# Cheat Sheet – Convolutional Neural Network

## Convolutional Neural Network:

The data gets into the CNN through the input layer and passes through various hidden layers before getting to the output layer. The output of the network is compared to the actual labels in terms of loss or error. The partial derivatives of this loss w.r.t the trainable weights are calculated, and the weights are updated through one of the various methods using backpropagation.
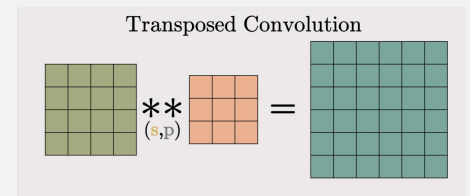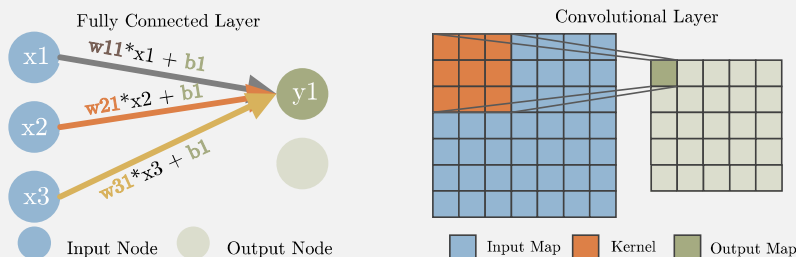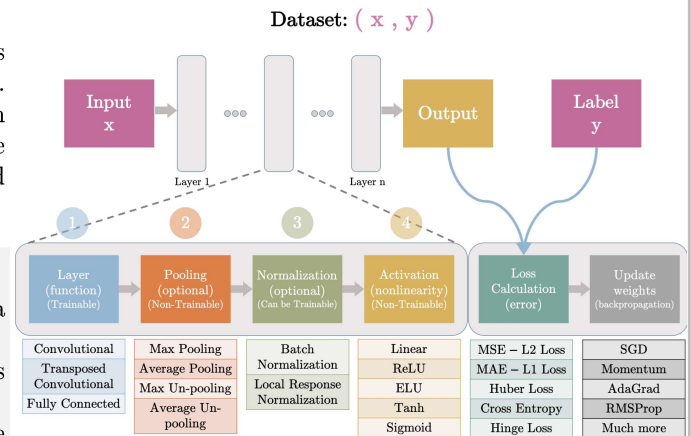
Dataset: ( x , y )



## CNN Template:

Most of the commonly used hidden layers (not all) follow a pattern

**1. Layer function:** Basic transforming function such as convolutional or fully connected layer.

**a. Fully Connected:** Linear functions between the input and the output.

**a. Convolutional Layers:** These layers are applied to 2D (3D) input feature maps. The trainable weights are a 2D (3D) kernel/filter that moves across the input feature map, generating dot products with the overlapping region of the input feature map.

**b. Transposed Convolutional (DeConvolutional) Layer:** Usually used to increase the size of the output feature map (Upsampling) The idea behind the transposed convolutional layer is to undo (not exactly) the convolutional layer



**2. Pooling:** Non-trainable layer to change the size of the feature map

**a. Max/Average Pooling:** Decrease the spatial size of the input layer based on selecting the maximum/average value in receptive field defined by the kernel

**b. UnPooling:** A non-trainable layer used to increase the spatial size of the input layer based on placing the input pixel at a certain index in the receptive field of the output defined by the kernel.



**3. Normalization:** Usually used just before the activation functions to limit the unbounded activation from increasing the output layer values too high

**a. Local Response Normalization LRN:** A **non-trainable layer** that square-normalizes the pixel values in a feature map within a local neighborhood.

**b. Batch Normalization:** A trainable approach to normalizing the data by learning scale and shift variable during training.

**3. Activation:** Introduce non-linearity so CNN can efficiently map non-linear complex mapping.

**a. Non-parametric/Static functions:** Linear, ReLU
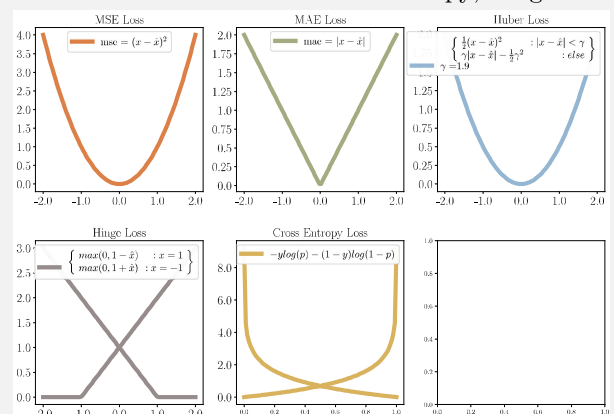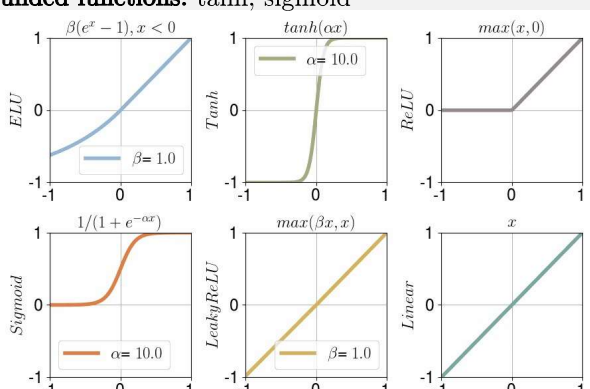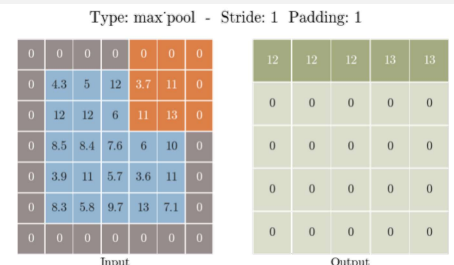
**b. Parametric functions:** ELU, tanh, sigmoid, Leaky ReLU

**c. Bounded functions:** tanh, sigmoid

**5. Loss function:** Quantifies how far off the CNN prediction is from the actual labels.

**a.** Regression Loss Functions: MAE, MSE, Huber loss

**b.** Classification Loss Functions: Cross entropy, Hinge loss





Source: https://www.cheatsheets.aqeel-anwar.com  Tutorial: Click here

# Cheat Sheet – Famous CNNs

## AlexNet – 2012

**Why:** AlexNet was born out of the need to improve the results of the ImageNet challenge.

**What:** The network consists of 5 Convolutional (CONV) layers and 3 Fully Connected (FC) layers. The activation used is the Rectified Linear Unit (ReLU).

**How:** Data augmentation is carried out to reduce over-fitting, Uses Local response localization.

## VGGNet – 2014

**Why:** VGGNet was born out of the need to reduce the # of parameters in the CONV layers and improve on training time

**What:** There are multiple variants of VGGNet (VGG16, VGG19, etc.)

**How:** The important point to note here is that all the conv kernels are of size 3x3 and maxpool kernels are of size 2x2 with a stride of two.

## ResNet – 2015

**Why:** Neural Networks are notorious for not being able to find a simpler mapping when it exists. ResNet solves that.

**What:** There are multiple versions of ResNetXX architectures where 'XX' denotes the number of layers. The most used ones are ResNet50 and ResNet101. Since the vanishing gradient problem was taken care of (more about it in the How part), CNN started to get deeper and deeper

**How:** ResNet architecture makes use of shortcut connections do solve the vanishing gradient problem. The basic building block of ResNet is a Residual block that is repeated throughout the network.
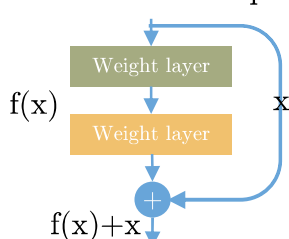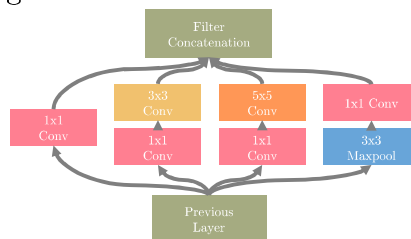
Figure 1 ResNet Block

Figure 2 Inception Block

## Inception – 2014

**Why:** Lager kernels are preferred for more global features, on the other hand, smaller kernels provide good results in detecting area-specific features. For effective recognition of such a variable-sized feature, we need kernels of different sizes. That is what Inception does.

**What:** The Inception network architecture consists of several inception modules of the following structure. Each inception module consists of four operations in parallel, 1x1 conv layer, 3x3 conv layer, 5x5 conv layer, max pooling

**How:** Inception increases the network space from which the best network is to be chosen via training. Each inception module can capture salient features at different levels.

### AlexNet Network - Structural Details

| Input | | | Output | | | Layer | Stride | Pad | Kernel size | | in | out | # of Param |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 227 | 227 | 3 | 55 | 55 | 96 | conv1 | 4 | 0 | 11 | 11 | 3 | 96 | 34944 |
| 55 | 55 | 96 | 27 | 27 | 96 | maxpool1 | 2 | 0 | 3 | 3 | 96 | 96 | 0 |
| 27 | 27 | 96 | 27 | 27 | 256 | conv2 | 1 | 2 | 5 | 5 | 96 | 256 | 614656 |
| 27 | 27 | 256 | 13 | 13 | 256 | maxpool2 | 2 | 0 | 3 | 3 | 256 | 256 | 0 |
| 13 | 13 | 256 | 13 | 13 | 384 | conv3 | 1 | 1 | 3 | 3 | 256 | 384 | 885120 |
| 13 | 13 | 384 | 13 | 13 | 384 | conv4 | 1 | 1 | 3 | 3 | 384 | 384 | 1327488 |
| 13 | 13 | 384 | 13 | 13 | 256 | conv5 | 1 | 1 | 3 | 3 | 384 | 256 | 884992 |
| 13 | 13 | 256 | 6 | 6 | 256 | maxpool5 | 2 | 0 | 3 | 3 | 256 | 256 | 0 |
| | | | | | | fc6 | | | 1 | 1 | 9216 | 4096 | 37752832 |
| | | | | | | fc7 | | | 1 | 1 | 4096 | 4096 | 16781312 |
| | | | | | | fc8 | | | 1 | 1 | 4096 | 1000 | 4097000 |
| | | | | | Total | | | | | | | | 62,378,344 |

### VGG16 - Structural Details

| # | Input Image | | | output | | | Layer | Stride | Kernel | | in | out | Param |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 224 | 224 | 3 | 224 | 224 | 64 | conv3-64 | 1 | 3 | 3 | 3 | 64 | 1792 |
| 2 | 224 | 224 | 64 | 224 | 224 | 64 | conv3064 | 1 | 3 | 3 | 64 | 64 | 36928 |
| | 224 | 224 | 64 | 112 | 112 | 64 | maxpool | 2 | 2 | 2 | 64 | 64 | 0 |
| 3 | 112 | 112 | 64 | 112 | 112 | 128 | conv3-128 | 1 | 3 | 3 | 64 | 128 | 73856 |
| 4 | 112 | 112 | 128 | 112 | 112 | 128 | conv3-128 | 1 | 3 | 3 | 128 | 128 | 147584 |
| | 112 | 112 | 128 | 56 | 56 | 128 | maxpool | 2 | 2 | 2 | 128 | 128 | 65664 |
| 5 | 56 | 56 | 128 | 56 | 56 | 256 | conv3-256 | 1 | 3 | 3 | 128 | 256 | 295168 |
| 6 | 56 | 56 | 256 | 56 | 56 | 256 | conv3-256 | 1 | 3 | 3 | 256 | 256 | 590080 |
| 7 | 56 | 56 | 256 | 56 | 56 | 256 | conv3-256 | 1 | 3 | 3 | 256 | 256 | 590080 |
| | 56 | 56 | 256 | 28 | 28 | 256 | maxpool | 2 | 2 | 2 | 256 | 256 | 0 |
| 8 | 28 | 28 | 256 | 28 | 28 | 512 | conv3-512 | 1 | 3 | 3 | 256 | 512 | 1180160 |
| 9 | 28 | 28 | 512 | 28 | 28 | 512 | conv3-512 | 1 | 3 | 3 | 512 | 512 | 2359808 |
| 10 | 28 | 28 | 512 | 28 | 28 | 512 | conv3-512 | 1 | 3 | 3 | 512 | 512 | 2359808 |
| | 28 | 28 | 512 | 14 | 14 | 512 | maxpool | 2 | 2 | 2 | 512 | 512 | 0 |
| 11 | 14 | 14 | 512 | 14 | 14 | 512 | conv3-512 | 1 | 3 | 3 | 512 | 512 | 2359808 |
| 12 | 14 | 14 | 512 | 14 | 14 | 512 | conv3-512 | 1 | 3 | 3 | 512 | 512 | 2359808 |
| 13 | 14 | 14 | 512 | 14 | 14 | 512 | conv3-512 | 1 | 3 | 3 | 512 | 512 | 2359808 |
| | 14 | 14 | 512 | 7 | 7 | 512 | maxpool | 2 | 2 | 2 | 512 | 512 | 0 |
| 14 | 1 | 1 | 25088 | 1 | 1 | 4096 | fc | | 1 | 1 | 25088 | 4096 | 102764544 |
| 15 | 1 | 1 | 4096 | 1 | 1 | 4096 | fc | | 1 | 1 | 4096 | 4096 | 16781312 |
| 16 | 1 | 1 | 4096 | 1 | 1 | 1000 | fc | | 1 | 1 | 4096 | 1000 | 4097000 |
| | | | | | | Total | | | | | | | 138,423,208 |

### ResNet18 - Structural Details

| # | Input Image | | | output | | | Layer | Stride | Pad | Kernel | | in | out | Param |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 227 | 227 | 3 | 112 | 112 | 64 | conv1 | 2 | 1 | 7 | 7 | 3 | 64 | 9472 |
| | 112 | 112 | 64 | 56 | 56 | 64 | maxpool | 2 | 0.5 | 3 | 3 | 64 | 64 | 0 |
| 2 | 56 | 56 | 64 | 56 | 56 | 64 | conv2-1 | 1 | 1 | 3 | 3 | 64 | 64 | 36928 |
| 3 | 56 | 56 | 64 | 56 | 56 | 64 | conv2-2 | 1 | 1 | 3 | 3 | 64 | 64 | 36928 |
| 4 | 56 | 56 | 64 | 56 | 56 | 64 | conv2-3 | 1 | 1 | 3 | 3 | 64 | 64 | 36928 |
| 5 | 56 | 56 | 64 | 56 | 56 | 64 | conv2-4 | 1 | 1 | 3 | 3 | 64 | 64 | 36928 |
| 6 | 56 | 56 | 64 | 28 | 28 | 128 | conv3-1 | 2 | 0.5 | 3 | 3 | 64 | 128 | 73856 |
| 7 | 28 | 28 | 128 | 28 | 28 | 128 | conv3-2 | 1 | 1 | 3 | 3 | 128 | 128 | 147584 |
| 8 | 28 | 28 | 128 | 28 | 28 | 128 | conv3-3 | 1 | 1 | 3 | 3 | 128 | 128 | 147584 |
| 9 | 28 | 28 | 128 | 28 | 28 | 128 | conv3-4 | 1 | 1 | 3 | 3 | 128 | 128 | 147584 |
| 10 | 28 | 28 | 128 | 14 | 14 | 256 | conv4-1 | 2 | 0.5 | 3 | 3 | 128 | 256 | 295168 |
| 11 | 14 | 14 | 256 | 14 | 14 | 256 | conv4-2 | 1 | 1 | 3 | 3 | 256 | 256 | 590080 |
| 12 | 14 | 14 | 256 | 14 | 14 | 256 | conv4-3 | 1 | 1 | 3 | 3 | 256 | 256 | 590080 |
| 13 | 14 | 14 | 256 | 14 | 14 | 256 | conv4-4 | 1 | 1 | 3 | 3 | 256 | 256 | 590080 |
| 14 | 14 | 14 | 256 | 7 | 7 | 512 | conv5-1 | 2 | 0.5 | 3 | 3 | 256 | 512 | 1180160 |
| 15 | 7 | 7 | 512 | 7 | 7 | 512 | conv5-2 | 1 | 1 | 3 | 3 | 512 | 512 | 2359808 |
| 16 | 7 | 7 | 512 | 7 | 7 | 512 | conv5-3 | 1 | 1 | 3 | 3 | 512 | 512 | 2359808 |
| 17 | 7 | 7 | 512 | 7 | 7 | 512 | conv5-4 | 1 | 1 | 3 | 3 | 512 | 512 | 2359808 |
| | 7 | 7 | 512 | 1 | 1 | 512 | avg pool | 7 | 0 | 7 | 7 | 512 | 512 | 0 |
| 18 | 1 | 1 | 512 | 1 | 1 | 1000 | fc | | | | | 512 | 1000 | 513000 |
| | | | | | | Total | | | | | | | | 11,511,784 |

### GoogLeNet - Structural Details

| | Input Image | | | output | | | Layer | Input Layer | Stride | Pad | Kernel | | in | out | Param |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 227 | 227 | 3 | 112 | 112 | 64 | conv1 | input | 2 | 1 | 7 | 7 | 3 | 64 | 9472 |
| | 56 | 56 | 64 | 56 | 56 | 64 | conv1x1 | maxpool1 | 2 | 0.5 | 3 | 3 | 64 | 64 | 4160 |
| | 56 | 56 | 64 | 56 | 56 | 192 | conv2-1 | | 1 | 1 | 3 | 3 | 64 | 192 | 110784 |
| | 56 | 56 | 192 | 28 | 28 | 192 | maxpool2 | | 2 | 0.5 | 3 | 3 | 192 | 192 | 0 |
| inception (3a) | 28 | 28 | 192 | 28 | 28 | 96 | conv1x1a | maxpool2 | 1 | 0 | 1 | 1 | 192 | 96 | 18528 |
| | 28 | 28 | 96 | 28 | 28 | 16 | conv1x1b | maxpool2 | 1 | 0 | 1 | 1 | 192 | 16 | 3088 |
| | 28 | 28 | 192 | 28 | 28 | 192 | maxpool-a | maxpool2 | 1 | 1 | 3 | 3 | 192 | 192 | 0 |
| | 28 | 28 | 192 | 28 | 28 | 64 | conv1x1c | maxpool2 | 1 | 0 | 1 | 1 | 192 | 64 | 12352 |
| | 28 | 28 | 96 | 28 | 28 | 128 | conv3-3 | conv1x1a | 1 | 1 | 3 | 3 | 96 | 128 | 12832 |
| | 28 | 28 | 16 | 28 | 28 | 32 | conv5x5 | conv1x1b | 1 | 2 | 5 | 5 | 16 | 32 | 12832 |
| | 28 | 28 | 192 | 28 | 28 | 32 | conv1x1d | maxpool-a | 1 | 0 | 1 | 1 | 192 | 32 | 6176 |
| | 28 | 28 | 256 | depth-concat | | | | | | | | | | | |
| inception (3b) | 28 | 28 | 256 | 28 | 28 | 128 | conv1x1a | depth-concat | 1 | 0 | 1 | 1 | 256 | 128 | 32896 |
| | 28 | 28 | 128 | 28 | 28 | 32 | conv1x1b | depth-concat | 1 | 0 | 1 | 1 | 256 | 32 | 8224 |
| | 28 | 28 | 256 | 28 | 28 | 256 | maxpool-a | depth-concat | 1 | 1 | 3 | 3 | 256 | 256 | 0 |
| | 28 | 28 | 192 | 28 | 28 | 128 | conv1x1c | depth-concat | 1 | 0 | 1 | 1 | 256 | 128 | 32896 |
| | 28 | 28 | 96 | 28 | 28 | 192 | conv3-3 | conv1x1a | 1 | 1 | 3 | 3 | 128 | 192 | 221376 |
| | 28 | 28 | 16 | 28 | 28 | 96 | conv5x5 | conv1x1b | 1 | 2 | 5 | 5 | 32 | 96 | 76896 |
| | 28 | 28 | 192 | 28 | 28 | 64 | conv1x1d | maxpool-a | 1 | 0 | 1 | 1 | 256 | 64 | 16448 |
| | 28 | 28 | 480 | depth-concat | | | | | | | | | | | |
| | 28 | 28 | 480 | 14 | 14 | 480 | maxpool3 | depth-concat | 2 | 0.5 | 3 | 3 | 480 | 480 | 0 |
| inception (4a) | 14 | 14 | 480 | 14 | 14 | 96 | conv1x1a | maxpool3 | 1 | 0 | 1 | 1 | 480 | 96 | 46176 |
| | 14 | 14 | 480 | 14 | 14 | 16 | conv1x1b | maxpool3 | 1 | 0 | 1 | 1 | 480 | 16 | 7696 |
| | 14 | 14 | 480 | 14 | 14 | 480 | maxpool-a | maxpool3 | 1 | 1 | 3 | 3 | 480 | 480 | 0 |
| | 14 | 14 | 480 | 14 | 14 | 192 | conv1x1c | maxpool3 | 1 | 0 | 1 | 1 | 480 | 192 | 92352 |
| | 14 | 14 | 96 | 14 | 14 | 208 | conv3-3 | conv1x1a | 1 | 1 | 3 | 3 | 96 | 208 | 179920 |
| | 14 | 14 | 16 | 14 | 14 | 48 | conv5x5 | conv1x1b | 1 | 2 | 5 | 5 | 16 | 48 | 19248 |
| | 14 | 14 | 192 | 14 | 14 | 64 | conv1x1d | maxpool-a | 1 | 0 | 1 | 1 | 480 | 64 | 30784 |
| | 14 | 14 | 512 | depth-concat | | | | | | | | | | | |
| inception (4b) | 14 | 14 | 512 | 14 | 14 | 112 | conv1x1a | depth-concat | 1 | 0 | 1 | 1 | 512 | 112 | 57456 |
| | 14 | 14 | 512 | 14 | 14 | 24 | conv1x1b | depth-concat | 1 | 0 | 1 | 1 | 64 | 24 | 1560 |
| | 14 | 14 | 512 | 14 | 14 | 64 | maxpool-a | depth-concat | 1 | 1 | 3 | 3 | 64 | 64 | 0 |
| | 14 | 14 | 512 | 14 | 14 | 160 | conv1x1c | depth-concat | 1 | 0 | 1 | 1 | 64 | 160 | 10400 |
| | 14 | 14 | 16 | 14 | 14 | 224 | conv3-3 | conv1x1a | 1 | 1 | 3 | 3 | 112 | 224 | 226016 |
| | 14 | 14 | 16 | 14 | 14 | 64 | conv5x5 | conv1x1b | 1 | 2 | 5 | 5 | 24 | 64 | 38464 |
| | 14 | 14 | 160 | 14 | 14 | 64 | conv1x1d | maxpool-a | 1 | 0 | 1 | 1 | 64 | 64 | 4160 |
| | 14 | 14 | 512 | depth-concat | | | | | | | | | | | |
| inception (4c) | 14 | 14 | 512 | 14 | 14 | 128 | conv1x1a | depth-concat | 1 | 0 | 1 | 1 | 512 | 128 | 65664 |
| | 14 | 14 | 512 | 14 | 14 | 24 | conv1x1b | depth-concat | 1 | 0 | 1 | 1 | 64 | 24 | 1560 |
| | 14 | 14 | 512 | 14 | 14 | 64 | maxpool-a | depth-concat | 1 | 1 | 3 | 3 | 64 | 64 | 0 |
| | 14 | 14 | 512 | 14 | 14 | 128 | conv1x1c | depth-concat | 1 | 0 | 1 | 1 | 64 | 128 | 8320 |
| | 14 | 14 | 96 | 14 | 14 | 256 | conv3-3 | conv1x1a | 1 | 1 | 3 | 3 | 128 | 256 | 38464 |
| | 14 | 14 | 16 | 14 | 14 | 64 | conv5x5 | conv1x1b | 1 | 2 | 5 | 5 | 24 | 64 | 38464 |
| | 14 | 14 | 128 | 14 | 14 | 64 | conv1x1d | maxpool-a | 1 | 0 | 1 | 1 | 64 | 64 | 4160 |
| | 14 | 14 | 512 | depth-concat | | | | | | | | | | | |
| inception (4d) | 14 | 14 | 512 | 14 | 14 | 144 | conv1x1a | depth-concat | 1 | 0 | 1 | 1 | 512 | 144 | 73872 |
| | 14 | 14 | 512 | 14 | 14 | 32 | conv1x1b | depth-concat | 1 | 0 | 1 | 1 | 64 | 32 | 2080 |
| | 14 | 14 | 512 | 14 | 14 | 64 | maxpool-a | depth-concat | 1 | 1 | 3 | 3 | 64 | 64 | 0 |
| | 14 | 14 | 512 | 14 | 14 | 112 | conv1x1c | depth-concat | 1 | 0 | 1 | 1 | 64 | 112 | 7280 |
| | 14 | 14 | 96 | 14 | 14 | 288 | conv3-3 | conv1x1a | 1 | 1 | 3 | 3 | 144 | 288 | 373536 |
| | 14 | 14 | 16 | 14 | 14 | 64 | conv5x5 | conv1x1b | 1 | 2 | 5 | 5 | 32 | 64 | 51264 |
| | 14 | 14 | 112 | 14 | 14 | 64 | conv1x1d | maxpool-a | 1 | 0 | 1 | 1 | 64 | 64 | 4160 |
| | 14 | 14 | 528 | depth-concat | | | | | | | | | | | |
| inception (4e) | 14 | 14 | 528 | 14 | 14 | 160 | conv1x1a | depth-concat | 1 | 0 | 1 | 1 | 528 | 160 | 84640 |
| | 14 | 14 | 528 | 14 | 14 | 32 | conv1x1b | depth-concat | 1 | 0 | 1 | 1 | 64 | 32 | 2080 |
| | 14 | 14 | 528 | 14 | 14 | 64 | maxpool-a | depth-concat | 1 | 1 | 3 | 3 | 64 | 64 | 0 |
| | 14 | 14 | 528 | 14 | 14 | 256 | conv1x1c | depth-concat | 1 | 0 | 1 | 1 | 64 | 256 | 16640 |
| | 14 | 14 | 96 | 14 | 14 | 320 | conv3-3 | conv1x1a | 1 | 1 | 3 | 3 | 160 | 320 | 461120 |
| | 14 | 14 | 16 | 14 | 14 | 128 | conv5x5 | conv1x1b | 1 | 2 | 5 | 5 | 32 | 128 | 102528 |
| | 14 | 14 | 256 | 14 | 14 | 128 | conv1x1d | maxpool-a | 1 | 0 | 1 | 1 | 64 | 128 | 8320 |
| | 14 | 14 | 832 | depth-concat | | | | | | | | | | | |
| | 14 | 14 | 832 | 7 | 7 | 832 | maxpool4 | depth-concat | 2 | 0.5 | 3 | 3 | 832 | 832 | 0 |
| inception (5a) | 7 | 7 | 832 | 7 | 7 | 160 | conv1x1a | maxpool4 | 1 | 0 | 1 | 1 | 832 | 160 | 133280 |
| | 7 | 7 | 832 | 7 | 7 | 32 | conv1x1b | maxpool4 | 1 | 0 | 1 | 1 | 832 | 32 | 26656 |
| | 7 | 7 | 832 | 7 | 7 | 832 | maxpool-a | maxpool4 | 1 | 1 | 3 | 3 | 832 | 832 | 0 |
| | 7 | 7 | 832 | 7 | 7 | 256 | conv1x1c | maxpool4 | 1 | 0 | 1 | 1 | 832 | 256 | 213248 |
| | 7 | 7 | 96 | 7 | 7 | 320 | conv3-3 | conv1x1a | 1 | 1 | 3 | 3 | 160 | 320 | 461120 |
| | 7 | 7 | 16 | 7 | 7 | 128 | conv5x5 | conv1x1b | 1 | 2 | 5 | 5 | 32 | 128 | 102528 |
| | 7 | 7 | 256 | 7 | 7 | 128 | conv1x1d | maxpool-a | 1 | 0 | 1 | 1 | 64 | 128 | 106624 |
| | 7 | 7 | 832 | depth-concat | | | | | | | | | | | |
| inception (5b) | 7 | 7 | 832 | 7 | 7 | 192 | conv1x1a | depth-concat | 1 | 0 | 1 | 1 | 832 | 192 | 159936 |
| | 7 | 7 | 832 | 7 | 7 | 48 | conv1x1b | depth-concat | 1 | 0 | 1 | 1 | 832 | 48 | 39984 |
| | 7 | 7 | 832 | 7 | 7 | 832 | maxpool-a | depth-concat | 1 | 1 | 3 | 3 | 832 | 832 | 0 |
| | 7 | 7 | 832 | 7 | 7 | 384 | conv1x1c | depth-concat | 1 | 0 | 1 | 1 | 832 | 384 | 319872 |
| | 7 | 7 | 96 | 7 | 7 | 384 | conv3-3 | conv1x1a | 1 | 1 | 3 | 3 | 192 | 384 | 663936 |
| | 7 | 7 | 16 | 7 | 7 | 128 | conv5x5 | conv1x1b | 1 | 2 | 5 | 5 | 48 | 128 | 153728 |
| | 7 | 7 | 384 | 7 | 7 | 128 | conv1x1d | maxpool-a | 1 | 0 | 1 | 1 | 64 | 128 | 16512 |
| | 7 | 7 | 1024 | depth-concat | | | | | | | | | | | |
| | 7 | 7 | 1024 | 1 | 1 | 1024 | avgpool | depth-concat | 1 | 0 | 7 | 7 | 1024 | 1024 | 0 |
| | 1 | 1 | 1024 | 1 | 1 | 1000 | fc | depth-concat | 1 | 0 | 1 | 1 | 1024 | 1000 | 1025000 |
| | | | | | | Total | | | | | | | | | 6,414,360 |

### Comparison

| Network | Year | Salient Feature | top5 accuracy | Parameters | FLOP |
|---|---|---|---|---|---|
| AlexNet | 2012 | Deeper | 84.70% | 62M | 1.5B |
| VGGNet | 2014 | Fixed-size kernels | 92.30% | 138M | 19.6B |
| Inception | 2014 | Wider - Parallel kernels | 93.30% | 6.4M | 2B |
| ResNet-152 | 2015 | Shortcut connections | 95.51% | 60.3M | 11B |