

MobileNet v2 Architecture

- MobileNet-v2 is a convolutional neural network that is 92 layers deep. The network has pretrained version which trained on more than a million images from the ImageNet database. The pretrained network can classify images into 1000 object categories, such as keyboard, mouse, pencil, and many animals. As a result, the network has learned rich feature representations for a wide range of images. The network has an image input size of 224-by-224.
- It's one of the applications of Depth-wise Separable Convolutional Neural Networks

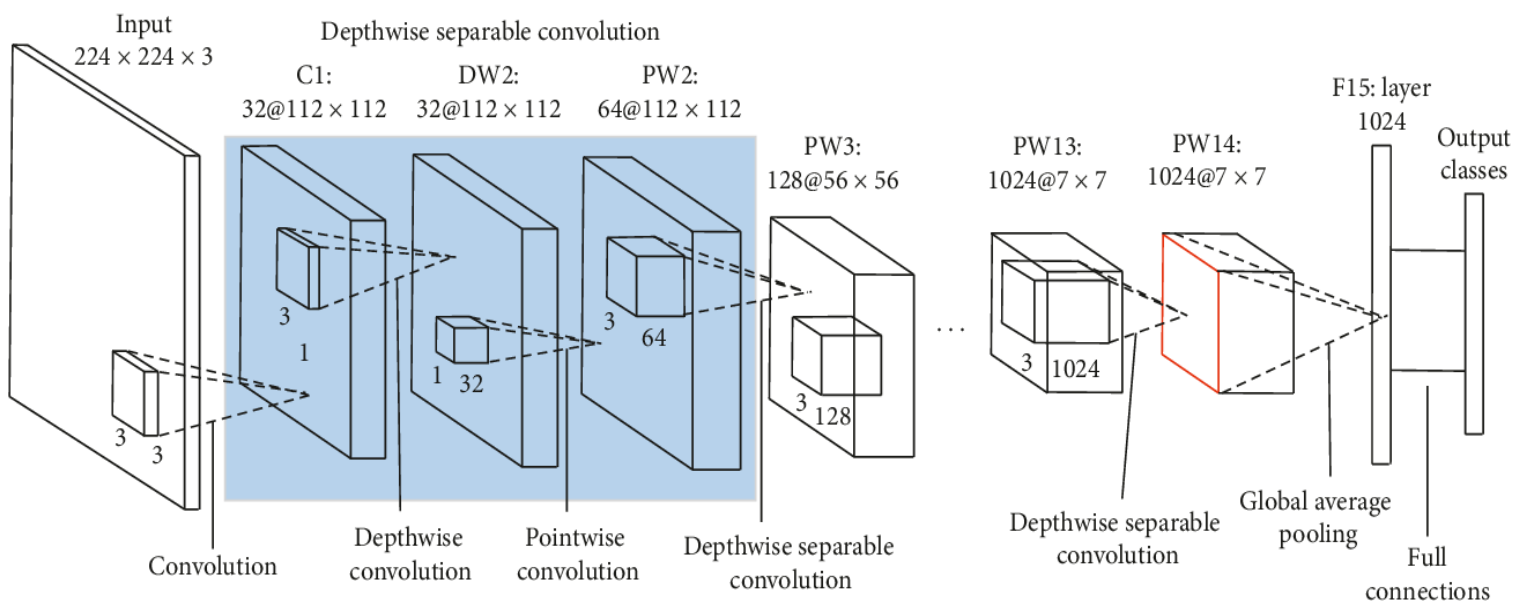
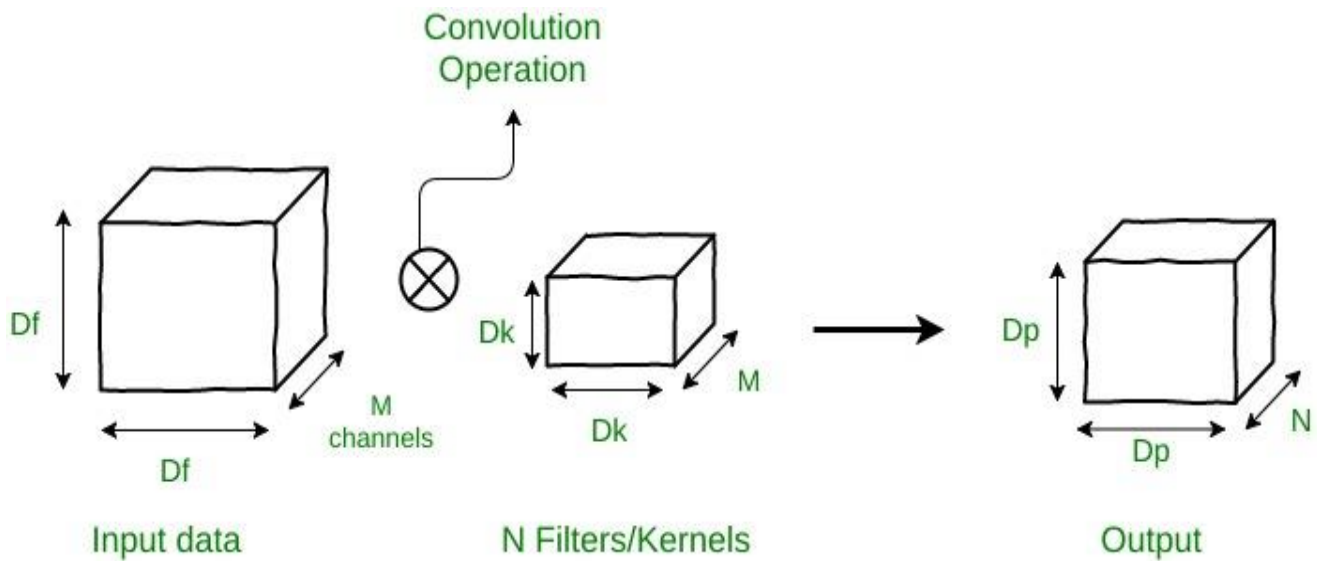


FIGURE 1: Architecture of MobileNet.

Understanding Normal Convolution operation



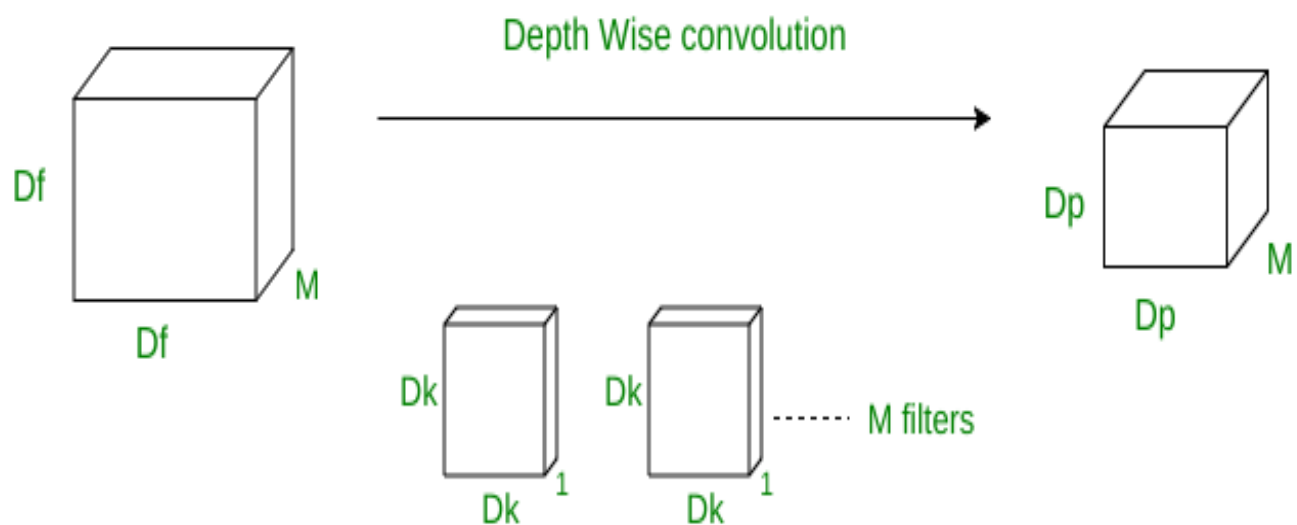
- Suppose there is an input data of size $D_f \times D_f \times M$, where $D_f \times D_f$ can be the image size and M is the number of channels (3 for an RGB image). Suppose there are N filters/kernels of size $D_k \times D_k \times M$. If a normal convolution operation is done, then, the output size will be $D_p \times D_p \times N$.
- *The number of multiplications in 1 convolution operation = size of filter = $D_k \times D_k \times M$*
- *the total number of multiplications become $N \times D_p \times D_p \times$ (Multiplications per convolution)*
- *Total no of multiplications = $N \times D_p^2 \times D_k^2 \times M$*

Depth-Wise Separable Convolutions

- This process is broken down into 2 operations :
 - Depth-wise convolutions
 - Point-wise convolutions

1-DEPTH WISE CONVOLUTION :

In **depth-wise operation**, convolution is applied to a **single channel** at a time unlike standard CNN's in which it is done for all the M channels. So here the filters/kernels will be of size **$D_k \times D_k \times 1$** . Given there are M channels in the input data, then M such filters are required. Output will be of size **$D_p \times D_p \times M$** .

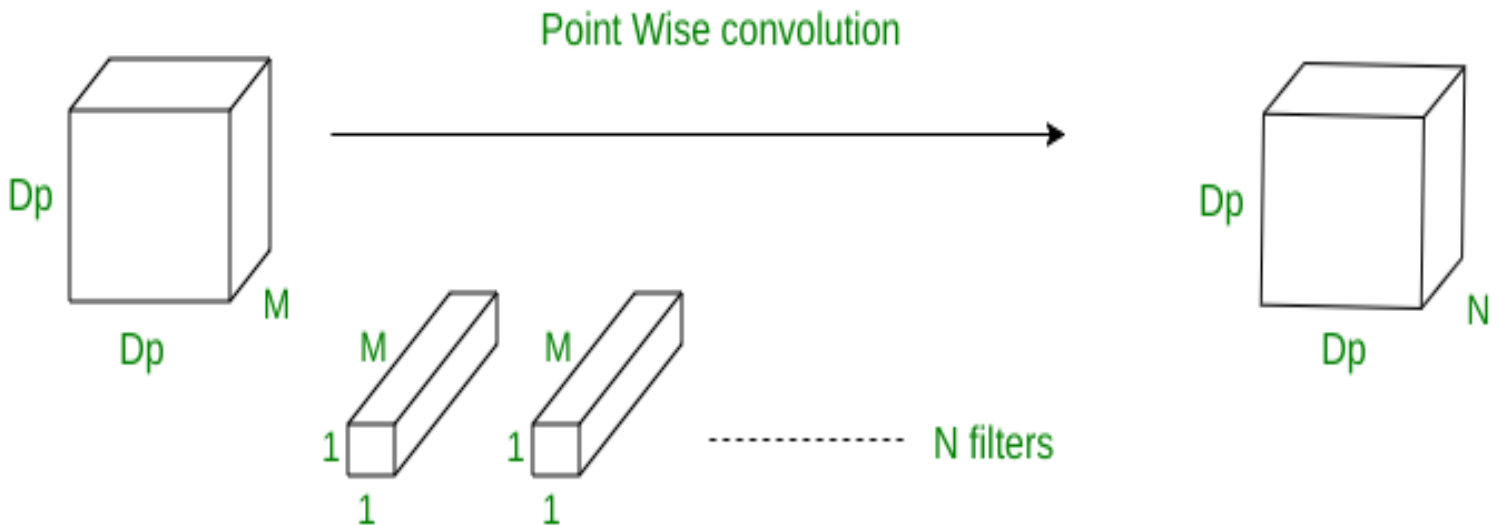


Cost of this operation:

- A single convolution operation require **$D_k \times D_k$** multiplications.
- Since the filter are slided by **$D_p \times D_p$** times across all the M channels, the total number of multiplications is equal to **$M \times D_p \times D_p \times D_k \times D_k$**
- So for depth wise convolution operation, **Total no of multiplications = $M \times D_k^2 \times D_p^2$**

2-POINT WISE CONVOLUTION :

In **point-wise operation**, a **1x1 convolution** operation is applied on the M channels. So the filter size for this operation will be **1 x 1 x M**. Say we use N such filters, the output size becomes **$D_p \times D_p \times N$** .



Cost of this operation:

- A single convolution operation require **1 x M** multiplications.
- Since the filter is being slided by **$D_p \times D_p$** times, the total number of multiplications is equal to **$M \times D_p \times D_p \times (\text{no. of filters})$**
- So for point wise convolution operation, **Total no of multiplications = $M \times D_p^2 \times N$**

Therefore, for overall operation:

- *Total multiplications = Depth wise conv. multiplications + Point wise conv. multiplications*
- *Total multiplications = $M * Dk^2 * Dp^2 + M * Dp^2 * N = M * Dp^2 * (Dk^2 + n)$*
- *So for depth wise separable convolution operation, Total no of multiplications = $M \times Dp^2 \times (Dk^2 + N)$*

Comparison between the complexities of these types of convolution operations

<i>Type of Convolution</i>	<i>Complexity</i>
<i>Standard</i>	<i>$N \times Dp^2 \times Dg^2 \times M$</i>
<i>Depth wise separable</i>	<i>$M \times Dp^2 \times (Dk^2 + N)$</i>

$$\frac{\text{Complexity of depth wise separable convolutions}}{\text{Complexity of standard convolution}} = \text{RATIO (R)}$$

$$\text{Ratio(R)} = 1/N + 1/Dk^2$$

As an example, consider $N = 100$ and $Dk = 512$. Then the ratio $R = 0.010004$

This means that the depth wise separable convolution network, in this example, performs 100 times lesser multiplications as compared to a standard constitutional neural network.

This implies that we can deploy faster convolution neural network models without losing much of the accuracy.