convolutional Layer



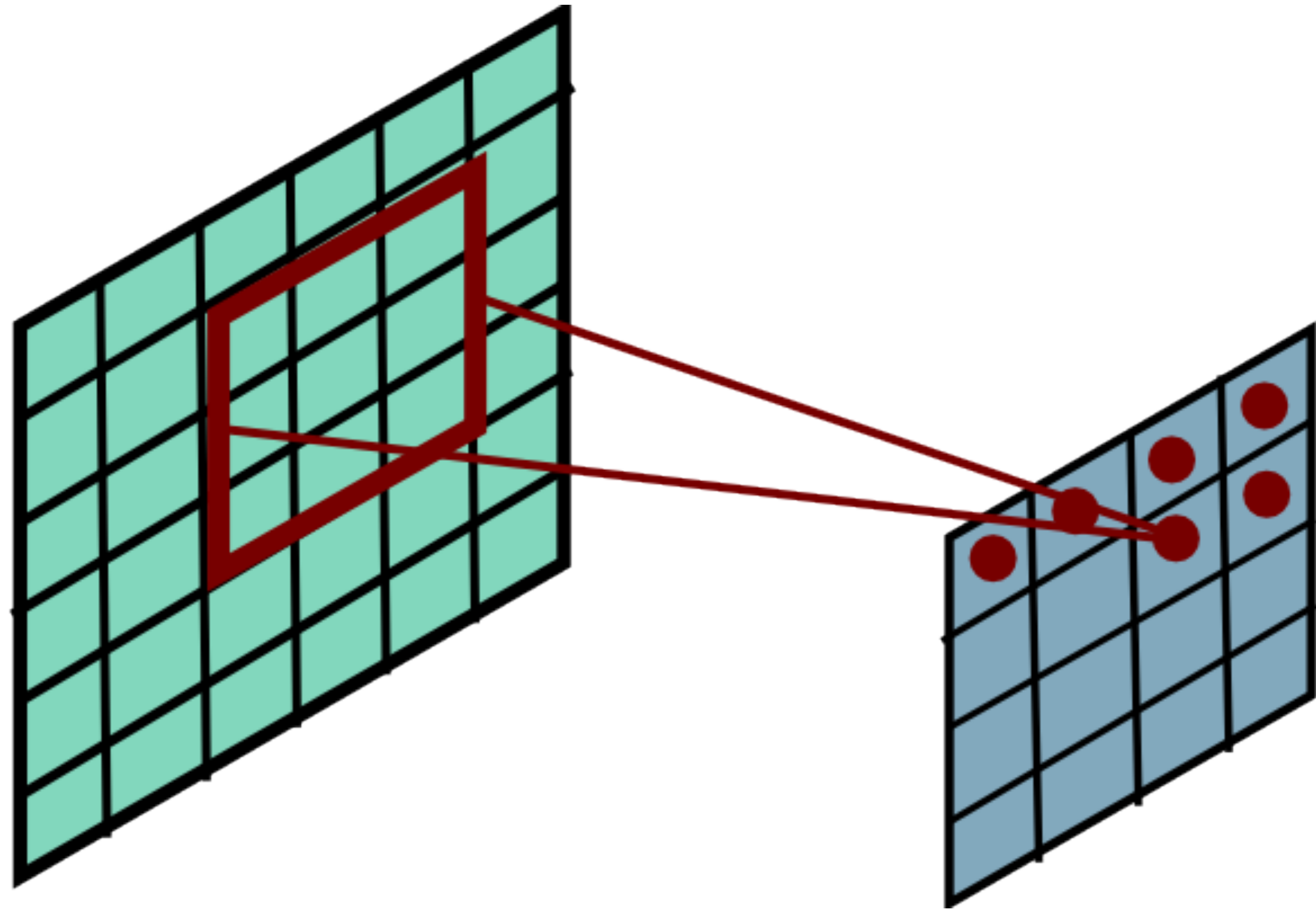
# Convolutional Layer



# Convolutional Layer

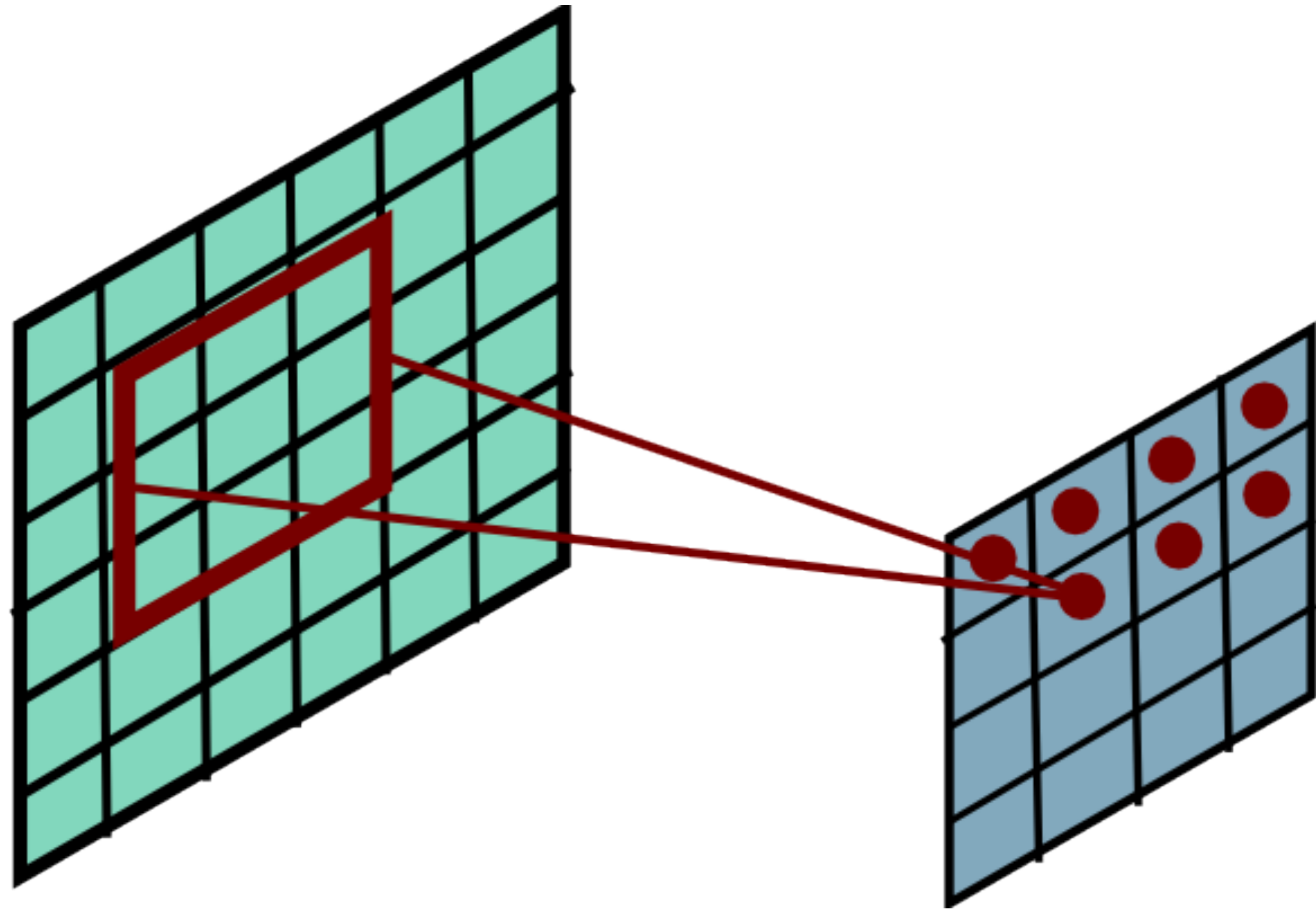


# Convolutional Layer

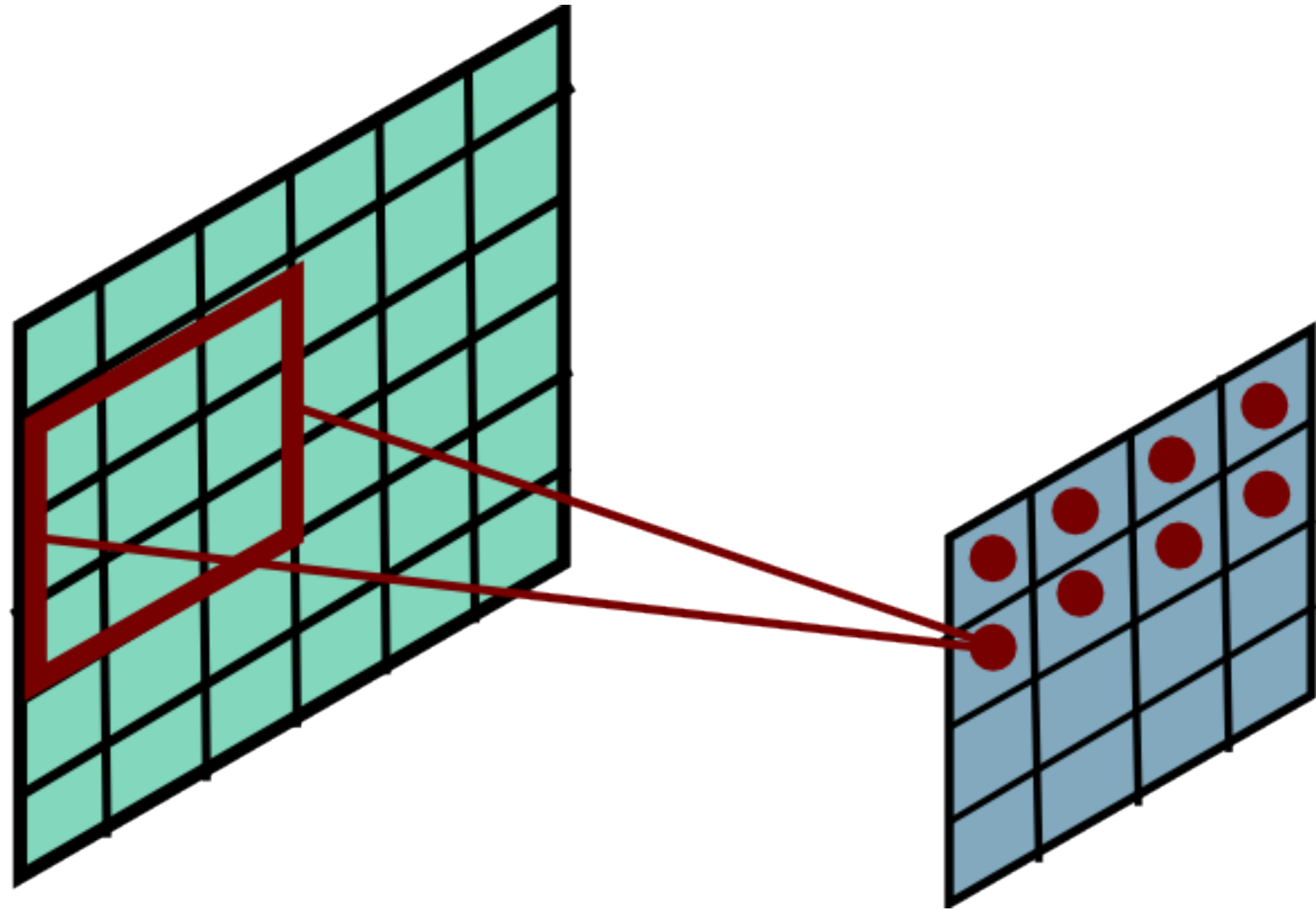




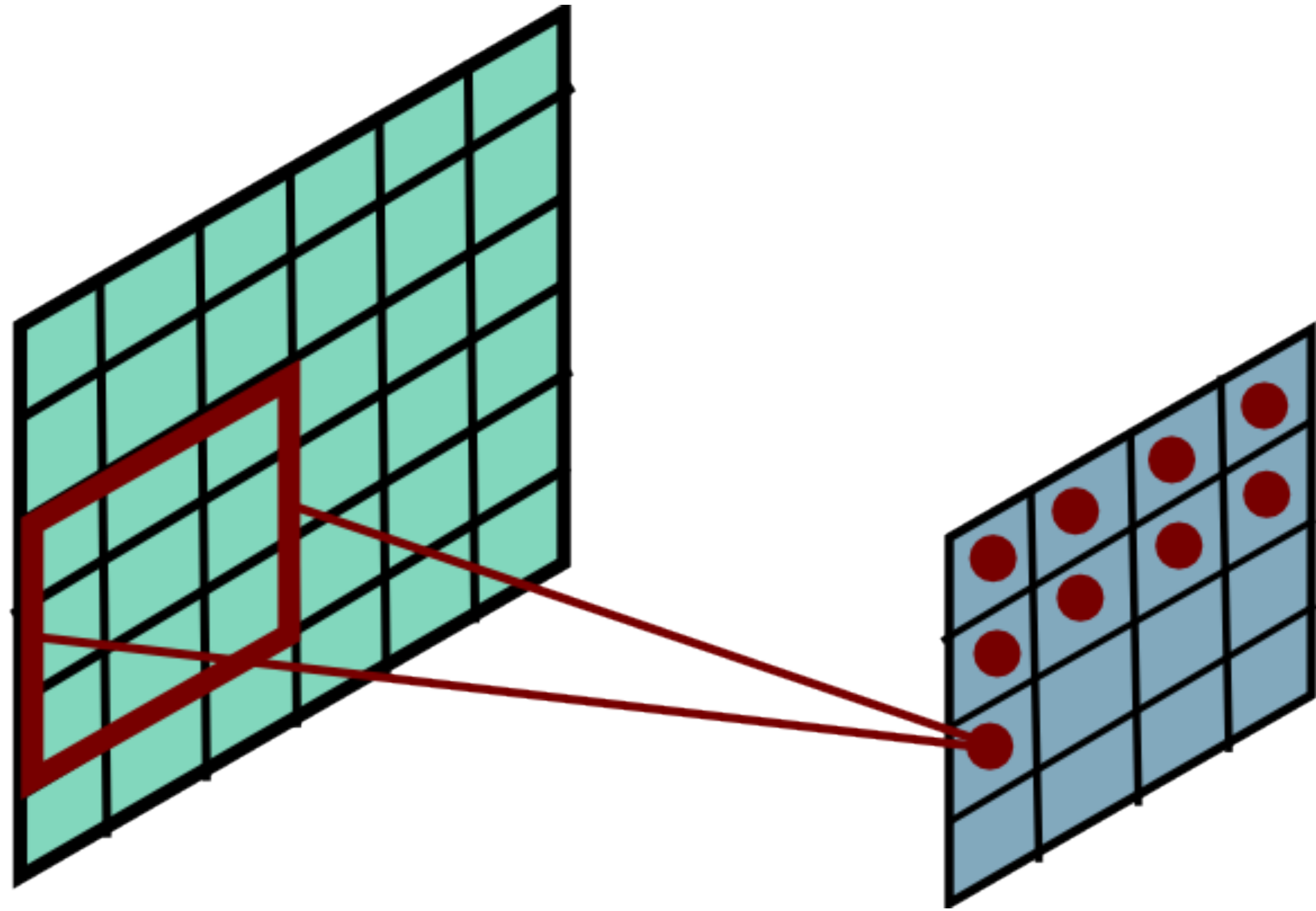
# Convolutional Layer



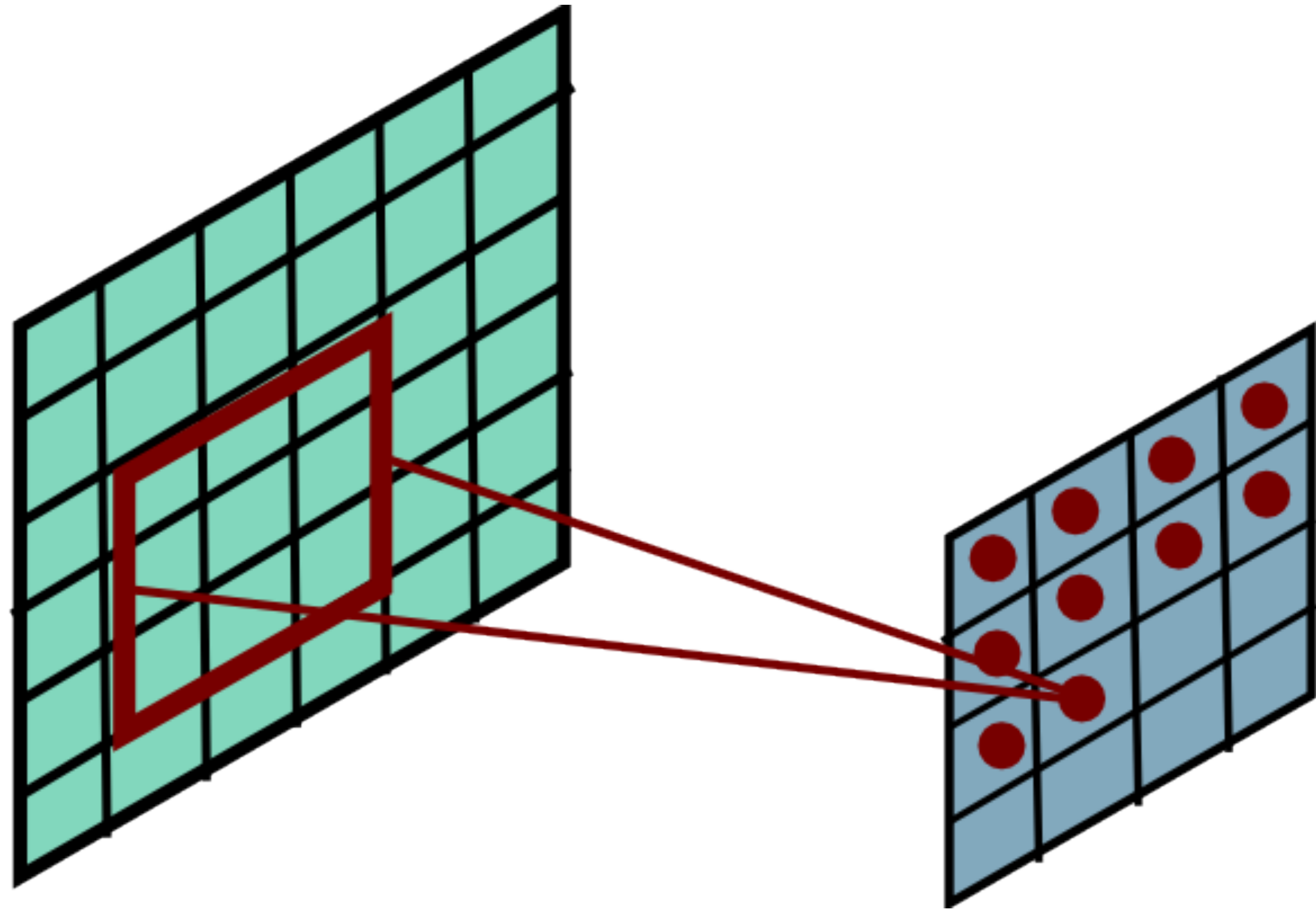
# Convolutional Layer



# Convolutional Layer

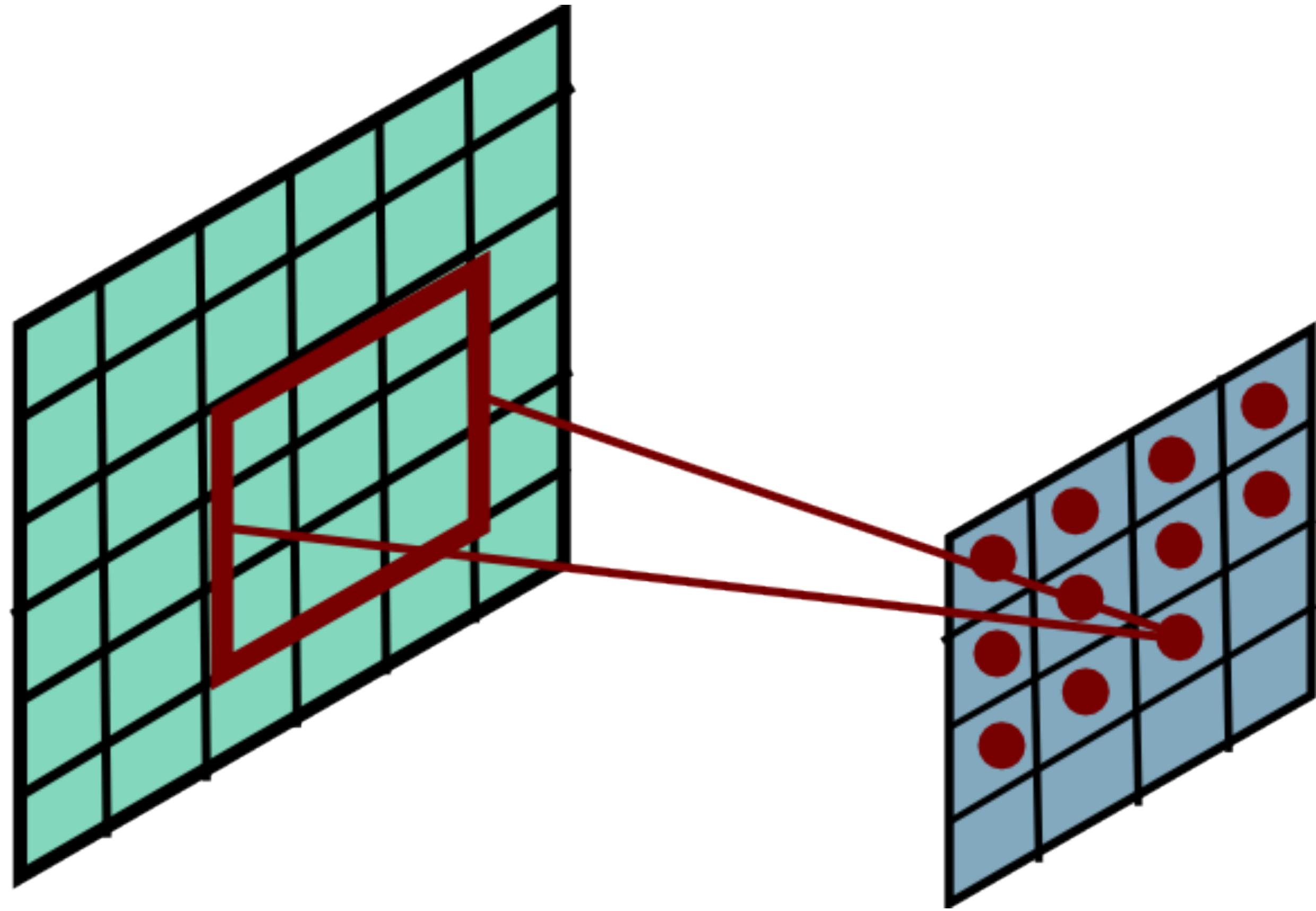


# Convolutional Layer

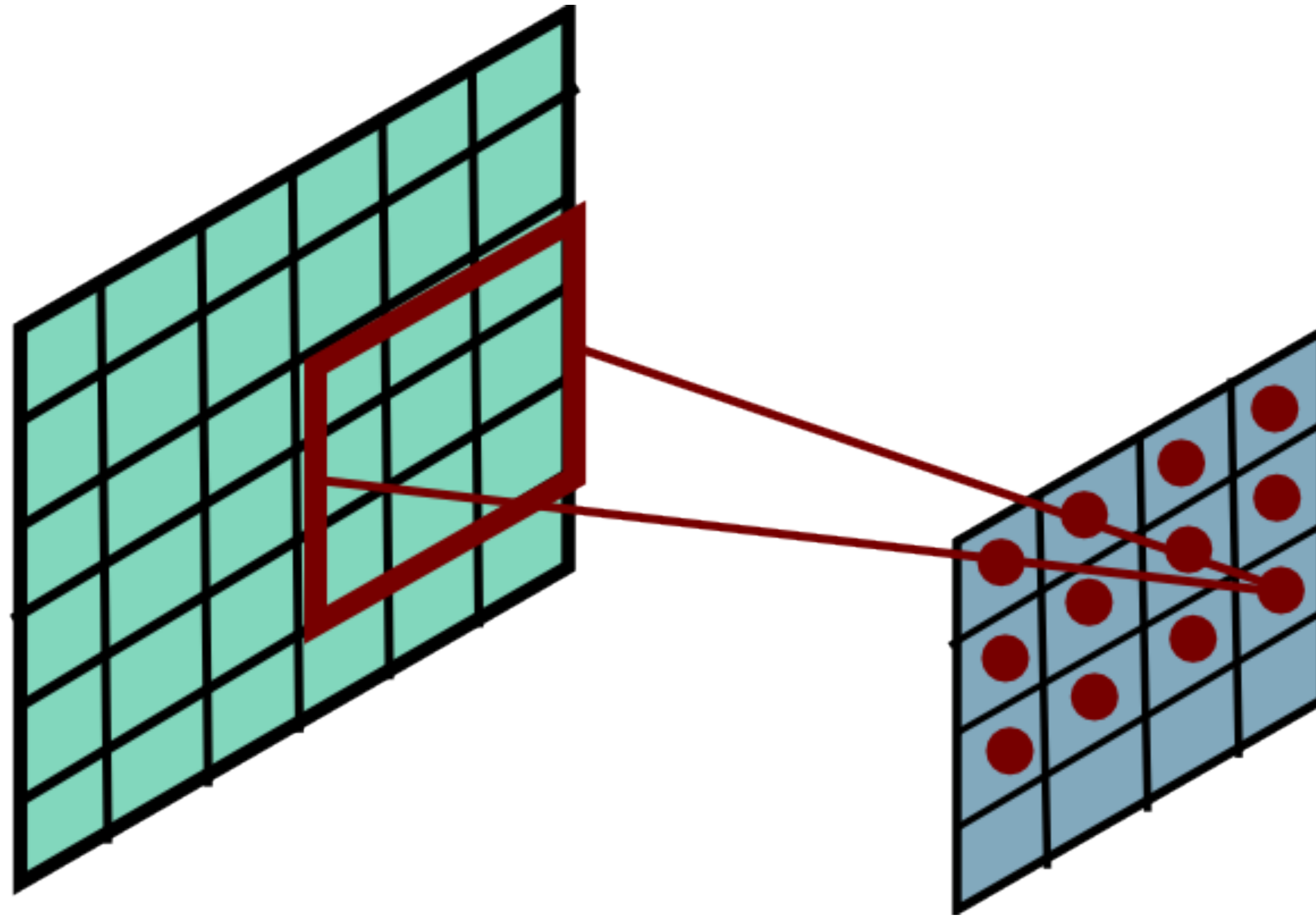




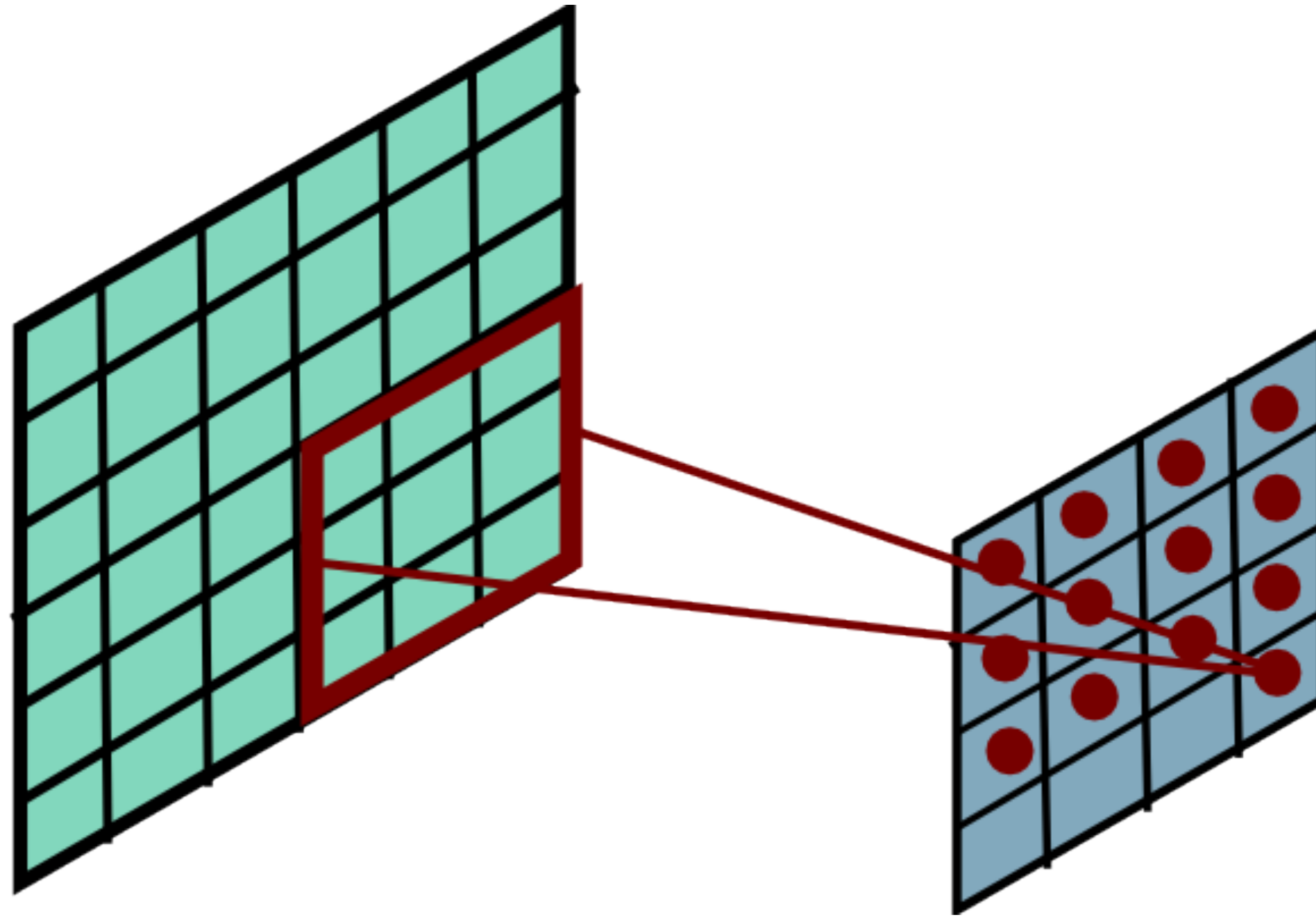
# Convolutional Layer



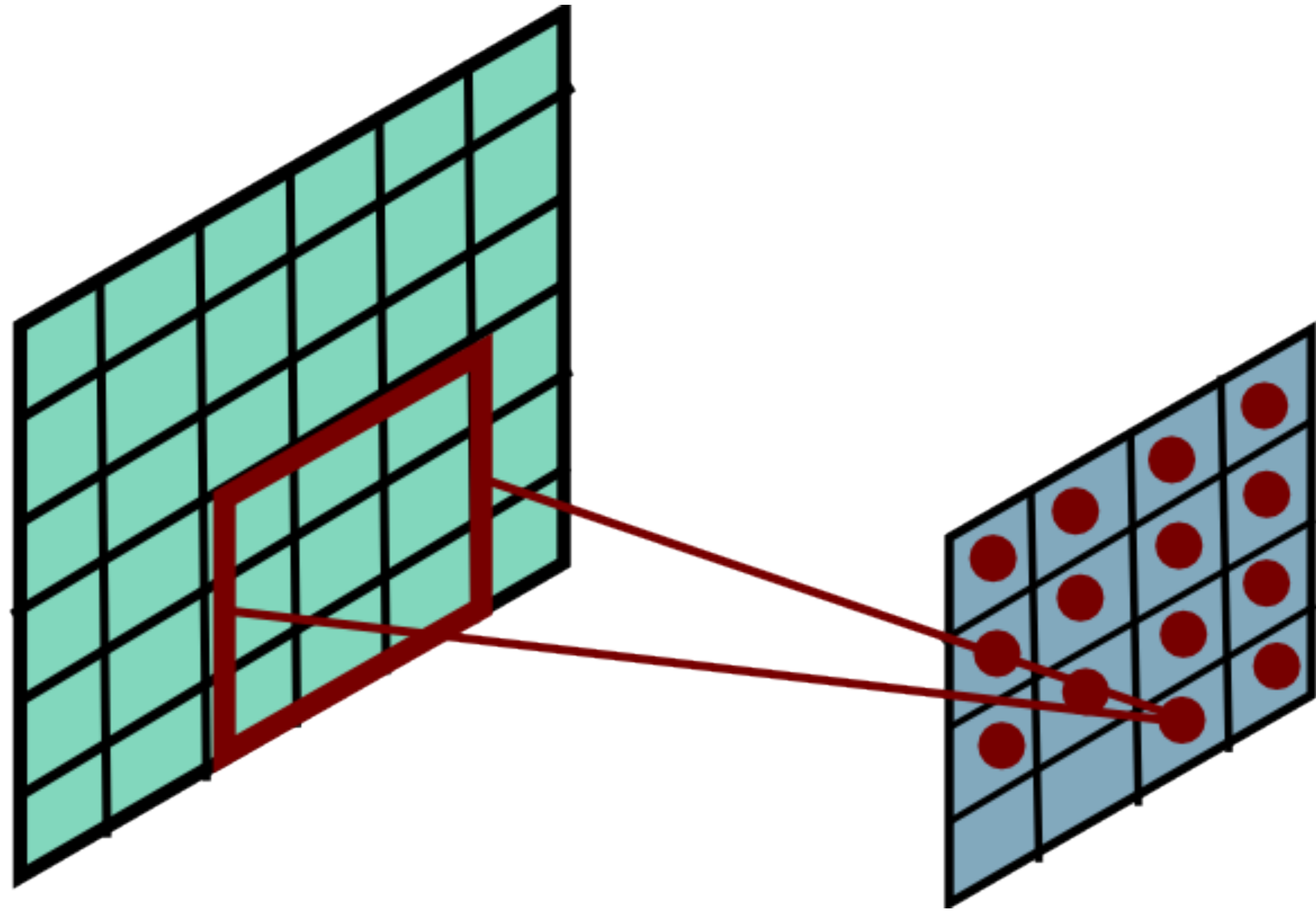
# Convolutional Layer



# Convolutional Layer

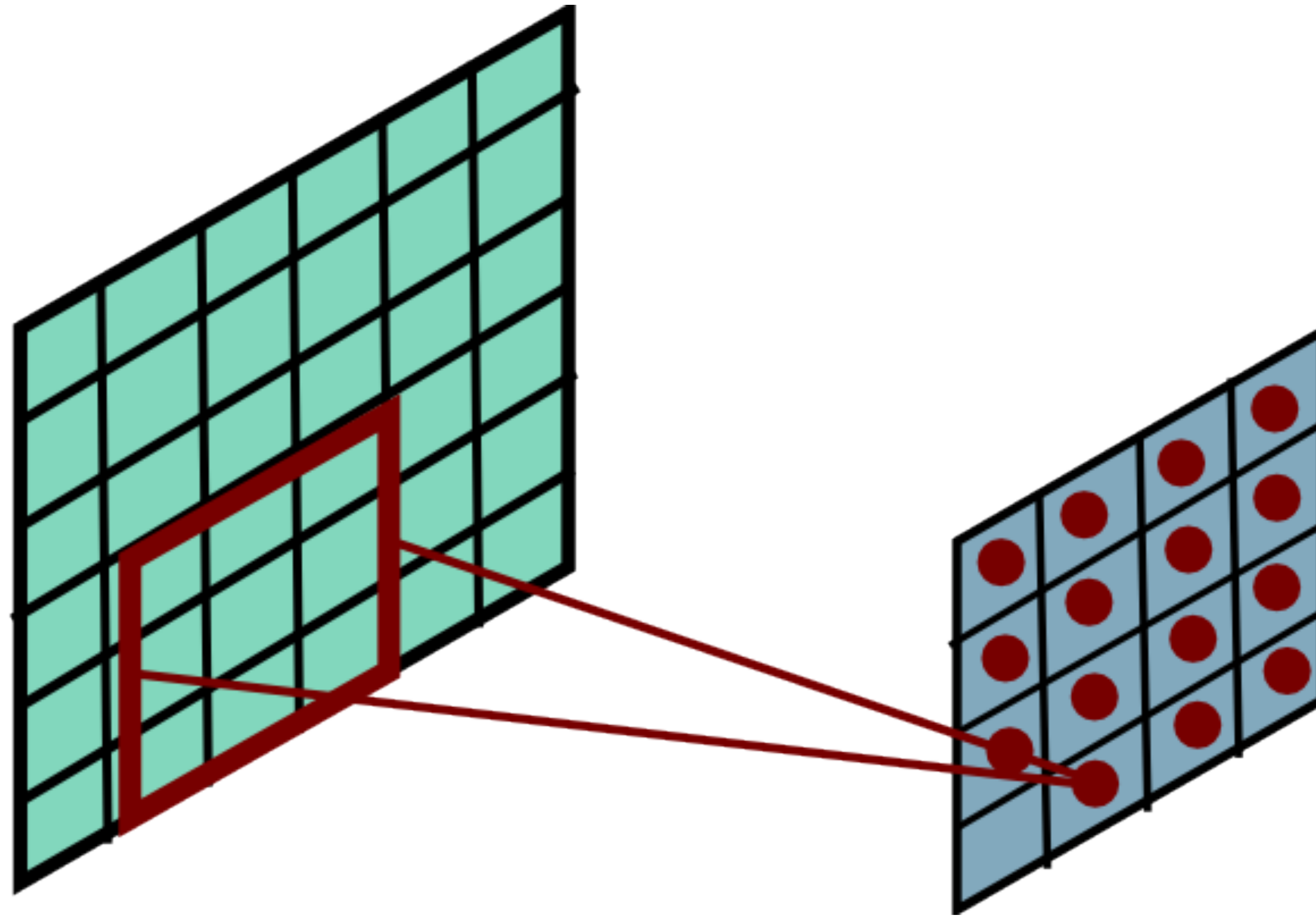


# Convolutional Layer

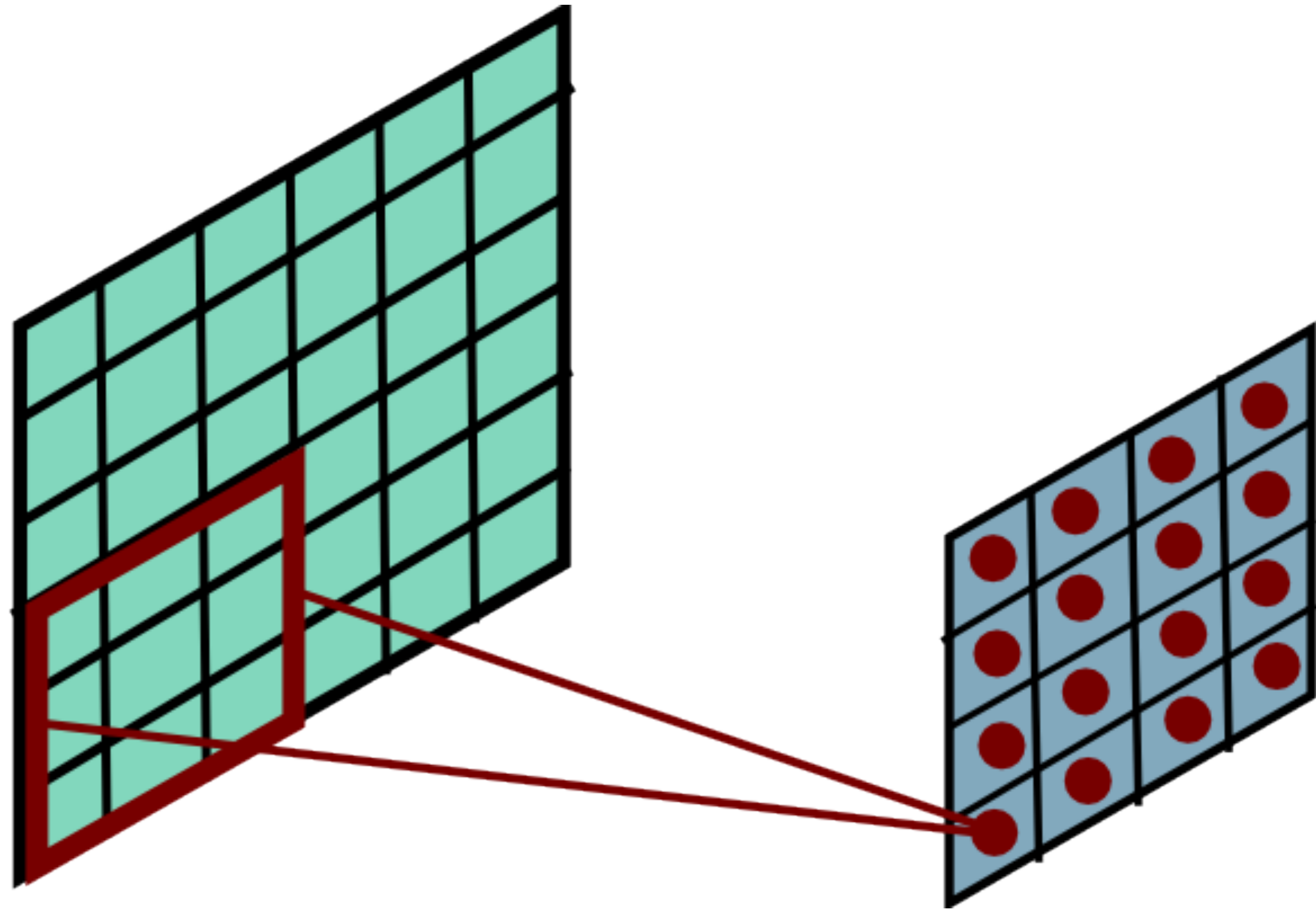




# Convolutional Layer



# Convolutional Layer

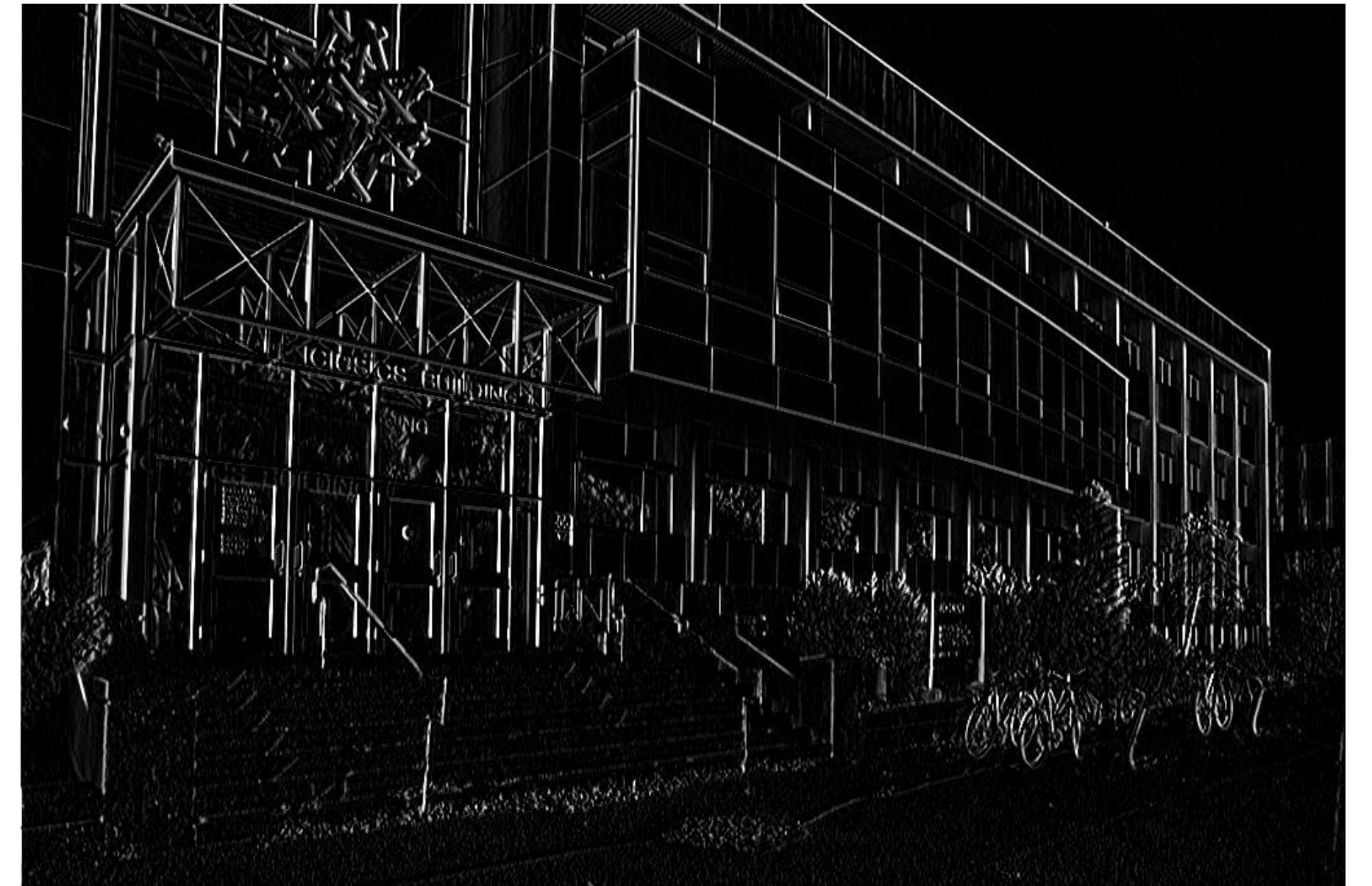




# Convolution Layer



$$\star \begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix} \longrightarrow$$





# Convolution Layer

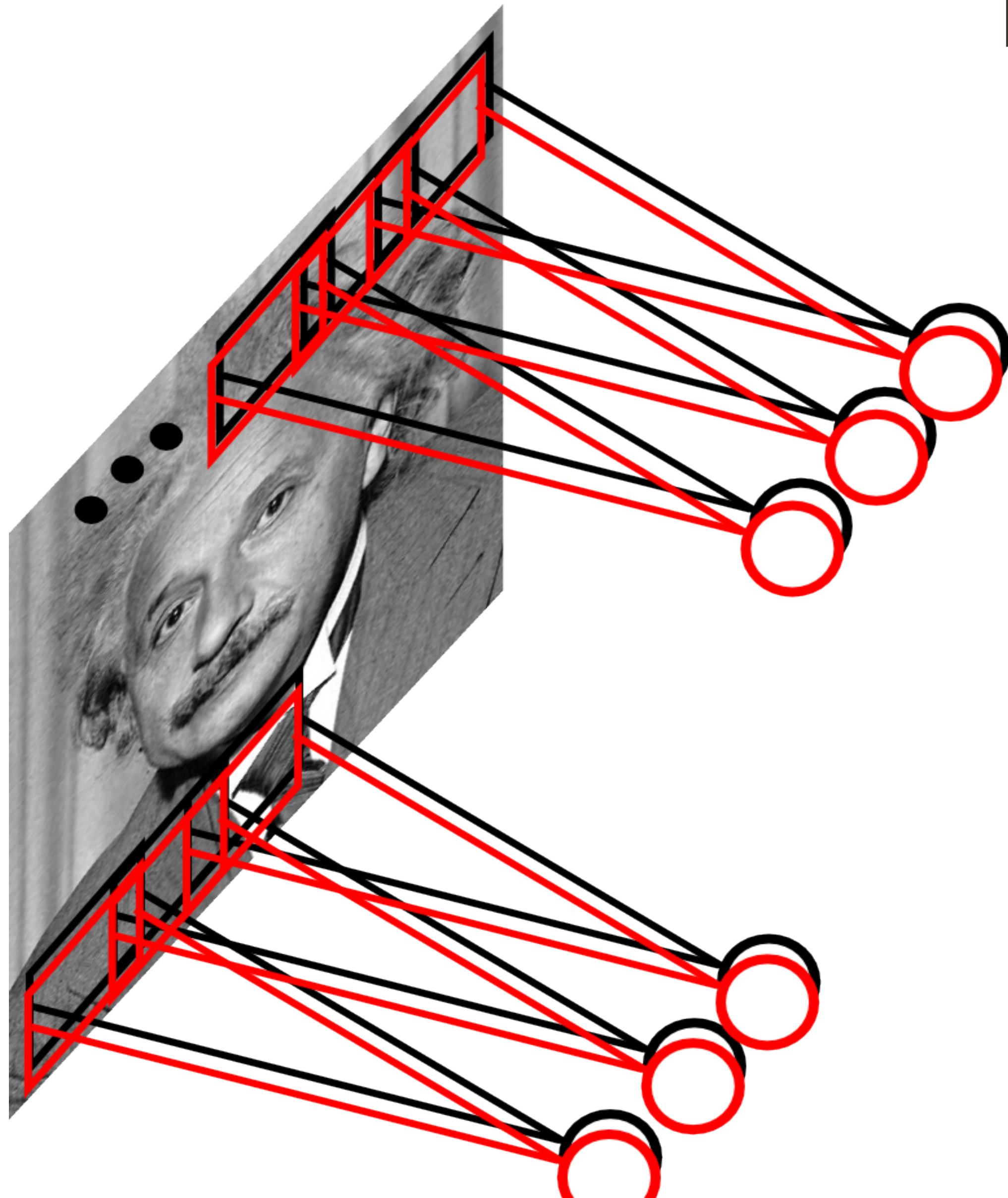


$$\star \begin{bmatrix} 0.11 & 0.11 & 0.11 \\ 0.11 & 0.11 & 0.11 \\ 0.11 & 0.11 & 0.11 \end{bmatrix} \rightarrow$$





# Convolutional Layer



**Example:** 200 x 200 image (small)  
x 40K hidden units

**Filter size:** 10 x 10

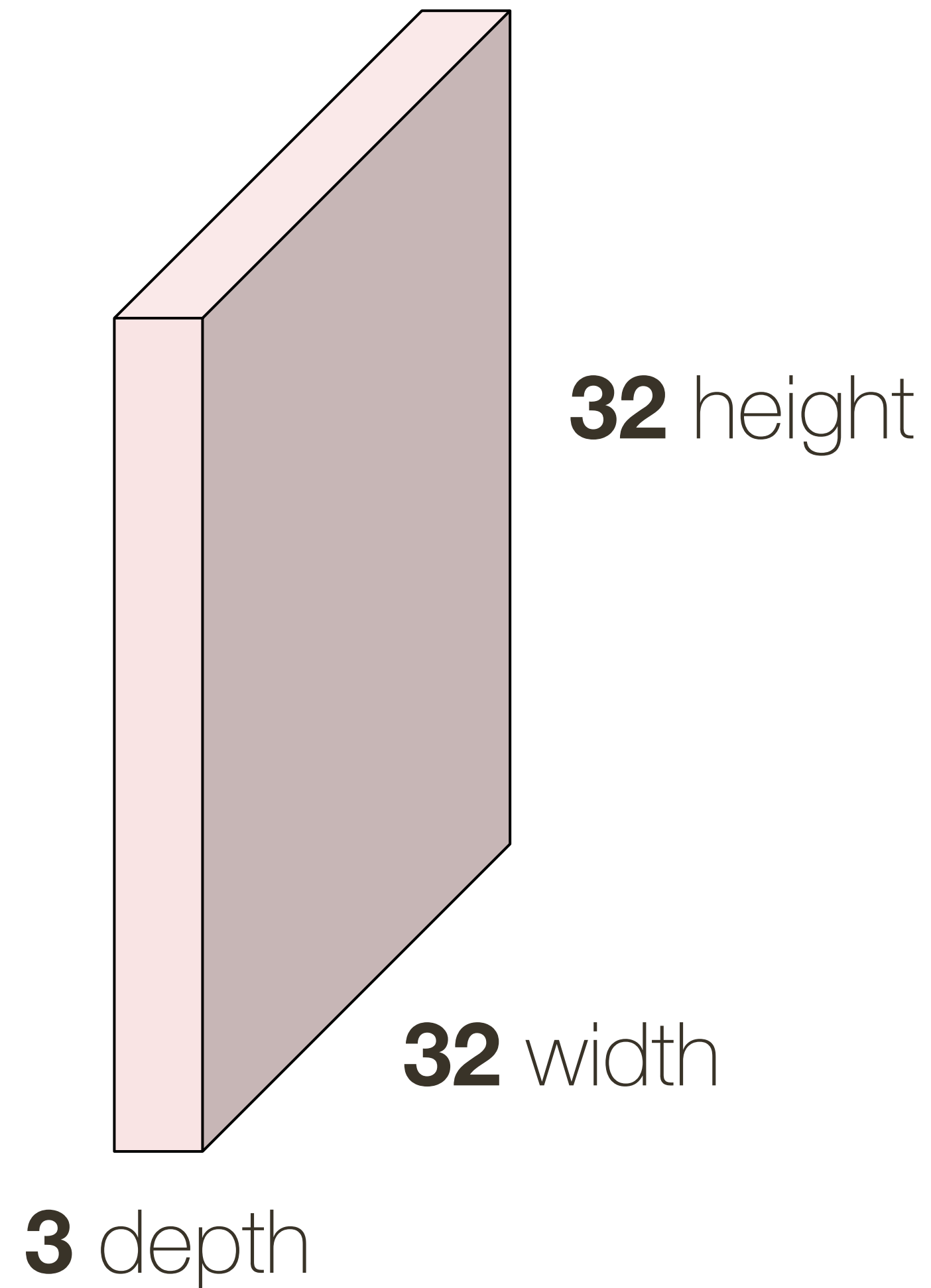
**# of filters:** 20

= 2000 parameters

**Learn multiple filters**

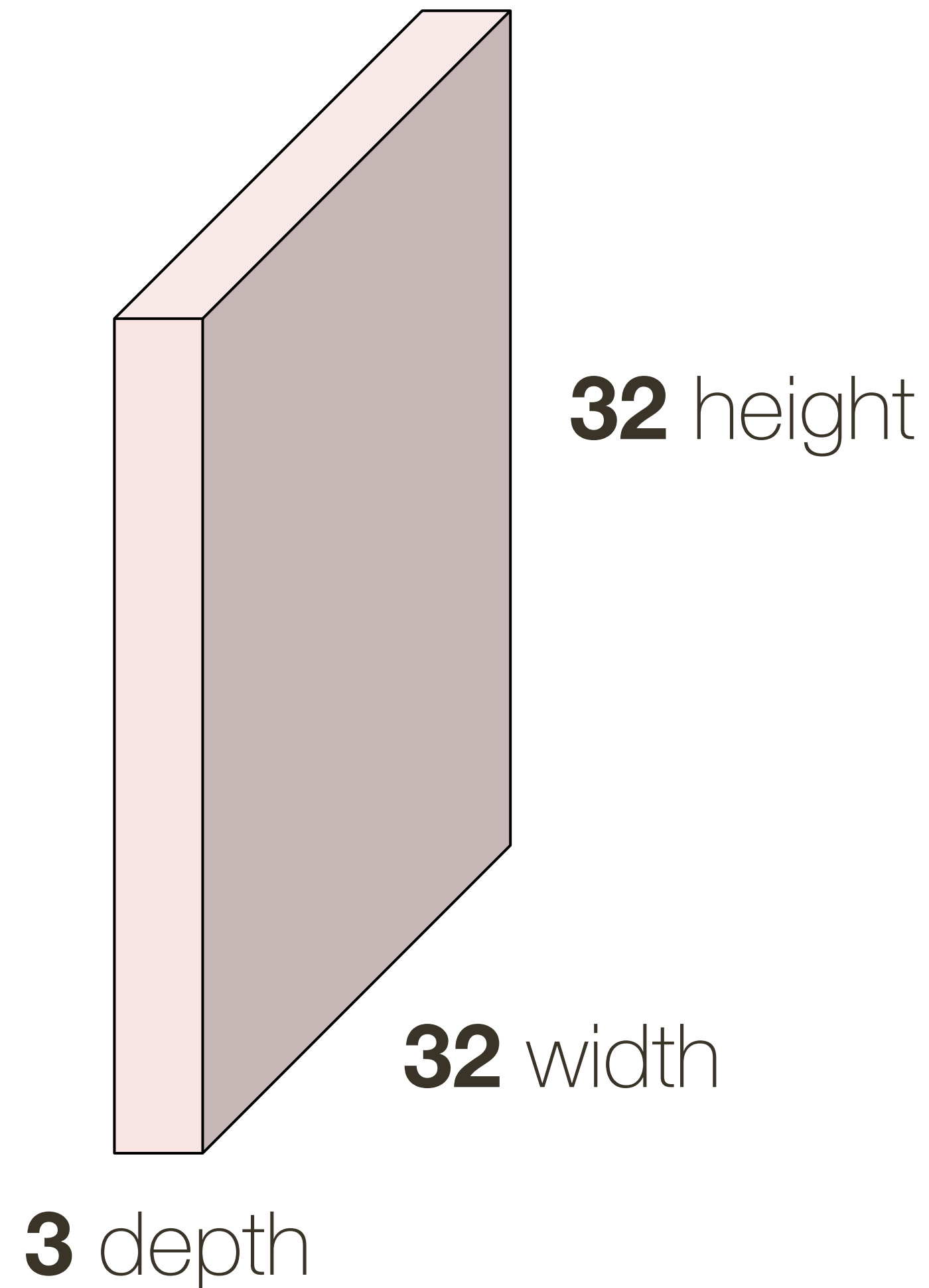
# Convolutional Layer

32 x 32 x 3 **image** (note the image preserves spatial structure)

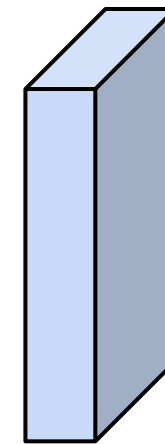


# Convolutional Layer

32 x 32 x 3 **image**



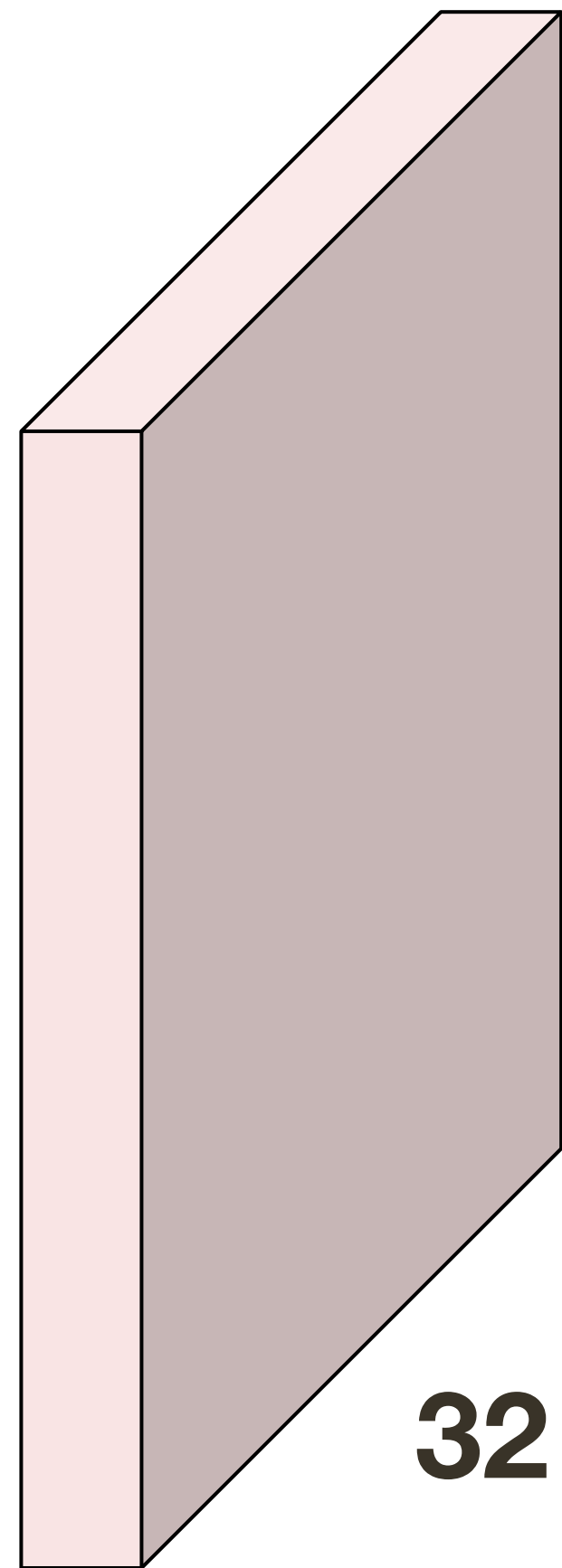
5 x 5 x 3 **filter**



**Convolve** the filter with the image (i.e., “slide over the image spatially, computing dot products”)

# Convolutional Layer

32 x 32 x **3** image



**32** height

**32** width

**3** depth

Filters always extend the full depth of the input volume

5 x 5 x **3** filter

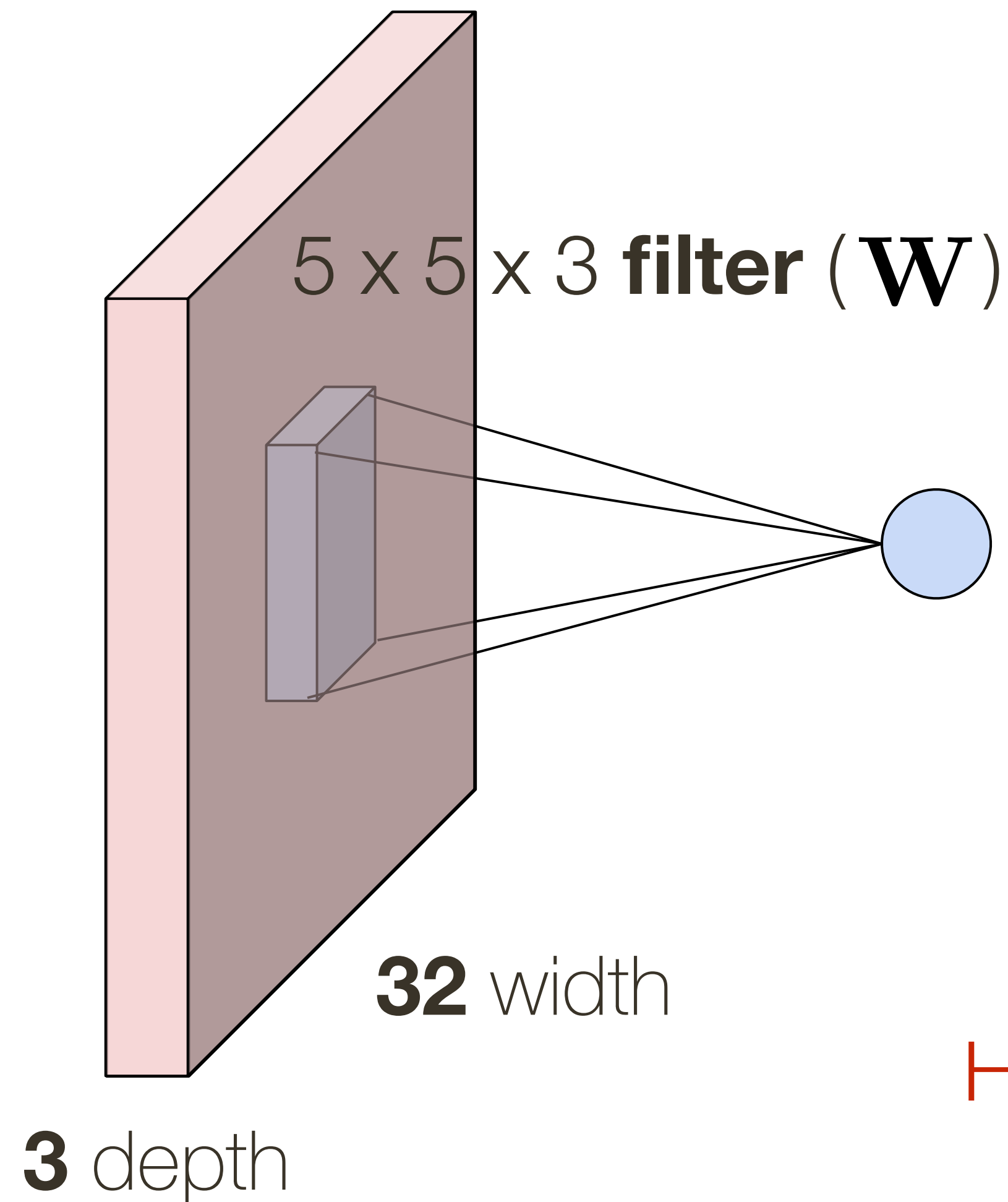


**Convolve** the filter with the image (i.e., “slide over the image spatially, computing dot products”)



# Convolutional Layer

32 x 32 x 3 **image**



**1 number:** the result of taking a dot product between the filter and a small 5 x 5 x 3 part of the image

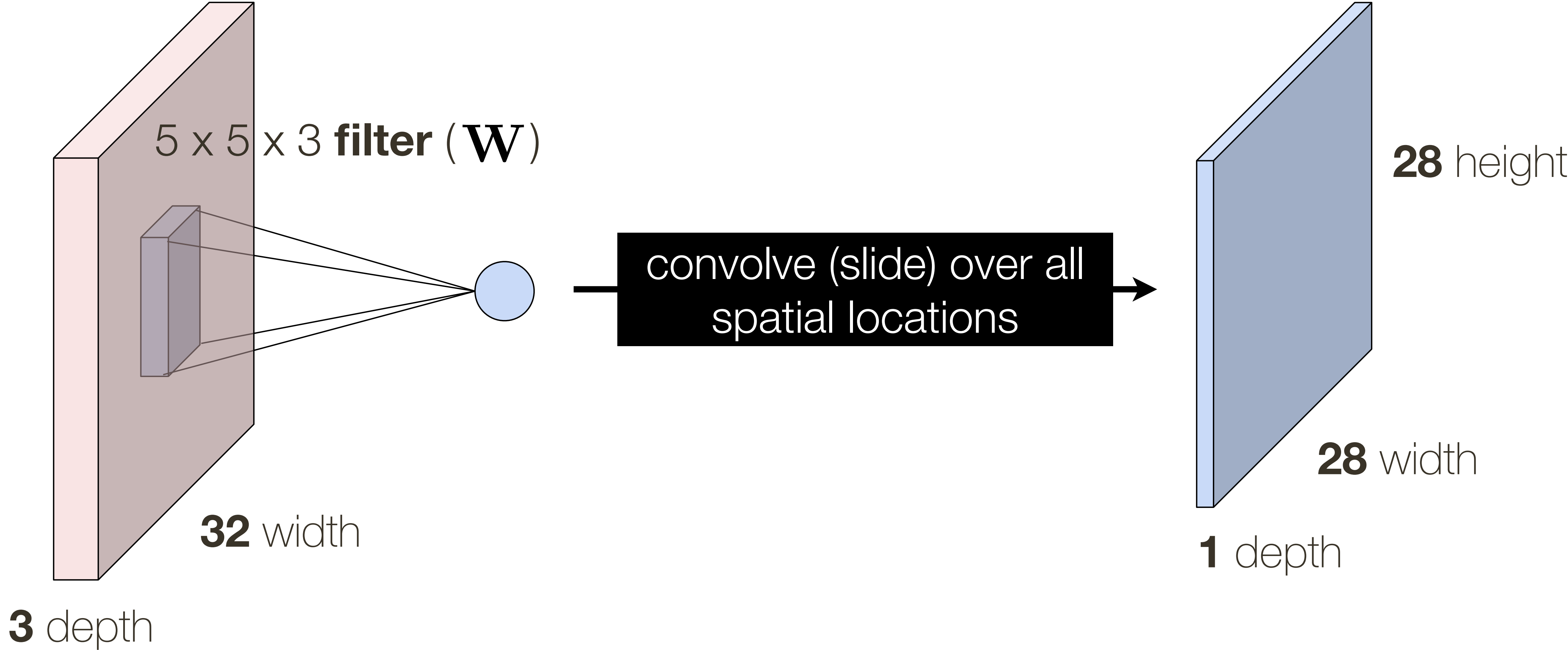
$$\mathbf{W}^T \mathbf{x} + b, \text{ where } \mathbf{W}, \mathbf{x} \in \mathbb{R}^{75}$$

How many **parameters** does the layer have? **76**

# Convolutional Layer

32 x 32 x 3 **image**

**activation** map

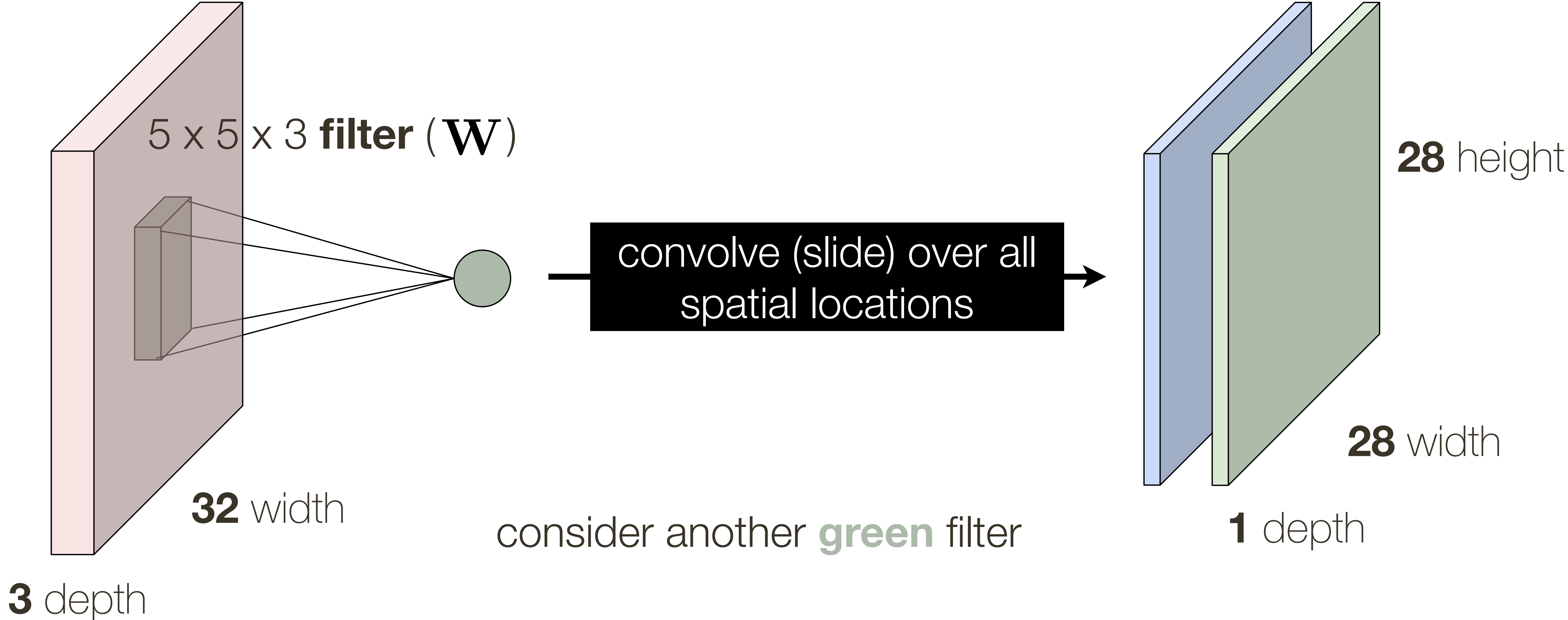


\* slide from Fei-Dei Li, Justin Johnson, Serena Yeung, **cs231n Stanford**

# Convolutional Layer

32 x 32 x 3 **image**

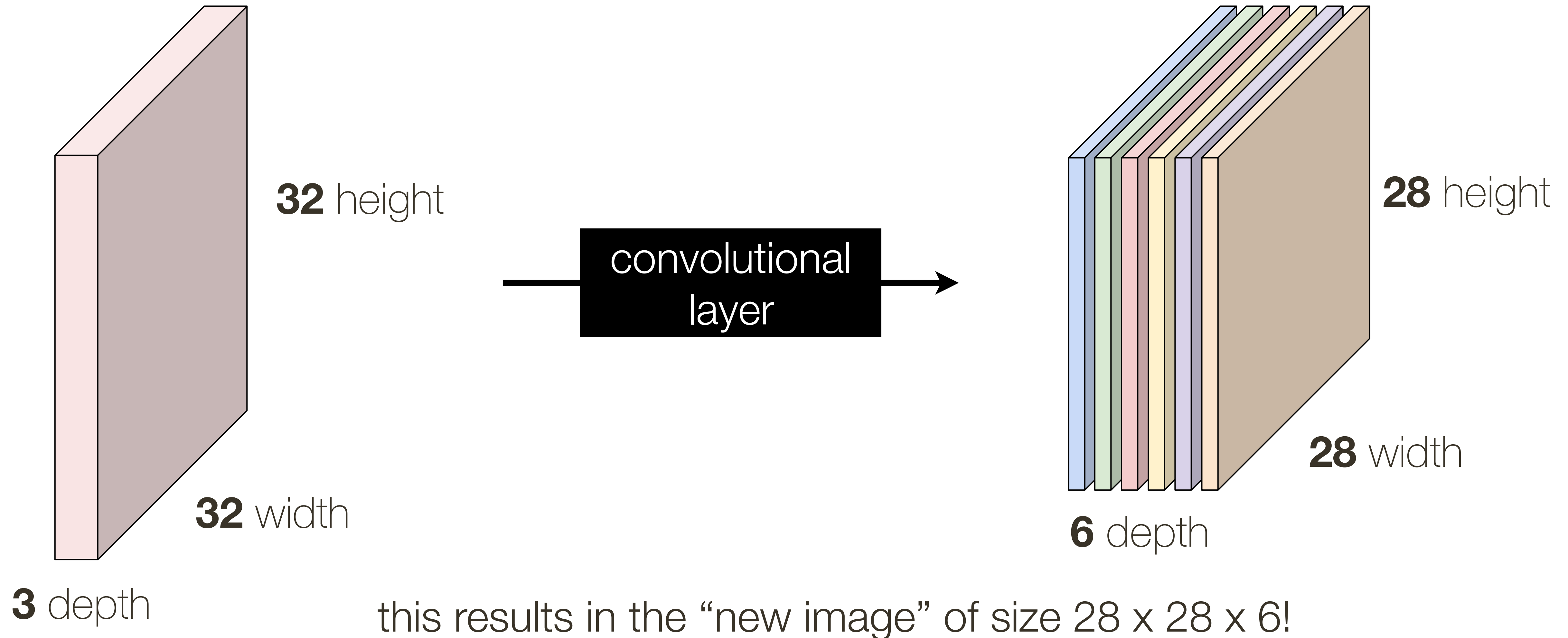
**activation** map



\* slide from Fei-Dei Li, Justin Johnson, Serena Yeung, **cs231n Stanford**

# Convolutional Layer

If we have 6 5x5 filter, we'll get 6 separate activation maps: **activation** map





# Convolutional Layer

The number of neurons in a layer is determined by depth and stride parameter  
— also affected by zero-padding

**Depth:** Controls number of neurons that connect to the same region of the input layer

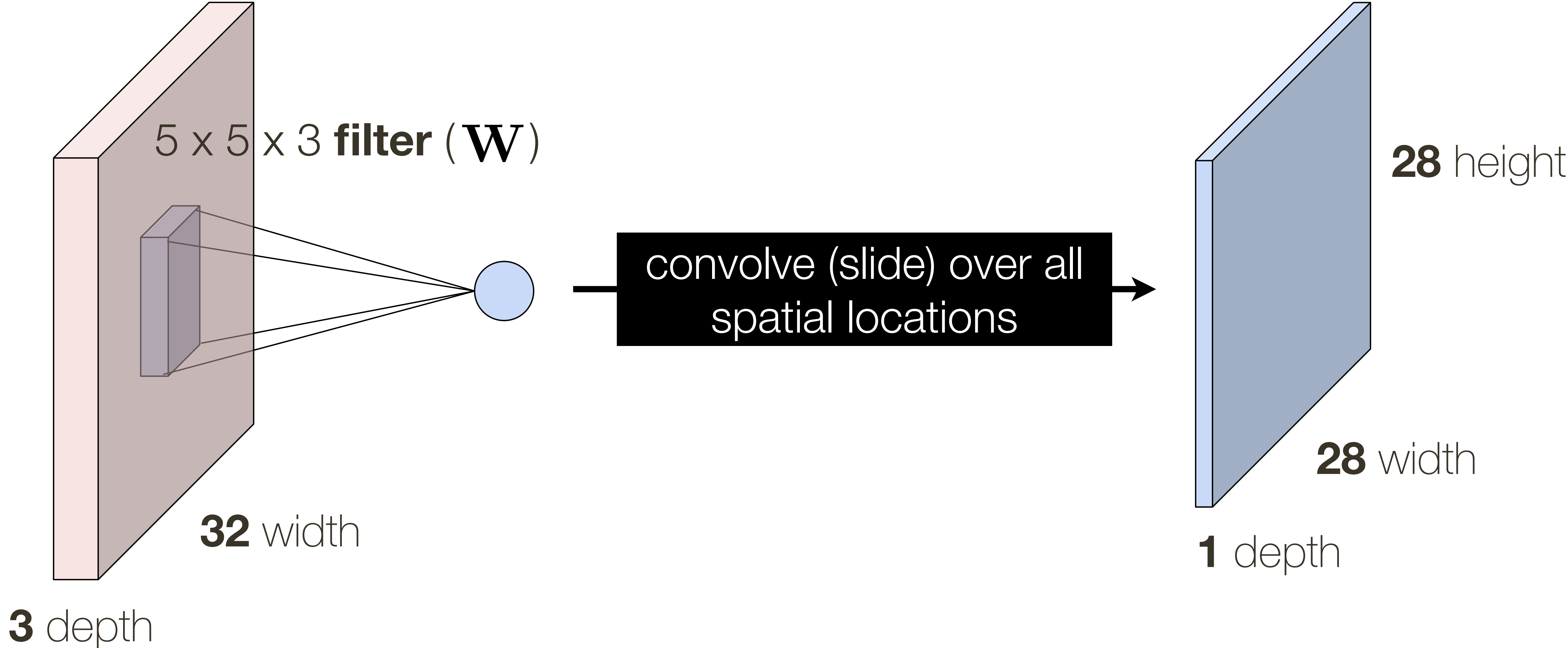
— a set of neurons connected to the same region is called a **depth column**

**Stride:** Controls spatial density. How far apart are depth columns?

# Convolutional Layer: Closer Look at **Spatial Dimensions**

32 x 32 x 3 **image**

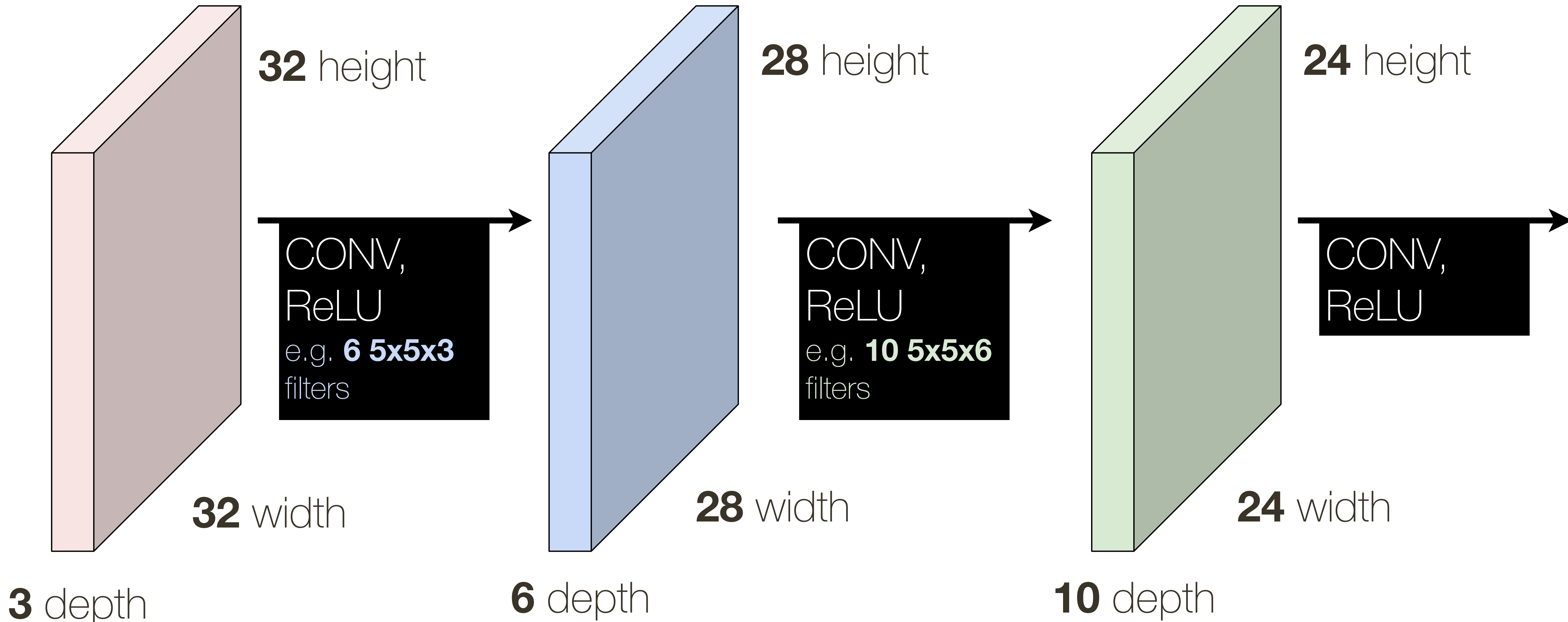
**activation** map



\* slide from Fei-Dei Li, Justin Johnson, Serena Yeung, **cs231n Stanford**

# Convolutional Neural Network (ConvNet)

With padding we can achieve no shrinking (32 -> 28 -> 24); shrinking quickly (which happens with larger filters) doesn't work well in practice



\* slide from Fei-Dei Li, Justin Johnson, Serena Yeung, **cs231n Stanford**

# Convolutional Neural Network (ConvNet)

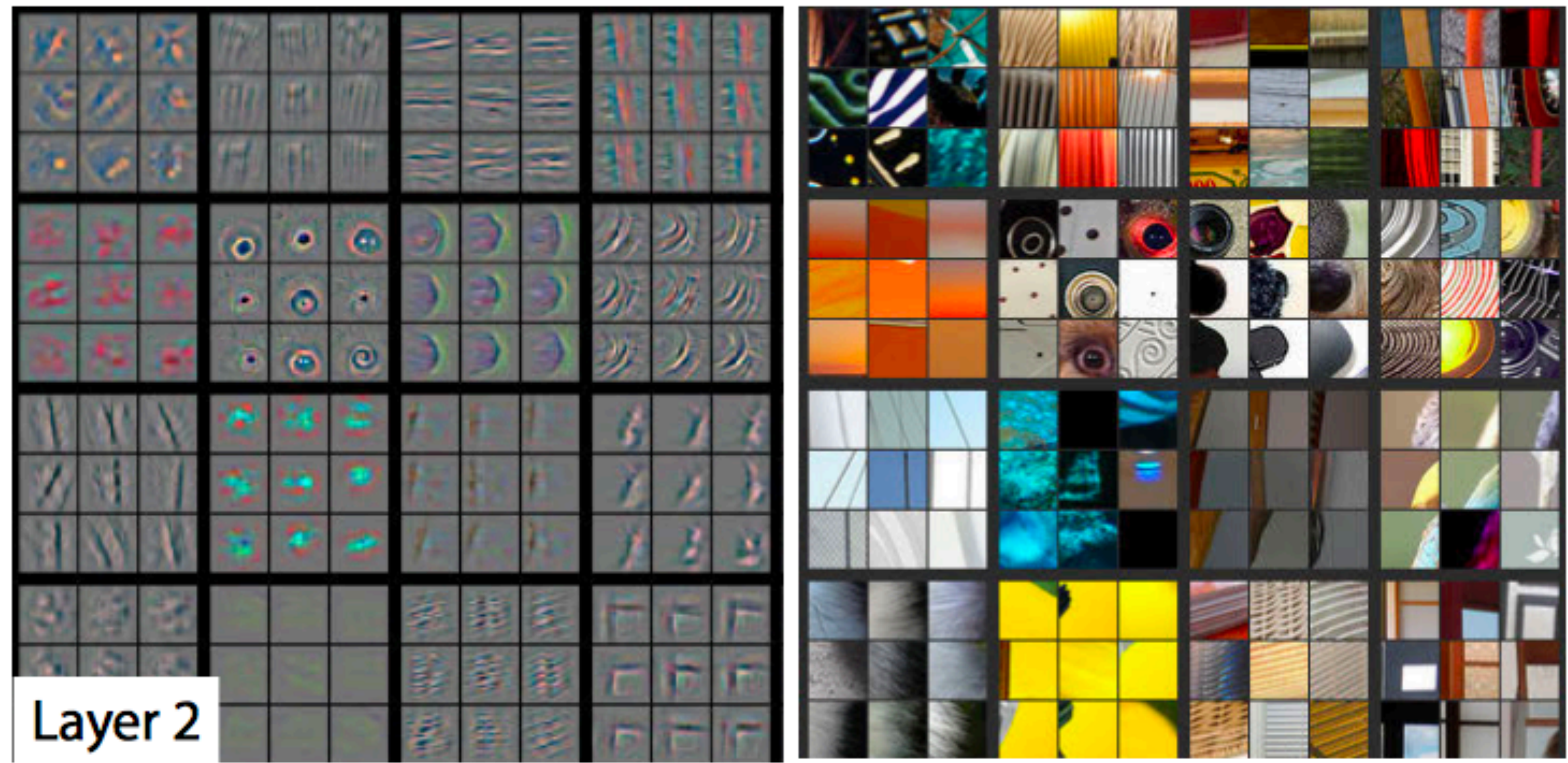
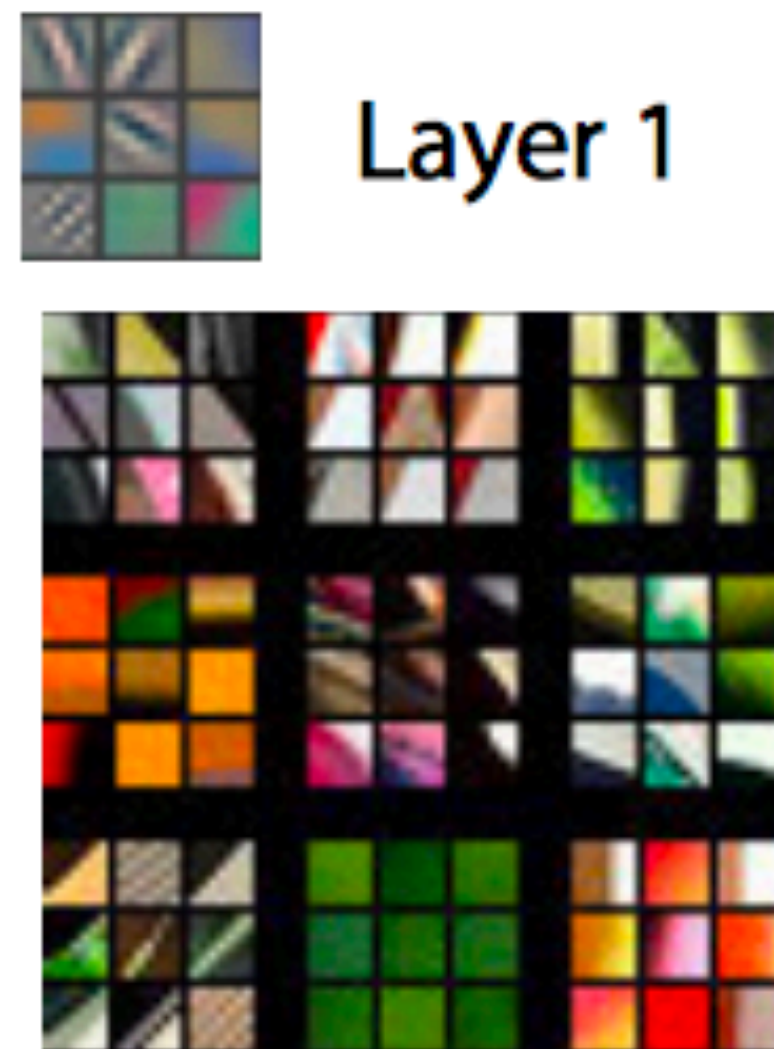
**Convolutional neural networks** can be seen as learning a hierarchy of filters.

As we go deeper in the network, filters learn and respond to increasingly specialized structures

— The first layers may contain simple orientation filters, middle layers may respond to common substructures, and final layers may respond to entire objects

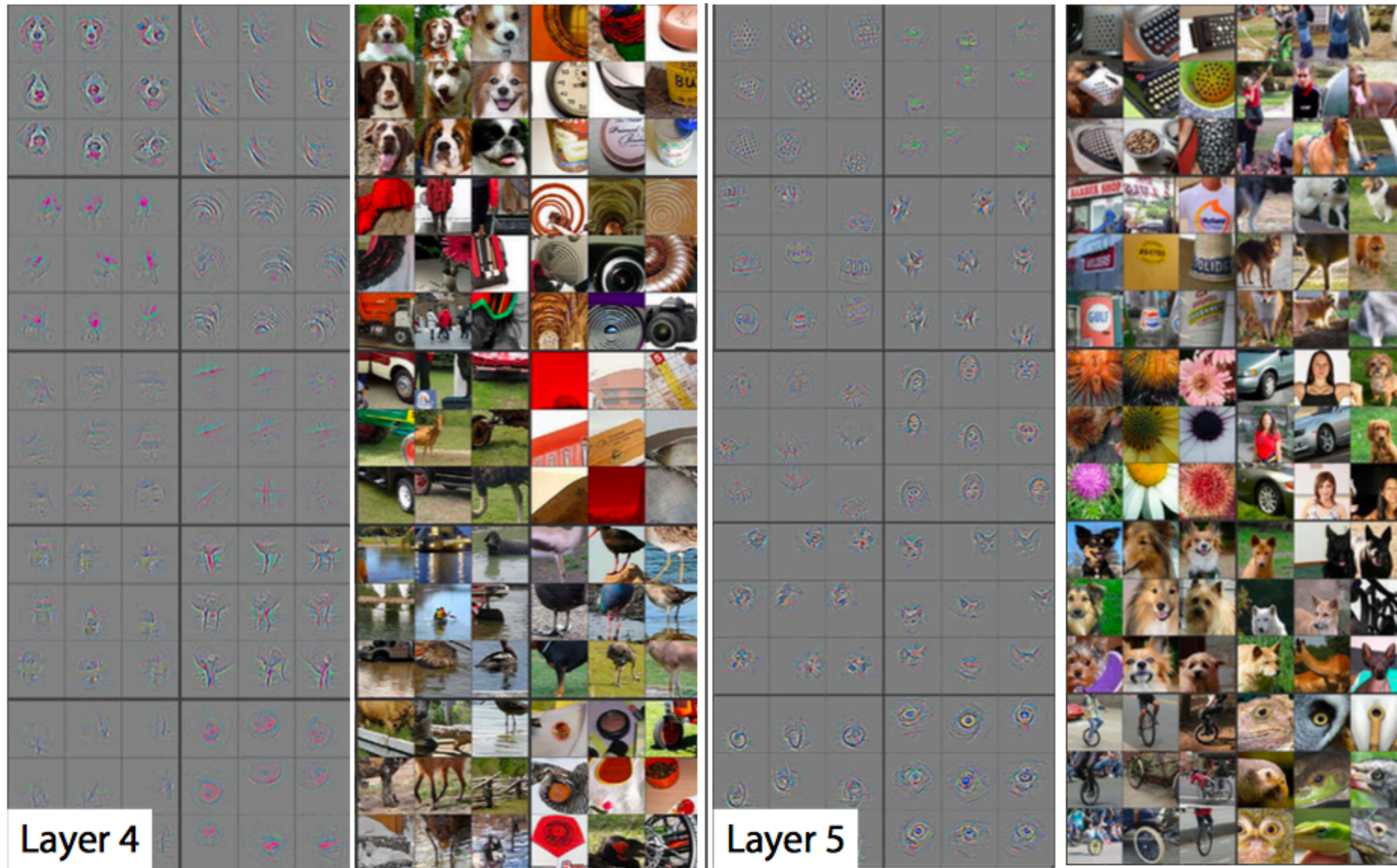


# What **filters** do networks learn?





# What **filters** do networks learn?



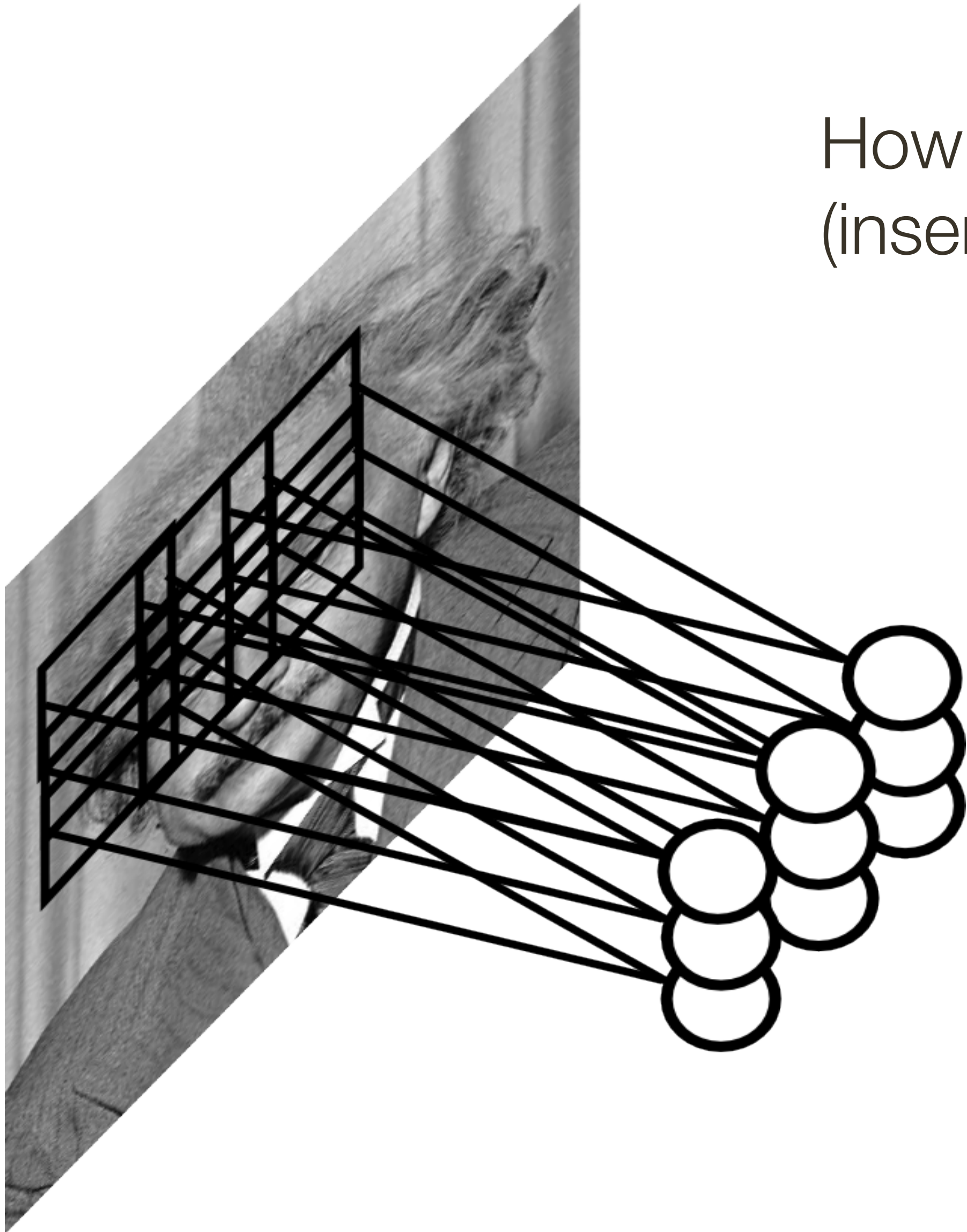
[ Zeiler and Fergus, 2013 ]



# Pooling Layer

Let us assume the filter is an “eye” detector

How can we make detection spatially invariant  
(insensitive to position of the eye in the image)

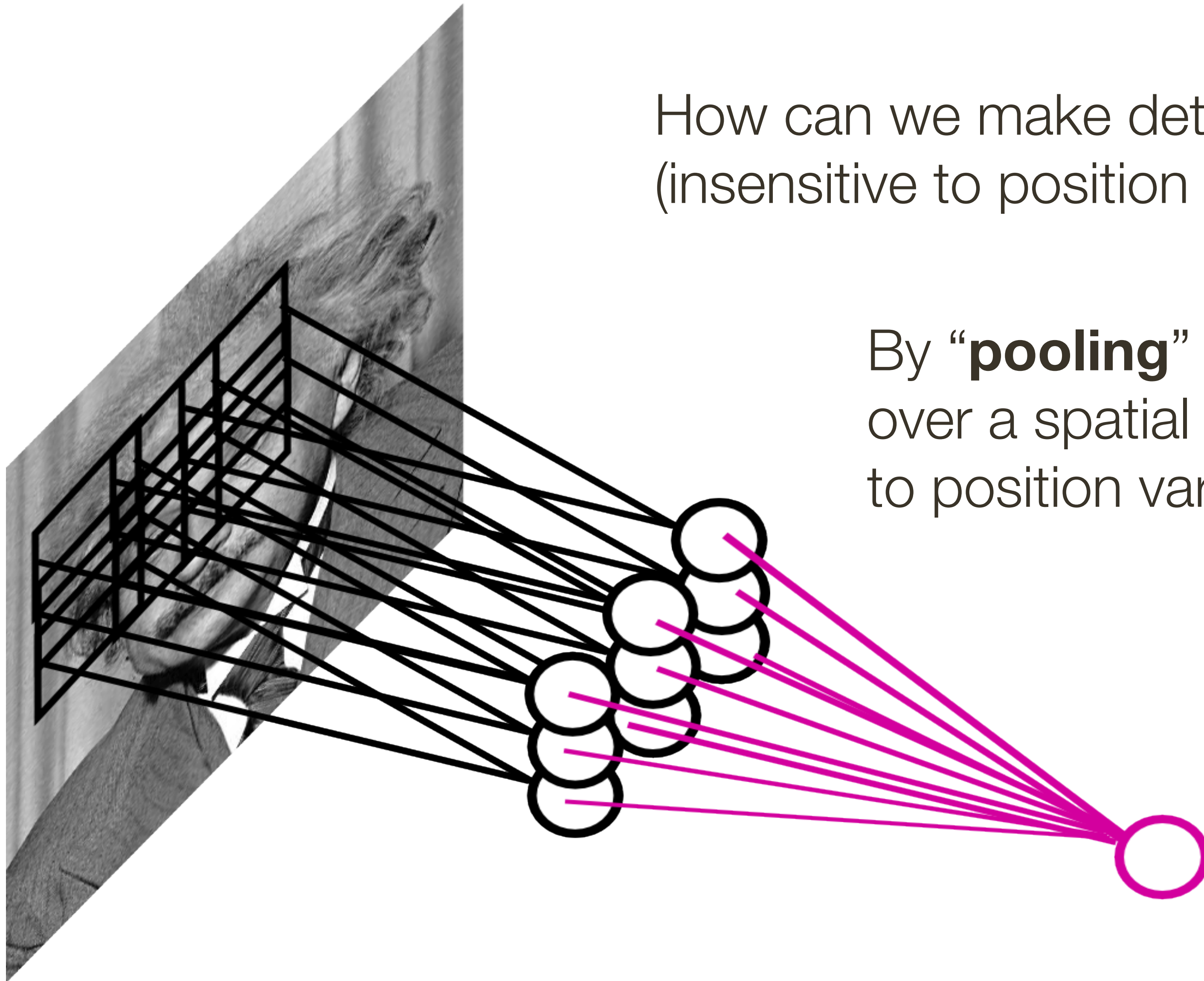


# Pooling Layer

Let us assume the filter is an “eye” detector

How can we make detection spatially invariant (insensitive to position of the eye in the image)

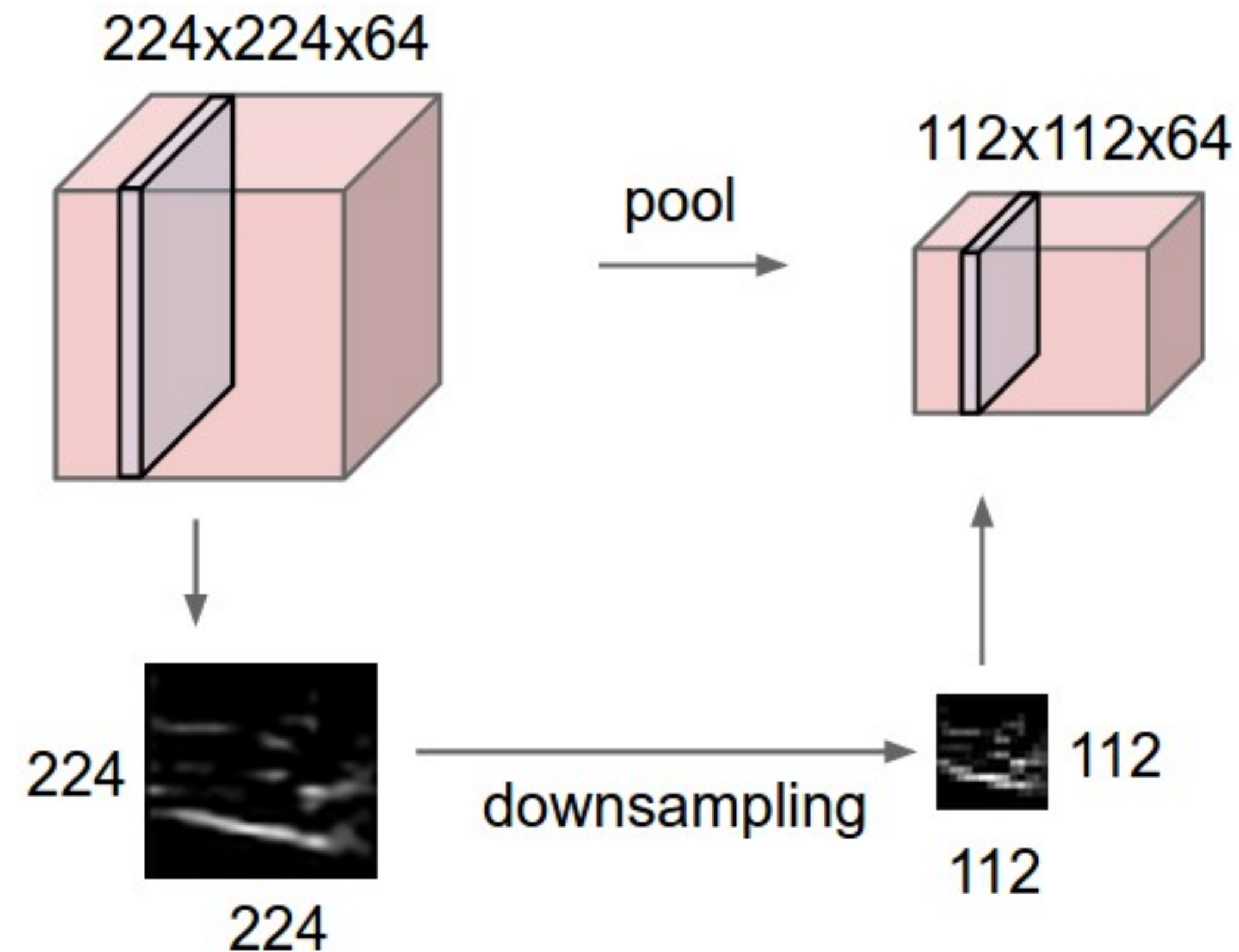
By “**pooling**” (e.g., taking a max) response over a spatial locations we gain robustness to position variations





# Pooling Layer

- Makes representation smaller, more manageable and spatially invariant
- Operates over each activation map independently



How many **parameters**?

**None!**

# Max Pooling

activation map

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

max pool with 2 x 2 filter  
and stride of 2

6	8
3	4

# Average Pooling

activation map

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

avg pool with 2 x 2 filter  
and stride of 2

3.25	5.25
2	2

# Object Classification

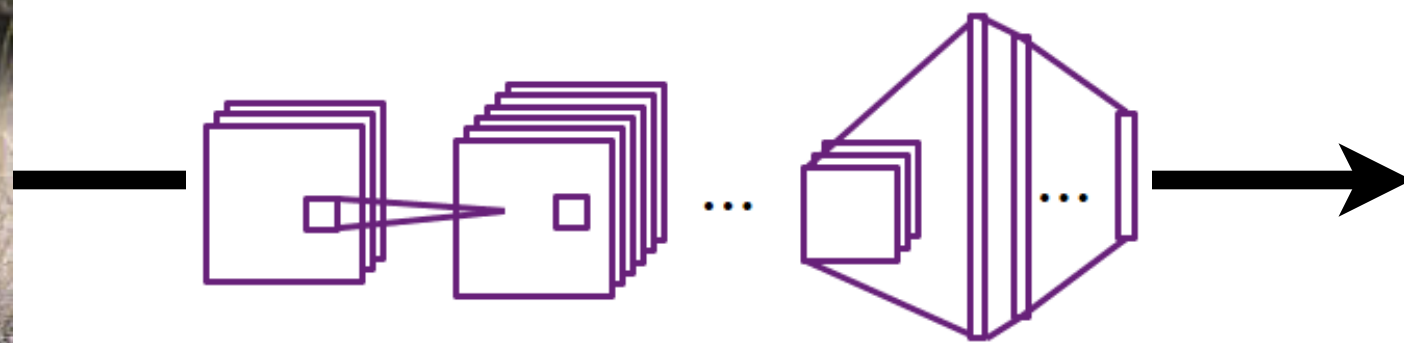


Category	Prediction
Dog	No
Cat	No
Couch	No
Flowers	No
Leopard	<b>Yes</b>
...	...

**Problem:** For each image predict which category it belongs to out of a fixed set



# Object Classification

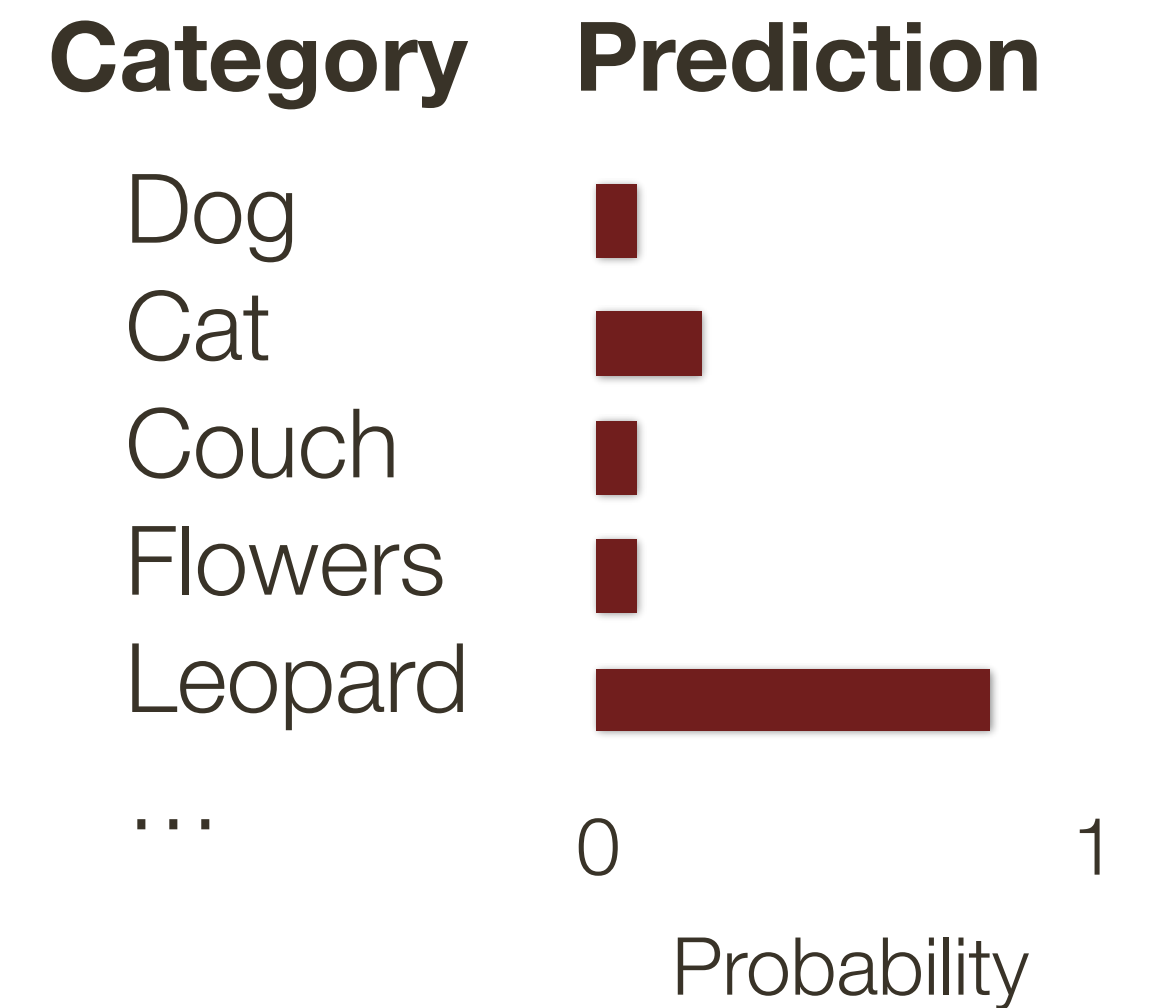
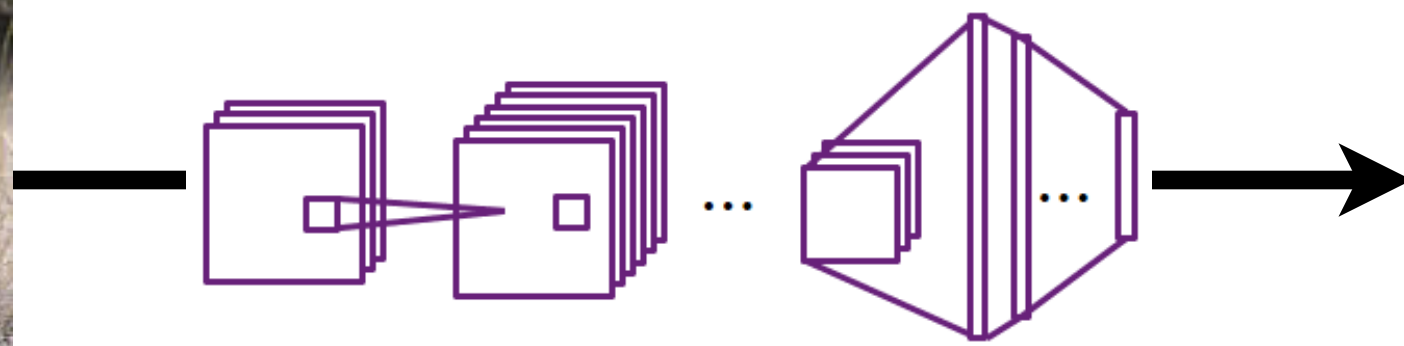


Category	Prediction
Dog	No
Cat	No
Couch	No
Flowers	No
Leopard	<b>Yes</b>
...	...

**Problem:** For each image predict which category it belongs to out of a fixed set



# Object Classification



**Problem:** For each image predict which category it belongs to out of a fixed set

# R-CNN

[ Girshick et al, CVPR 2014 ]



Input **Image**

\* image from Ross Girshick



# R-CNN

[ Girshick et al, CVPR 2014 ]



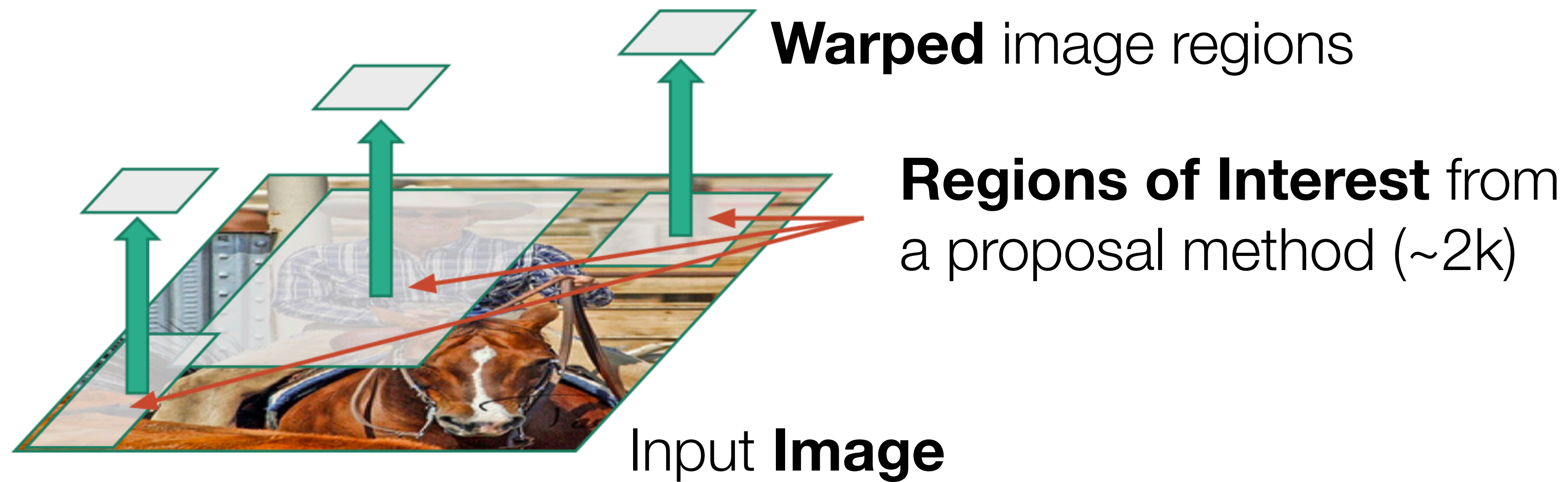
**Regions of Interest** from  
a proposal method (~2k)

Input **Image**



# R-CNN

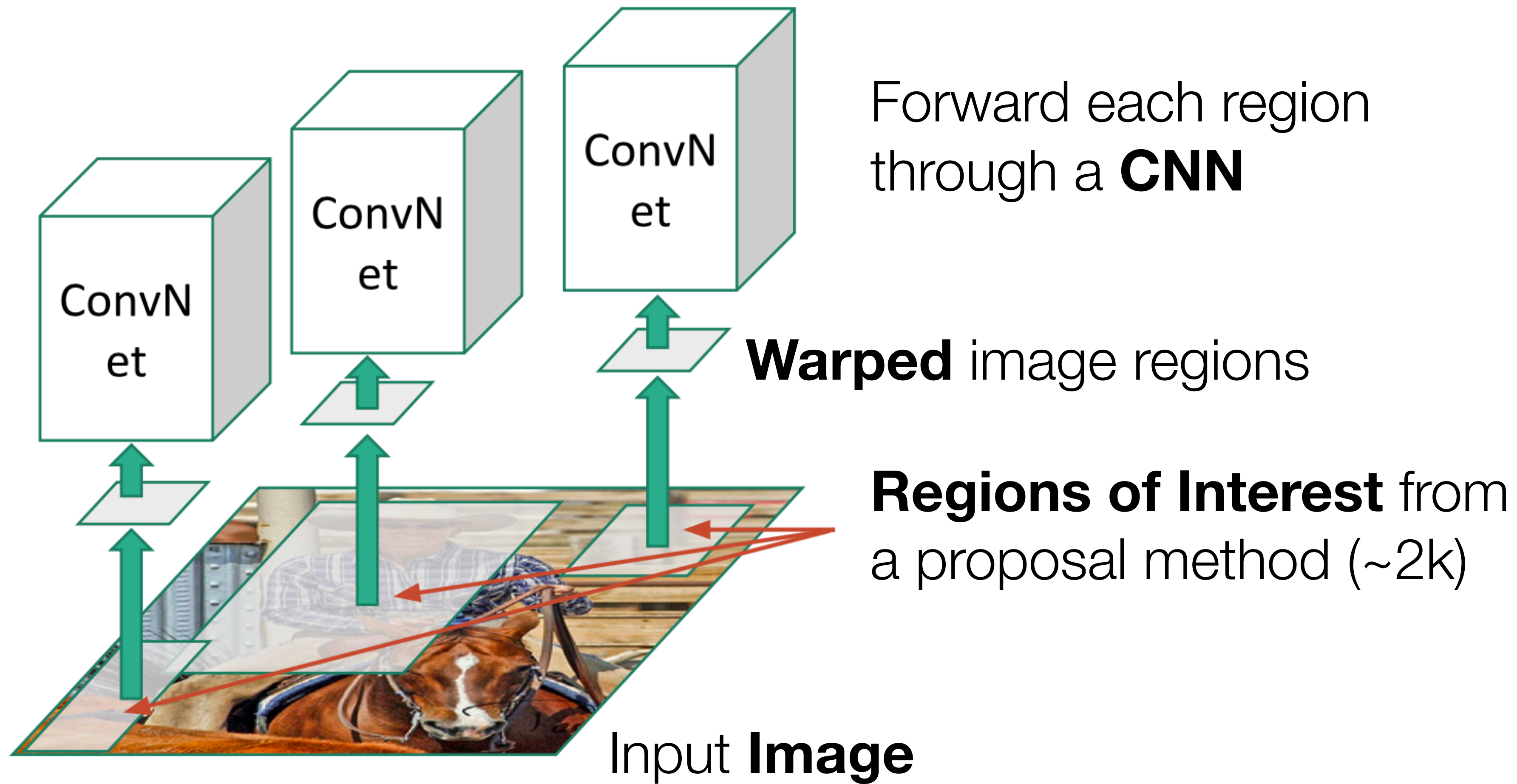
[ Girshick et al, CVPR 2014 ]



\* image from Ross Girshick

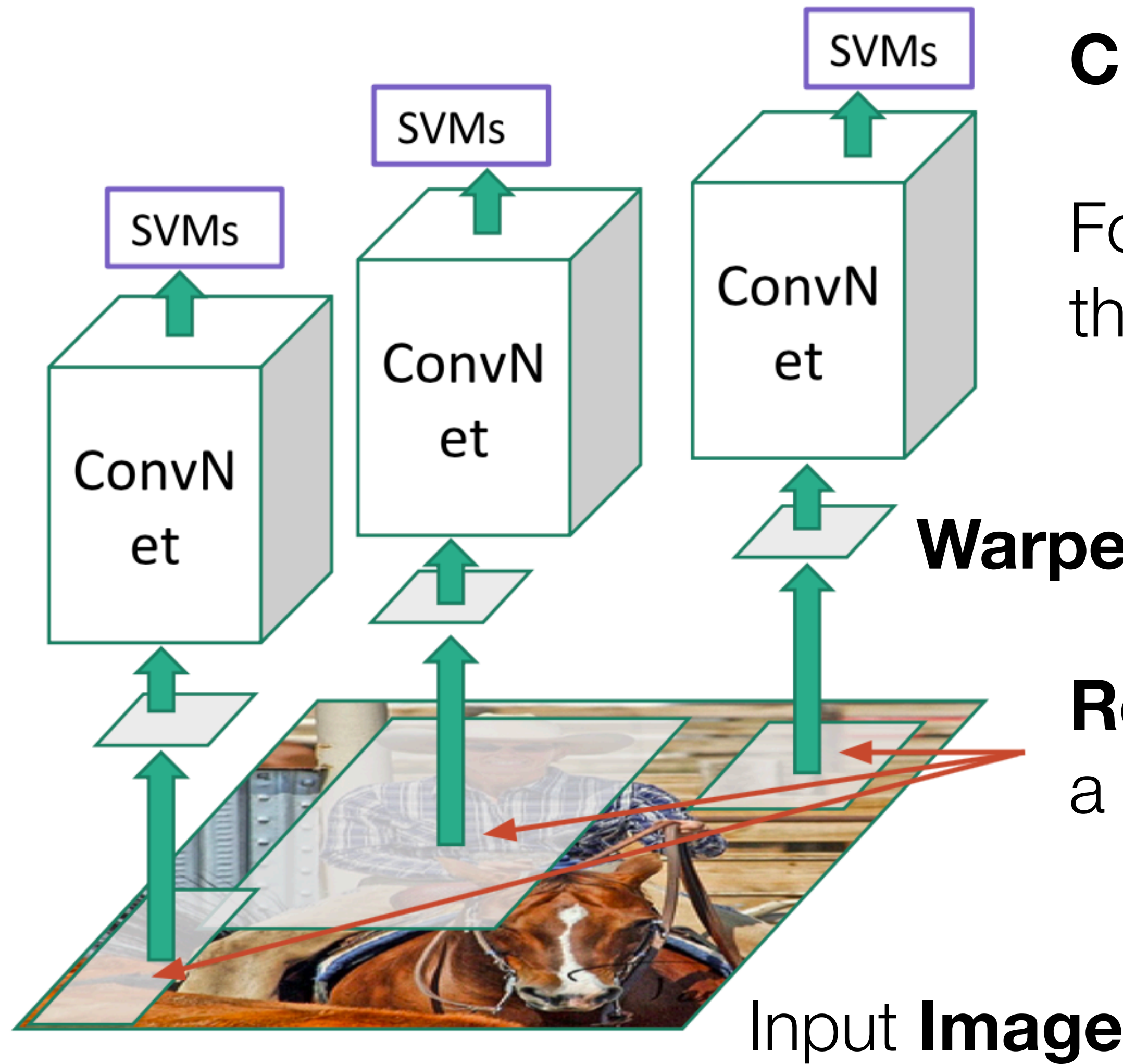
# R-CNN

[ Girshick et al, CVPR 2014 ]



# R-CNN

[ Girshick et al, CVPR 2014 ]



**Classify** regions with SVM

Forward each region through a **CNN**

**Warped** image regions

**Regions of Interest** from a proposal method (~2k)

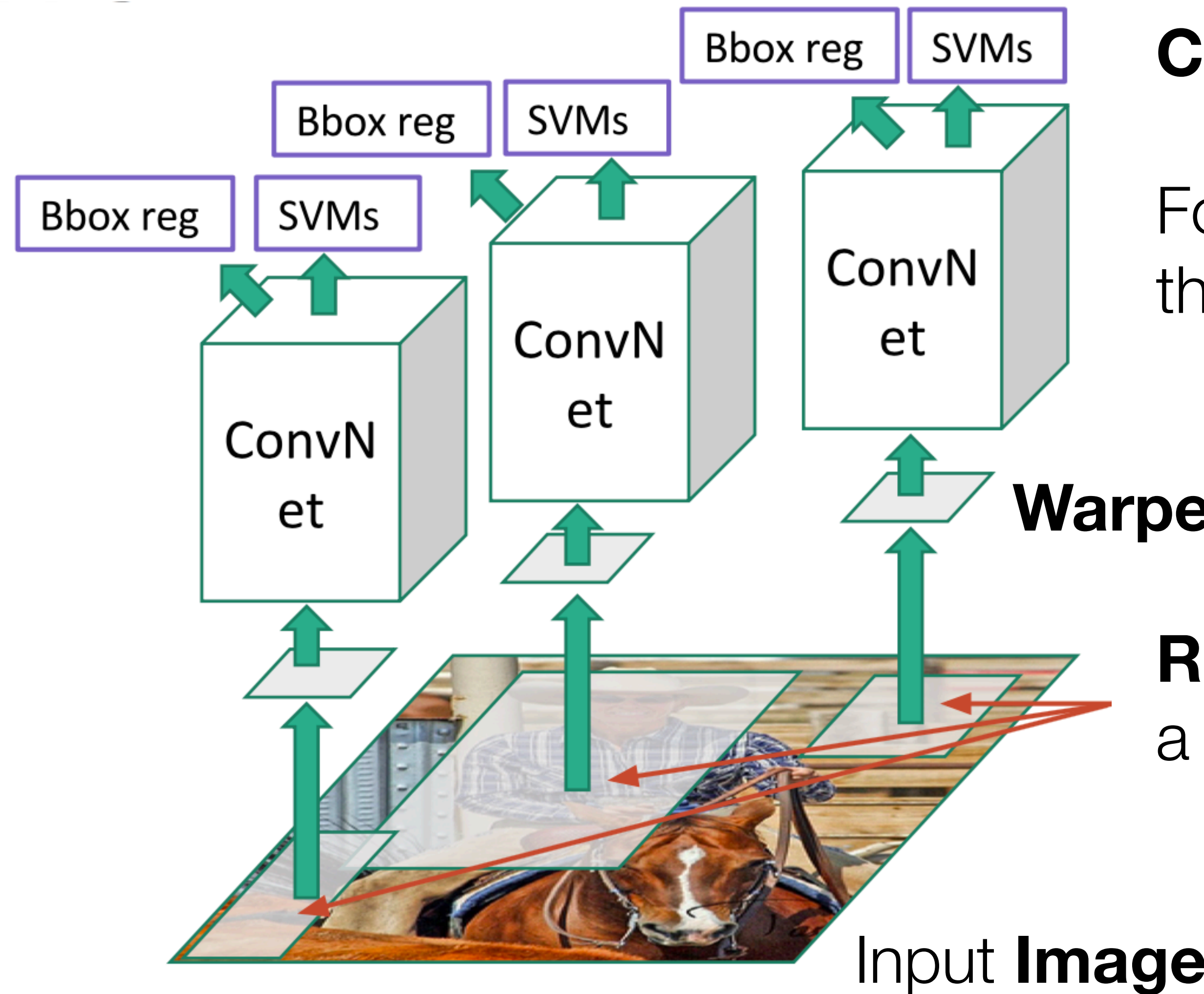
Input **Image**



# R-CNN

**Linear Regression** for bounding box offsets

[ Girshick et al, CVPR 2014 ]



**Classify** regions with SVM

Forward each region through a **CNN**

**Warped** image regions

**Regions of Interest** from a proposal method (~2k)



# R-CNN

R-CNN (Regions with CNN features) algorithm:

- Extract promising candidate regions using an object proposals algorithm
- Resize each proposal window to the size of the input layer of a trained convolutional neural network
- Input each resized image patch to the convolutional neural network

**Implementation detail:** Instead of using the classification scores of the network directly, the output of the final fully-connected layer can be used as an input feature to a trained support vector machine (SVM)

# Summary

The parameters of a neural network are learned using **backpropagation**, which computes gradients via recursive application of the chain rule

A **convolutional neural network** assumes inputs are images, and constrains the network architecture to reduce the number of parameters

A **convolutional layer** applies a set of learnable filters

A **pooling layer** performs spatial downsampling

A **fully-connected** layer is the same as in a regular neural network

Convolutional neural networks can be seen as learning a hierarchy of filters