> ## Other Visual Tasks
>
> There was stunning progress as well in other visual tasks such as object detection and localization, and image segmentation. In object detection and localization, the neural network typically outputs a sequence of bounding boxes around various objects in the image. For example, see Maxine Oquab et al.'s 2015 paper that outputs a heat map for each object class, or Russell Stewart et al.'s 2015 paper that uses a combination of a CNN to detect faces and a recurrent neural network to output a sequence of bounding boxes around them. In image segmentation, the net outputs an image (usually of the same size as the input) where each pixel indicates the class of the object to which the corresponding input pixel belongs. For example, check out Evan Shelhamer et al.'s 2016 paper.

# LeNet-5

The LeNet-5 architecture is perhaps the most widely known CNN architecture. As mentioned earlier, it was created by Yann LeCun in 1998 and widely used for hand-written digit recognition (MNIST). It is composed of the layers shown in Table 13-1.

*Table 13-1. LeNet-5 architecture*

| Layer | Type | Maps | Size | Kernel size | Stride | Activation |
|-------|------|------|------|-------------|--------|------------|
| Out | Fully Connected | – | 10 | – | – | RBF |
| F6 | Fully Connected | – | 84 | – | – | tanh |
| C5 | Convolution | 120 | $1 \times 1$ | $5 \times 5$ | 1 | tanh |
| S4 | Avg Pooling | 16 | $5 \times 5$ | $2 \times 2$ | 2 | tanh |
| C3 | Convolution | 16 | $10 \times 10$ | $5 \times 5$ | 1 | tanh |
| S2 | Avg Pooling | 6 | $14 \times 14$ | $2 \times 2$ | 2 | tanh |
| C1 | Convolution | 6 | $28 \times 28$ | $5 \times 5$ | 1 | tanh |
| In | Input | 1 | $32 \times 32$ | – | – | – |

There are a few extra details to be noted:

- MNIST images are $28 \times 28$ pixels, but they are zero-padded to $32 \times 32$ pixels and normalized before being fed to the network. The rest of the network does not use any padding, which is why the size keeps shrinking as the image progresses through the network.

- The average pooling layers are slightly more complex than usual: each neuron computes the mean of its inputs, then multiplies the result by a learnable coefficient (one per map) and adds a learnable bias term (again, one per map), then finally applies the activation function.

- Most neurons in C3 maps are connected to neurons in only three or four S2 maps (instead of all six S2 maps). See table 1 in the original paper for details.

- The output layer is a bit special: instead of computing the dot product of the inputs and the weight vector, each neuron outputs the square of the Euclidian distance between its input vector and its weight vector. Each output measures how much the image belongs to a particular digit class. The cross entropy cost function is now preferred, as it penalizes bad predictions much more, producing larger gradients and thus converging faster.

Yann LeCun's website ("LENET" section) features great demos of LeNet-5 classifying digits.

## AlexNet

The *AlexNet* CNN architecture[9] won the 2012 ImageNet ILSVRC challenge by a large margin: it achieved 17% top-5 error rate while the second best achieved only 26%! It was developed by Alex Krizhevsky (hence the name), Ilya Sutskever, and Geoffrey Hinton. It is quite similar to LeNet-5, only much larger and deeper, and it was the first to stack convolutional layers directly on top of each other, instead of stacking a pooling layer on top of each convolutional layer. Table 13-2 presents this architecture.

*Table 13-2. AlexNet architecture*

| Layer | Type | Maps | Size | Kernel size | Stride | Padding | Activation |
|---|---|---|---|---|---|---|---|
| Out | Fully Connected | – | 1,000 | – | – | – | Softmax |
| F9 | Fully Connected | – | 4,096 | – | – | – | ReLU |
| F8 | Fully Connected | – | 4,096 | – | – | – | ReLU |
| C7 | Convolution | 256 | $13 \times 13$ | $3 \times 3$ | 1 | SAME | ReLU |
| C6 | Convolution | 384 | $13 \times 13$ | $3 \times 3$ | 1 | SAME | ReLU |
| C5 | Convolution | 384 | $13 \times 13$ | $3 \times 3$ | 1 | SAME | ReLU |
| S4 | Max Pooling | 256 | $13 \times 13$ | $3 \times 3$ | 2 | VALID | – |
| C3 | Convolution | 256 | $27 \times 27$ | $5 \times 5$ | 1 | SAME | ReLU |
| S2 | Max Pooling | 96 | $27 \times 27$ | $3 \times 3$ | 2 | VALID | – |
| C1 | Convolution | 96 | $55 \times 55$ | $11 \times 11$ | 4 | SAME | ReLU |
| In | Input | 3 (RGB) | $224 \times 224$ | – | – | – | – |

To reduce overfitting, the authors used two regularization techniques we discussed in previous chapters: first they applied dropout (with a 50% dropout rate) during training to the outputs of layers F8 and F9. Second, they performed data augmentation by

---

9 "ImageNet Classification with Deep Convolutional Neural Networks," A. Krizhevsky et al. (2012).