

# MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications

Andrew G. Howard Menglong Zhu Bo Chen Dmitry Kalenichenko  
Weijun Wang Tobias Weyand Marco Andreetto Hartwig Adam  
Google Inc.

Speaker : Wade  
Date : 2018/08/23

# Outline

- Relate work
- Main idea
- Experiments

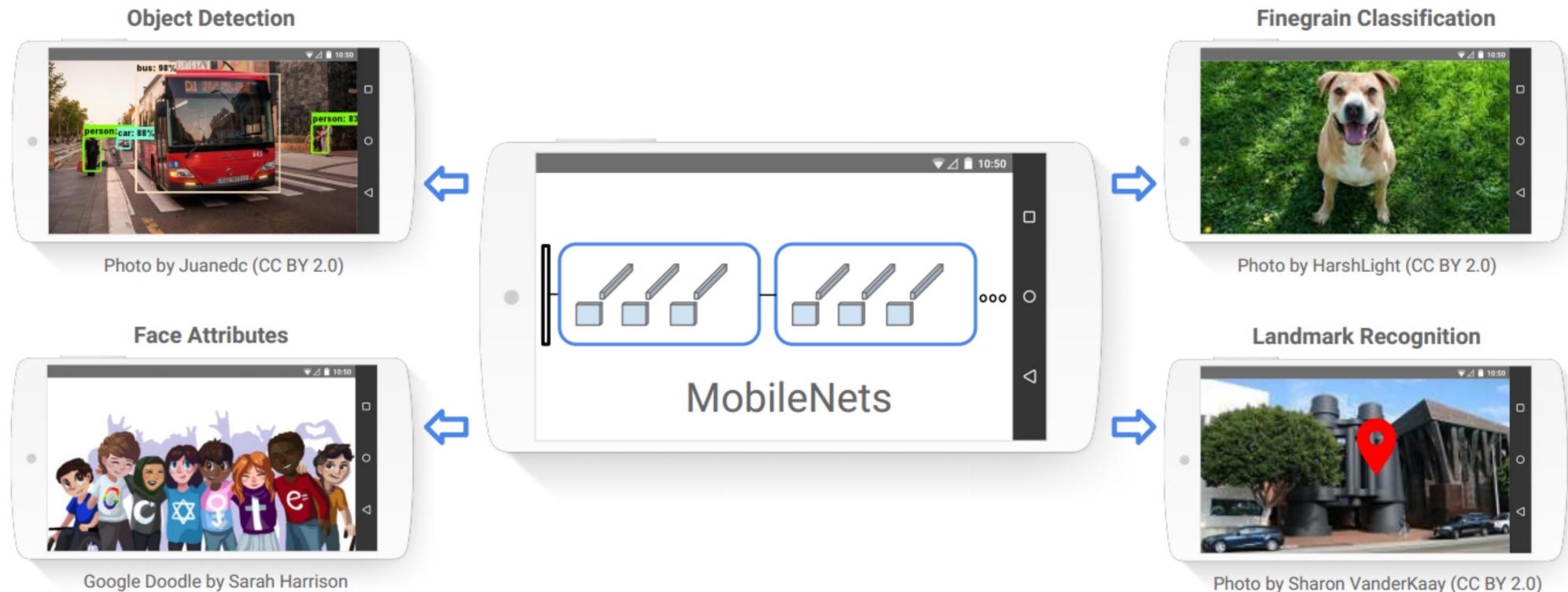


Figure 1. MobileNet models can be applied to various recognition tasks for efficient on device intelligence.

# What's the “Right” Neural Network for Use in a Gadget?

- Sufficiently high accuracy
- Low computational complexity
- Low energy usage
- Small model size

# Related Work

- Quantization, pruning, decomposition and distillation
- Small network, SqueezeNet, Xception network

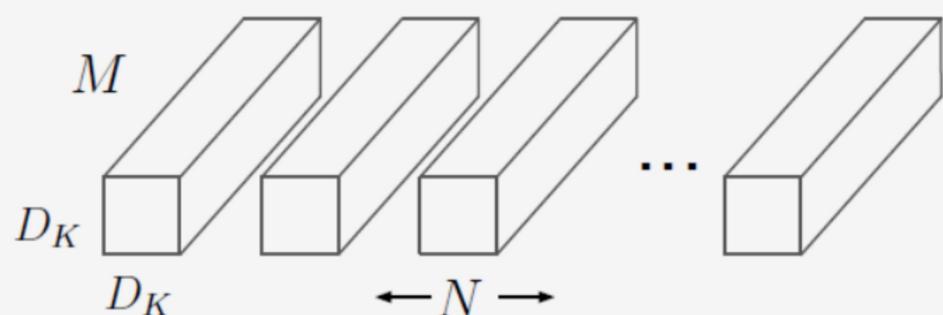
# Techniques for Small Deep Neural Networks

- Remove Fully-Connected Layers
- Kernel Reduction (  $3 \times 3 \rightarrow 1 \times 1$  )
- Channel Reduction
- Depthwise Separable Convolutions

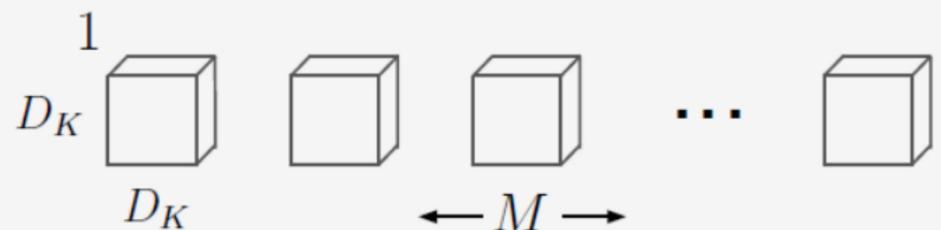
- Key Idea : Depthwise Separable Convolution!

- The MobileNet model is based on depthwise separable convolutions which is a form of factorized convolutions which factorize a standard convolution into a **depthwise convolution** and a  $1 \times 1$  convolution called a **pointwise convolution**

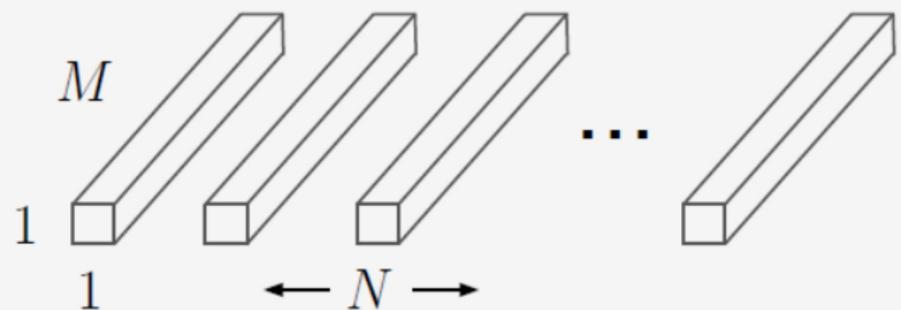
# Standard Convolution vs Depthwise Separable Convolution



(a) Standard Convolution Filters

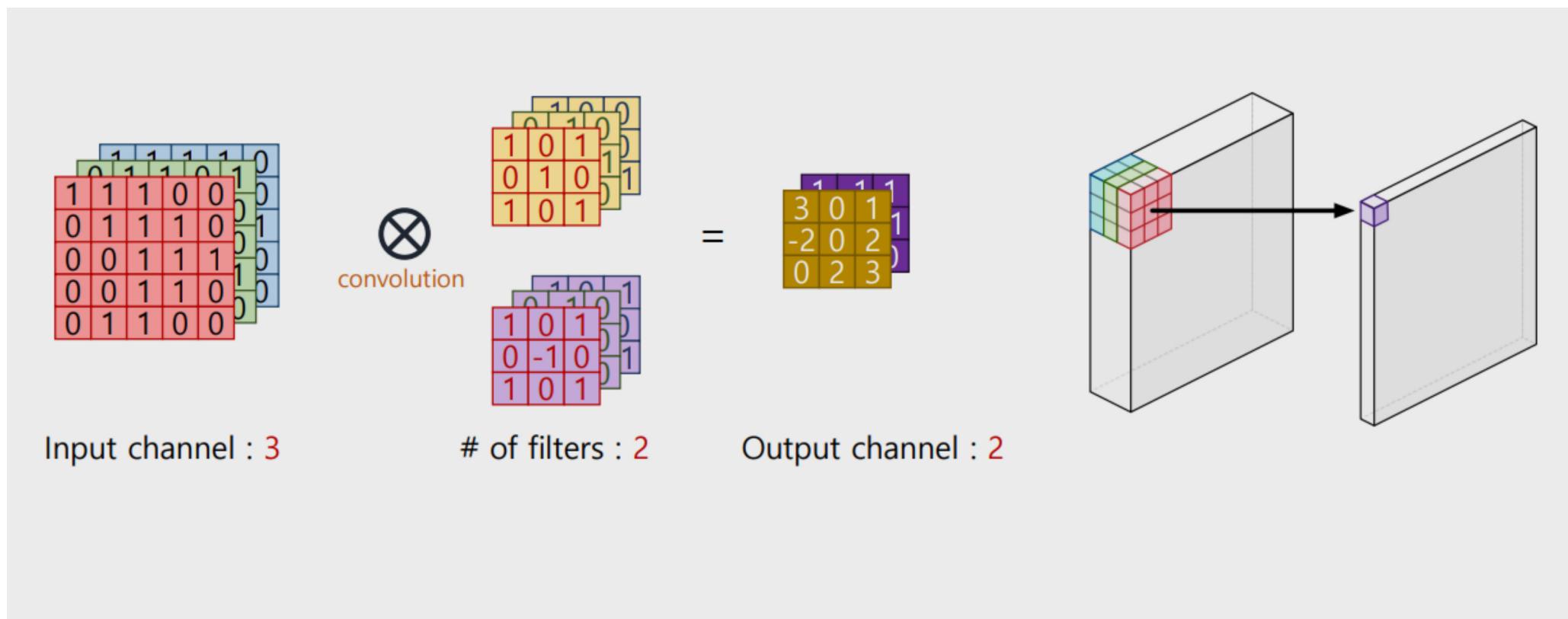


(b) Depthwise Convolutional Filters



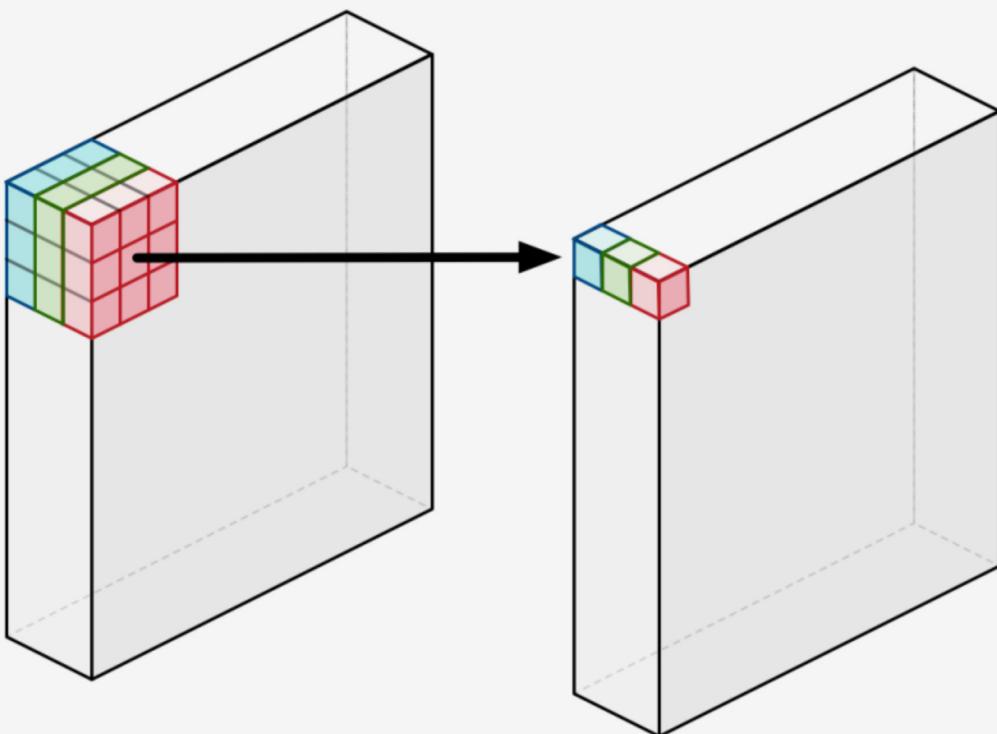
(c)  $1 \times 1$  Convolutional Filters called Pointwise Convolution in the context of Depthwise Separable Convolution

# Recap – Convolution Operation

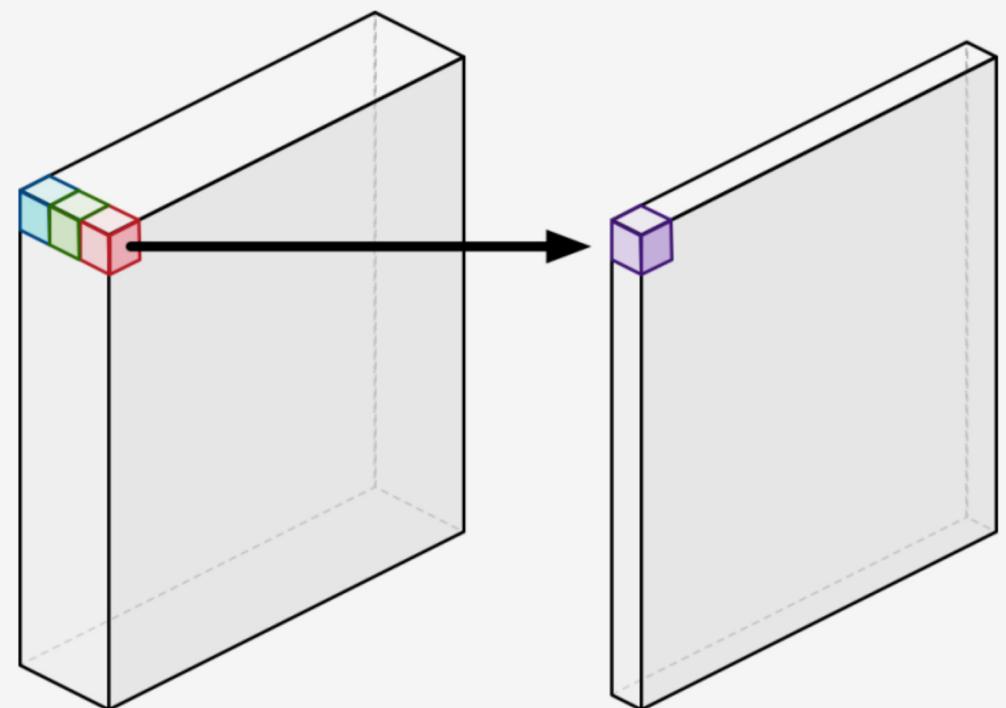


# Depthwise Separable Convolution

- Depthwise Convolution + Pointwise Convolution( $1 \times 1$  convolution)



**Depthwise convolution**



**Pointwise convolution**

Figures from <http://machinethink.net/blog/googles-mobile-net-architecture-on-iphone/>

# Standard Convolution vs Depthwise Separable Convolution

- Standard convolutions have the computational cost of
  - $D_K \times D_K \times M \times N \times D_F \times D_F$
- Depthwise separable convolutions cost
  - $D_K \times D_K \times M \times D_F \times D_F + M \times N \times D_F \times D_F$
- Reduction in computations
  - $1/N + 1/D_K^2$
  - If we use  $3 \times 3$  depthwise separable convolutions, we get between 8 to 9 times less computations

$D_K$  : width/height of filters

$D_F$  : width/height of feature maps

M : number of input channels

N : number of output channels(number of filters)

# Depthwise Separable Convolutions

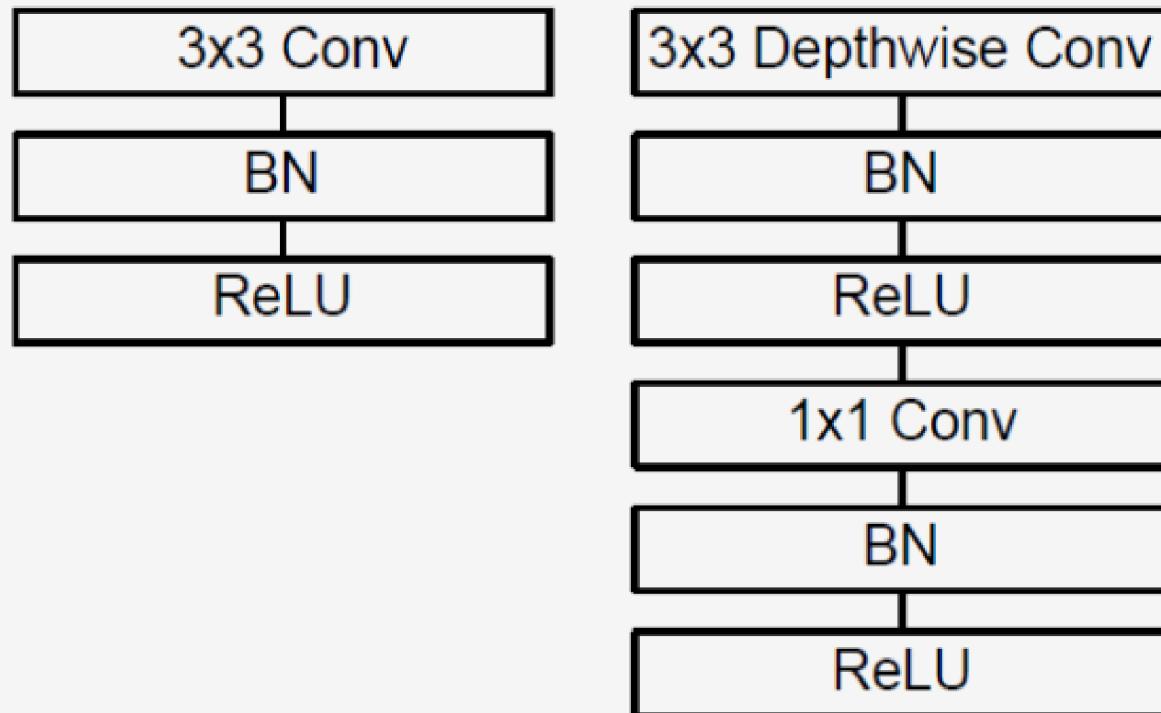


Figure 3. Left: Standard convolutional layer with batchnorm and ReLU. Right: Depthwise Separable convolutions with Depthwise and Pointwise layers followed by batchnorm and ReLU.

# Model Structure

Table 1. MobileNet Body Architecture

Type / Stride	Filter Shape	Input Size
Conv / s2	$3 \times 3 \times 3 \times 32$	$224 \times 224 \times 3$
Conv dw / s1	$3 \times 3 \times 32$ dw	$112 \times 112 \times 32$
Conv / s1	$1 \times 1 \times 32 \times 64$	$112 \times 112 \times 32$
Conv dw / s2	$3 \times 3 \times 64$ dw	$112 \times 112 \times 64$
Conv / s1	$1 \times 1 \times 64 \times 128$	$56 \times 56 \times 64$
Conv dw / s1	$3 \times 3 \times 128$ dw	$56 \times 56 \times 128$
Conv / s1	$1 \times 1 \times 128 \times 128$	$56 \times 56 \times 128$
Conv dw / s2	$3 \times 3 \times 128$ dw	$56 \times 56 \times 128$
Conv / s1	$1 \times 1 \times 128 \times 256$	$28 \times 28 \times 128$
Conv dw / s1	$3 \times 3 \times 256$ dw	$28 \times 28 \times 256$
Conv / s1	$1 \times 1 \times 256 \times 256$	$28 \times 28 \times 256$
Conv dw / s2	$3 \times 3 \times 256$ dw	$28 \times 28 \times 256$
Conv / s1	$1 \times 1 \times 256 \times 512$	$14 \times 14 \times 256$
5× Conv dw / s1	$3 \times 3 \times 512$ dw	$14 \times 14 \times 512$
	$1 \times 1 \times 512 \times 512$	$14 \times 14 \times 512$
	$3 \times 3 \times 512$ dw	$14 \times 14 \times 512$
	$1 \times 1 \times 512 \times 1024$	$7 \times 7 \times 512$
	$3 \times 3 \times 1024$ dw	$7 \times 7 \times 1024$
Conv / s1	$1 \times 1 \times 1024 \times 1024$	$7 \times 7 \times 1024$
Avg Pool / s1	Pool $7 \times 7$	$7 \times 7 \times 1024$
FC / s1	$1024 \times 1000$	$1 \times 1 \times 1024$
Softmax / s1	Classifier	$1 \times 1 \times 1000$

Table 2. Resource Per Layer Type

Type	Mult-Adds	Parameters
Conv $1 \times 1$	94.86%	74.59%
Conv DW $3 \times 3$	3.06%	1.06%
Conv $3 \times 3$	1.19%	0.02%
Fully Connected	0.18%	24.33%

# Width Multiplier & Resolution Multiplier

- Width Multiplier – Thinner Models
  - For a given layer and width multiplier  $\alpha$ , the number of input channels  $M$  becomes  $\alpha M$  and the number of output channels  $N$  becomes  $\alpha N$  – where  $\alpha$  with typical settings of 1, 0.75, 0.6 and 0.25
- Resolution Multiplier – Reduced Representation
  - The second hyper-parameter to reduce the computational cost of a neural network is a resolution multiplier  $\rho$
  - $0 < \rho \leq 1$ , which is typically set implicitly so that input resolution of network is 224, 192, 160 or 128 ( $\rho = 1, 0.857, 0.714, 0.571$ )
- Computational cost:  
$$D_K \times D_K \times \alpha M \times \rho D_F \times \rho D_F + \alpha M \times \alpha N \times \rho D_F \times \rho D_F$$

# Width Multiplier & Resolution Multiplier

- input feature map of size  $14 \times 14 \times 512$  with a kernel K of size  $3 \times 3 \times 512 \times 512$
- 

Table 3. Resource usage for modifications to standard convolution. Note that each row is a cumulative effect adding on top of the previous row. This example is for an internal MobileNet layer with  $D_K = 3$ ,  $M = 512$ ,  $N = 512$ ,  $D_F = 14$ .

Layer/Modification	Million	Million
	Mult-Adds	Parameters
Convolution	462	2.36
Depthwise Separable Conv	52.3	0.27
$\alpha = 0.75$	29.6	0.15
$\rho = 0.714$	15.1	0.15

# Experiments

Table 4. Depthwise Separable vs Full Convolution MobileNet

Model	ImageNet Accuracy	Million Mult-Adds	Million Parameters
Conv MobileNet	71.7%	4866	29.3
MobileNet	70.6%	569	4.2

Table 5. Narrow vs Shallow MobileNet

Model	ImageNet Accuracy	Million Mult-Adds	Million Parameters
0.75 MobileNet	68.4%	325	2.6
Shallow MobileNet	65.3%	307	2.9

Table 6. MobileNet Width Multiplier

Width Multiplier	ImageNet Accuracy	Million Mult-Adds	Million Parameters
1.0 MobileNet-224	70.6%	569	4.2
0.75 MobileNet-224	68.4%	325	2.6
0.5 MobileNet-224	63.7%	149	1.3
0.25 MobileNet-224	50.6%	41	0.5

Table 7. MobileNet Resolution

Resolution	ImageNet Accuracy	Million Mult-Adds	Million Parameters
1.0 MobileNet-224	70.6%	569	4.2
1.0 MobileNet-192	69.1%	418	4.2
1.0 MobileNet-160	67.2%	290	4.2
1.0 MobileNet-128	64.4%	186	4.2

Shallow or deep. The full network has a group of 5 layers in the middle that you can leave out.