

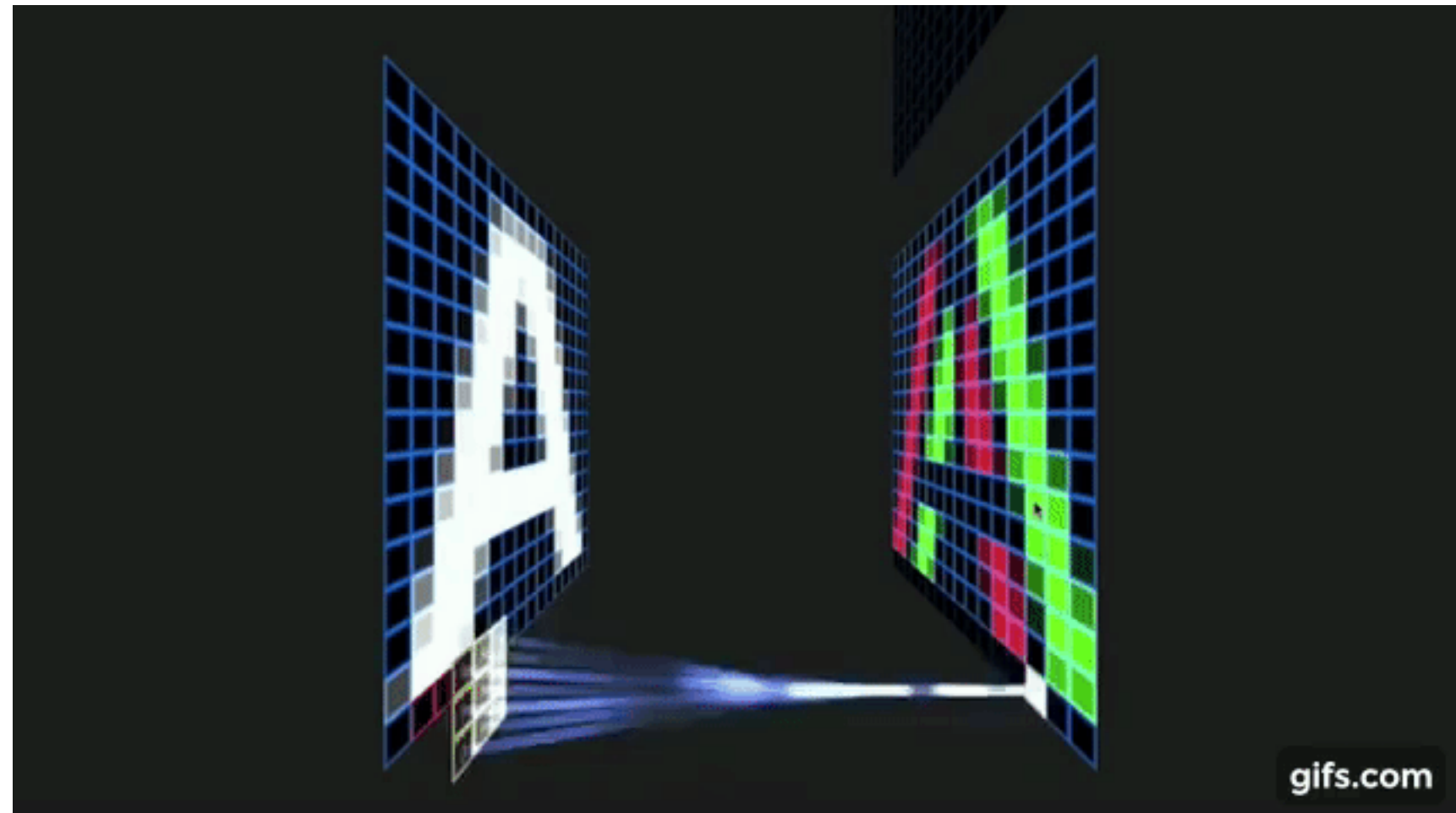
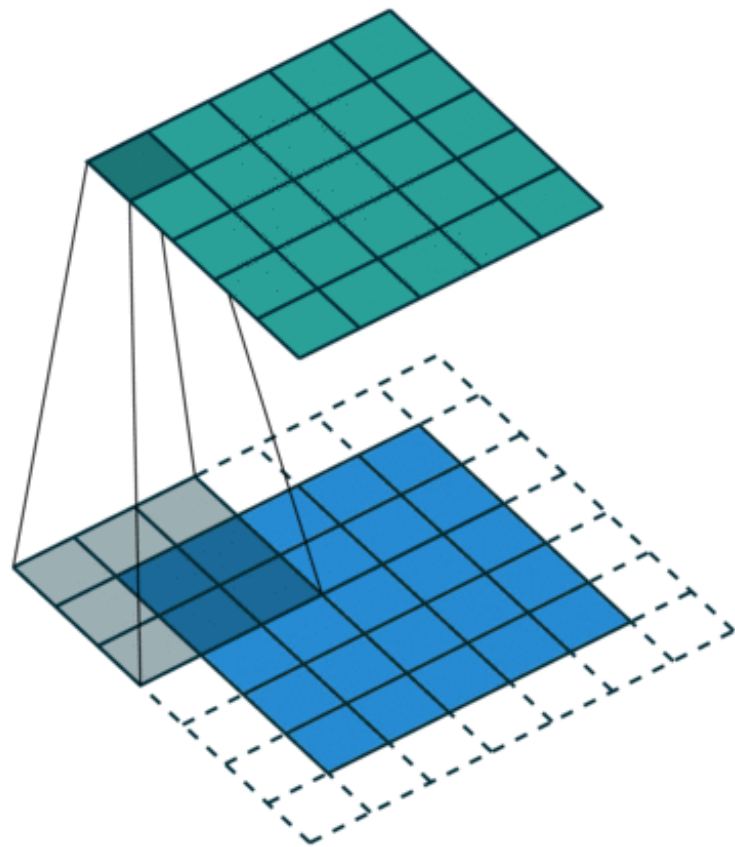
# ENM 360: Introduction to Data-driven Modeling

## *Lecture #17: Convolutional neural networks*

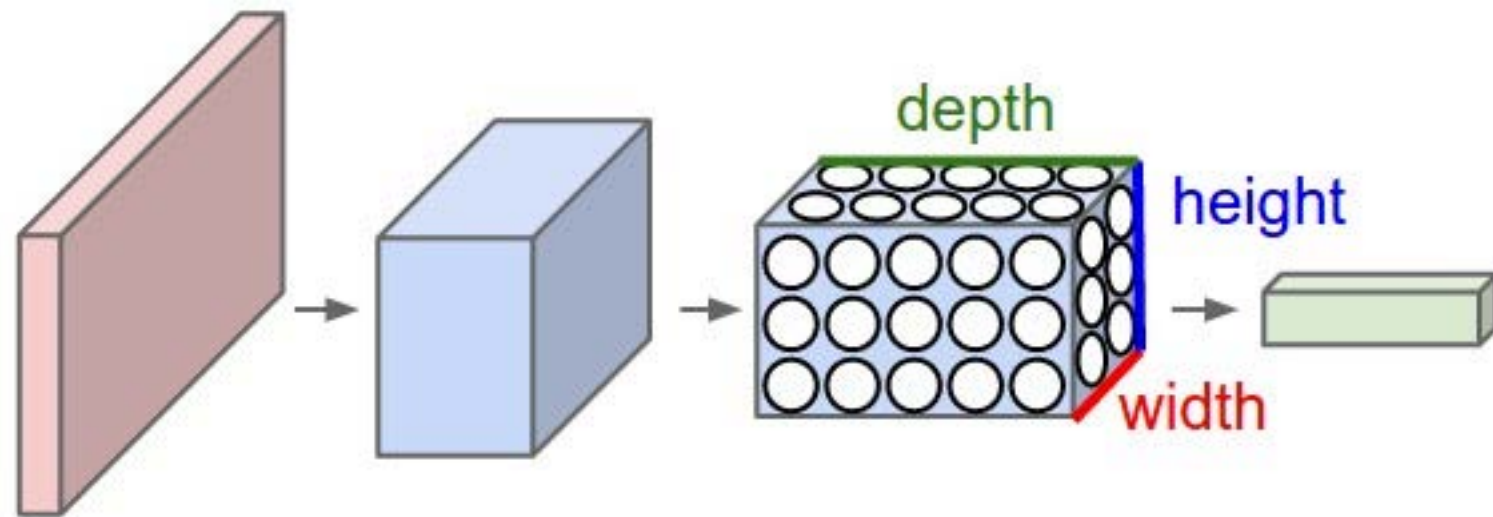
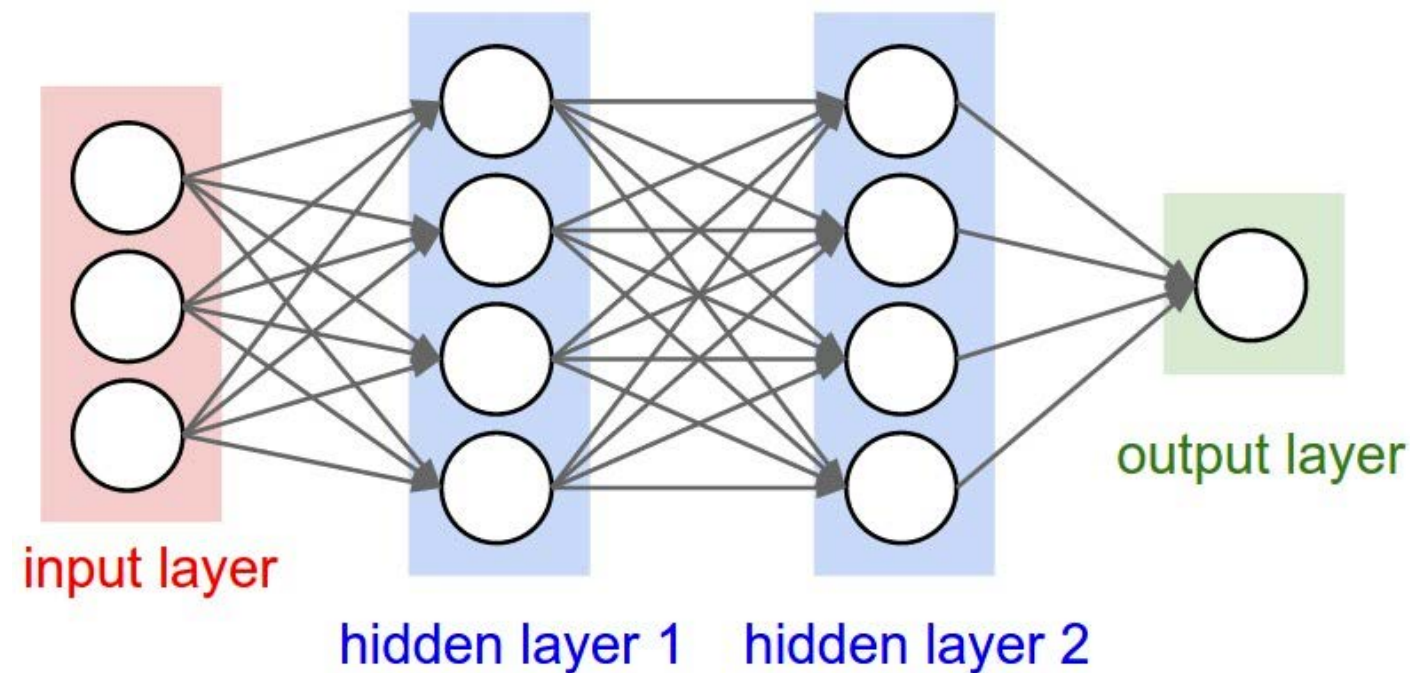
Paris Perdikaris  
October 27, 2019



# ConvNets for pattern recognition



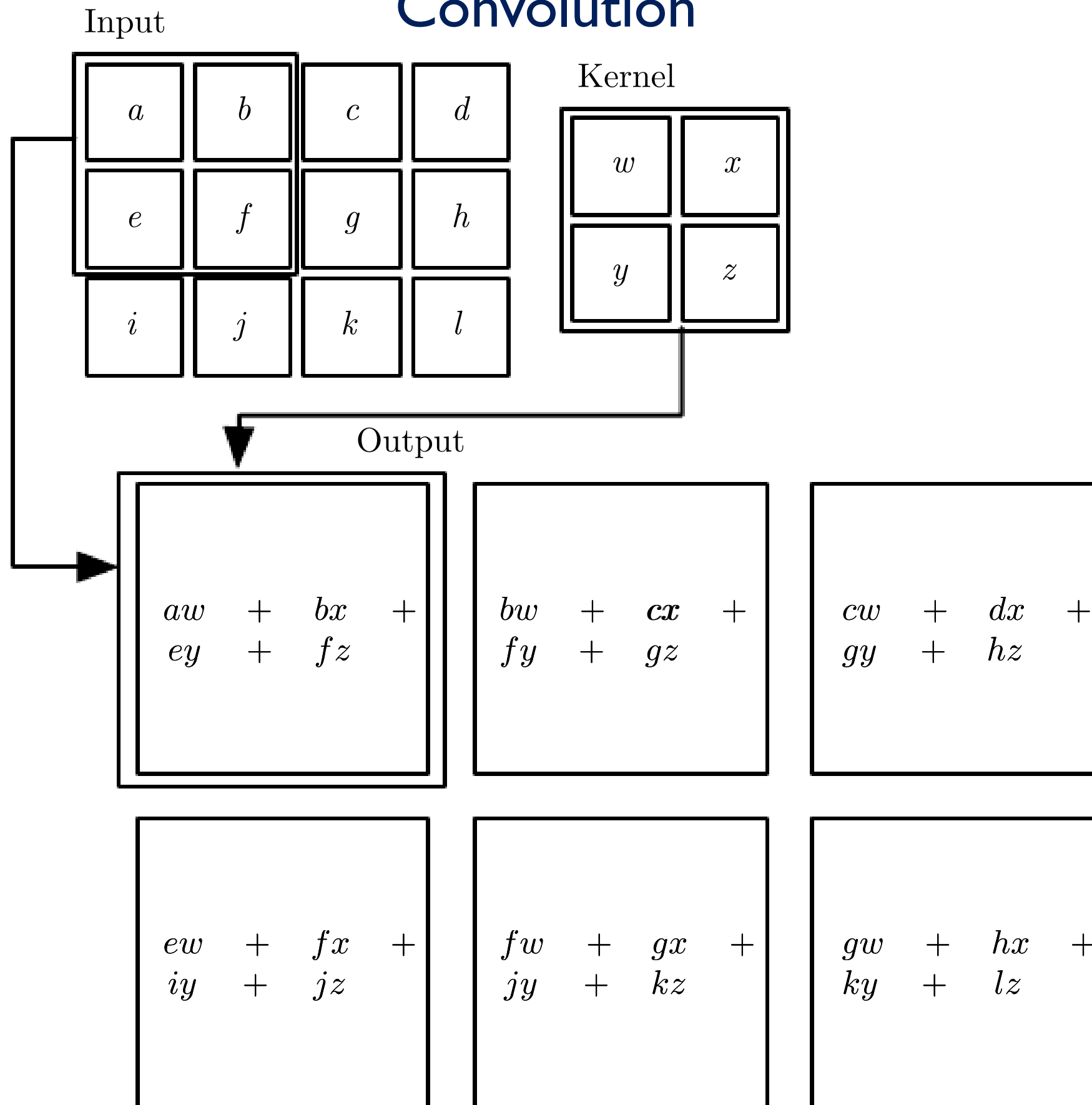
# From MLPs to ConvNets



**Top:** A regular 3-layer Neural Network.

**Bottom:** A ConvNet arranges its neurons in three dimensions (width, height, depth), as visualized in one of the layers. Every layer of a ConvNet transforms the 3D input volume to a 3D output volume of neuron activations. In this example, the red input layer holds the image, so its width and height would be the dimensions of the image, and the depth would be 3 (Red, Green, Blue channels).

# Convolution



# Convolution

Input Volume (+pad 1) (7x7x3)

$x[:, :, 0]$

0	0	0	0	0	0	0
0	0	0	1	0	2	0
0	1	0	2	0	1	0

0	1	0	2	2	0	0
0	2	0	0	2	0	0
0	2	1	2	2	0	0
0	0	0	0	0	0	0

$x[:, :, 1]$

0	0	0	0	0	0	0
0	2	1	2	1	1	0
0	2	1	2	0	1	0

0	0	2	1	0	1	0
0	1	2	2	2	2	0
0	0	1	2	0	1	0
0	0	0	0	0	0	0

$x[:, :, 2]$

0	0	0	0	0	0	0
0	2	1	1	2	0	0
0	1	0	0	1	0	0

0	0	1	0	0	0	0
0	1	0	2	1	0	0
0	2	2	1	1	1	0
0	0	0	0	0	0	0

Filter W0 (3x3x3)

$w0[:, :, 0]$

-1	0	1
0	0	1
1	-1	1

$w0[:, :, 1]$

-1	0	1
1	-1	1
0	1	0

$w0[:, :, 2]$

-1	1	1
1	1	0
0	-1	0

Bias b0 (1x1x1)

$b0[:, :, 0]$

1
---

Filter W1 (3x3x3)

$w1[:, :, 0]$

0	1	-1
0	-1	0
0	-1	1

$w1[:, :, 1]$

-1	0	0
1	-1	0
1	-1	0

$w1[:, :, 2]$

-1	1	-1
0	-1	-1
1	0	0

Bias b1 (1x1x1)

$b1[:, :, 0]$

0
---

Output Volume (3x3x2)

$o[:, :, 0]$

2	3	3
3	7	3
8	10	-3

$o[:, :, 1]$

-8	-8	-3
-3	1	0
-3	-8	-5

toggle movement



## Sparse connectivity

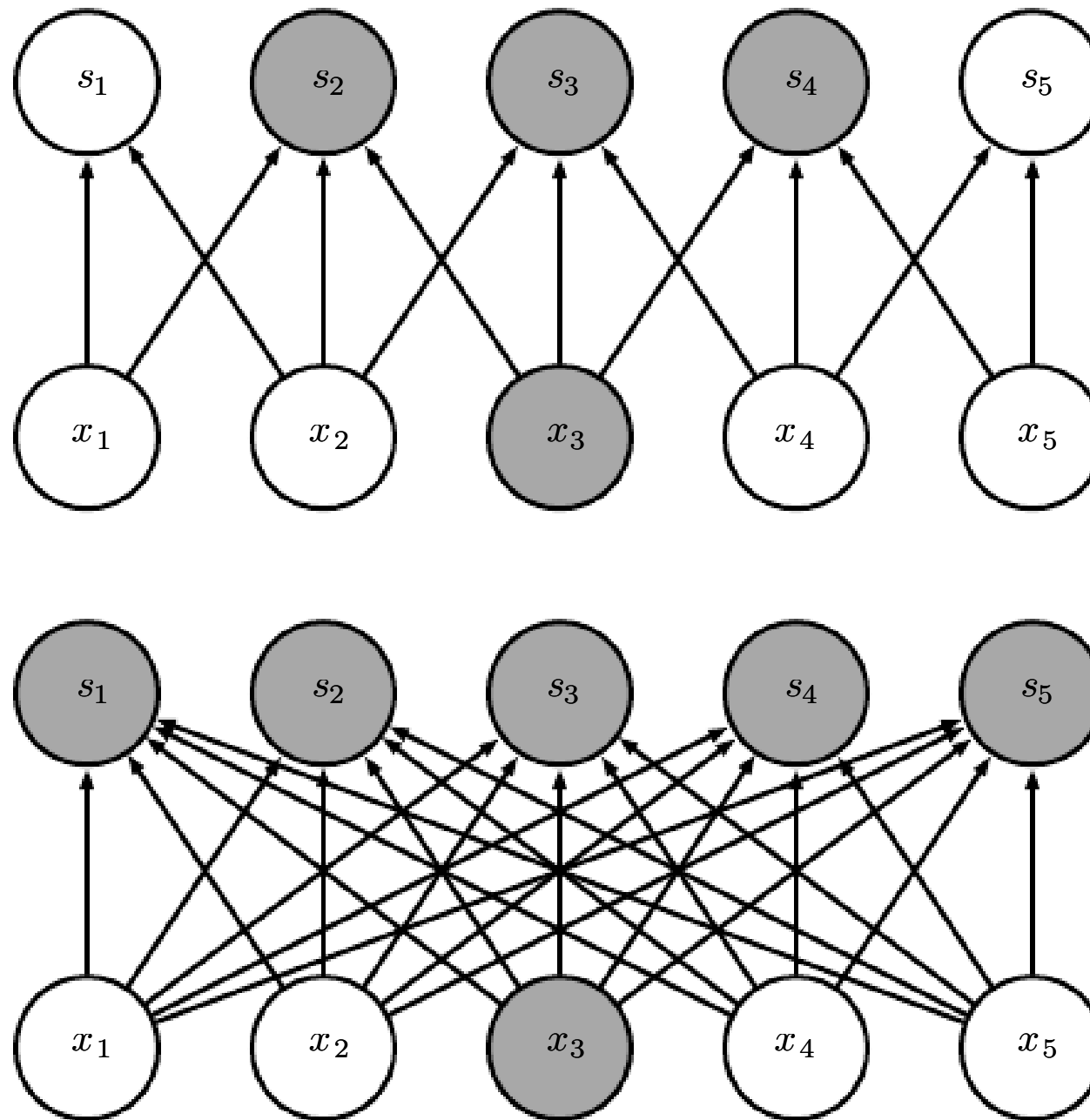


Figure 9.2: *Sparse connectivity, viewed from below:* We highlight one input unit,  $x_3$ , and also highlight the output units in  $\mathbf{s}$  that are affected by this unit. (*Top*) When  $\mathbf{s}$  is formed by convolution with a kernel of width 3, only three outputs are affected by  $\mathbf{x}$ . (*Bottom*) When  $\mathbf{s}$  is formed by matrix multiplication, connectivity is no longer sparse, so all of the outputs are affected by  $x_3$ .

## Sparse connectivity

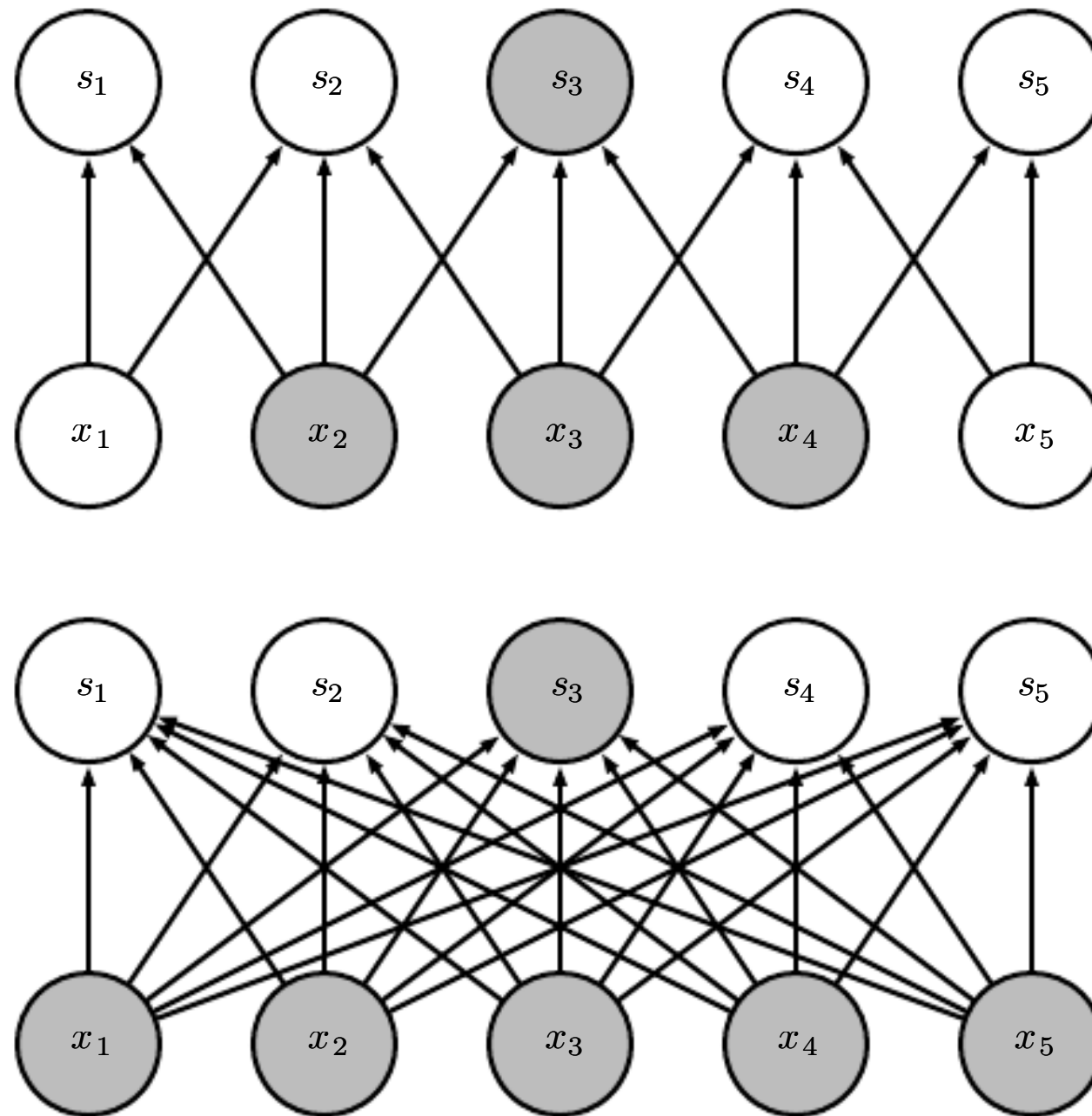


Figure 9.3: *Sparse connectivity, viewed from above*: We highlight one output unit,  $s_3$ , and also highlight the input units in  $\mathbf{x}$  that affect this unit. These units are known as the **receptive field** of  $s_3$ . (*Top*) When  $\mathbf{s}$  is formed by convolution with a kernel of width 3, only three inputs affect  $s_3$ . (*Bottom*) When  $\mathbf{s}$  is formed by matrix multiplication, connectivity is no longer sparse, so all of the inputs affect  $s_3$ .

# Sparse connectivity

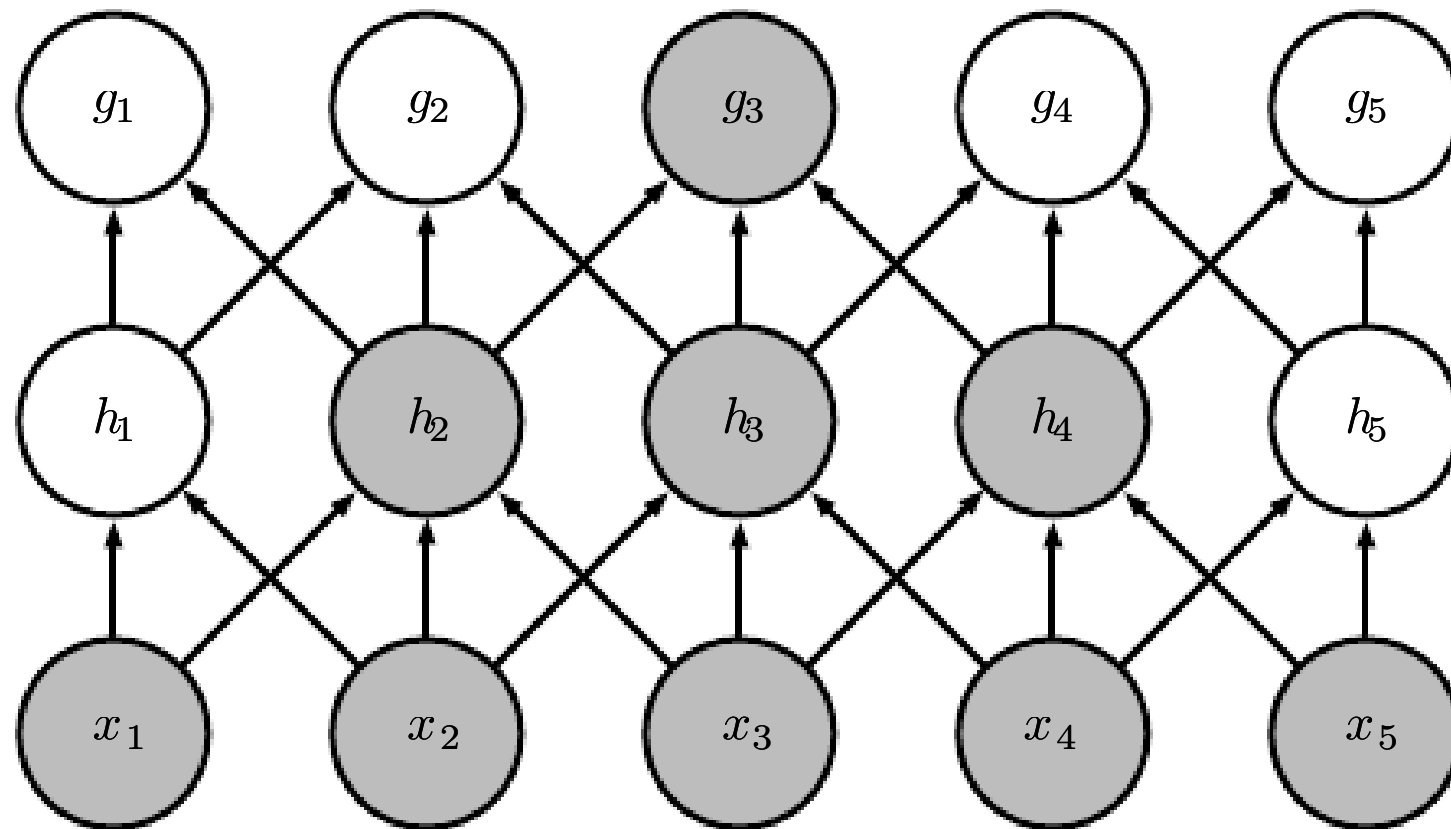


Figure 9.4: The receptive field of the units in the deeper layers of a convolutional network is larger than the receptive field of the units in the shallow layers. This effect increases if the network includes architectural features like strided convolution (figure 9.12) or pooling (section 9.3). This means that even though *direct* connections in a convolutional net are very sparse, units in the deeper layers can be *indirectly* connected to all or most of the input image.



# Parameter sharing

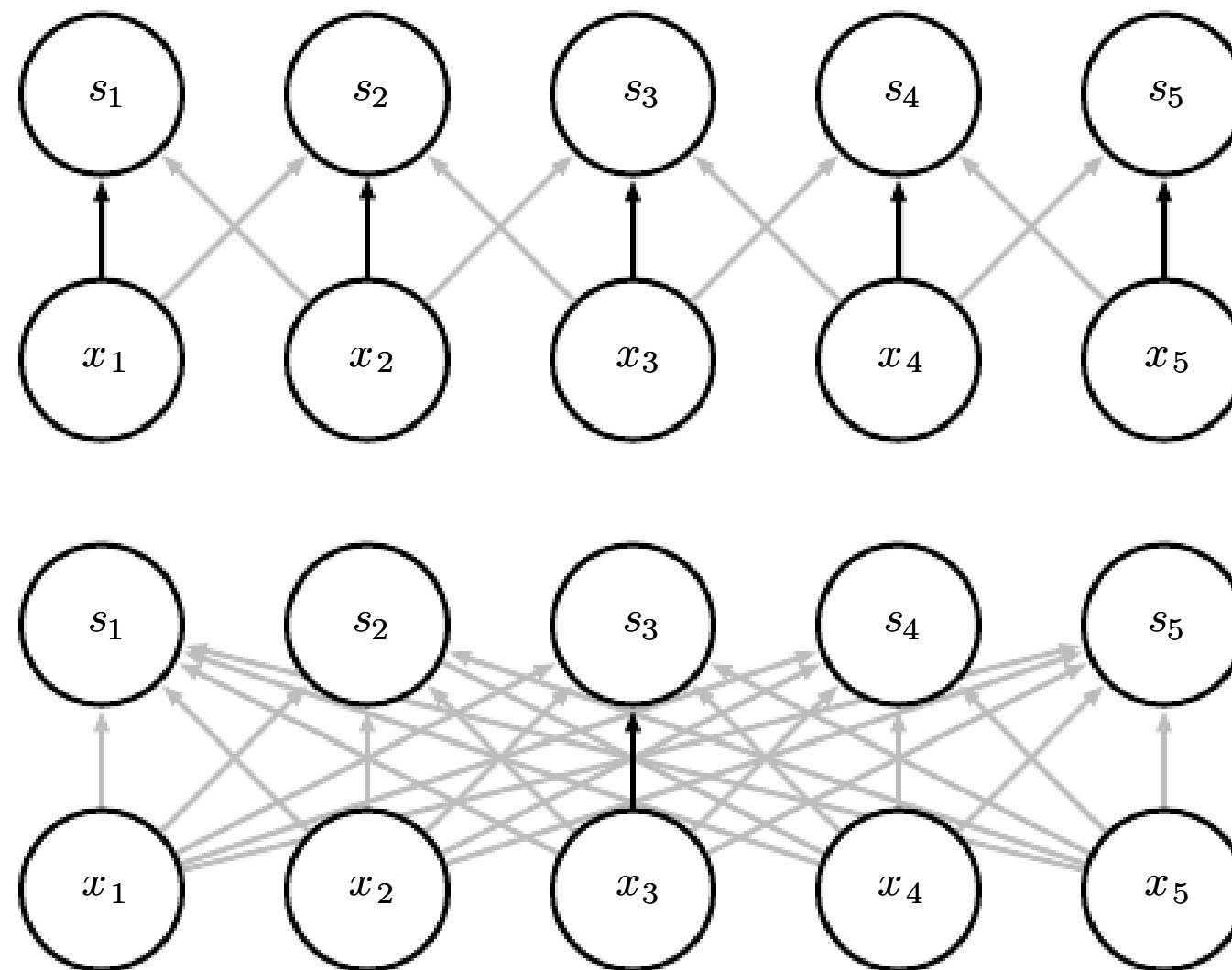
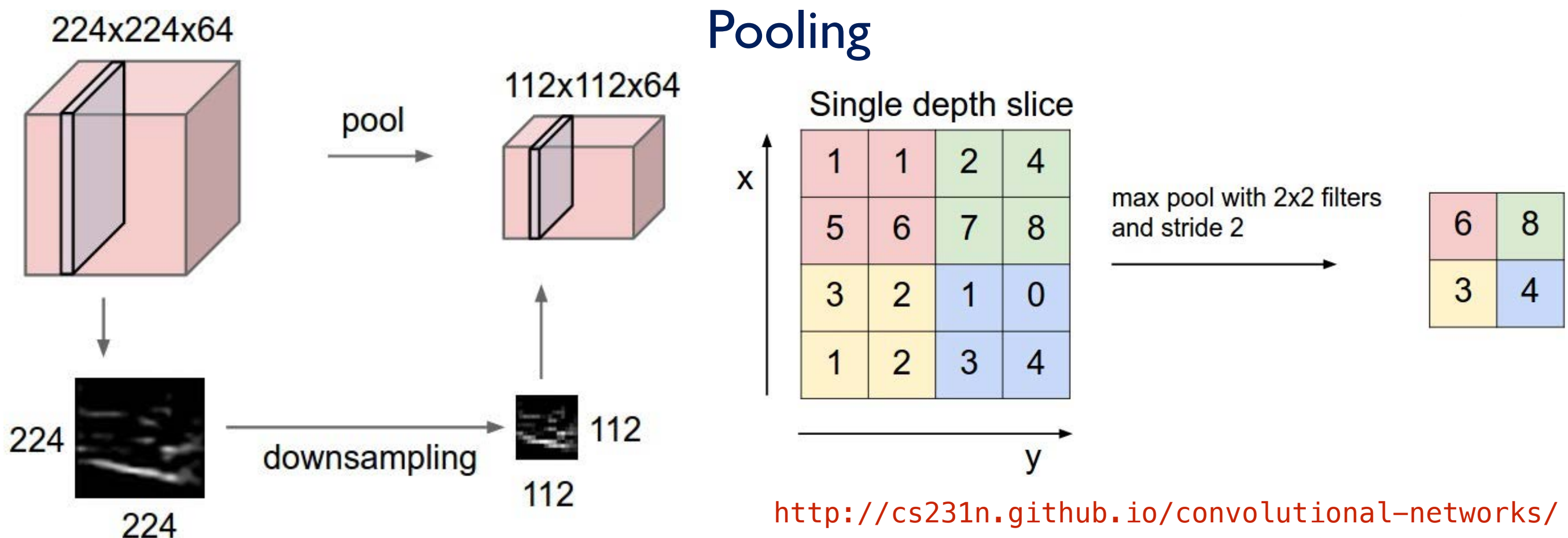


Figure 9.5: Parameter sharing: Black arrows indicate the connections that use a particular parameter in two different models. *(Top)* The black arrows indicate uses of the central element of a 3-element kernel in a convolutional model. Due to parameter sharing, this single parameter is used at all input locations. *(Bottom)* The single black arrow indicates the use of the central element of the weight matrix in a fully connected model. This model has no parameter sharing so the parameter is used only once.

# Equivariance to translation

In the case of convolution, the particular form of parameter sharing causes the layer to have a property called **equivariance** to translation. To say a function is equivariant means that if the input changes, the output changes in the same way. Specifically, a function  $f(x)$  is equivariant to a function  $g$  if  $f(g(x)) = g(f(x))$ . In the case of convolution, if we let  $g$  be any function that translates the input, i.e., shifts it, then the convolution function is equivariant to  $g$ . For example, let  $I$  be a function giving image brightness at integer coordinates. Let  $g$  be a function mapping one image function to another image function, such that  $I' = g(I)$  is the image function with  $I'(x, y) = I(x - 1, y)$ . This shifts every pixel of  $I$  one unit to the right. If we apply this transformation to  $I$ , then apply convolution, the result will be the same as if we applied convolution to  $I'$ , then applied the transformation  $g$  to the output.

Convolution is not naturally equivariant to some other transformations, such as changes in the scale or rotation of an image. Other mechanisms are necessary for handling these kinds of transformations.



Pooling layer downsamples the volume spatially, independently in each depth slice of the input volume. Left: In this example, the input volume of size  $[224 \times 224 \times 64]$  is pooled with filter size 2, stride 2 into output volume of size  $[112 \times 112 \times 64]$ . Notice that the volume depth is preserved. Right: The most common downsampling operation is max, giving rise to max pooling, here shown with a stride of 2. That is, each max is taken over 4 numbers (little  $2 \times 2$  square).

In all cases, pooling helps to make the representation become approximately **invariant** to small translations of the input. Invariance to translation means that if we translate the input by a small amount, the values of most of the pooled outputs do not change. See figure 9.8 for an example of how this works. *Invariance to local translation can be a very useful property if we care more about whether some feature is present than exactly where it is.*

# Pooling

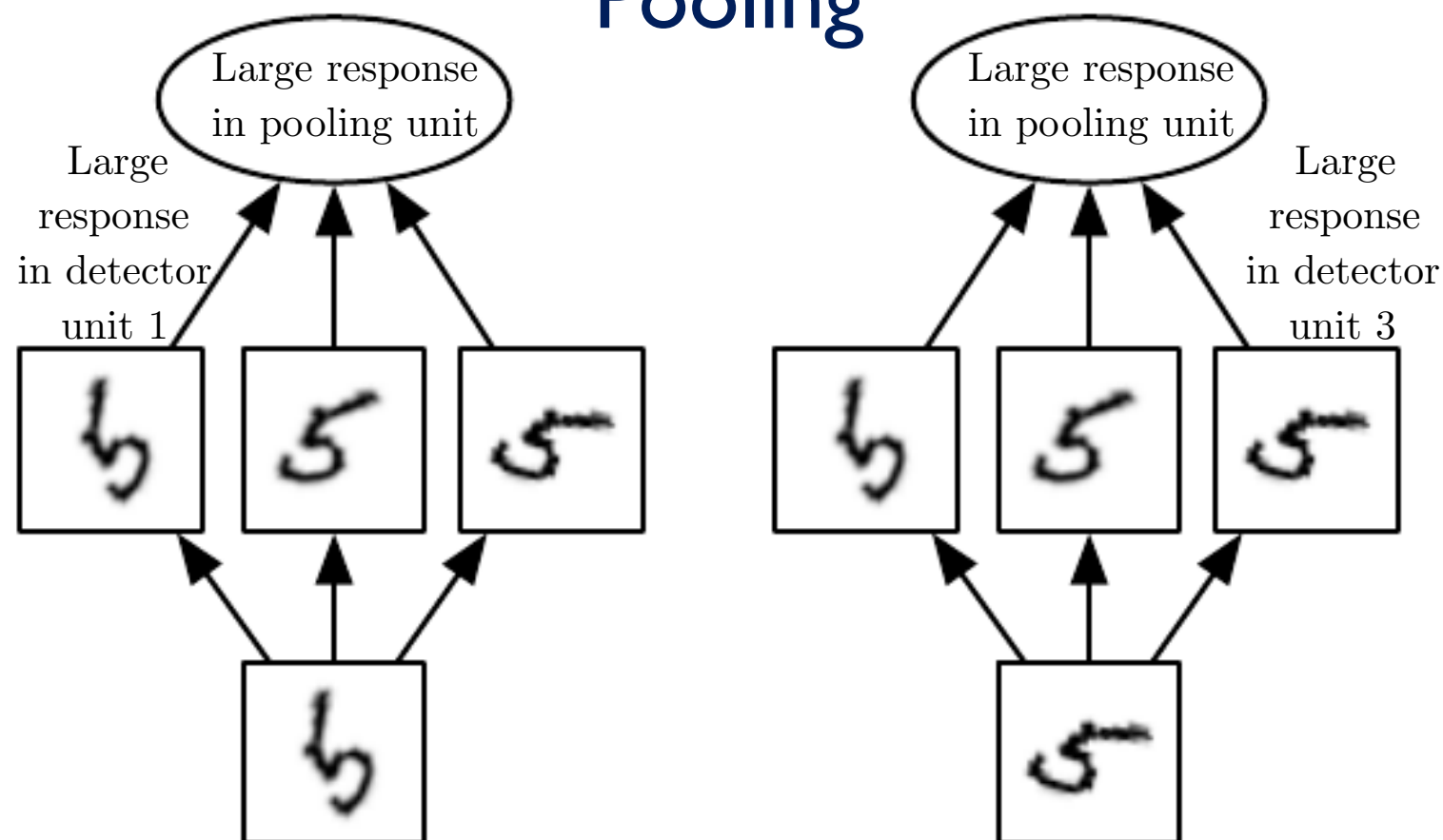


Figure 9.9: *Example of learned invariances:* A pooling unit that pools over multiple features that are learned with separate parameters can learn to be invariant to transformations of the input. Here we show how a set of three learned filters and a max pooling unit can learn to become invariant to rotation. All three filters are intended to detect a hand-written 5. Each filter attempts to match a slightly different orientation of the 5. When a 5 appears in the input, the corresponding filter will match it and cause a large activation in a detector unit. The max pooling unit then has a large activation regardless of which detector unit was activated. We show here how the network processes two different inputs, resulting in two different detector units being activated. The effect on the pooling unit is roughly the same either way. This principle is leveraged by maxout networks ([Goodfellow et al., 2013a](#)) and other convolutional networks. Max pooling over spatial positions is naturally invariant to translation; this multi-channel approach is only necessary for learning other transformations.

# Pooling

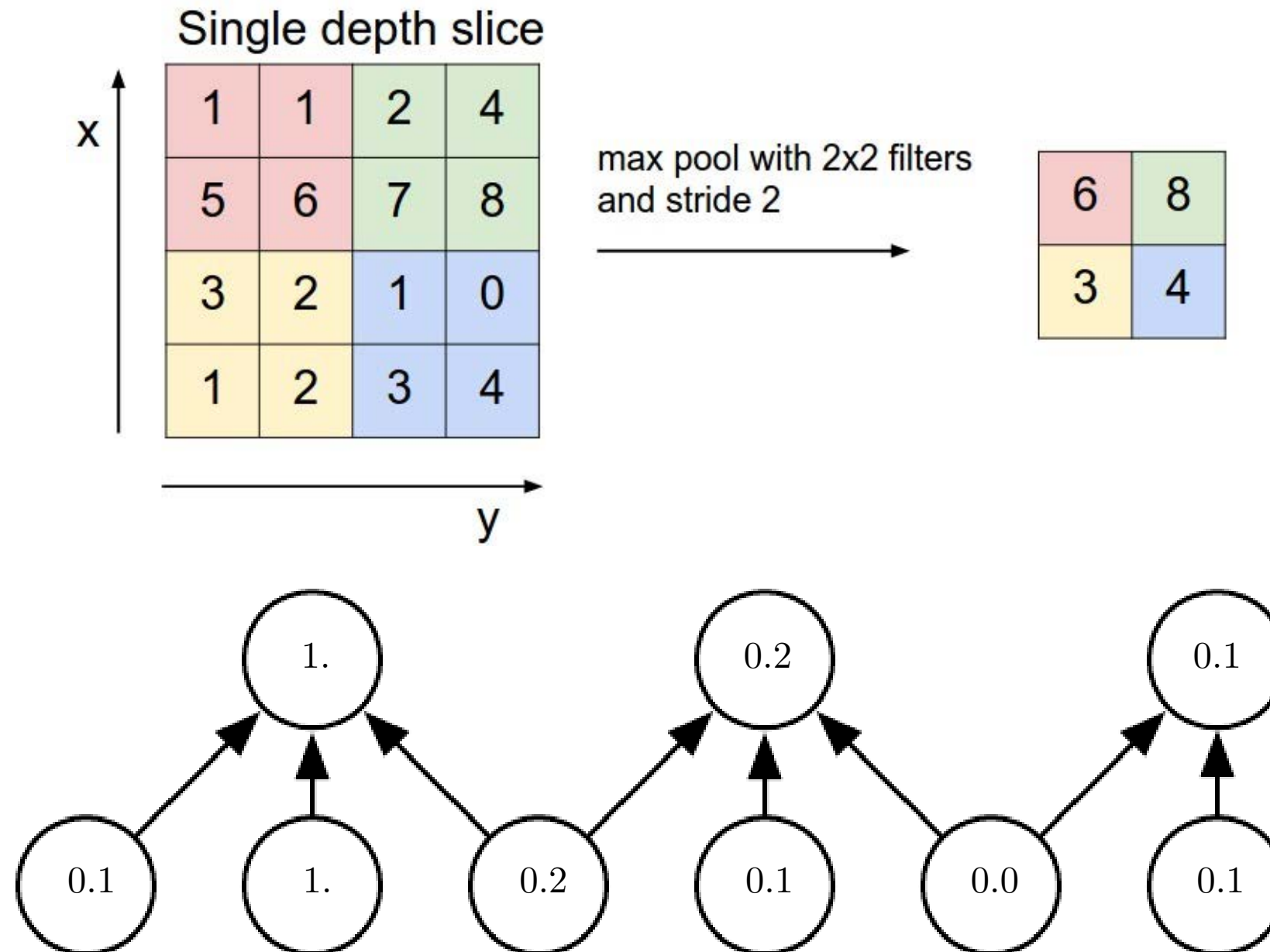
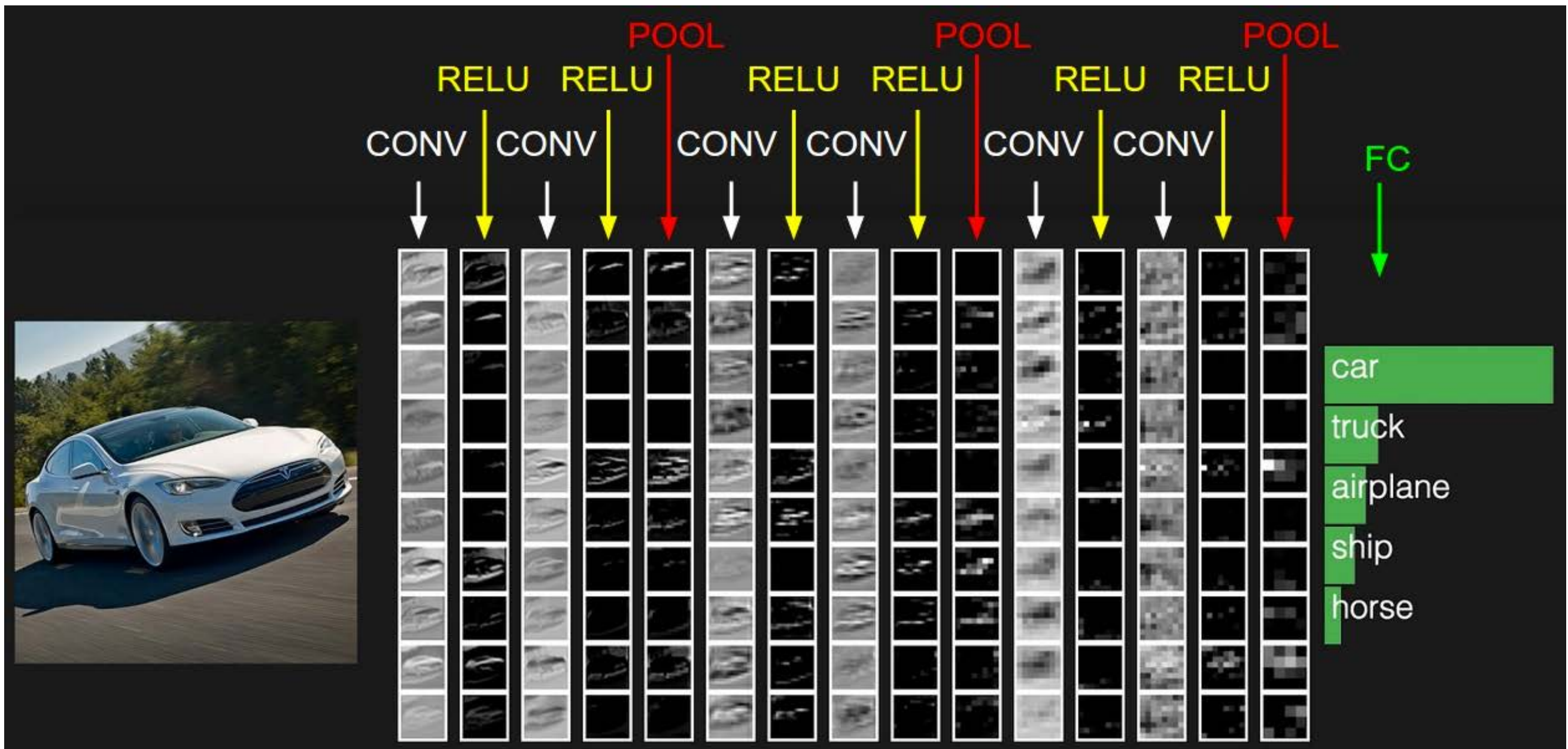


Figure 9.10: *Pooling with downsampling*. Here we use max-pooling with a pool width of three and a stride between pools of two. This reduces the representation size by a factor of two, which reduces the computational and statistical burden on the next layer. Note that the rightmost pooling region has a smaller size, but must be included if we do not want to ignore some of the detector units.



# Image classification using ConvNets



The activations of an example ConvNet architecture. The initial volume stores the raw image pixels (left) and the last volume stores the class scores (right). Each volume of activations along the processing path is shown as a column. Since it's difficult to visualize 3D volumes, we lay out each volume's slices in rows. The last layer volume holds the scores for each class, but here we only visualize the sorted top 5 scores, and print the labels of each one.



# Layers used to build ConvNets

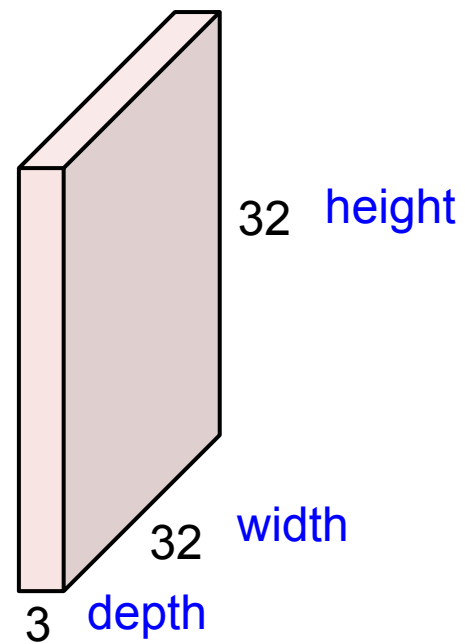
As we described above, a simple ConvNet is a sequence of layers, and every layer of a ConvNet transforms one volume of activations to another through a differentiable function. We use three main types of layers to build ConvNet architectures: **Convolutional Layer**, **Pooling Layer**, and **Fully-Connected Layer** (exactly as seen in regular Neural Networks). We will stack these layers to form a full ConvNet **architecture**.

*Example Architecture: Overview.* We will go into more details below, but a simple ConvNet for CIFAR-10 classification could have the architecture [INPUT - CONV - RELU - POOL - FC]. In more detail:

- INPUT [32x32x3] will hold the raw pixel values of the image, in this case an image of width 32, height 32, and with three color channels R,G,B.
- CONV layer will compute the output of neurons that are connected to local regions in the input, each computing a dot product between their weights and a small region they are connected to in the input volume. This may result in volume such as [32x32x12] if we decided to use 12 filters.
- RELU layer will apply an elementwise activation function, such as the  $\max(0, x)$  thresholding at zero. This leaves the size of the volume unchanged ([32x32x12]).
- POOL layer will perform a downsampling operation along the spatial dimensions (width, height), resulting in volume such as [16x16x12].
- FC (i.e. fully-connected) layer will compute the class scores, resulting in volume of size [1x1x10], where each of the 10 numbers correspond to a class score, such as among the 10 categories of CIFAR-10. As with ordinary Neural Networks and as the name implies, each neuron in this layer will be connected to all the numbers in the previous volume.

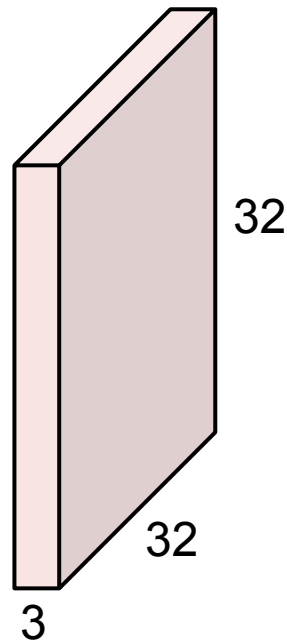
# The input layer

32x32x3 image



# The convolution layer

32x32x3 image

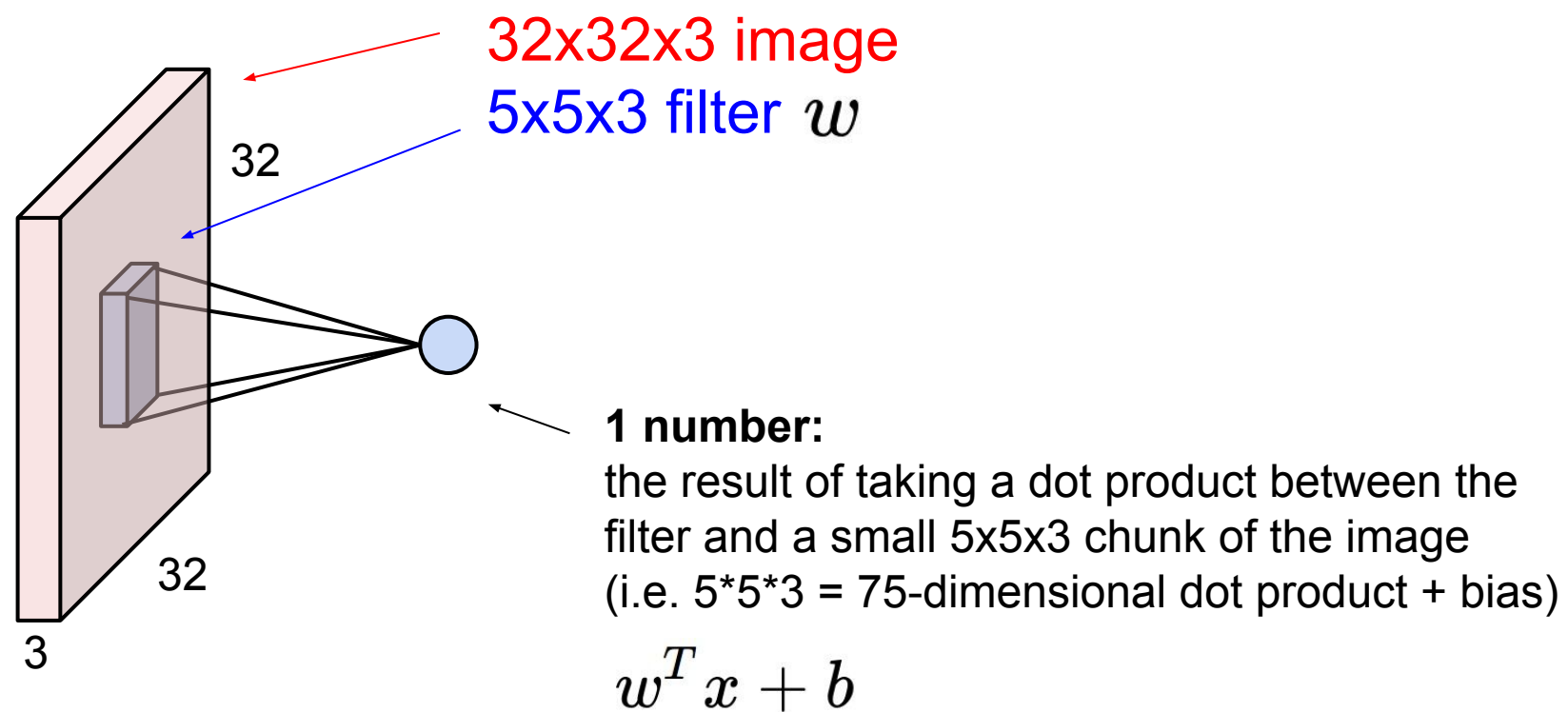


5x5x3 filter

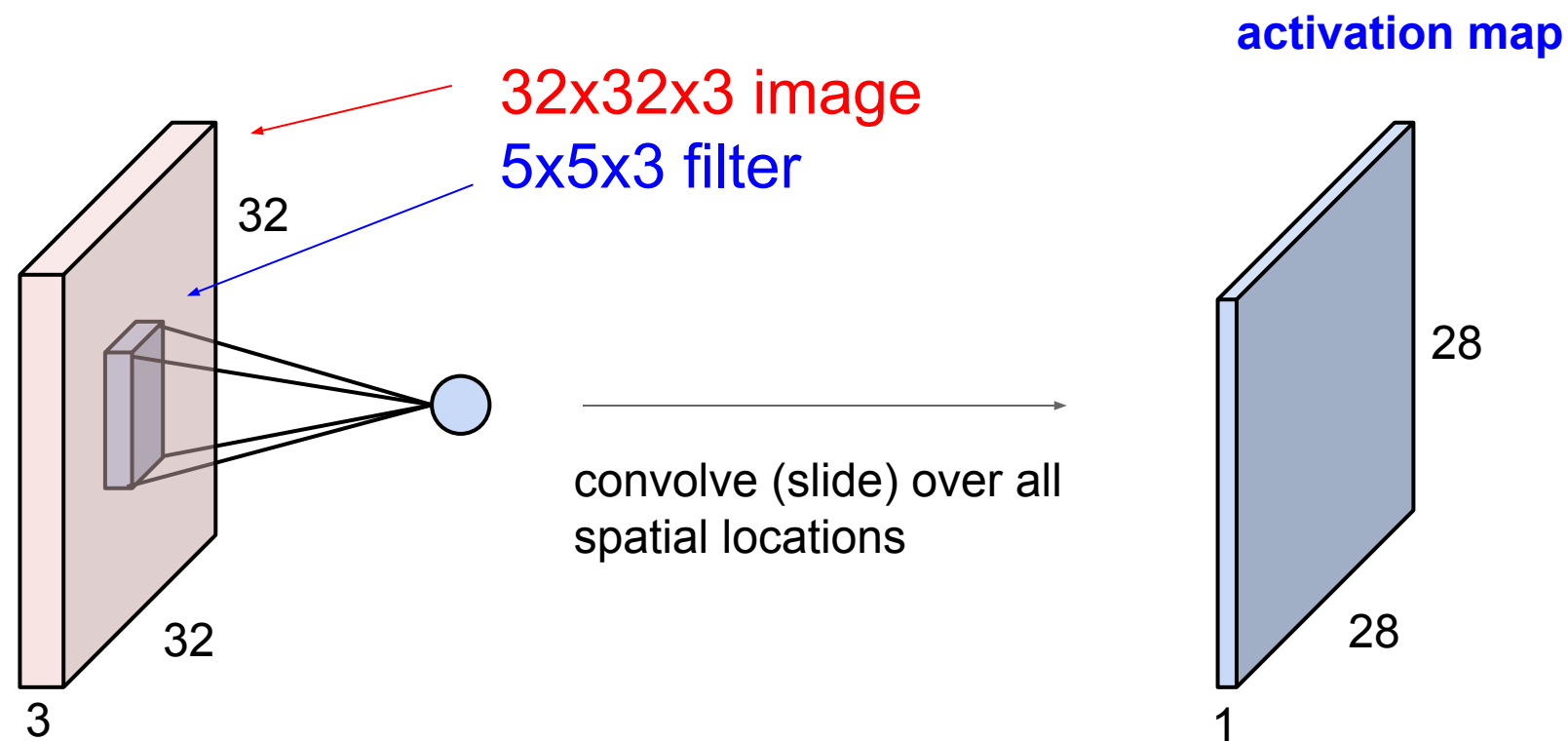


**Convolve** the filter with the image  
i.e. “slide over the image spatially,  
computing dot products”

# The convolution layer



# The convolution layer





# The convolution layer

In practice: Common to zero pad the border

0	0	0	0	0	0			
0								
0								
0								
0								

e.g. input 7x7

**3x3** filter, applied with **stride 1**

**pad with 1 pixel** border => what is the output?

**7x7 output!**

in general, common to see CONV layers with stride 1, filters of size  $F \times F$ , and zero-padding with  $(F-1)/2$ . (will preserve size spatially)

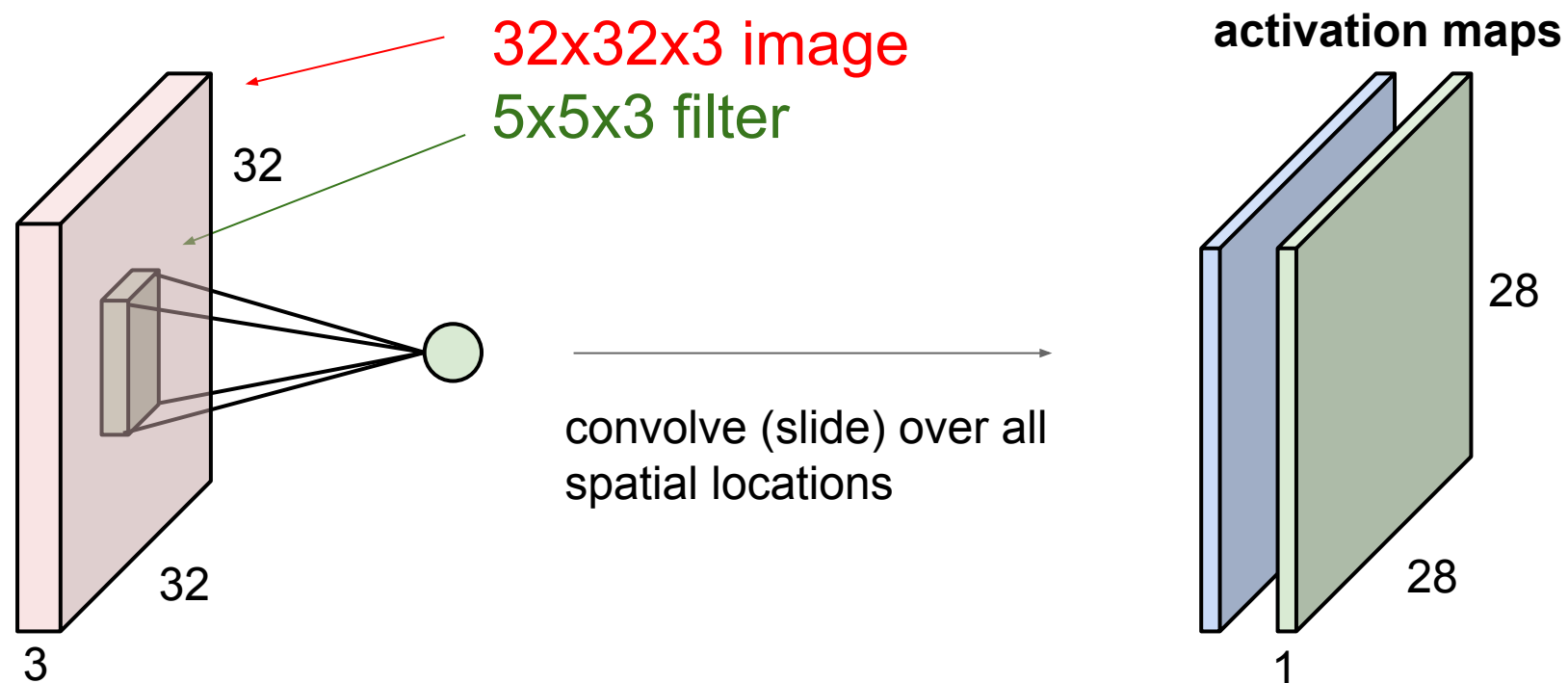
e.g.  $F = 3 \Rightarrow$  zero pad with 1

$F = 5 \Rightarrow$  zero pad with 2

$F = 7 \Rightarrow$  zero pad with 3

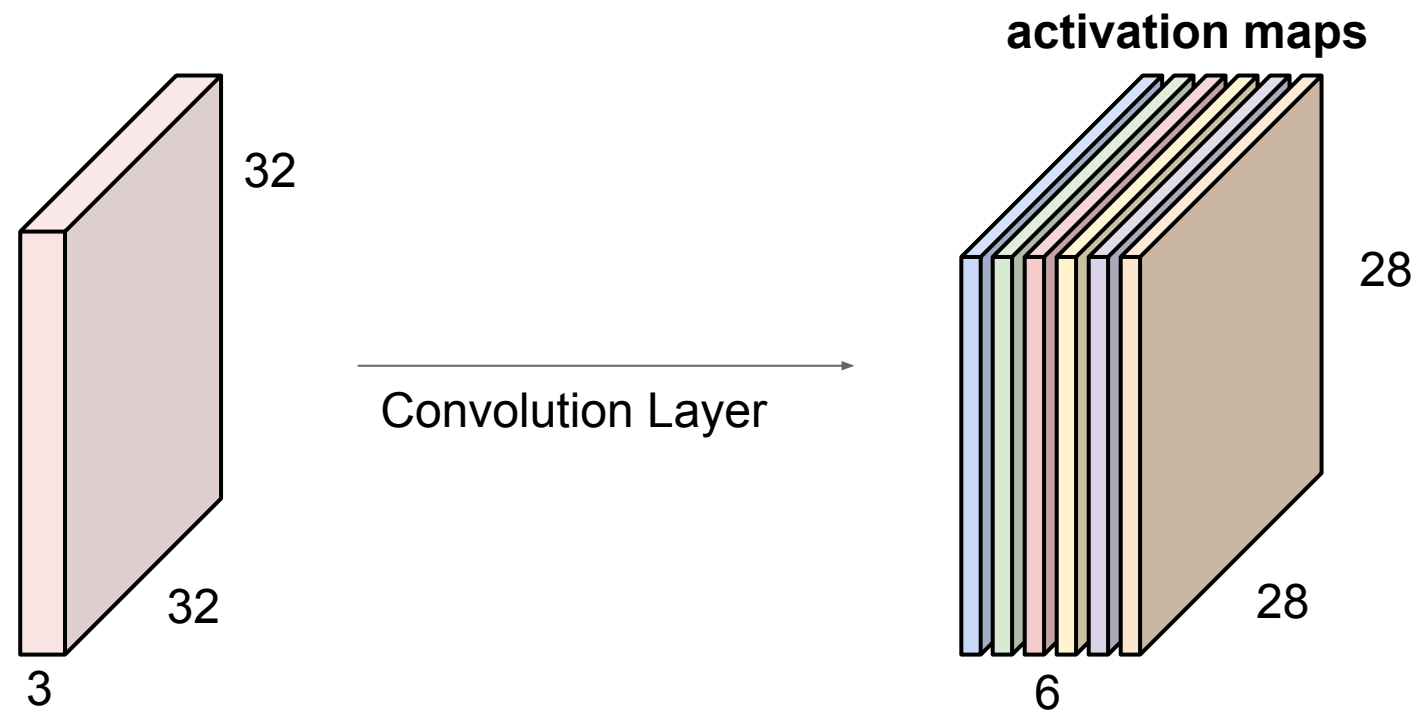


# The convolution layer



# The convolution layer

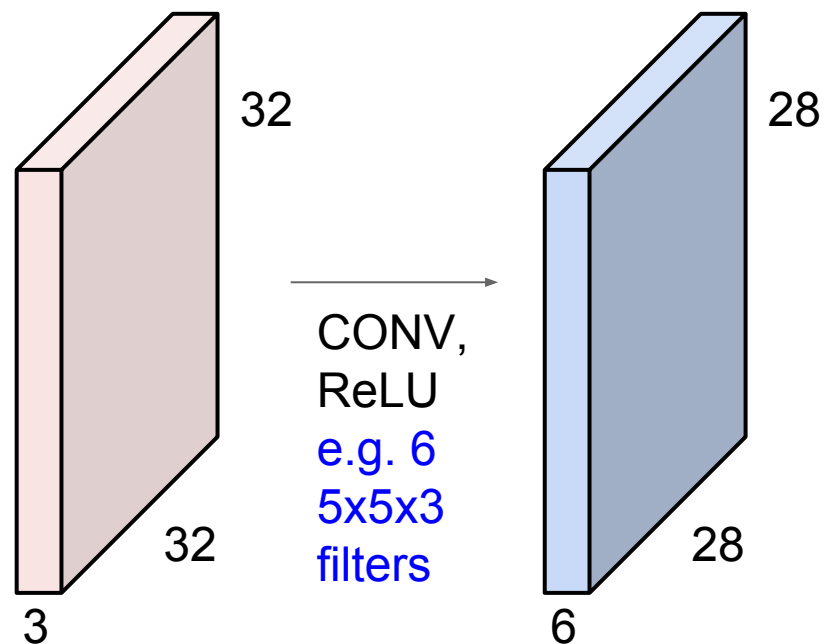
For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:



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

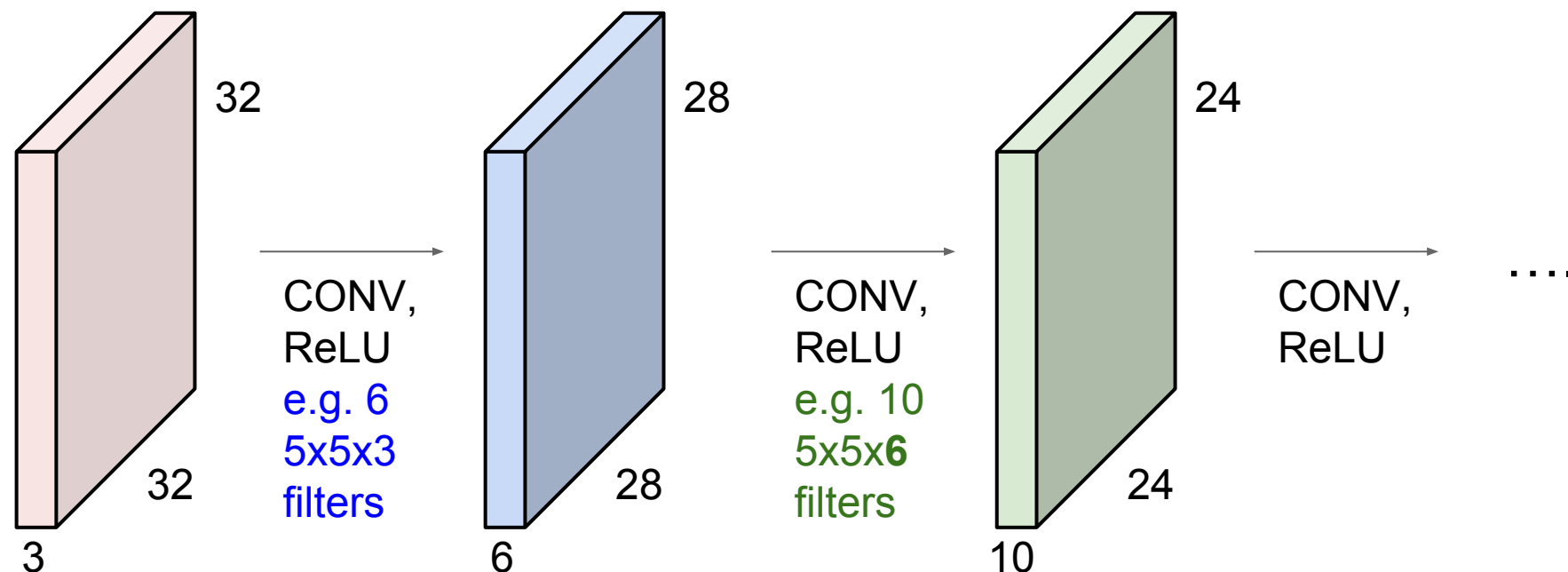
# The convolution layer

**Preview:** ConvNet is a sequence of Convolution Layers, interspersed with activation functions



# The convolution layer

**Preview:** ConvNet is a sequence of Convolutional Layers, interspersed with activation functions



# The convolution layer

**Summary.** To summarize, the Conv Layer:

- Accepts a volume of size  $W_1 \times H_1 \times D_1$
- Requires four hyperparameters:
  - Number of filters  $K$ ,
  - their spatial extent  $F$ ,
  - the stride  $S$ ,
  - the amount of zero padding  $P$ .
- Produces a volume of size  $W_2 \times H_2 \times D_2$  where:
  - $W_2 = (W_1 - F + 2P)/S + 1$
  - $H_2 = (H_1 - F + 2P)/S + 1$  (i.e. width and height are computed equally by symmetry)
  - $D_2 = K$
- With parameter sharing, it introduces  $F \cdot F \cdot D_1$  weights per filter, for a total of  $(F \cdot F \cdot D_1) \cdot K$  weights and  $K$  biases.
- In the output volume, the  $d$ -th depth slice (of size  $W_2 \times H_2$ ) is the result of performing a valid convolution of the  $d$ -th filter over the input volume with a stride of  $S$ , and then offset by  $d$ -th bias.

Common settings:

- $K =$  (powers of 2, e.g. 32, 64, 128, 512)
- $F = 3, S = 1, P = 1$
  - $F = 5, S = 1, P = 2$
  - $F = 5, S = 2, P = ?$  (whatever fits)
  - $F = 1, S = 1, P = 0$

A common setting of the hyperparameters is  $F = 3, S = 1, P = 1$ . However, there are common conventions and rules of thumb that motivate these hyperparameters. See the [ConvNet architectures](http://cs231n.github.io/convolutional-networks/) section below.



# The pooling layer

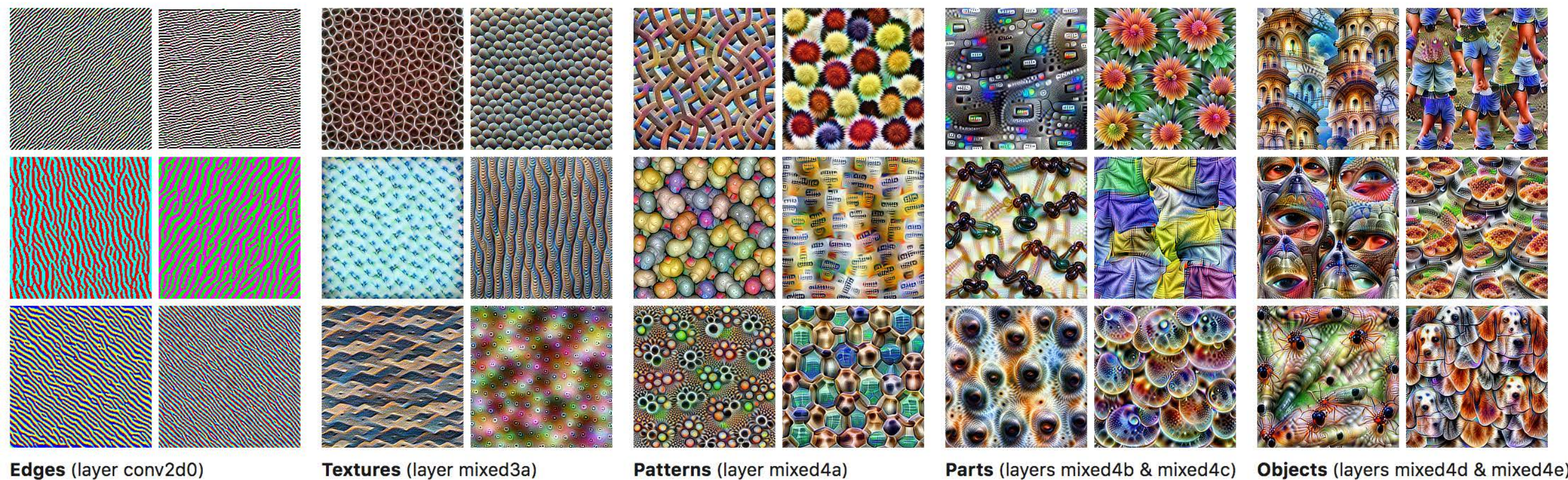
It is common to periodically insert a Pooling layer in-between successive Conv layers in a ConvNet architecture. Its function is to progressively reduce the spatial size of the representation to reduce the amount of parameters and computation in the network, and hence to also control overfitting. The Pooling Layer operates independently on every depth slice of the input and resizes it spatially, using the MAX operation. The most common form is a pooling layer with filters of size  $2 \times 2$  applied with a stride of 2 downsamples every depth slice in the input by 2 along both width and height, discarding 75% of the activations. Every MAX operation would in this case be taking a max over 4 numbers (little  $2 \times 2$  region in some depth slice). The depth dimension remains unchanged. More generally, the pooling layer:

- Accepts a volume of size  $W_1 \times H_1 \times D_1$
- Requires two hyperparameters:
  - their spatial extent  $F$ ,
  - the stride  $S$ ,
- Produces a volume of size  $W_2 \times H_2 \times D_2$  where:
  - $W_2 = (W_1 - F)/S + 1$
  - $H_2 = (H_1 - F)/S + 1$
  - $D_2 = D_1$
- Introduces zero parameters since it computes a fixed function of the input
- Note that it is not common to use zero-padding for Pooling layers

It is worth noting that there are only two commonly seen variations of the max pooling layer found in practice: A pooling layer with  $F = 3, S = 2$  (also called overlapping pooling), and more commonly  $F = 2, S = 2$ . Pooling sizes with larger receptive fields are too destructive. <http://cs231n.github.io/convolutional-networks/>



# Feature visualization



Feature visualization allows us to see how GoogLeNet [1], trained on the ImageNet [2] dataset, builds up its understanding of images over many layers. Visualizations of all channel are available in the [appendix](#).



# Evolution of architectures

LeNet architecture representation from the original 1998 LeCun et al. paper.

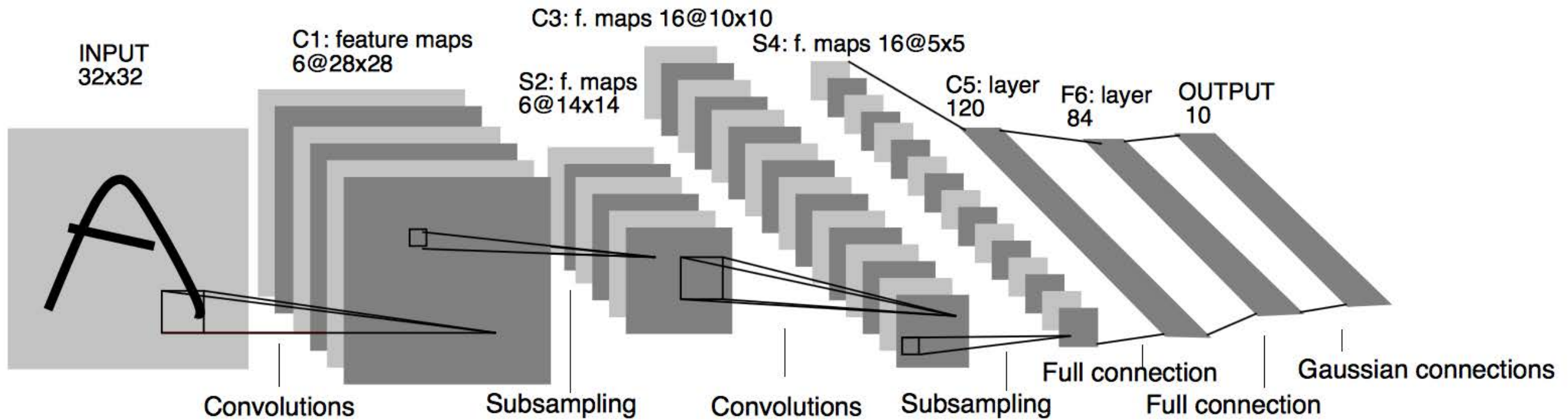
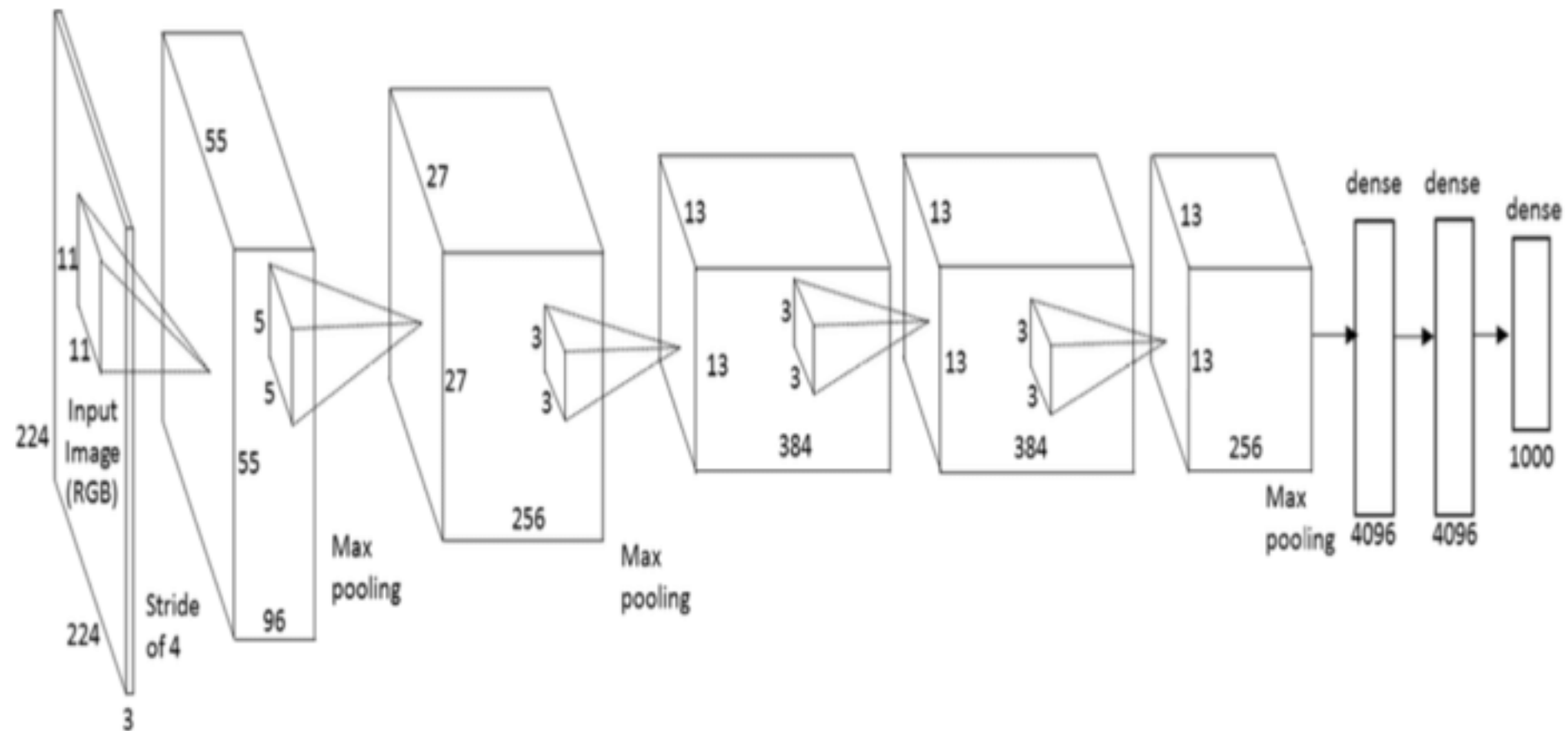


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

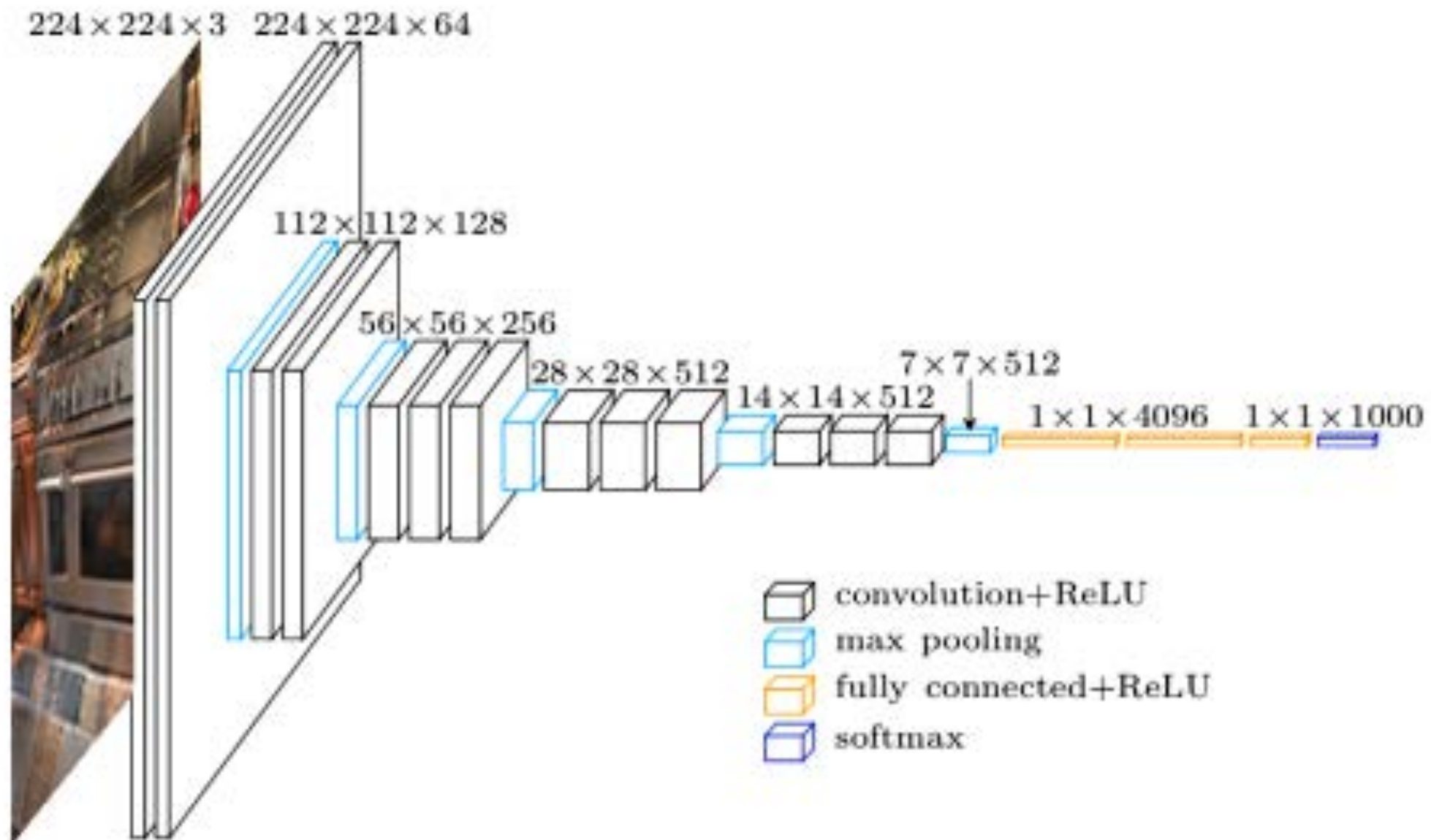
# Evolution of architectures

AlexNet architecture representation from the original 2012 Krizhevsky et al. paper.



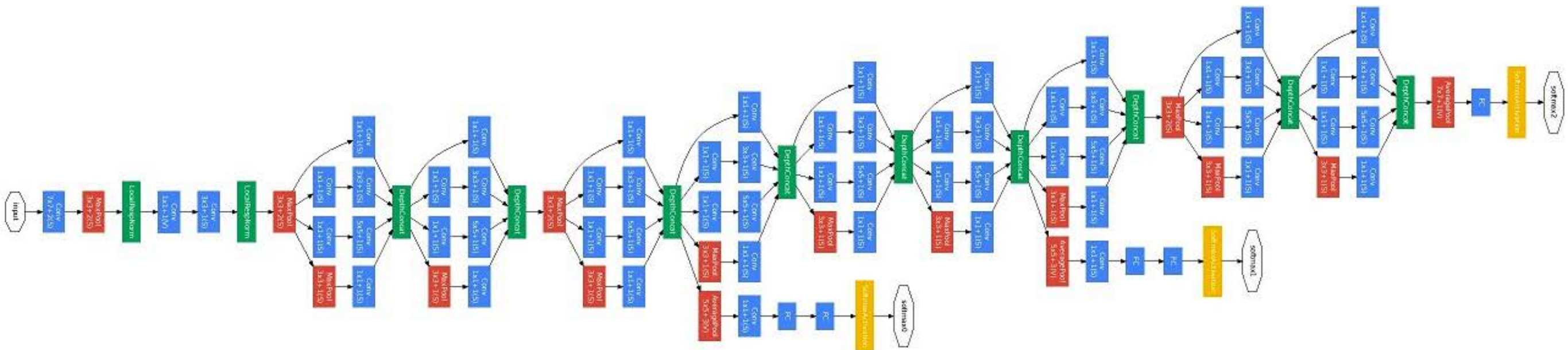
# Evolution of architectures

VGGI6 architecture representation from the original 2014 Simonyan et al. paper.



# Evolution of architectures

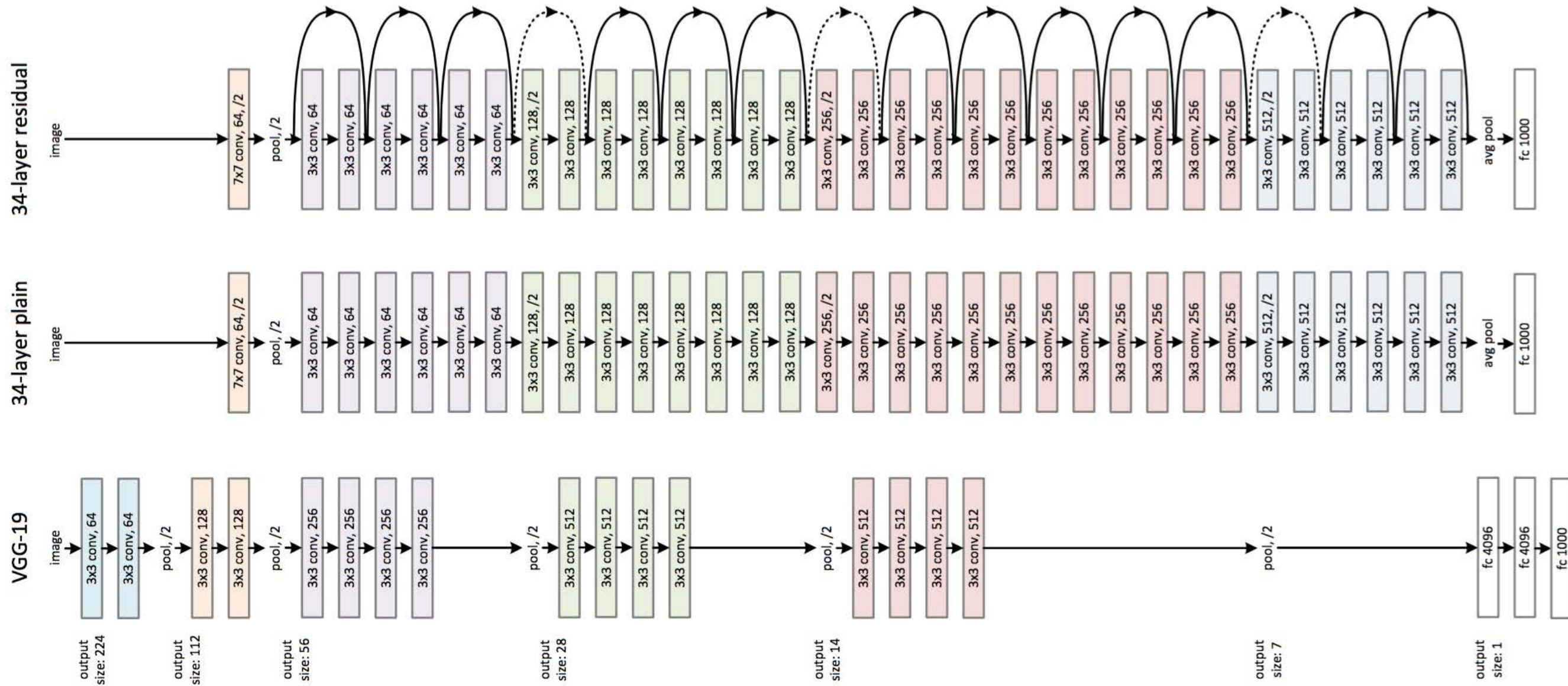
GoogLeNet architecture representation from the original 2014 Szegedy et al. paper.





# Evolution of architectures

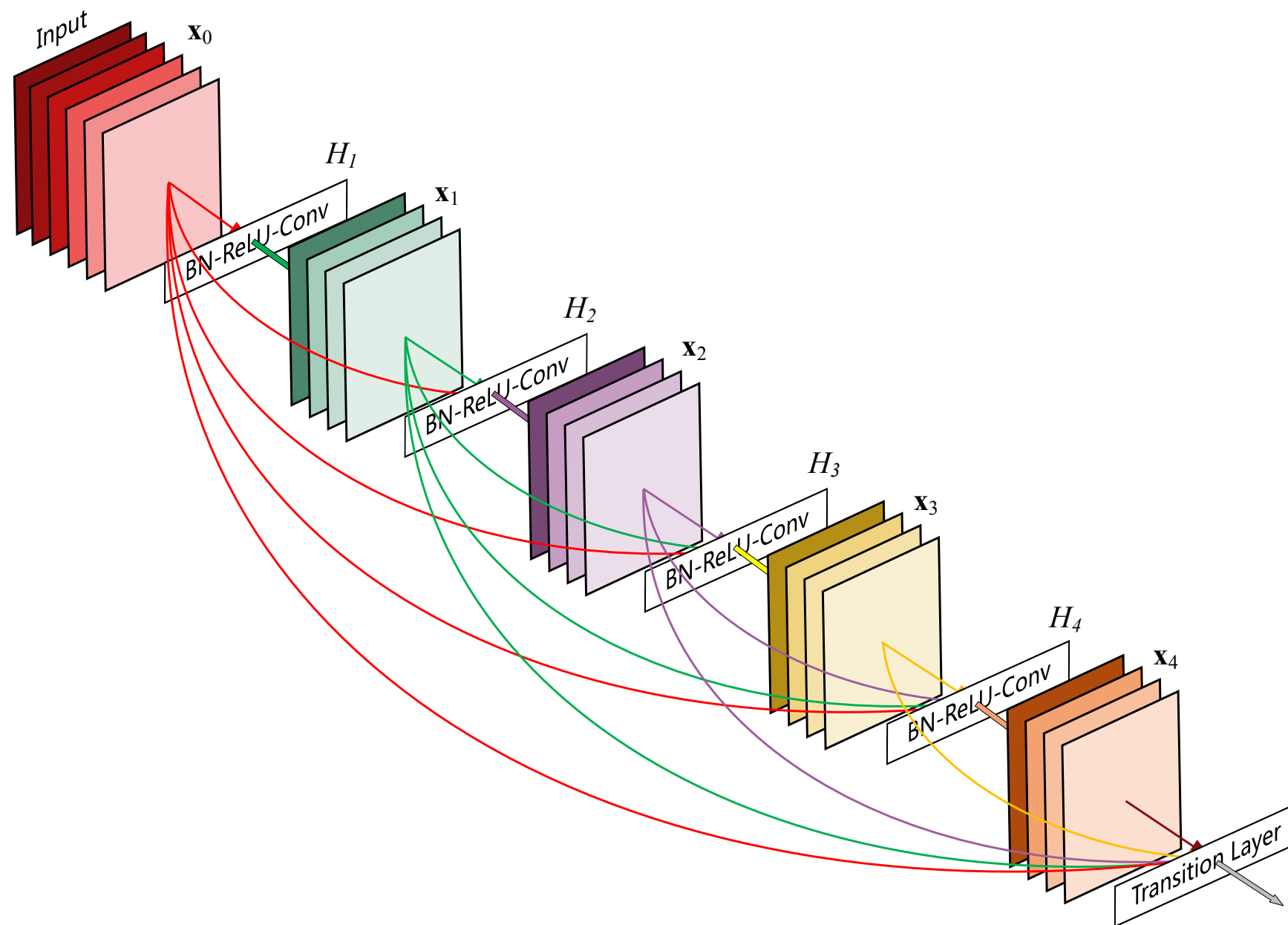
# ResNet architecture representation from the original 2015 He et al. paper.





# Evolution of architectures

DenseNet architecture representation from the original 2016 Huang et al. paper.



**Figure 1:** A 5-layer dense block with a growth rate of  $k = 4$ . Each layer takes all preceding feature-maps as input.

# Data augmentation

In many applications, the acquisition of labeled data is very expensive and tedious. In order to make the most of our few training examples, we will "augment" them via a number of random transformations, so that our model would never see twice the exact same picture. This helps prevent overfitting and helps the model generalize better.

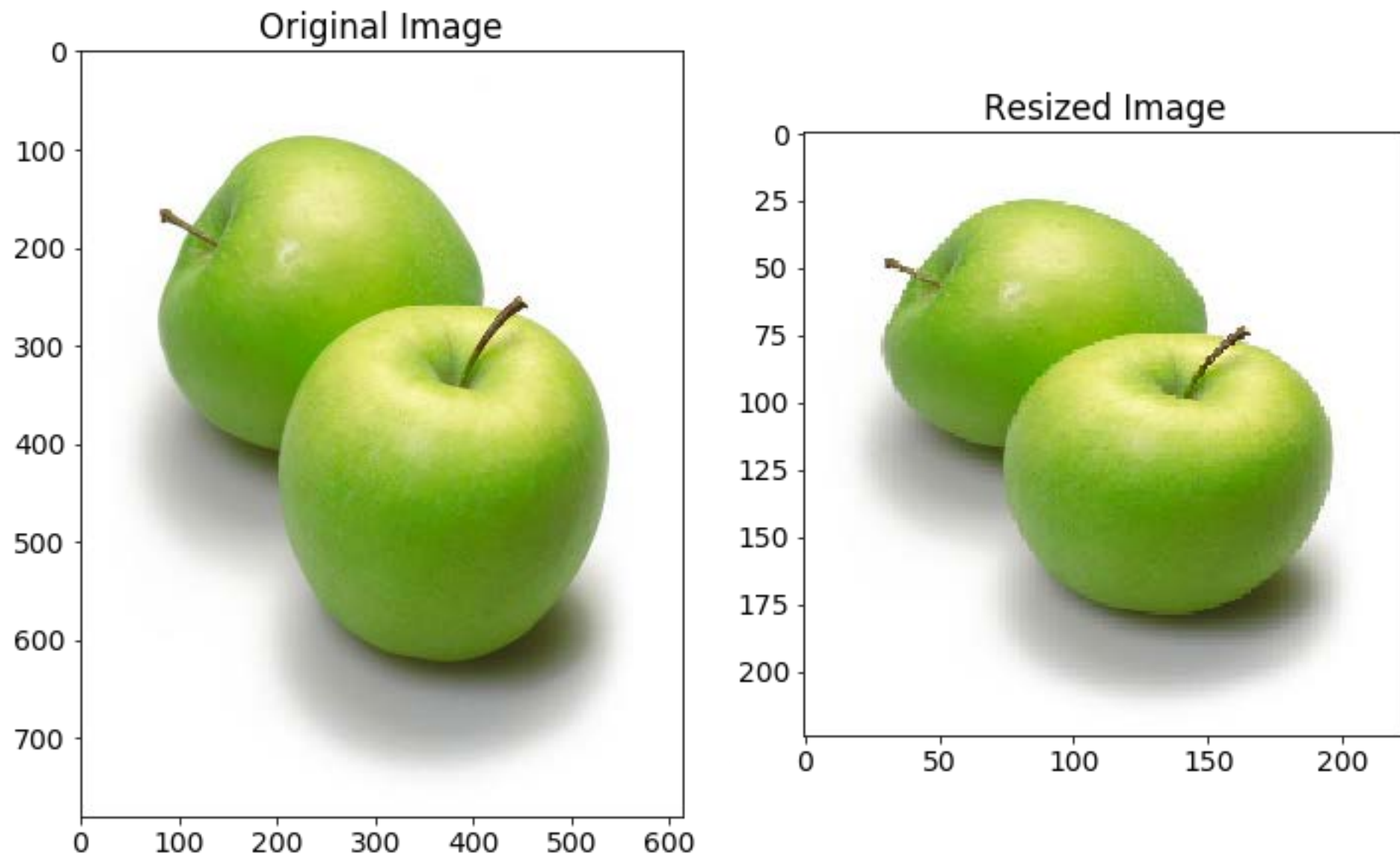
Here are the most common transformations used in practice:

- Rotation
- Horizontal or vertical translations
- Rescaling
- Shearing transformations
- Zooming
- Mirroring
- Adding noise
- Tweaking the lighting conditions
- Elastic distortions

\*Some of these transformations may alter the image by creating new pixel regions. This should be addressed in order not to introduce bias.

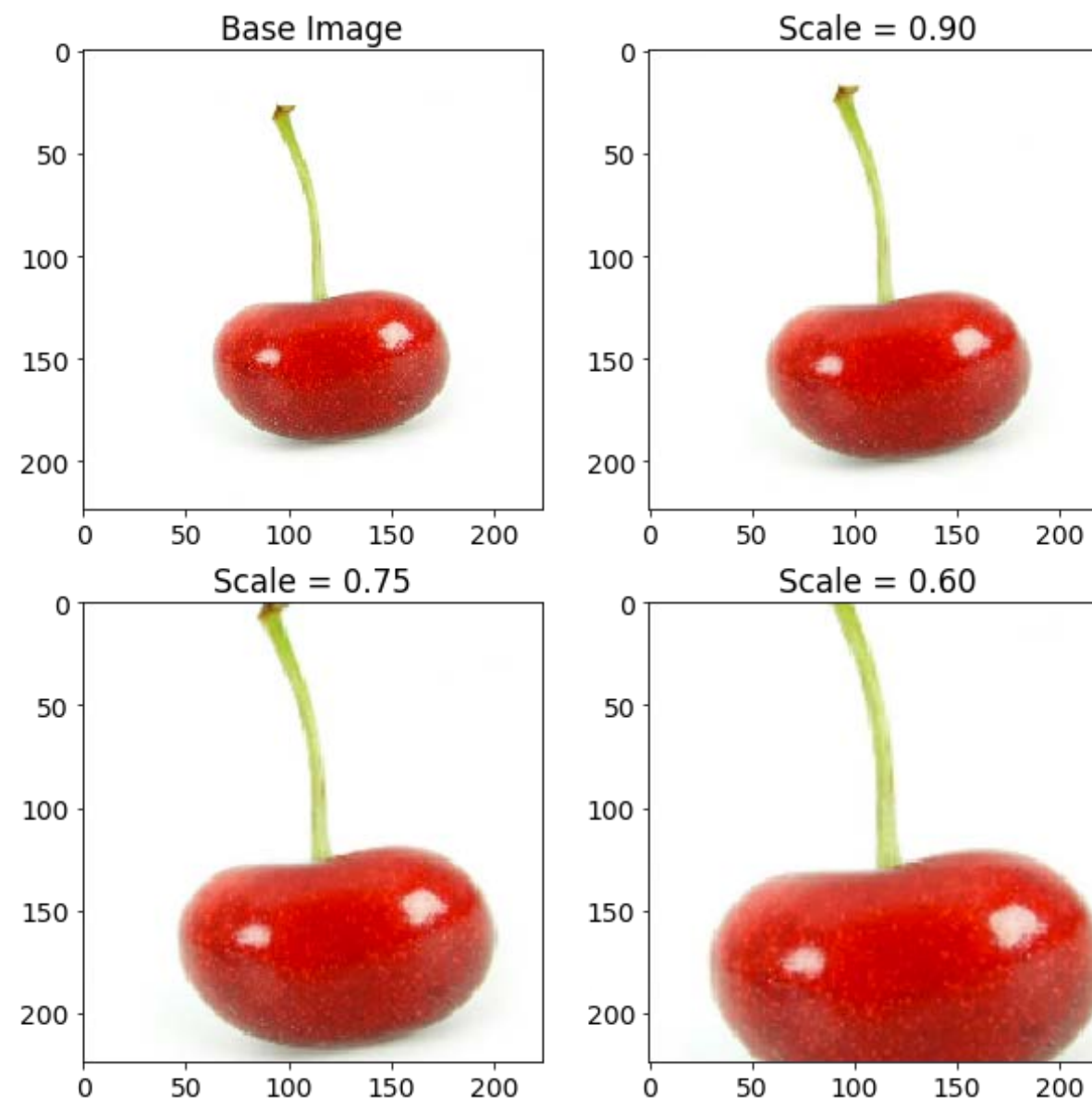
# Resizing

Training data may have varying sizes. Due to presence of fully connected layers in most of the neural networks, the images being fed to network will be required of a fixed size. Because of this, before the image augmentation happens, we have to preprocess the images to the size which our network needs. With the fixed sized image, we get the benefits of processing them in batches.



# Scaling

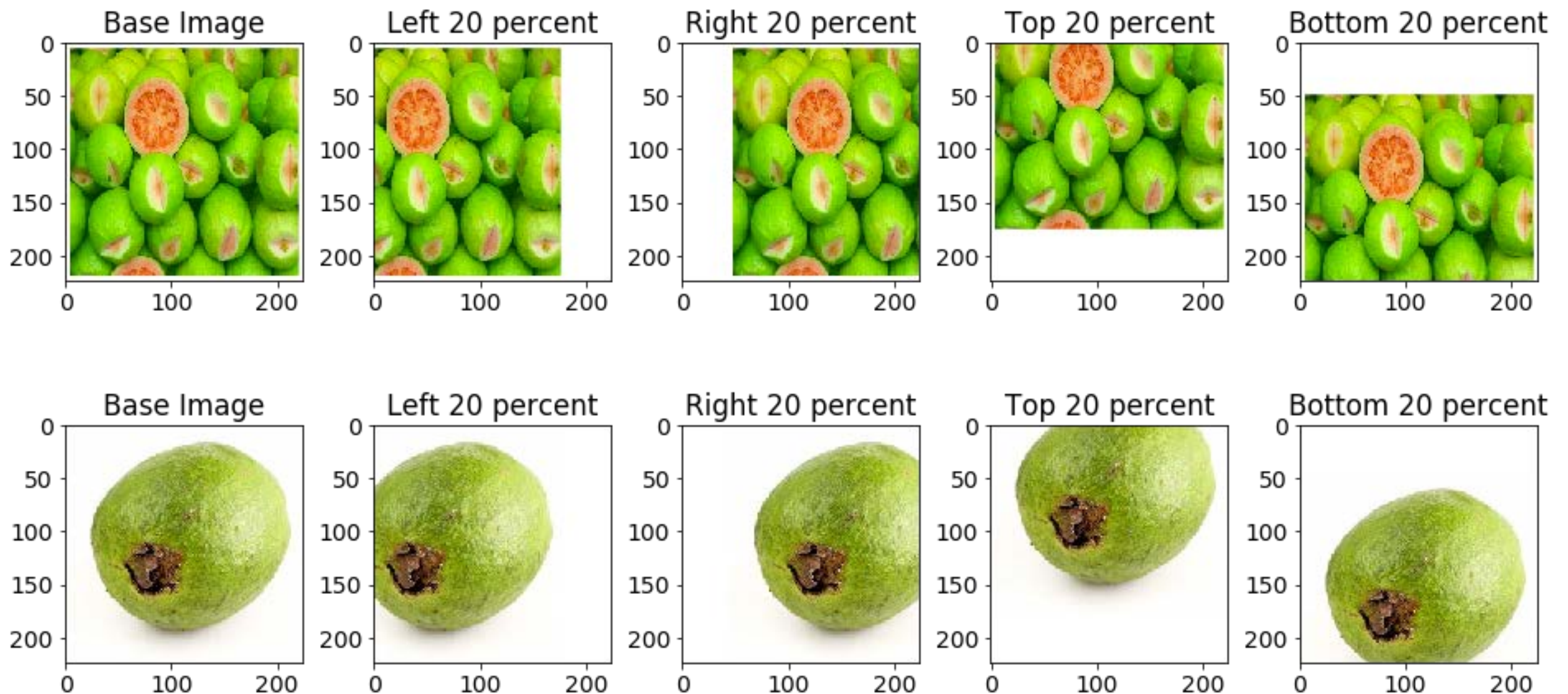
Having differently scaled object of interest in the images is the most important aspect of image diversity. When your network is in hands of real users, the object in the image can be tiny or large. Also, sometimes, object can cover the entire image and yet will not be present totally in image (i.e cropped at edges of object).





# Translation

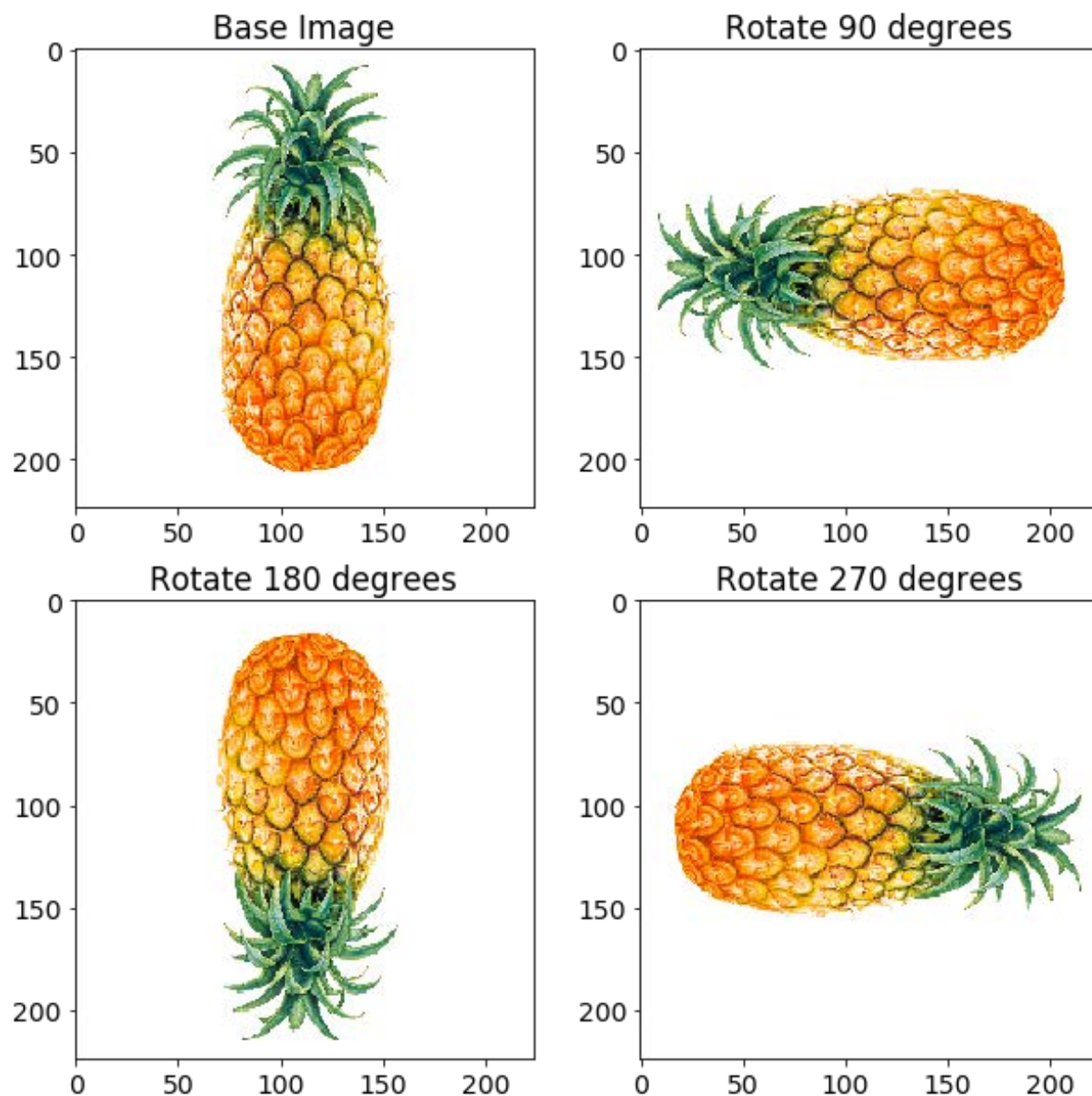
We would like our network to recognize the object present in any part of the image. Also, the object can be present partially in the corner or edges of the image. For this reason, we shift the object to various parts of the image. This may also result in addition of a background noise.



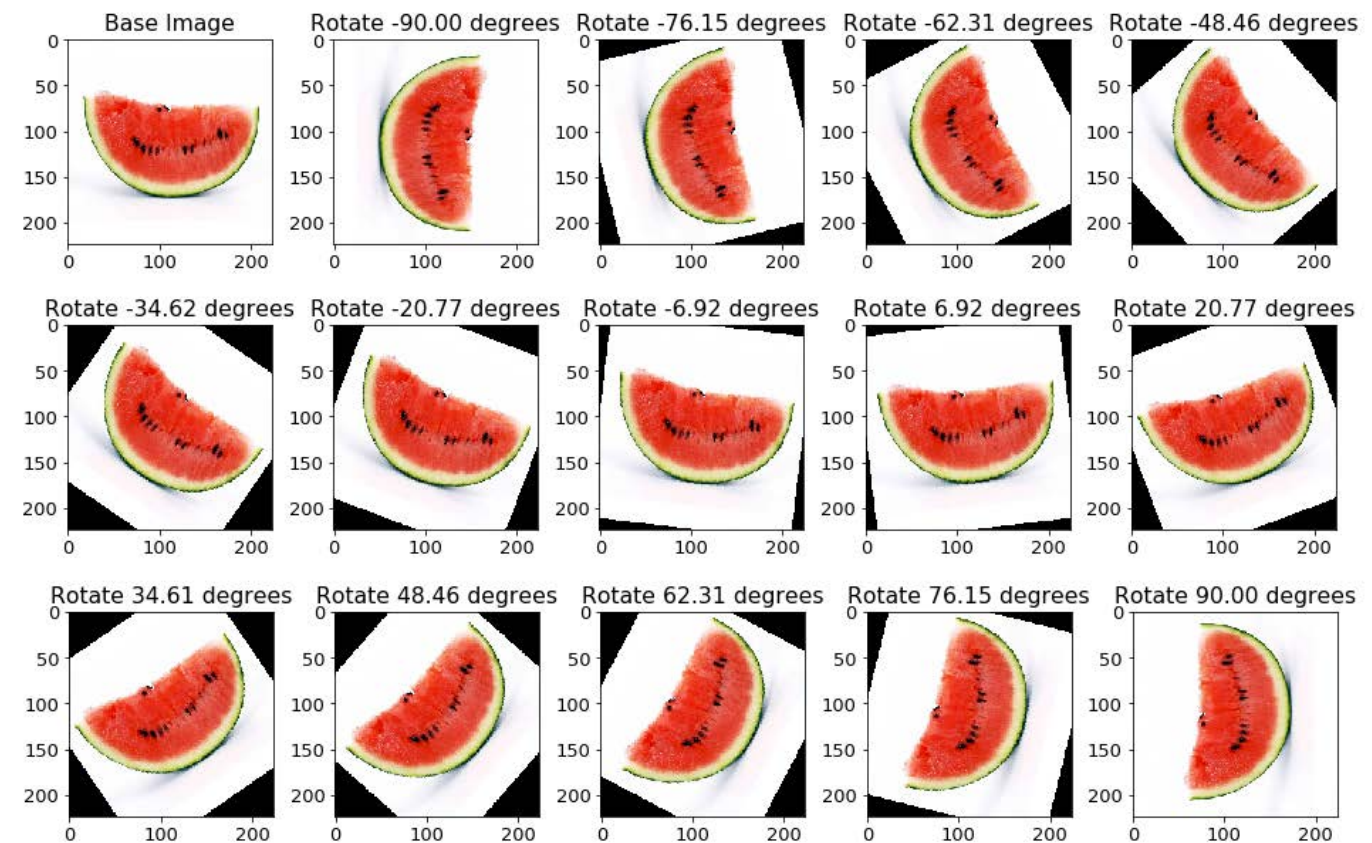


# Rotation

The network has to recognize the object present in any orientation. Assuming the image is square, rotating the image at 90 degrees will not add any background noise in the image.

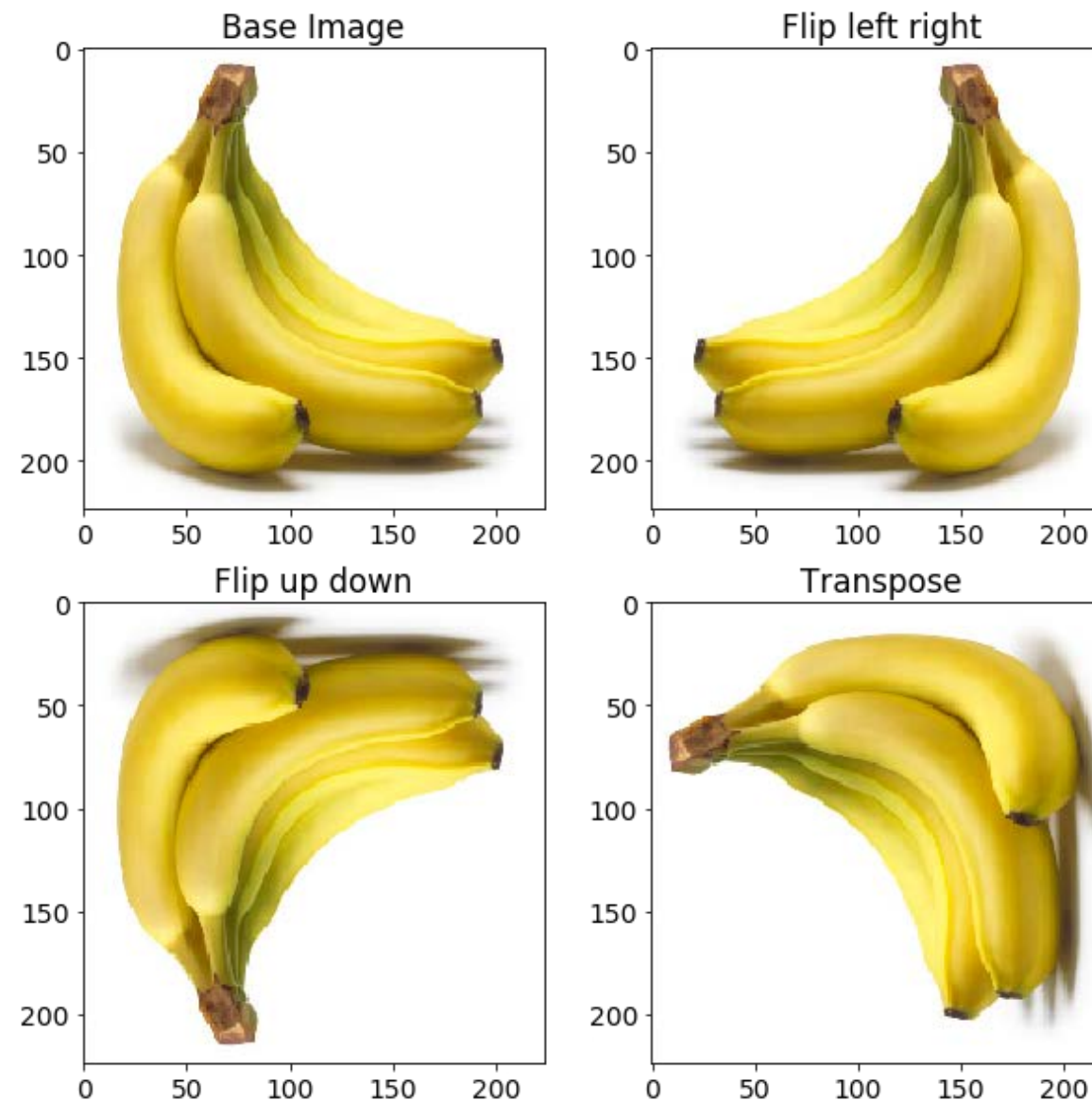


Depending upon the requirement, there maybe a necessity to orient the object at minute angles. However problem with this approach is, it will add background noise. If the background in image is of a fixed color (say white or black), the newly added background can blend with the image. However, if the newly added background color doesn't blend, the network may consider it as to be a feature and learn unnecessary features.



# Mirroring

This scenario is more important for network to remove biasness of assuming certain features of the object is available in only a particular side. Consider the case shown the image below. You don't want network to learn that tilt of banana happens only in right side as observed in the base image. Also notice that mirroring produces different set of images from rotation at multiple of 90 degrees.

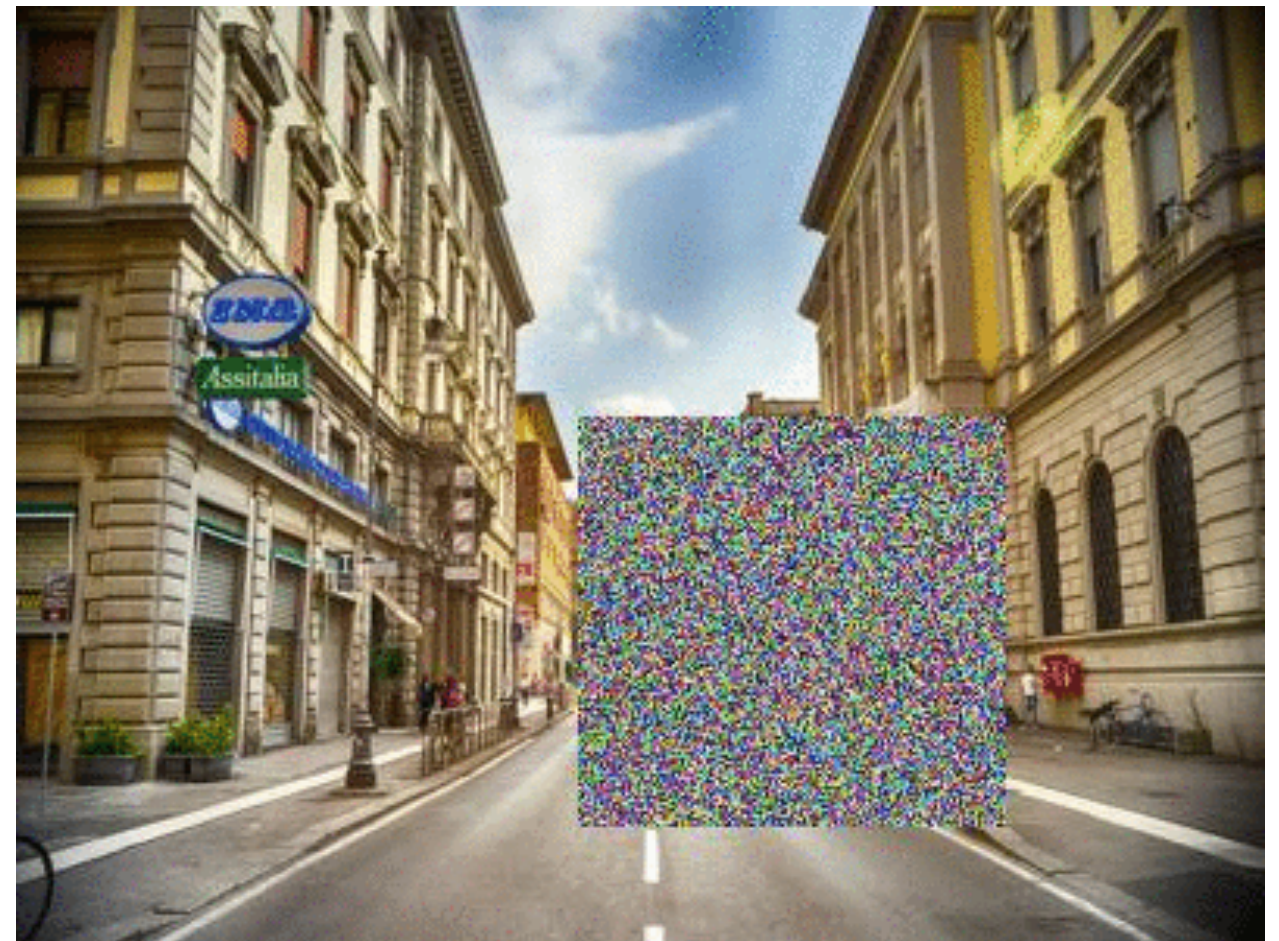
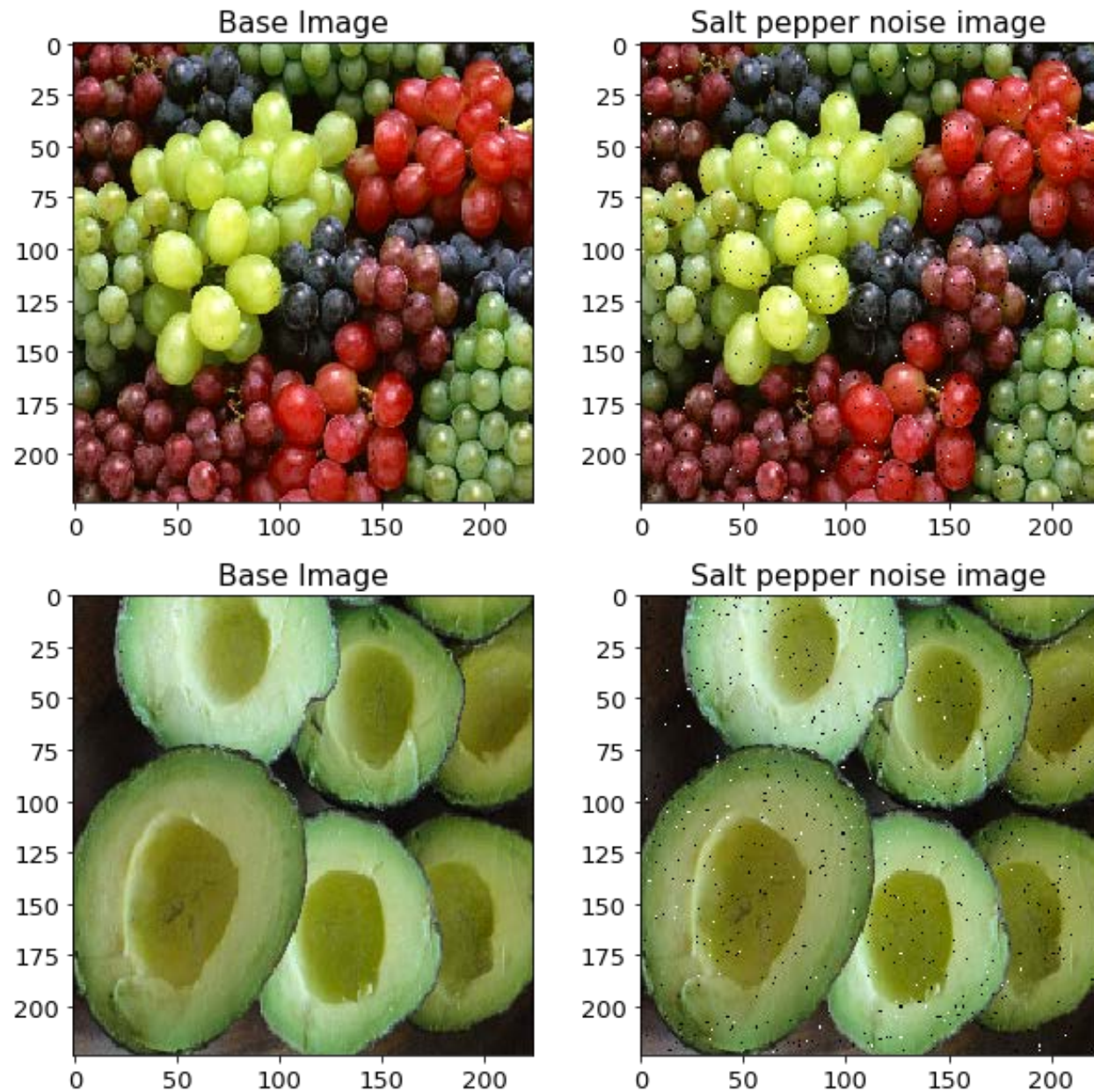




# Adding noise

Though this may seem unnecessary, it may be helpful to safeguard against corrupted data.

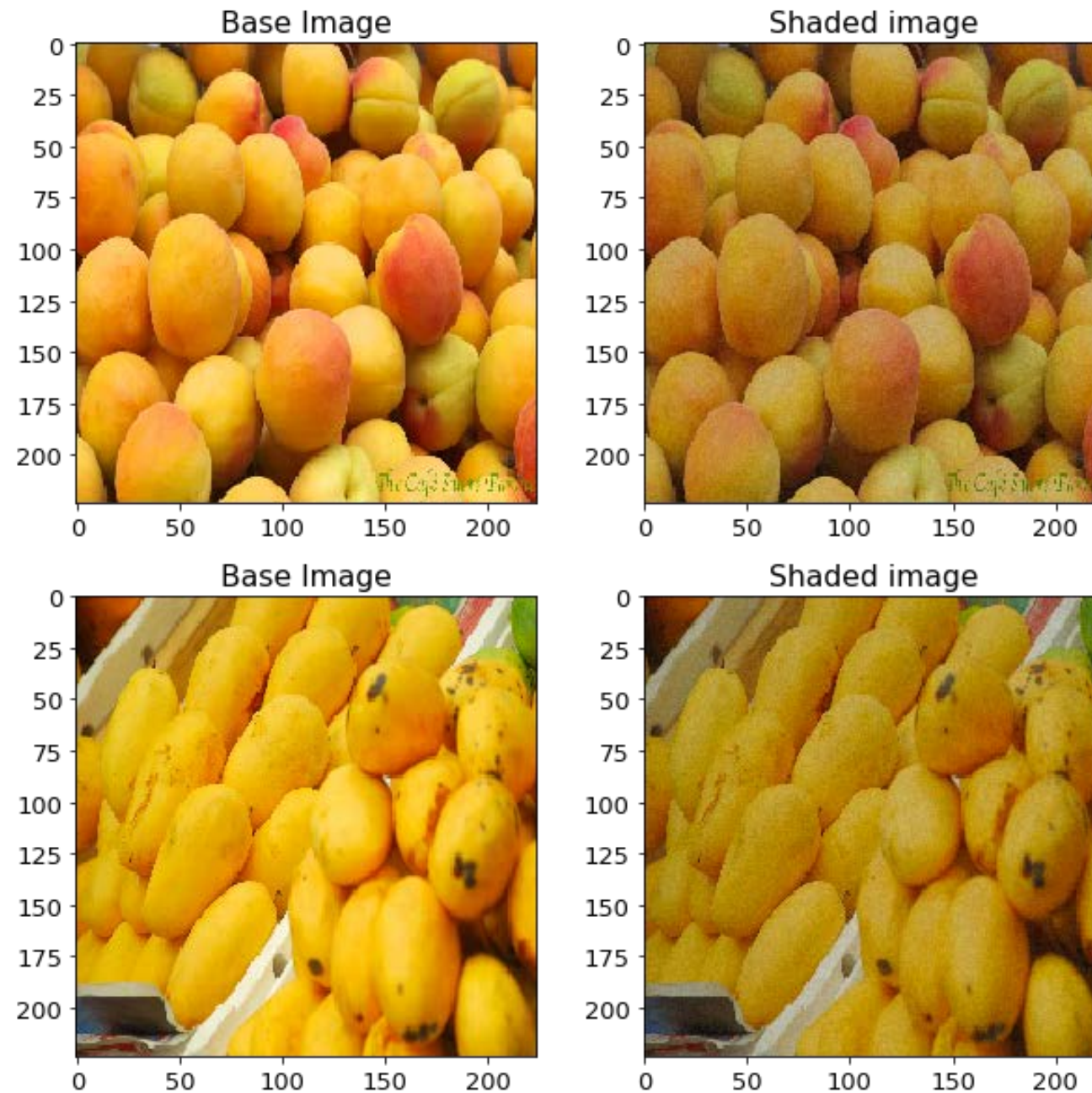
Removing entire regions can be used to make models robust to occlusion. This may be useful for training neural networks used in object detection in navigation scenarios, for example.





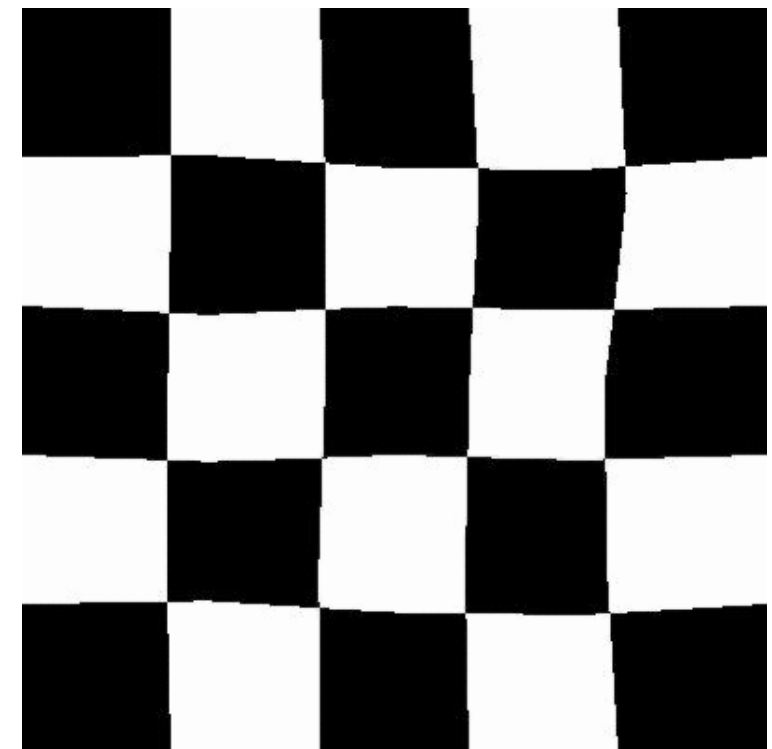
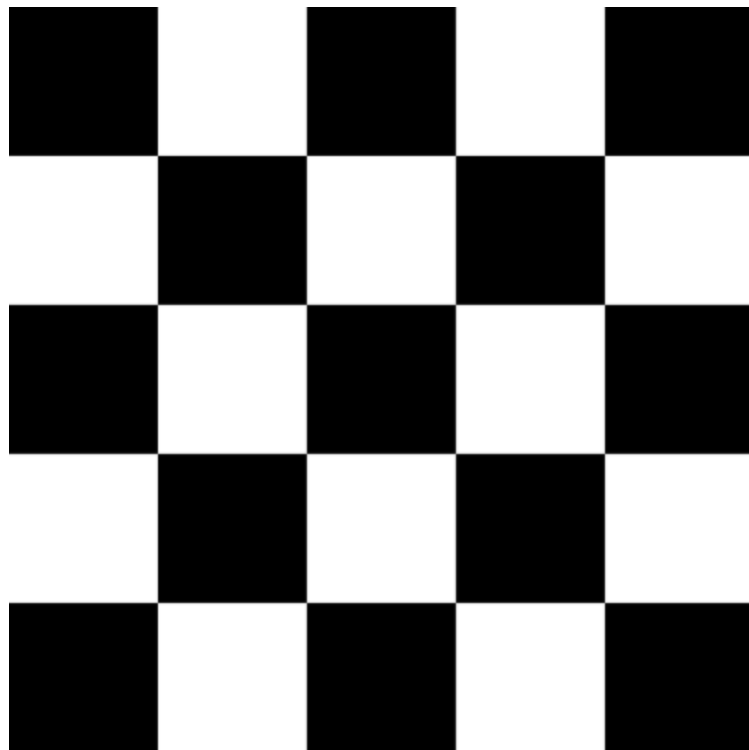
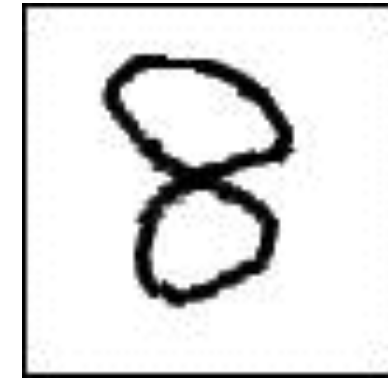
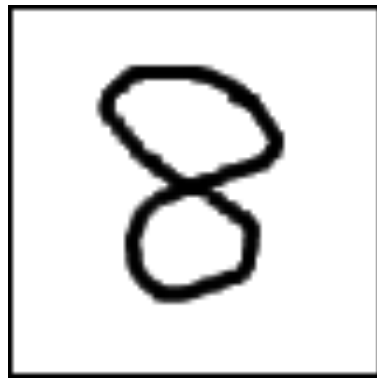
# Lighting conditions

This is a very important type of diversity needed in the image dataset not only for the network to learn properly the object of interest but also to simulate the practical scenario of images being taken by the user. The lighting condition of the images are varied by adding Gaussian noise in the image.



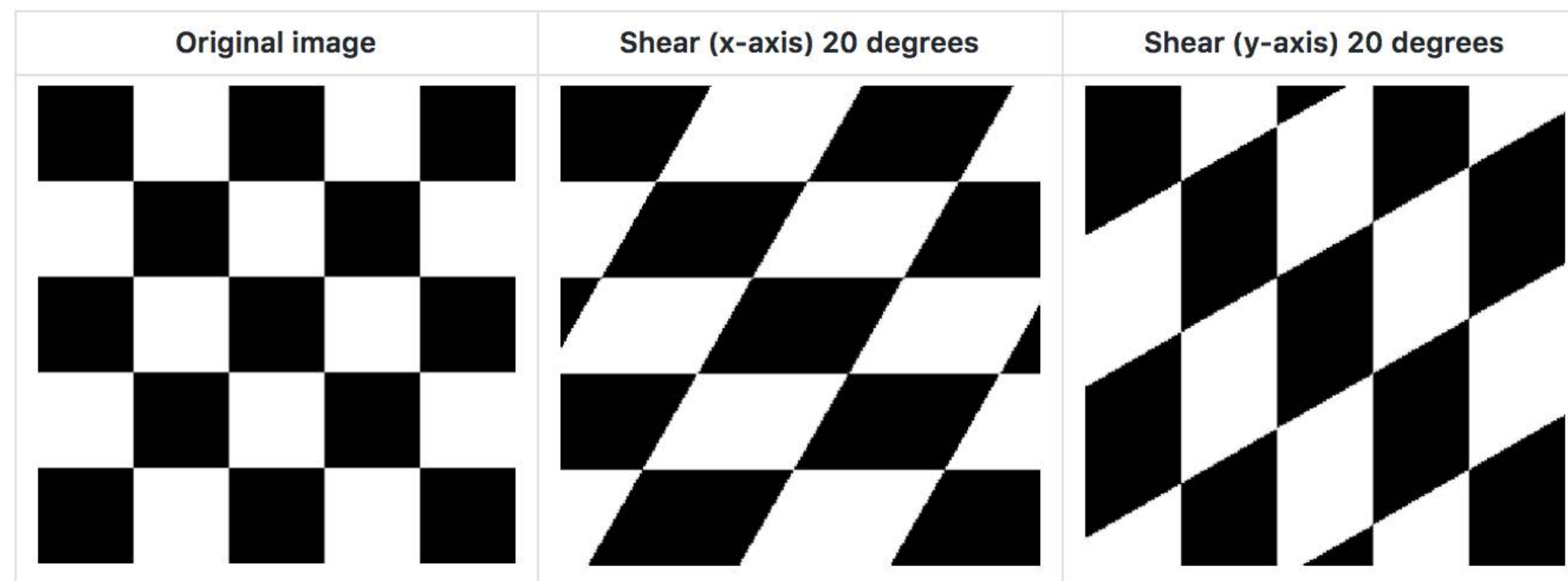
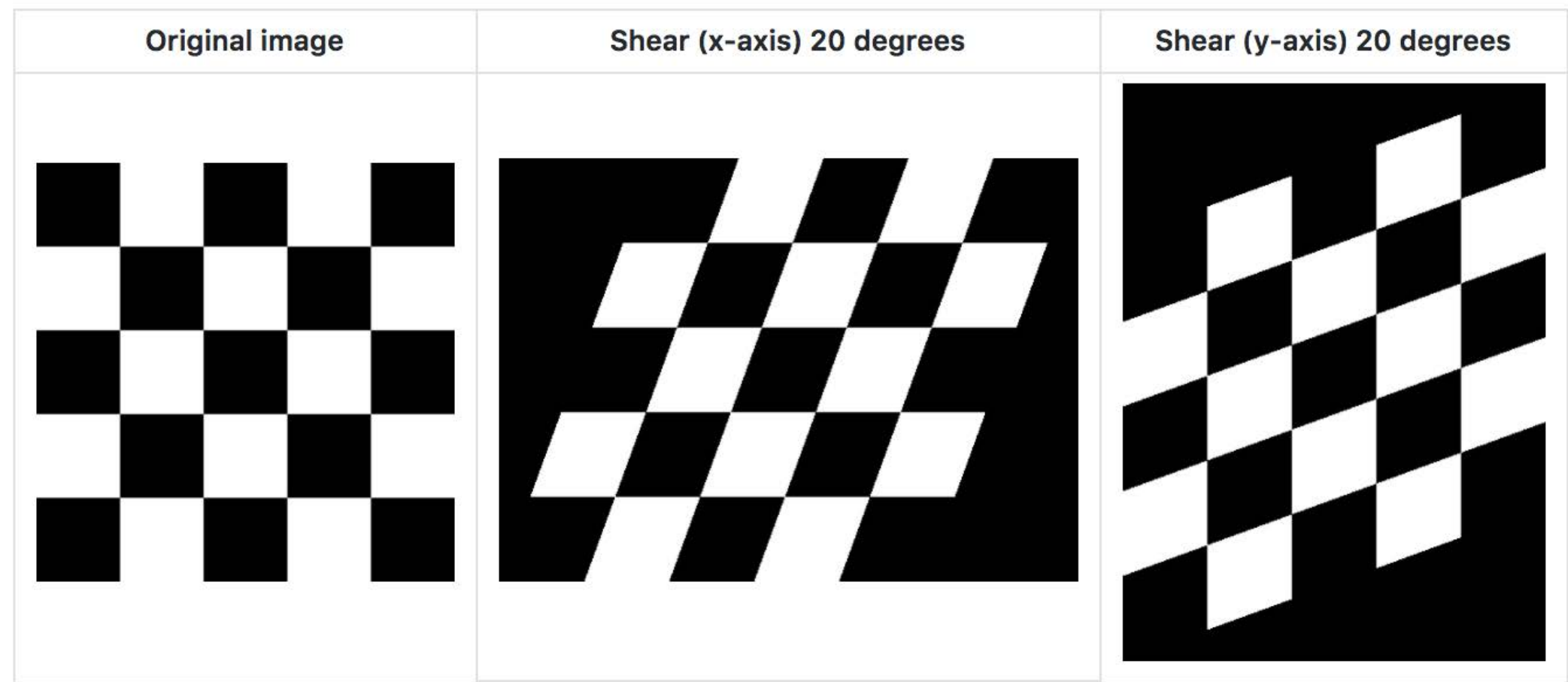
# Elastic Distortions

Using elastic distortions, one can generate many “artificial” images that are label preserving.





# Size Preserving Shearing



# Classification of hand-written digits

