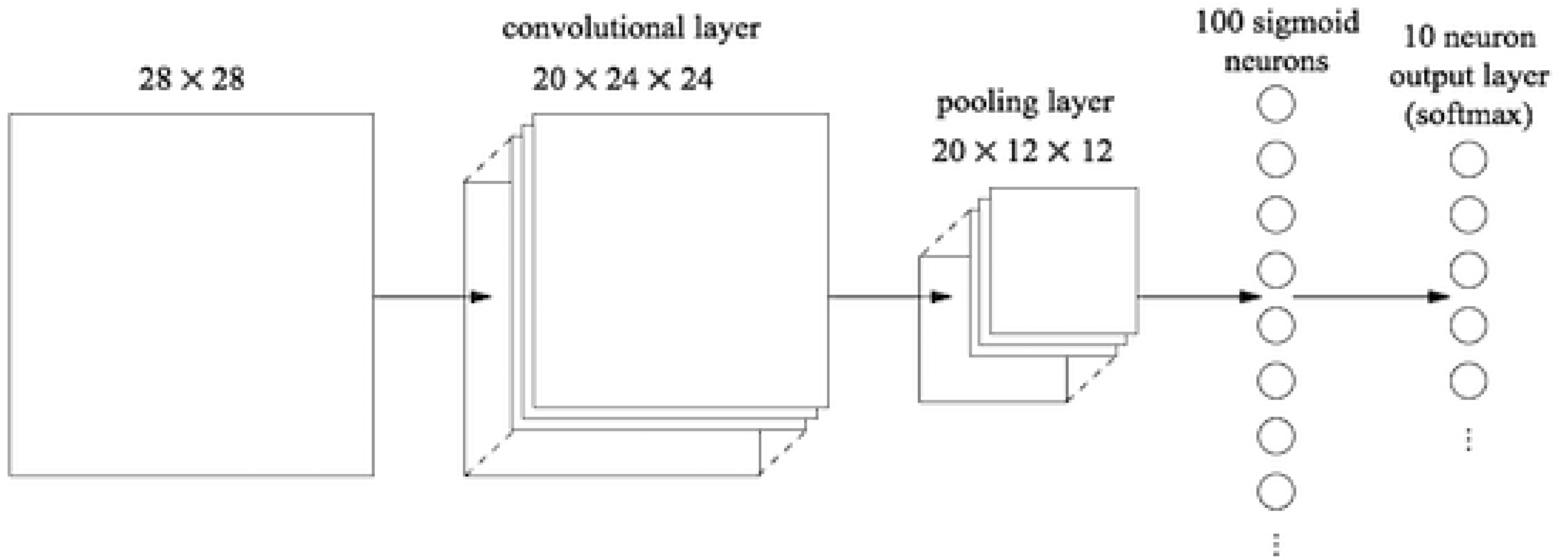


Conv Nets

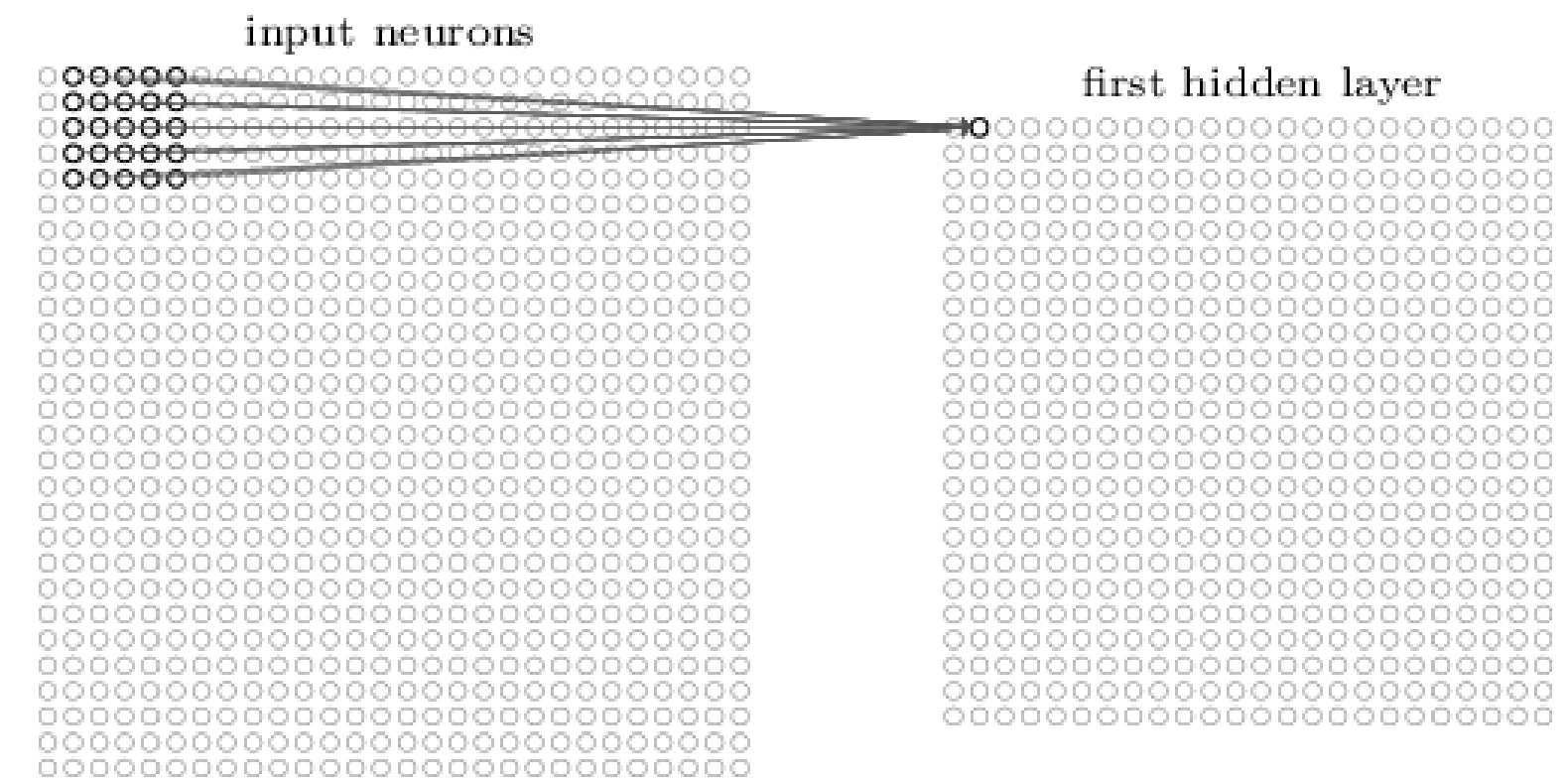


A Convolutional Neural Network (CNN) is a type of deep learning algorithm designed to analyze visual data, such as images or videos. It has proven to be highly effective in tasks like image classification, object detection, and image recognition.

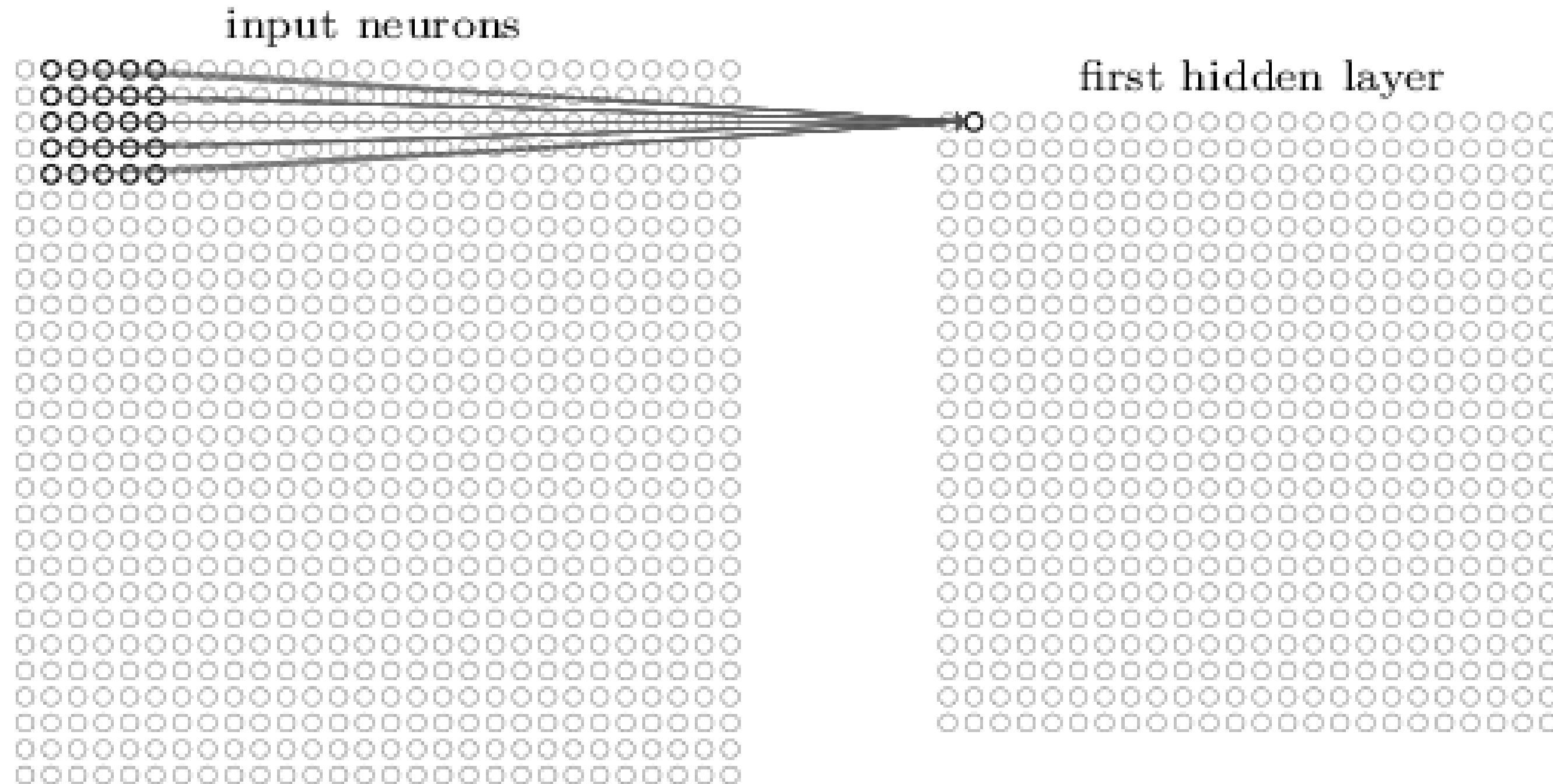
Convolutional neural networks use three basic ideas: local receptive fields, shared weights, and pooling. Let's look at each of these ideas in turn.

1. Local receptive fields: In the fully-connected layers shown earlier, the inputs were depicted as a vertical line of neurons. In a convolutional net, the inputs are a $m \times n$ matrix of neurons, where m, n are the size of the images.

Again in the FC layers, each hidden neuron was connected to each and every neuron in the previous layer. But in Convnet, we won't connect every input pixel to every hidden neuron. Instead, we only make connections in small, localized regions of the input image. To be more precise, each neuron in the first hidden layer will be connected to a small region of the input neurons, say, for example, a 5×5 region, corresponding to 25 input pixels. That region in the input image is called the local receptive field for the hidden neuron.



Then we slide the local receptive field over by one pixel to the right (i.e., by one neuron), to connect to a second hidden neuron:



And so on, building up the first hidden layer. Note that if we have a 28×28 input image, and 5×5 local receptive fields, then there will be 24×24 neurons in the hidden layer.

2. Shared weights and biases: Each hidden neuron has a bias and 5×5 weights connected to its local receptive field. Moreover, we're going to use the same weights and bias for each of the 24×24 hidden neurons.

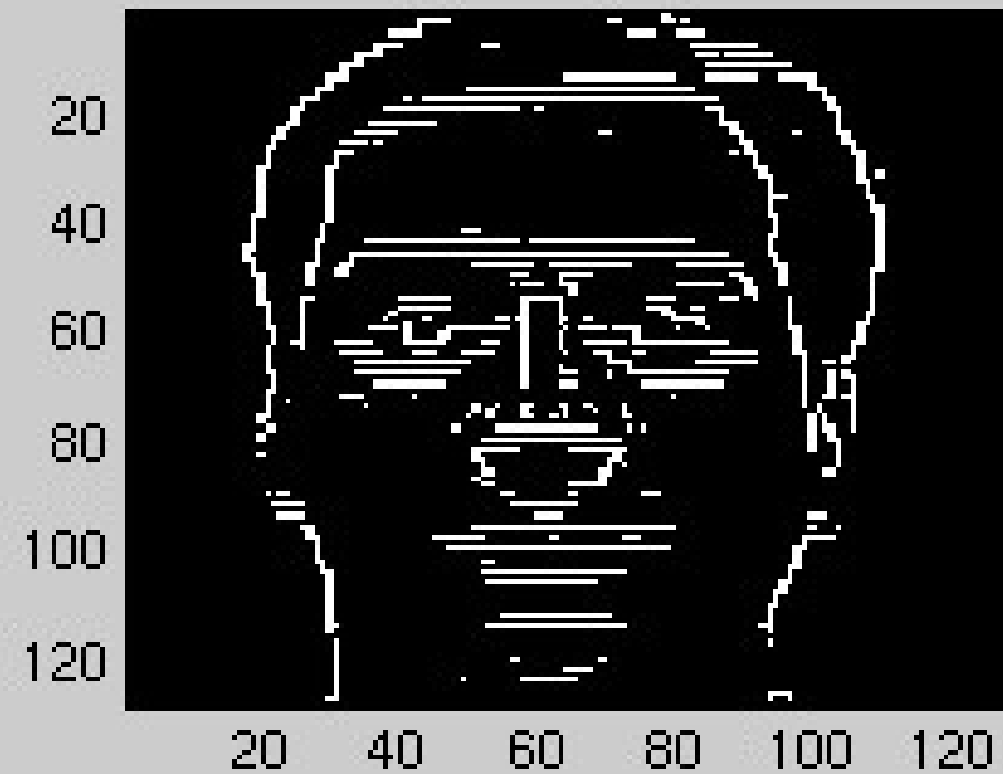
This means that all the neurons in the first hidden layer detect exactly the same feature, just at different locations in the input image. To see why this makes sense, suppose the weights and bias are such that the hidden neuron can pick out, say, a vertical edge in a particular local receptive field. That ability is also likely to be useful at other places in the image. And so it is useful to apply the same feature detector everywhere in the image. For this reason, we sometimes call the map from the input layer to the hidden layer a feature map. We call the weights defining the feature map the shared weights. And we call the bias defining the feature map in this way the shared bias. The shared weights and bias are often said to define a kernel or filter.

Advantage of sharing weights

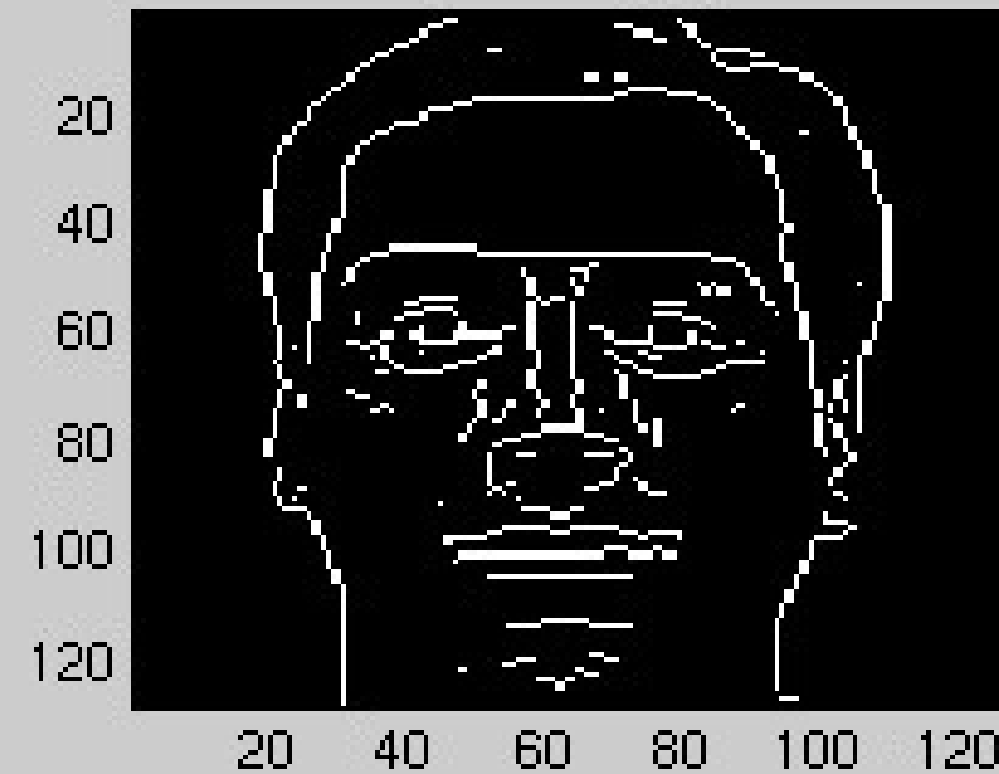
A big advantage of sharing weights and biases is that it greatly reduces the number of parameters involved in a convolutional network. For each feature map we need $25=5\times 5$ shared weights, plus a single shared bias. So each feature map requires 26 parameters. If we have 20 feature maps that's a total of $20\times 26=520$ parameters defining the convolutional layer. By comparison, suppose we had a fully connected first layer, with $784=28\times 28$ input neurons, and a relatively modest 30 hidden neurons, as we used in many of the examples earlier in the book. That's a total of 784×30 weights, plus an extra 30 biases, for a total of 23,550 parameters. In other words, the fully-connected layer would have more than 40 times as many parameters as the convolutional layer.

Output image on application of different filters/kernels

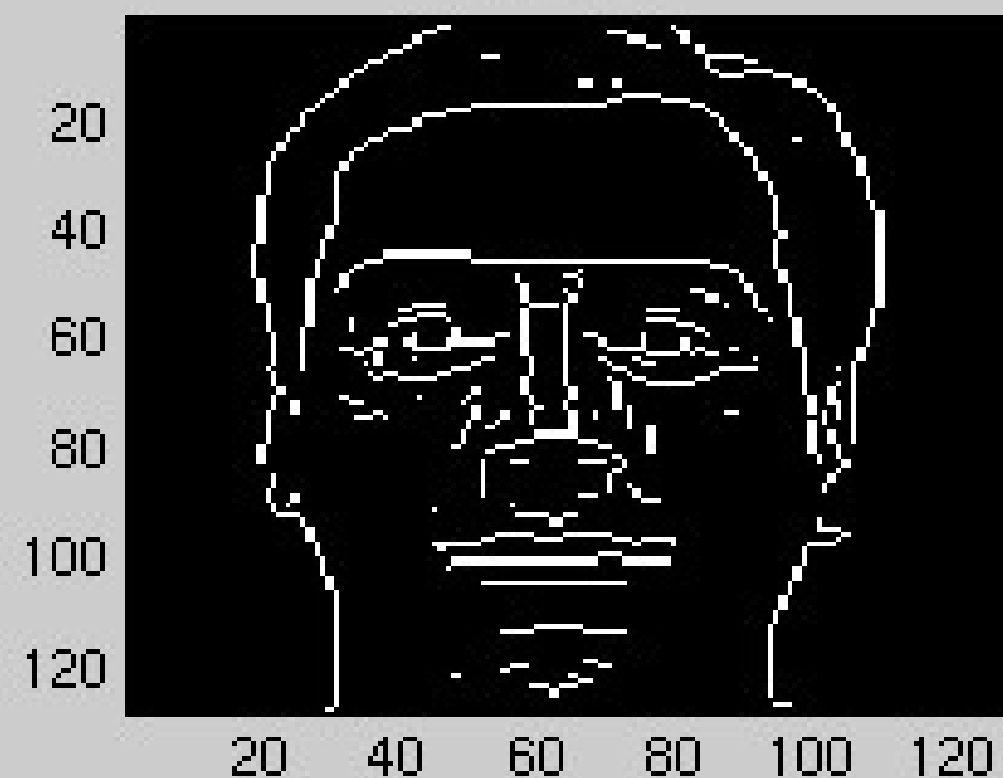
Roberts



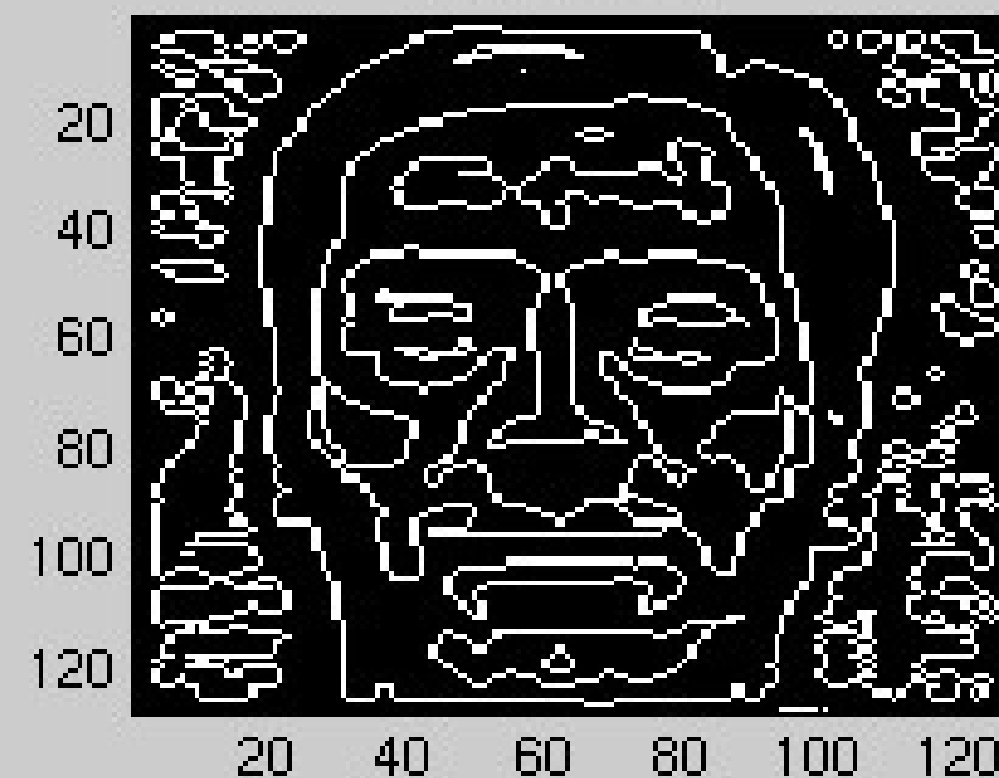
Sobel



Prewitt



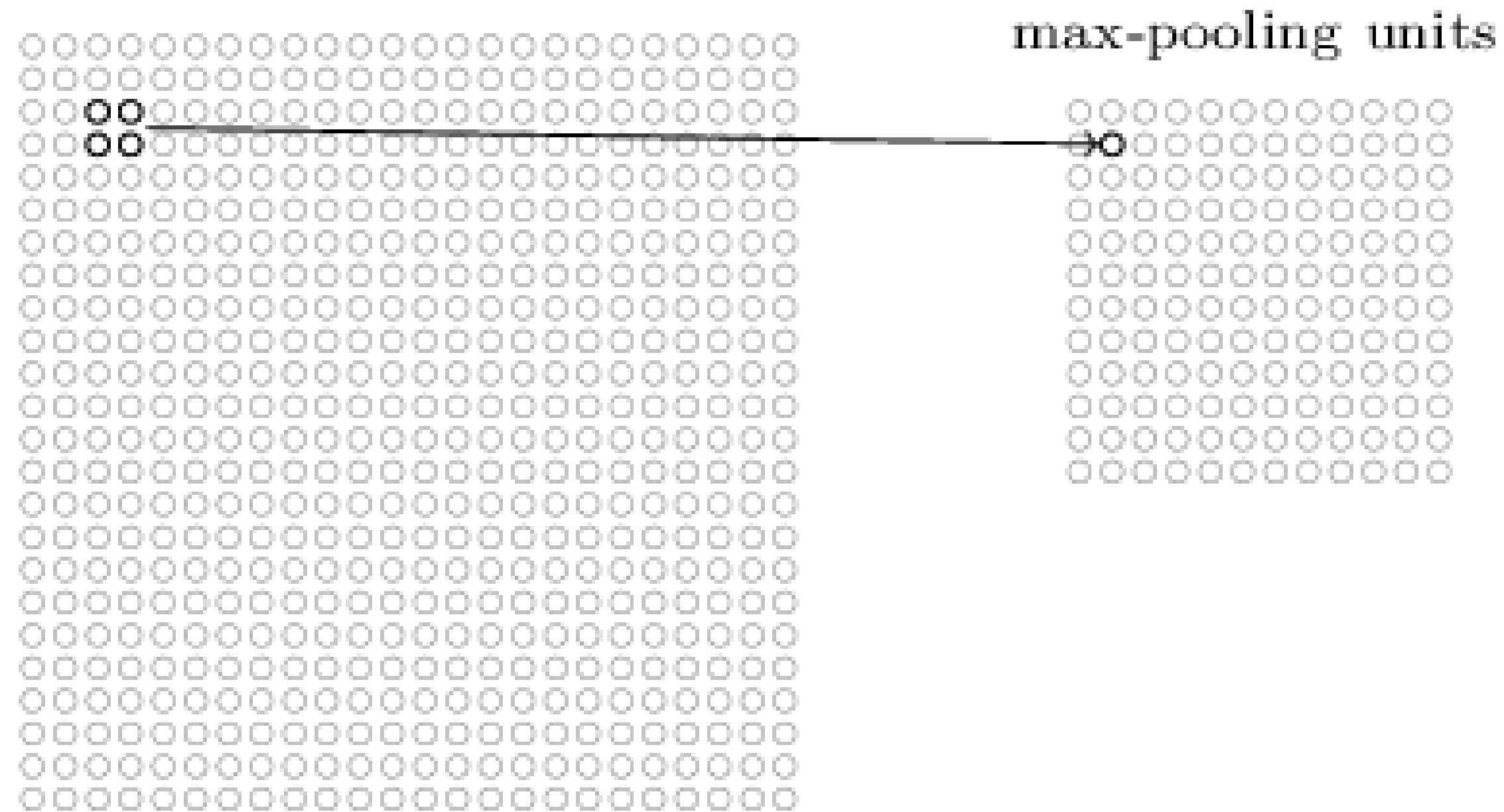
Marr-Hildreth



Pooling layers: In addition to the convolutional layers just described, convolutional neural networks also contain pooling layers. Pooling layers are usually used immediately after convolutional layers. What the pooling layers do is simplify the information in the output from the convolutional layer.

In detail, a pooling layer takes each feature map, output from the convolutional layer and prepares a condensed feature map. For instance, each unit in the pooling layer may summarize a region of (say) 2×2 neurons in the previous layer. As a concrete example, one common procedure for pooling is known as max-pooling. In max-pooling, a pooling unit simply outputs the maximum activation in the 2×2 input region, as illustrated in the following diagram:

hidden neurons (output from feature map)

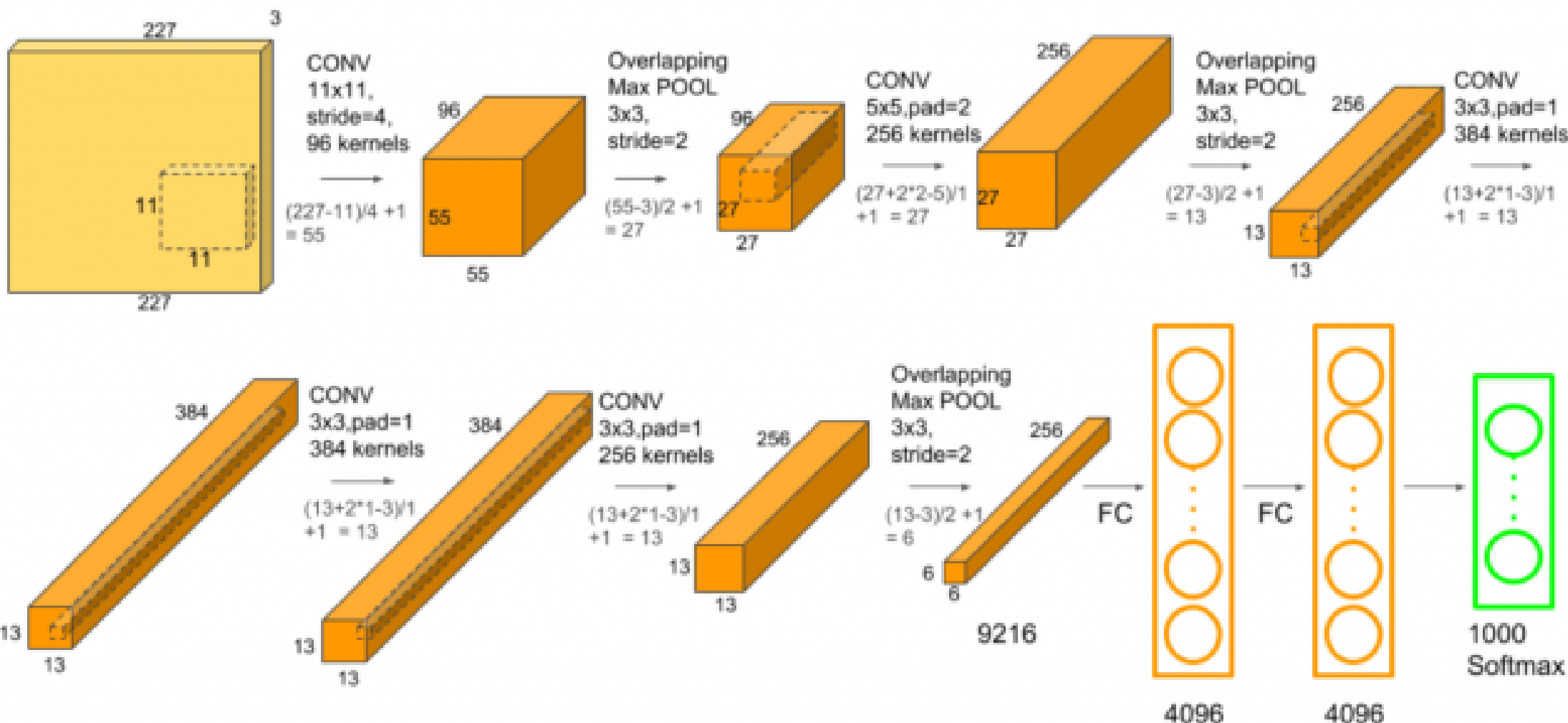
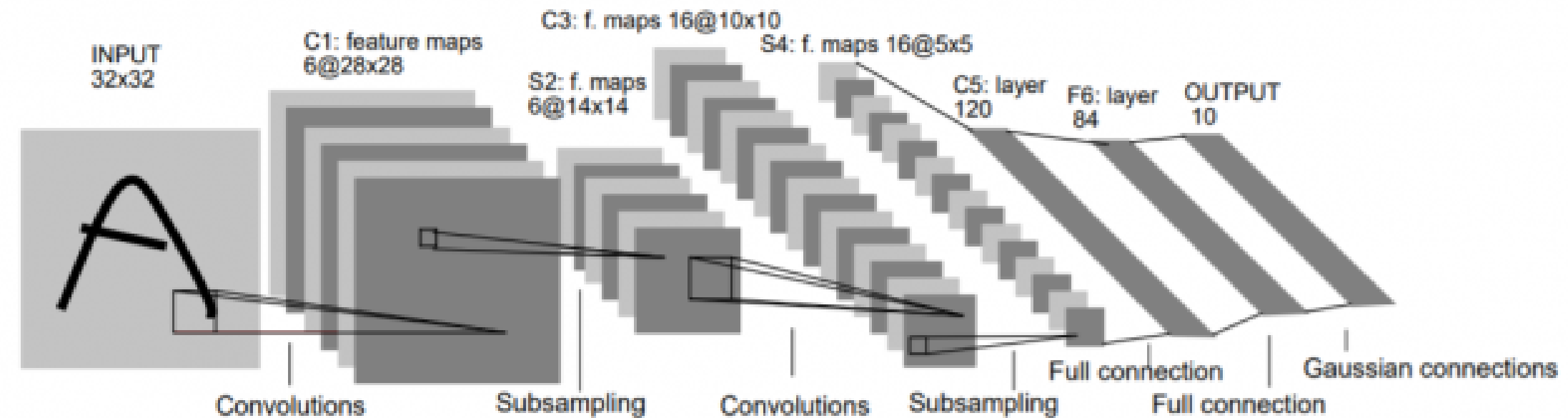


Max-pooling isn't the only technique used for pooling. Another common approach is known as L2 pooling. Here, instead of taking the maximum activation of a 2×2 region of neurons, we take the square root of the sum of the squares of the activations in the $2 \times 2 \times 2$ region.

Some Common CNN Architectures

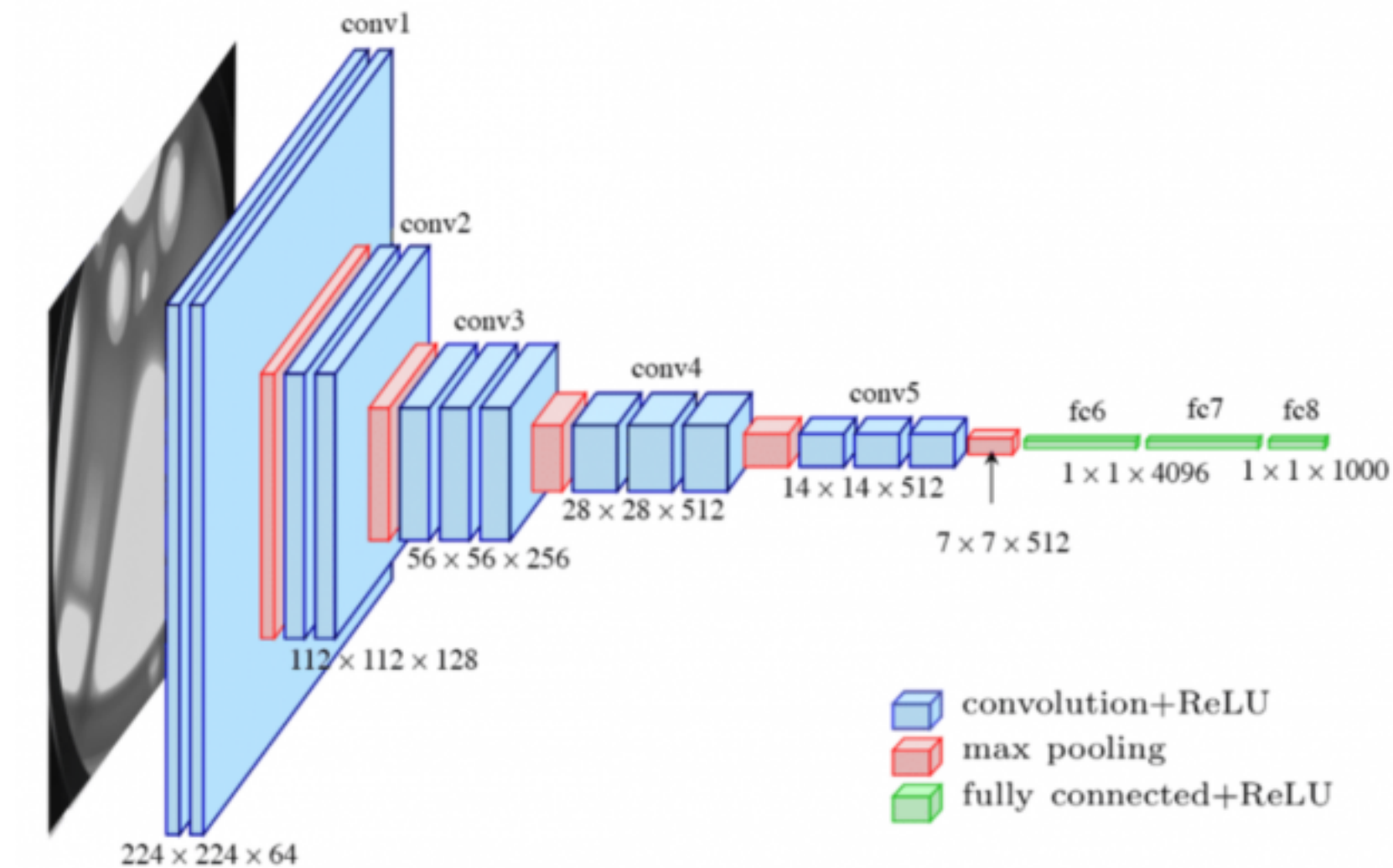
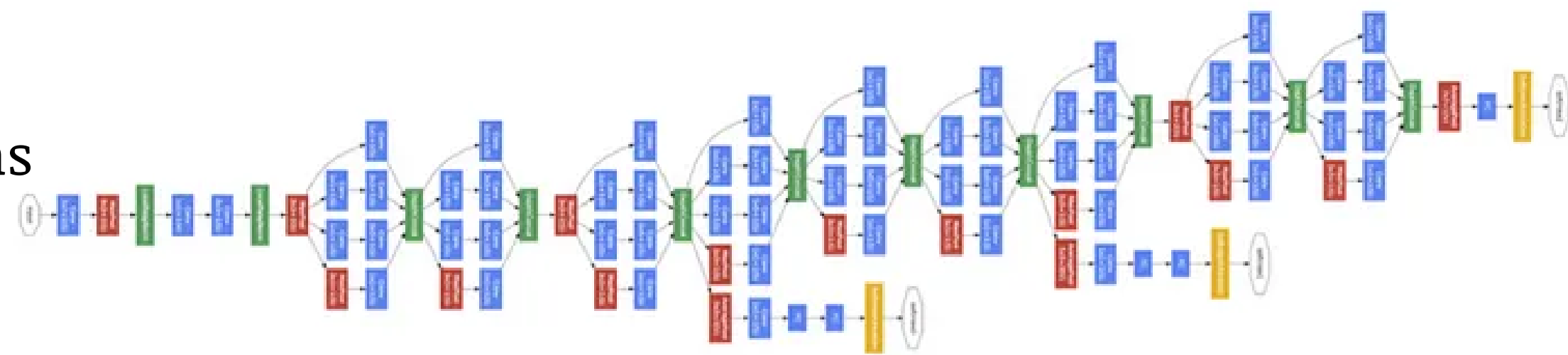
For more info

1. **LeNet-5**: LeNet is the first CNN architecture. The model has five convolution layers followed by two fully connected layers. It takes input images of shape (32,32).



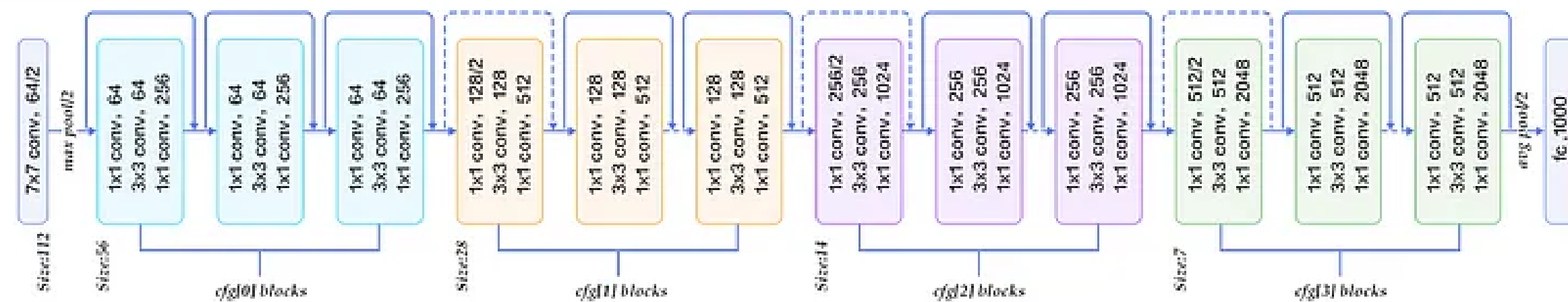
2. **AlexNet**: 5 convolutional layers with a combination of max-pooling layers, 3 fully connected layers, and 2 dropout layers. The activation function used in all layers is Relu. The activation function used in the output layer is Softmax. It takes input images of shape (227,227,3).

3. GoogLeNet / Inception : Implemented a novel element called inception module. It used batch normalization, image distortions and RMSprop optimizer. Their architecture consisted of a 22 layer deep CNN with 4 million parameters.



4. VGGNet : VGGNet is a 16-layer CNN with up to 95 million parameters and trained on over one billion images (1000 classes). It can take large input images of 224 x 224-pixel size for which it has 4096 convolutional features.

5. ResNet : The network has 152 layers and over one million parameters. Mostly used for image classification tasks with 1000 classes,



The handwritten digit classification task :

- The network begins with 28×28 input neurons, which encode the MNIST image's pixel intensities.
- This is followed by a convolutional layer using a 5×5 local receptive field and 6 feature maps. The result is a layer of $6 \times 24 \times 24$ hidden feature neurons.
- The next step is a max-pooling layer, applied to 2×2 regions across each of the 6 feature maps. The result is a layer of $6 \times 12 \times 12$ hidden feature neurons.
- Followed by another convolutional layer of 5×5 local receptive field and 16 feature maps, and then a max-pooling layer, resulting in outputs of shape $(4, 4, 16)$.
- The final layer of connections in the network is fully-connected. This layer connects every neuron from the max-pooled layer to every one of the 10 output neurons.

<https://colab.research.google.com/drive/1ywUwvrOc9N6ctr84ke0bFs5-XrA7XFQx?usp=sharing>