

# 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.

{howarda, menglong, bochen, dkalenichenko, weijunw, weyand, anm, hadam}@google.com

## Abstract

We present a class of efficient models (called MobileNets) for mobile and embedded vision applications. MobileNets are based on a streamlined architecture that uses depthwise separable convolutions to build light weight deep neural networks. We introduce two simple global hyper-parameters that efficiently trade off between latency and accuracy. These hyper-parameters allow the model builder to choose the right sized model for their application based on the constraints of the problem. We present extensive experiments on resource and accuracy tradeoffs and show strong performance compared to other popular models on ImageNet classification. We then demonstrate the effectiveness of MobileNets across a wide range of applications and use cases including object detection, finegrain classification, face attributes and large scale geo-localization.

## 1. Introduction

Convolutional neural networks have become ubiquitous in computer vision ever since AlexNet [19] popularized deep convolutional neural networks by winning the ImageNet Challenge: ILSVRC 2012 [24]. The general trend has been to make deeper and more complicated networks in order to achieve higher accuracy [27, 31, 29, 8]. However, these advances to improve accuracy are not necessarily making networks more efficient with respect to size and speed. In many real world applications (such as robotics, self-driving car and augmented reality), the recognition tasks need to be carried out in a timely fashion on a computationally limited platform. This paper describes an efficient network architecture and a set of two hyper-parameters in order to build very small, low latency models that can be easily matched to the design requirements for mobile and embedded vision applications. Section 2 reviews prior work in building small

models. Section 3 describes the MobileNet architecture and two hyper-parameters width multiplier and resolution multiplier to define smaller and more efficient MobileNets. Section 4 describes experiments on ImageNet as well a variety of different applications and use cases. Section 5 closes with a summary and conclusion. ↗ 암드로니 쓰리 편

## 2. Prior Work

There has been rising interest in building small and efficient neural networks in the recent literature, e.g. [16, 34, 12, 36, 22]. Many different approaches can be generally categorized into either compressing pretrained networks or training small networks directly. This paper proposes a class of network architectures that allows a model developer to specifically choose a small network that matches the resource restrictions (latency, size) for their application. MobileNets primarily focus on optimizing for latency but also yield small networks. Many papers on small networks focus only on size but do not consider speed. ↗ MobileNet은 시간과 크기 둘 고려!

MobileNets are built primarily from depthwise separable convolutions initially introduced in [26] and subsequently used in Inception models [13] to reduce the computation in the first few layers. Flattened networks [16] build a network out of fully factorized convolutions and showed the potential of extremely factorized networks. Independent of this current paper, Factorized Networks [34] introduces a similar factorized convolution as well as the use of topological connections. Subsequently, the Xception network [3] demonstrated how to scale up depthwise separable filters to outperform Inception V3 networks. Another small network is SqueezeNet [12] which uses a bottleneck approach to design a very small network. Other reduced computation networks include structured transform networks [28] and deep fried convnets [37]. ↗ 1. depthwise separable conv 2. Inception model

A different approach for obtaining small networks is shrinking, factorizing or compressing pretrained networks. Compression based on product quantization [36], hashing

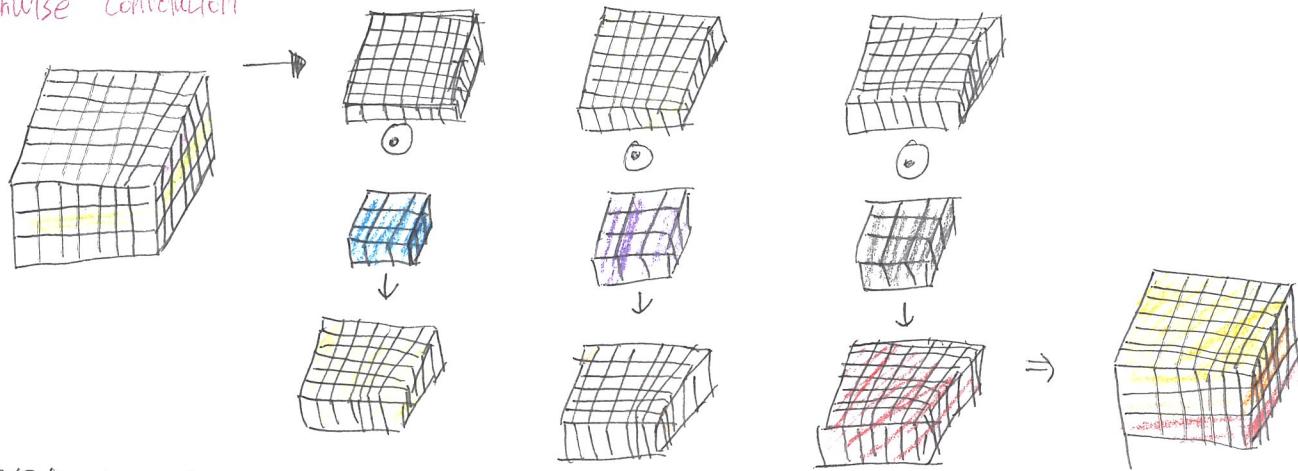
### \* factorization

기존의  $5 \times 5$  CONV 연산을 2번의  $3 \times 3$  CONV 연산으로 교체

### \* hashing

: 반복 비교 X, 특정 계산만으로 바로 자료 저장 주소 알아내는 단축법

## \* Depthwise Convolution



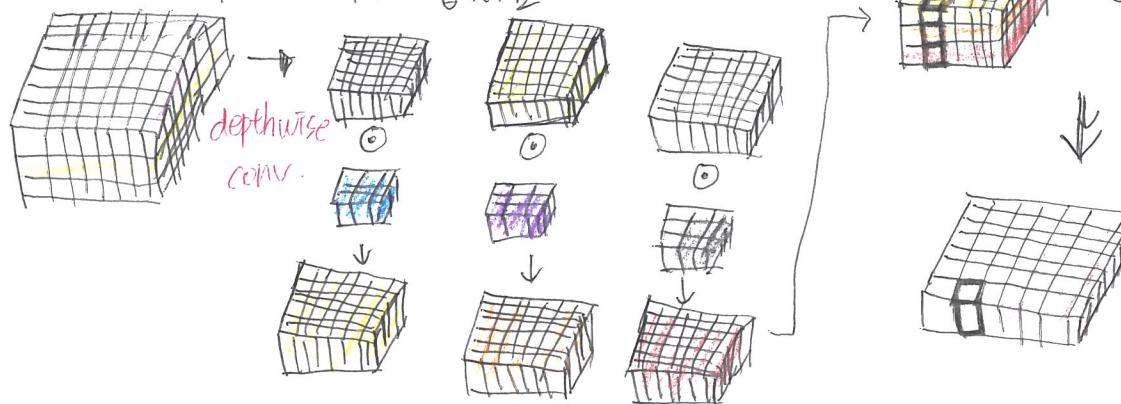
$\Rightarrow 8 \times 8 \times 3$ 의 matrix를 depthwise convolution 하기 위해  $3 \times 3 \times 3$ 의 커널 사용.

이때,  $3 \times 3 \times 3$  커널은 각 채널별로 분리되어  $3 \times 3$ 의 2차원 커널을  $8 \times 8$  크기의 각 분리된 matrix에  
붙여 가는 convolution 진행하고 다시 합쳐지는 구조.

: 채널 방향의 convolution은 하지 않고 공간 방향의 convolution만 진행.

## \* Depthwise Separable convolution.

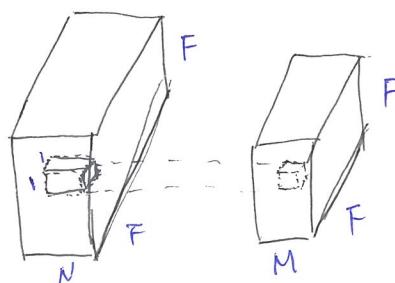
: 각 channel output이 하나로 합쳐짐



depth wise convolution을 진행한 결과물에 각 채널을 1개의 채널로 압축할 수 있는 추가 conv. 진행

## \* Pointwise convolution (1x1 conv)

: depthwise conv. depthwise separable conv는 공간 방향의 conv 진행 후, channel 방향에 conv. 진행함/하지  
pointwise conv. 는 공간 방향의 conv는 진행 X. 채널방향의 conv. 진행



1x1x $C$ 의 크기 kernel을 사용해 input의 특징을 1개의 채널로  
압축시키는 효과 O  
input의 채널을 압축시키는 효과  $\rightarrow$  연산속도 ↑  
압축하면서 없어지는 데이터는 연산 속도와 trade-off 관계. 어느정도 조정

## \* distillation (증류)

- 여러 갈 학습된 큰 네트워크 (teacher Network)에
- 자신을 실제로 사용하고자 하는 작은 네트워크 (student network)에게 전달하는 것.

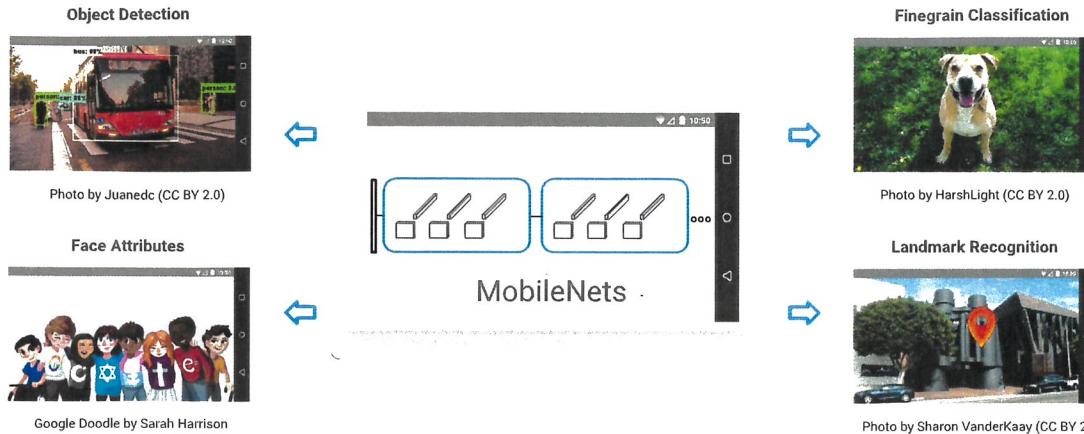


Figure 1. MobileNet models can be applied to various recognition tasks for efficient on device intelligence.

small network 얻기  
1. shrinking [2], and pruning, vector quantization and Huffman coding [5] have been proposed in the literature. Additionally various factorizations have been proposed to speed up pre-trained networks [14, 20]. Another method for training small networks is distillation [9] which uses a larger network to teach a smaller network. It is complementary to our approach and is covered in some of our use cases in section 4. Another emerging approach is low bit networks [4, 22, 11]. **small network : shrinking, factorizing, compressing, distillation.**

### 3. MobileNet Architecture

**core layer**  
In this section we first describe the core layers that MobileNet is built on which are depthwise separable filters. We then describe the MobileNet network structure and conclude with descriptions of the two model shrinking hyperparameters width multiplier and resolution multiplier.  $\rightarrow$  (설명 필요)

#### 3.1. 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.) For MobileNets the depthwise convolution applies a single filter to each input channel. The pointwise convolution then applies a  $1 \times 1$  convolution to combine the outputs of the depthwise convolution. A standard convolution both filters and combines inputs into a new set of outputs in one step. The depthwise separable convolution splits this into two layers, a separate layer for filtering and a separate layer for combining. This factorization has the effect of drastically reducing computation and model size. Figure 2 shows how a standard convolution 2(a) is factorized into a depthwise convolution 2(b) and a  $1 \times 1$  pointwise convolution 2(c).  $\rightarrow$  depthwise separable conv. depthwise conv. point wise, standard conv.

A standard convolutional layer takes as input a  $D_F \times D_F$  depthwise separable conv.  $\leftarrow$  layer for filtering  $\leftarrow$  layer for combining.  $\rightarrow$  standard conv.  $\rightarrow$  depthwise conv. & pointwise conv.

: 각 채널에 하나의 필터 적용

\* pointwise conv.  
 $\therefore$  1 conv 적용.  $\therefore$  depth wise conv. 복잡한 계산은 없어짐.

$D_F \times M$  feature map  $F$  and produces a  $D_F \times D_F \times N$  feature map  $G$  where  $D_F$  is the spatial width and height of a square input feature map,  $M$  is the number of input channels (input depth),  $D_G$  is the spatial width and height of a square output feature map and  $N$  is the number of output channels (output depth).  $\rightarrow$  standard conv. (기본)

The standard convolutional layer is parameterized by convolution kernel  $K$  of size  $D_K \times D_K \times M \times N$  where  $D_K$  is the spatial dimension of the kernel (assumed to be square) and  $M$  is number of input channels and  $N$  is the number of output channels (as defined previously).  $\rightarrow$  standard conv. (기본) (K)의 설명.

The output feature map for standard convolution assuming stride one and padding is computed as:  $\mathbf{F} = D_F \times D_F \times N$

$$G_{k,l,n} = \sum_{i,j,m} K_{i,j,m,n} \cdot F_{k+i-1, l+j-1, m} \quad (1)$$

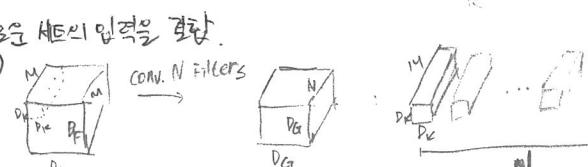
Standard convolutions have the computational cost of:

$$D_K \cdot D_K \cdot M \cdot N / D_F \cdot D_F \quad (2)$$

where the computational cost depends multiplicatively on the number of input channels  $M$ , the number of output channels  $N$ , the kernel size  $D_K \times D_K$  and the feature map size  $D_F \times D_F$ . MobileNet models address each of these terms and their interactions. First it uses depthwise separable convolutions to break the interaction between the number of output channels and the size of the kernel.  $\rightarrow$  standard conv. 이용.  $\rightarrow$  MN 상호작용 step!

The standard convolution operation has the effect of filtering features based on the convolutional kernels and combining features in order to produce a new representation. The filtering and combination steps can be split into two steps via the use of factorized convolutions (called depthwise

<sup>1</sup>We assume that the output feature map has the same spatial dimensions as the input and both feature maps are square. Our model shrinking results generalize to feature maps with arbitrary sizes and aspect ratios.



## \* Computation cost

: filter size  $\times$  Input size (가로, 세로)

4 dm

separable convolutions for substantial reduction/in computational cost.  $\rightarrow$  Standard conv  $\geq$  depthwise separable conv로 2단계로 나누기 가능

Depthwise separable convolution are made up of two layers: depthwise convolutions and pointwise convolutions. We use depthwise convolutions (to apply a single filter per each input channel) (input depth). Pointwise convolution, (a simple  $1 \times 1$  convolution), is then used to create a linear combination of the output of the depthwise layer. MobileNets use both batchnorm and ReLU nonlinearities for both layers.  $\rightarrow$  depthwise separable conv  $\leq$  depthwise conv + pointwise conv (각각 BN, ReLU 사용)

Depthwise convolution (with one filter per input channel) (input depth) can be written as:

$$\hat{G}_{k,l,m} = \sum_{i,j} \hat{K}_{i,j,m} \cdot F_{k+i-1, l+j-1, m} \quad (3)$$

where ( $\hat{K}$ ) is the depthwise convolutional kernel of size  $D_K \times D_K \times M$  (where the  $m$ th filter (in  $\hat{K}$ ) is applied to the  $m$ th channel (in  $F$ ) to produce the  $m$ th channel of the filtered output feature map  $\hat{G}$ ).  $\rightarrow$  depthwise conv. size

Depthwise convolution has a computational cost of:

$$\text{cost (co-computation)} = \text{filter size} \times \text{input size} \cdot (\text{width} \times \text{height}) \quad D_K \cdot D_K \cdot M \cdot D_F \cdot D_F \quad (4)$$

Depthwise convolution is extremely efficient relative to standard convolution. However it only filters input channels, it does not combine them to create new features. So an additional layer that computes a linear combination of the output of depthwise convolution via  $1 \times 1$  convolution is needed in order to generate these new features.

The combination of depthwise convolution and  $1 \times 1$  (pointwise) convolution is called depthwise separable convolution which was originally introduced in [26].

Depthwise separable convolutions cost:

$$D_K \cdot D_K \cdot M \cdot D_F \cdot D_F + M \cdot N \cdot D_F \cdot D_F \quad (5)$$

which is the sum of the depthwise and  $1 \times 1$  pointwise convolutions.

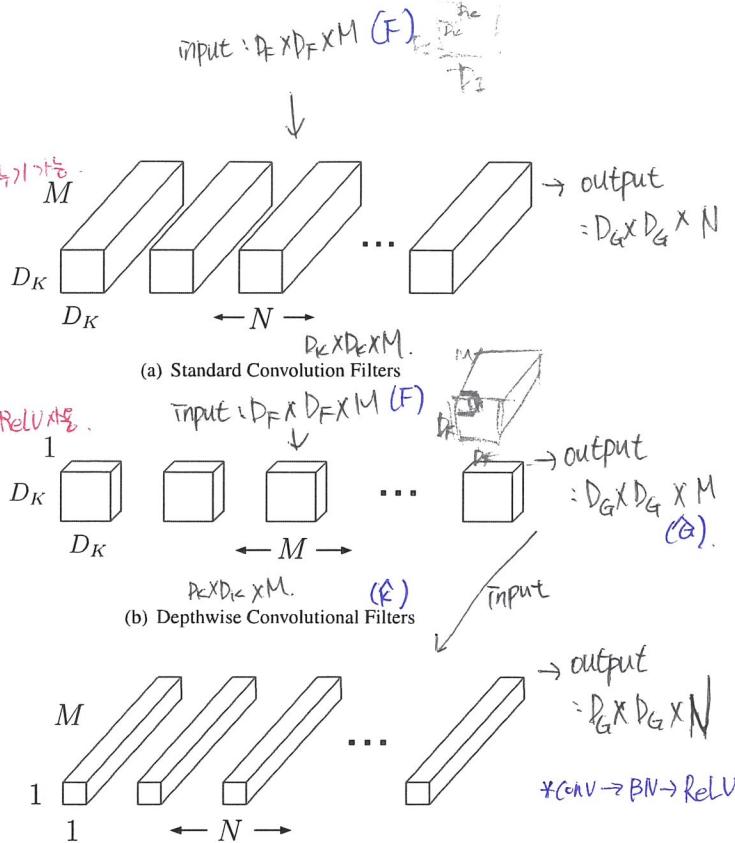
By expressing convolution as a two step process of filtering and combining we get a reduction in computation of:

$$\begin{aligned} \text{two-steps} &\rightarrow D_K \cdot D_K \cdot M \cdot D_F \cdot D_F + M \cdot N \cdot D_F \cdot D_F \\ \text{Standard conv} &\rightarrow \frac{D_K \cdot D_K \cdot M \cdot N \cdot D_F \cdot D_F}{D_K \cdot D_K \cdot M \cdot N \cdot D_F \cdot D_F} \\ &= \frac{1}{N} + \frac{1}{D_K^2} \end{aligned}$$

MobileNet uses  $3 \times 3$  depthwise separable convolutions which uses between 8 to 9 times less computation than standard convolutions (at only a small reduction in accuracy) as seen in Section 4.  $\rightarrow$  MobileNet에서 depthwise conv 사용됨

Additional factorization (in spatial dimension, such as in [16, 31]) does not save much additional computation (as very little computation is spent in depthwise convolutions)

$\rightarrow$  추가 분해는 별 효과 X.



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

Figure 2. The standard convolutional filters in (a) are replaced by two layers: depthwise convolution in (b) and pointwise convolution in (c) to build a depthwise separable filter.

Mobile Net :  
[1단계 레이어에 depthwise conv 사용]  
[2단계, full conv 일 때는 1st layer는?] ]

## 3.2. Network Structure and Training

The MobileNet structure is built on depthwise separable convolutions as mentioned in the previous section except for the first layer (which is a full convolution). By defining the network in such simple terms we are able to easily explore network topologies to find a good network. The MobileNet architecture is defined in Table 1. All layers are followed by a batchnorm [13] and ReLU nonlinearity with the exception of the final fully connected layer which has no nonlinearity and feeds into a softmax layer for classification. Figure 3 contrasts a layer with regular convolutions, batchnorm and ReLU nonlinearity to the factorized layer with depthwise convolution,  $1 \times 1$  pointwise convolution as well as batchnorm and ReLU after each convolutional layer. Down sampling is handled with strided convolution in the depthwise convolutions as well as in the first layer. A final average pooling reduces the spatial resolution to 1 before the fully connected layer. Counting depthwise and pointwise convolutions as separate layers, MobileNet has 28 layers.

It is not enough to simply define networks in terms of a small number of Mult-Adds. It is also important to make sure these operations can be efficiently implementable. For

Figure 3  
regular conv  
(BN, ReLU)  
factorized layer  
(depthwise conv)  
 $1 \times 1$  pointwise conv  
(BN, ReLU)

구현 가능한

## \* Strided convolution

: Stride  $\neq 1$  in conv.

\* down sampling  
: 이미지 사이즈 작게 줄이기

\* fully conv. Network (FCN)  
: 각 층은 마지막 레이어에 FC layer 포함해 softmax.  
FCN layer 값이 FC layer에서 위치 정렬 X.  
FCN의 input size가 고정되어야 함.  
FCN은 FC layer 사용 X, 모든 레이어 conv. layer로 바꿀

## \*GEMM

(General Matrix Multiplication) ↗ 복잡한 계산.  
= 행렬곱셈 나타내는 % \* %.

depthwise et pointwise conv를 다른 층으로 인식하면  
MobileNet은 28층!

↑  
마지막 FC layer 642 + BN, ReLU 따라다니

↗ 마지막 FC는 softmax 역할을.

\*down sampling

: 1st layer 포함

strided conv. 사용

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	$1 \times 1 \times 256 \times 512$	$14 \times 14 \times 256$
5x Conv dw / s1	$3 \times 3 \times 512$ dw	$14 \times 14 \times 512$
Conv / s1	$1 \times 1 \times 512 \times 512$	$14 \times 14 \times 512$
Conv dw / s2	$3 \times 3 \times 512$ dw	$14 \times 14 \times 512$
Conv / s1	$1 \times 1 \times 512 \times 1024$	$7 \times 7 \times 512$
Conv dw / s2	$3 \times 3 \times 1024$ 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% 거의 모든 parameter.

Conv  $1 \times 1$  빠짐.

↳ computation time.

and width multiplier  $\alpha$ , the number of input channels  $M$  becomes  $\alpha M$  and the number of output channels  $N$  becomes  $\alpha N$ . 더 좋은 모델위해 MobileNet 단축화면서 파라미터 X 줄임.

The computational cost of a depthwise separable convolution with width multiplier  $\alpha$  is:

$$D_K \cdot D_K \cdot \alpha M \cdot D_F \cdot D_F + \alpha M \cdot \alpha N \cdot D_F \cdot D_F \quad (6)$$

where  $\alpha \in (0, 1]$  (with typical settings of 1, 0.75, 0.5 and 0.25).  $\alpha = 1$  is the baseline MobileNet and  $\alpha < 1$  are reduced MobileNets. Width multiplier has the effect of reducing computational cost and the number of parameters quadratically by roughly  $\alpha^2$ . Width multiplier can be applied to any model structure to define a new smaller model with a reasonable accuracy, latency and size trade off. It is used to define a new reduced structure that needs to be trained from scratch. ↗ width multiplier는 다른 모델 적용 가능

### 3.4. Resolution Multiplier: Reduced Representation ↗ 해상도

The second hyper-parameter to reduce the computational cost of a neural network is a resolution multiplier  $\rho$ . We ap-

MobileNet 반복은 더 작고 빠르게! → 더 작고 더 짧은 계산이용 모델우선

즉 multiplier X 파라미터 설정

↓

네트워크를 각 층에서 고밀하게 않게 하기.

$\alpha \in [0, 1]$  사이의 값.

$\alpha_1$  : base line MobileNet  
 $\alpha_2$  : reduced MobileNet

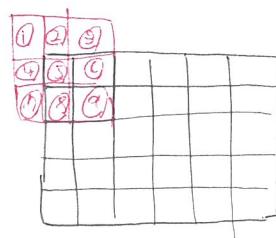
## \* im2col

: CNN은 1989년 고안 but 빙동화X. (why? Conv는 특정 size의 window를 입력 이미지 각 pixel마다 size의 자료 개수만큼 곱셉 연산 수행 → 대소의 Feature Map)  
 ↗ 학습해야 할 parameter 적음, 연산시간 많.

빠르게 conv 연산 위해 이미지 데이터를 conv하기 좋은 연속된 데이터로 변환 (im2col)

ex) 3x3 window 사용해 conv 수행.

source image



→ col

: 입력데이터 개수가 커지게 됨.  
 conv 연산 빠른 대신 연산 메모리↑

## \* RMSProp

: AdaGrad는 최소값에 도달하기 전에 학습률을 0에 수렴하게 만들 수 O.

AdaGrad가 간단한 convex function에서는 잘 작동 but, 복잡한 다차원 곡면함수를 잘 탐색X.  
 기울기의 달성을 누락만으로 충분X.

⇒ RMSProp : 기울기 달성을 누적X. 지수 가중 이동 평균 (Exponentially weighted moving average)  
 사용해 최신 기울기들이 더 크게 반영.

$$h_i \leftarrow p h_{i-1} + (1-p) \frac{\partial L_i}{\partial w} \odot \frac{\partial L_i}{\partial w}$$

이전 AdaGrad의 h에 새로운 hyper parameter p 추가.

$h_i$ 가 무한히 커지지X. p가 작을수록 가장 최신의 기울기 더 크게 반영.

## \* 분산 병렬 학습 구조.

진행 속도가

parameter server로부터 항상 일관된 최신 global-parameter  
 계산. 다음 단계 학습 진행

(동기식) (Synchronous) : 상대적으로 빠른 worker들이 모든 worker의 학습 진도에 맞춰 학습.  
 (비동기식) (Asynchronous) : worker들이 독립적 학습 진행. parameter server로부터 동일 X global parameter  
 제공받아 다음의 학습 과정 진행.

## \* label smoothing

: 이진 분류 문제에서 label은 0 또는 1.  
 But, 반드시 그걸지 않을 수 있음

why? 사람이 annotation하기 때문에 실수 가능성 존재

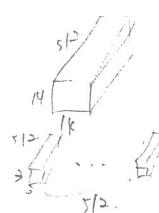
(label smoothing) X : 이때 발생할 수 있는 잘못된 label의 영향을 줄이기 위해  
 label을 0 또는 1이 아니라 smooth하게 부여.

→ label smoothing은 mislabeling 된 데이터 고려 위해 사용  
 추가적으로 regularization에 도움.

model generalization, calibration(수정)에 도움

full conv. layer :

Input feature map:  $14 \times 14 \times 512$   
kernel =  $3 \times 3 \times 512 \times 512$



계산량 & # of parameters.

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 Mult-Adds	Million 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

\*resolution multiplier:  $\rho \leftarrow$  원본 이미지 크기 / 2<sup>정</sup>에 입력 해상도 설정하여 같이 보정됨.  
ply this to the input image and the internal representation of every layer is subsequently reduced by the same multiplier. In practice we implicitly set  $\rho$  by setting the input resolution. → input image의 기본값은 커버리 파라미터  $\rho$

We can now express the computational cost for the core layers of our network as depthwise separable convolutions with width multiplier  $\alpha$  and resolution multiplier  $\rho$ :

$$D_K \cdot D_K \cdot @M \cdot @D_F \cdot @D_F + @M \cdot @N \cdot @D_F \cdot @D_F \quad (7)$$

Input resolution: where  $\rho \in (0, 1]$  which is typically set implicitly so that the input resolution of the network is 224, 192, 160 or 128.  $\rho = 1$  is the baseline MobileNet and  $\rho < 1$  are reduced computation MobileNets. Resolution multiplier has the effect of reducing computational cost by  $\rho^2$ . → ρ에 대한 설명

As an example (we can look at a typical layer) in MobileNet and see how depthwise separable convolutions, width multiplier and resolution multiplier reduce the cost and parameters. Table 3 shows the computation and number of parameters for a layer as architecture shrinking methods are sequentially applied to the layer. The first row shows the Mult-Adds and parameters for a full convolutional layer with an input feature map of size  $14 \times 14 \times 512$  with a kernel  $K$  of size  $3 \times 3 \times 512 \times 512$ . We will look in detail in the next section at the trade offs between resources and accuracy. → table 3 살펴

## 4. Experiments

In this section we first investigate the effects of depthwise convolutions as well as the choice of shrinking by reducing the width of the network rather than the number of layers. We then show the trade offs of reducing the network based on the two hyper-parameters: width multiplier and resolution multiplier and compare results to a number of popular models. We then investigate MobileNets applied to a number of different applications. → 모델의 구조

### 4.1. Model Choices

First we show results for MobileNet with depthwise separable convolutions compared to a model built with full convolutions. In Table 4 we see that using depthwise separable convolutions compared to full convolutions only reduces

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

→ 적은 PAZ

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

→ thin model  
→ 5 layer of separable filters

Table 6. MobileNet Width Multiplier ( $\alpha$ )

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

→ Smooth drop

Table 7. MobileNet Resolution ( $\rho$ )

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

→ Smooth drop

accuracy by 1% on ImageNet was saving tremendously on mult-adds and parameters. → full conv. 모델 VS MobileNet

We next show results comparing thinner models (with width multiplier) to shallower models using less layers. To make MobileNet shallower, the 5 layers of separable filters (with feature size  $14 \times 14 \times 512$  in Table 1) are removed. Table 5 shows that at similar computation and number of parameters, that making MobileNets thinner is 3% better than making them shallower. → thin VS shallow Model

### 4.2. Model Shrinking Hyperparameters

Table 6 shows the accuracy, computation and size trade offs of shrinking the MobileNet architecture with the width multiplier  $\alpha$ . Accuracy drops off smoothly until the architecture is made too small at  $\alpha = 0.25$ . → α에 따른 경계값

Table 7 shows the accuracy, computation and size trade offs for different resolution multipliers (by training MobileNets with reduced input resolutions). Accuracy drops off smoothly across resolution. → ρ에 따른 경계값

Figure 4 shows the trade off between ImageNet Accuracy and computation for the 16 models made from the cross product of width multiplier  $\alpha \in \{1, 0.75, 0.5, 0.25\}$  and resolutions  $\{224, 192, 160, 128\}$ . Results are log linear with a jump when models get very small at  $\alpha = 0.25$ .

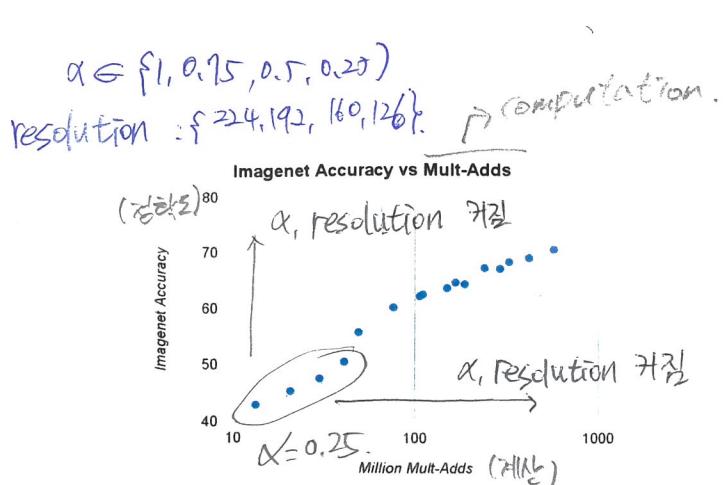


Figure 4. This figure shows the trade off between computation (Mult-Adds) and accuracy on the ImageNet benchmark. Note the log linear dependence between accuracy and computation.

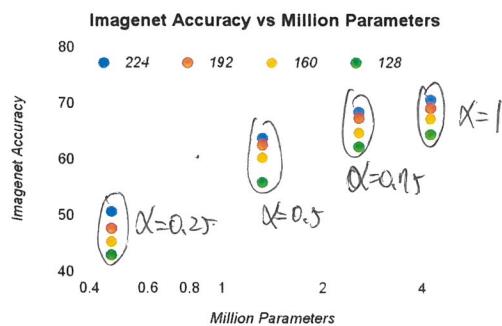


Figure 5. This figure shows the trade off between the number of parameters and accuracy on the ImageNet benchmark. The colors encode input resolutions. The number of parameters do not vary based on the input resolution.

Figure 5 shows the trade off between ImageNet Accuracy and number of parameters for the 16 models made from the cross product of width multiplier  $\alpha \in \{1, 0.75, 0.5, 0.25\}$  and resolutions  $\{224, 192, 160, 128\}$ .

Table 8 compares full MobileNet to the original GoogleNet [30] and VGG16 [27]. MobileNet is nearly as accurate as VGG16 while being 32 times smaller and 27 times less compute intensive. It is more accurate than GoogleNet while being smaller and more than 2.5 times less computation.

Table 9 compares a reduced MobileNet with width multiplier  $\alpha = 0.5$  and reduced resolution  $160 \times 160$ . Reduced MobileNet is 4% better than AlexNet [19] while being 45× smaller and 9.4× less compute than AlexNet. It is also 4% better than SqueezeNet [12] at about the same size and 22× less computation.

\*fine-grained

-잘게 쪼개 늘는 것.

$\hookleftarrow$  coarse-grained

: 템어리로 작업

Model	ImageNet Accuracy	Million Mult-Adds	Million Parameters
1.0 MobileNet-224	70.6%	569	4.2
GoogleNet	69.8%	1550	6.8
VGG 16	71.5%	15300	138

Table 8. MobileNet Comparison to Popular Models

Model	ImageNet Accuracy	Million Mult-Adds	Million Parameters
0.50 MobileNet-160	60.2%	76	1.32
SqueezeNet	57.5%	1700	1.25
AlexNet	57.2%	720	60

Table 9. Smaller MobileNet Comparison to Popular Models

Model	ImageNet Accuracy	Million Mult-Adds	Million Parameters
Inception V3 [18]	84%	5000	23.2
1.0 MobileNet-224	83.3%	569	3.3
0.75 MobileNet-224	81.9%	325	1.9
1.0 MobileNet-192	81.9%	418	3.3
0.75 MobileNet-192	80.5%	239	1.9

Table 10. MobileNet for Stanford Dogs

Scale	Im2GPS [7]	PlaNet [35]	PlaNet MobileNet
Continent (2500 km)	51.9%	77.6%	79.3%
Country (750 km)	35.4%	64.0%	60.3%
Region (200 km)	32.1%	51.1%	45.2%
City (25 km)	21.9%	31.7%	31.7%
Street (1 km)	2.5%	11.0%	11.4%

#### 4.3. Fine Grained Recognition

We train MobileNet for fine grained recognition on the Stanford Dogs dataset [17]. We extend the approach of [18] and collect an even larger but noisy training set than [18] from the web. We use the noisy web data to pretrain a fine grained dog recognition model and then fine tune the model on the Stanford Dogs training set. Results on Stanford Dogs test set are in Table 10. MobileNet can almost achieve the state of the art results from [18] at greatly reduced computation and size.

#### 4.4. Large Scale Geolocalization

PlaNet [35] casts the task of determining where on earth a photo was taken as a classification problem. The approach divides the earth into a grid of geographic cells that serve as the target classes and trains a convolutional neural network

행성 분류 문제 : 지구 상 어디에서 사진 찍었나?

지구를 grid(지역)로 나누고 수백만장의

geo-tag 된 사진을 CNN으로 훈련

## \* Pre training

: 선형 학습, 사전 훈련, 전처리 과정

Multi layered Perceptron (MLP)에서 weight와 bias를 잘 초기화 시키는 방법

⇒ 효과적으로 layer 쌓아 여러개의 hidden layer도 효율적 훈련하는 목적.

또한, unsupervised learning이 가능함

(물론 가중치, 평균 초기화 끝나면 레이블 된 데이터로 supervised learning 해야함)

레이블 X인 큰 데이터 넣어 훈련 시킬 수 있다.

Drop out, Mini-batch 방식으로 Pre-training 생각 가능.

## \* Fine-tuning

: 기존에 학습되어져 있는 모델 기반으로 아키텍처를 새로운 목적 (내 이미지 데이터에 맞게) 변형하고 이미 학습된 모델 weight로부터 학습 업데이트 하는 방법.

- 모델의 파라미터 미세하게 조절하는 행위

: 특히 DL에서 이미 존재하는 모델에 추가 데이터를 투입해 파라미터 업데이트

⇒ "정교한" "파라미터" 튜닝.

ex) 고양이, 개 분류기 VGG16 모델.

1000개의 카테고리 학습 → 고양이, 개 2개의 카테고리만 필요.

(모든 레이어 끝 필요 X)

데이터를 해당 모델로 예측해 bottleneck feature만 뽑아내고, 이를 이용해 어파인 레이어 (fully-connected layer)만 학습 시켜 사용하는 방법

→ 파인튜닝 X. 피처 추출 레이어 파라미터 업데이트 안함.

(파인튜닝) → 기존에 학습이 된 레이어 내 데이터를 추가로 학습시켜 파라미터 업뎃.  
단, 정교 해야함.

일반적인 피처 학습한

완전히 랜덤한 초기 파라미터 쓰거나 가장 아래 레이어의 모든 추상화된  
파라미터 학습시 오비피팅 or 전체 파라미터 막가짐.

PlaNet ← localizing 흐도 성능.

Im2GPS 보다 높은 성능

52M parameter

5.74 B multi-adds

on millions of geo-tagged photos. PlaNet has been shown to successfully localize a large variety of photos and to outperform Im2GPS [6, 7] (that addresses the same task.)

We re-train PlaNet using the MobileNet architecture on the same data. While the full PlaNet model based on the Inception V3 architecture [31] has 52 million parameters and 5.74 billion multi-adds. The MobileNet model has only 13 million parameters with the usual 3 million for the body and 10 million for the final layer and 0.58 Million multi-adds. As shown in Tab. 11, the MobileNet version delivers only slightly decreased performance compared to PlaNet despite being much more compact. Moreover, it still outperforms Im2GPS by a large margin.

#### 4.5. Face Attributes

Another use-case for MobileNet is compressing large systems with unknown or esoteric training procedures. In a face attribute classification task, we demonstrate a synergistic relationship between MobileNet and distillation [9], a knowledge transfer technique for deep networks. We seek to reduce a large face attribute classifier with 75 million parameters and 1600 million Mult-Adds. The classifier is trained on a multi-attribute dataset similar to YFCC100M [32].

We distill a face attribute classifier using the MobileNet architecture. Distillation [9] works by training the classifier to emulate the outputs of a larger model instead of the ground-truth labels, hence enabling training from large (and potentially infinite) unlabeled datasets. Marrying the scalability of distillation training and the parsimonious parameterization of MobileNet, the end system not only requires no regularization (e.g. weight-decay and early-stopping), but also demonstrates enhanced performances. It is evident from Tab. 12 that the MobileNet-based classifier is resilient to aggressive model shrinking: it achieves a similar mean average precision across attributes (mean AP) as the in-house while consuming only 1% the Multi-Adds.

#### 4.6. Object Detection

MobileNet can also be deployed as an effective base network in modern object detection systems. We report results for MobileNet trained for object detection on COCO data based on the recent work that won the 2016 COCO challenge [10]. In table 13, MobileNet is compared to VGG and Inception V2 [14] under both Faster-RCNN [23] and SSD [21] framework. In our experiments, SSD is evaluated with 300 input resolution (SSD 300) and Faster-RCNN is compared with both 300 and 600 input resolution (Faster-RCNN 300, Faster-RCNN 600). The Faster-RCNN model evaluates 300 RPN proposal boxes per image. The models are trained on COCO train+val excluding 8k minival images

<sup>2</sup>The emulation quality is measured by averaging the per-attribute cross-entropy over all attributes.

\* Ground-truth.

: 어느 한 장소에서 수집된 정보.

\* Parsimonious model

: 최소의 데이터를 이용한 최대의 정보

MobileNet

13. M Parameter  
0.58 M multi-adds

\* Framework

: 프로그램 기본 구조

기능적 비슷한지.

Table 12. Face attribute classification using the MobileNet architecture. Each row corresponds to a different hyper-parameter setting (width multiplier  $\alpha$  and image resolution).

Width Multiplier / Resolution	Mean AP	Million Mult-Adds	Million Parameters
1.0 MobileNet-224	88.7%	568	3.2
0.5 MobileNet-224	88.1%	149	0.8
0.25 MobileNet-224	87.2%	45	0.2
1.0 MobileNet-128	88.1%	185	3.2
0.5 MobileNet-128	87.7%	48	0.8
0.25 MobileNet-128	86.4%	15	0.2
Baseline	86.9%	1600	7.5

작은 네트워크에 대한

비교

Table 13. COCO object detection results comparison using different frameworks and network architectures. mAP is reported with COCO primary challenge metric (AP at IoU=0.50:0.05:0.95)

Framework * Resolution	Model	mAP	Billion Mult-Adds	Million Parameters
	deeplab-VGG	21.1%	34.9	33.1
Faster-RCNN 300	Inception V2	22.0%	3.8	13.7
	MobileNet	19.3%	1.2	6.8
	VGG	22.9%	64.3	138.5
Faster-RCNN 600	Inception V2	15.4%	118.2	13.3
	MobileNet	16.4%	25.2	6.1
	VGG	25.7%	149.6	138.5
SSD 300	Inception V2	21.9%	129.6	13.3
	MobileNet	19.8%	30.5	6.1



Figure 6. Example object detection results using MobileNet SSD.

and evaluated on minival. For both frameworks, MobileNet achieves comparable results to other networks with only a fraction of computational complexity and model size.

#### 4.7. Face Embeddings

The FaceNet model is a state-of-the-art face recognition model [25]. It builds face embeddings based on the triplet loss. To build a mobile FaceNet model we use distillation to train by minimizing the squared differences of the output

distillation

FaceNet과

MobileNet 출격

차이 적용 최소화

+

→

\* Coco data

: detection, segmentation, captioning 데이터 전처리

Table 14. MobileNet Distilled from FaceNet

Model	1e-4 Accuracy	Million Mult-Adds	Million Parameters
FaceNet [25]	83%	1600	7.5
1.0 MobileNet-160	79.4%	286	4.9
1.0 MobileNet-128	78.3%	185	5.5
0.75 MobileNet-128	75.2%	166	3.4
0.75 MobileNet-128	72.5%	108	3.8

of FaceNet and MobileNet on the training data. Results for very small MobileNet models can be found in table 14.

## 5. Conclusion

We proposed a new model architecture called MobileNets based on depthwise separable convolutions. We investigated some of the important design decisions leading to an efficient model. We then demonstrated how to build smaller and faster MobileNets using width multiplier and resolution multiplier by trading off a reasonable amount of accuracy to reduce size and latency. We then compared different MobileNets to popular models demonstrating superior size, speed and accuracy characteristics. We concluded by demonstrating MobileNet's effectiveness when applied to a wide variety of tasks. As a next step to help adoption and exploration of MobileNets, we plan on releasing models in Tensor Flow.

## References

- [1] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, et al. Tensorflow: Large-scale machine learning on heterogeneous systems, 2015. *Software available from tensorflow.org*, 1, 2015. 4
- [2] W. Chen, J. T. Wilson, S. Tyree, K. Q. Weinberger, and Y. Chen. Compressing neural networks with the hashing trick. *CoRR, abs/1504.04788*, 2015. 2
- [3] F. Chollet. Xception: Deep learning with depthwise separable convolutions. *arXiv preprint arXiv:1610.02357v2*, 2016. 1
- [4] M. Courbariaux, J.-P. David, and Y. Bengio. Training deep neural networks with low precision multiplications. *arXiv preprint arXiv:1412.7024*, 2014. 2
- [5] S. Han, H. Mao, and W. J. Dally. Deep compression: Compressing deep neural network with pruning, trained quantization and huffman coding. *CoRR, abs/1510.00149*, 2, 2015. 2
- [6] J. Hays and A. Efros. IM2GPS: estimating geographic information from a single image. In *Proceedings of the IEEE International Conference on Computer Vision and Pattern Recognition*, 2008. 7
- [7] J. Hays and A. Efros. Large-Scale Image Geolocation. In J. Choi and G. Friedland, editors, *Multimodal Location Estimation of Videos and Images*. Springer, 2014. 6, 7
- [8] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. *arXiv preprint arXiv:1512.03385*, 2015. 1
- [9] G. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. 2, 7
- [10] J. Huang, V. Rathod, C. Sun, M. Zhu, A. Korattikara, A. Fathi, I. Fischer, Z. Wojna, Y. Song, S. Guadarrama, et al. Speed/accuracy trade-offs for modern convolutional object detectors. *arXiv preprint arXiv:1611.10012*, 2016. 7
- [11] I. Hubara, M. Courbariaux, D. Soudry, R. El-Yaniv, and Y. Bengio. Quantized neural networks: Training neural networks with low precision weights and activations. *arXiv preprint arXiv:1609.07061*, 2016. 2
- [12] F. N. Iandola, M. W. Moskewicz, K. Ashraf, S. Han, W. J. Dally, and K. Keutzer. SqueezeNet: Alexnet-level accuracy with 50x fewer parameters and 1mb model size. *arXiv preprint arXiv:1602.07360*, 2016. 1, 6
- [13] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015. 1, 3, 7
- [14] M. Jaderberg, A. Vedaldi, and A. Zisserman. Speeding up convolutional neural networks with low rank expansions. *arXiv preprint arXiv:1405.3866*, 2014. 2
- [15] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. *arXiv preprint arXiv:1408.5093*, 2014. 4
- [16] J. Jin, A. Dundar, and E. Culurciello. Flattened convolutional neural networks for feedforward acceleration. *arXiv preprint arXiv:1412.5474*, 2014. 1, 3
- [17] A. Khosla, N. Jayadevaprakash, B. Yao, and L. Fei-Fei. Novel dataset for fine-grained image categorization. In *First Workshop on Fine-Grained Visual Categorization, IEEE Conference on Computer Vision and Pattern Recognition*, Colorado Springs, CO, June 2011. 6
- [18] J. Krause, B. Sapp, A. Howard, H. Zhou, A. Toshev, T. Duerig, J. Philbin, and L. Fei-Fei. The unreasonable effectiveness of noisy data for fine-grained recognition. *arXiv preprint arXiv:1511.06789*, 2015. 6
- [19] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012. 1, 6
- [20] V. Lebedev, Y. Ganin, M. Rakhuba, I. Oseledets, and V. Lempitsky. Speeding-up convolutional neural networks using fine-tuned cp-decomposition. *arXiv preprint arXiv:1412.6553*, 2014. 2
- [21] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, and S. Reed. Ssd: Single shot multibox detector. *arXiv preprint arXiv:1512.02325*, 2015. 7
- [22] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi. Xnor-net: Imagenet classification using binary convolutional neural networks. *arXiv preprint arXiv:1603.05279*, 2016. 1, 2
- [23] S. Ren, K. He, R. Girshick, and J. Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015. 7