

## Abstract

Convolutional Neural Networks (ConvNets) are commonly developed at a fixed resource budget, and then scaled up for better accuracy if more resources are available. In this paper, we systematically study model scaling and identify that carefully balancing network depth, width, and resolution can lead to better performance. Based on this observation, we propose a new scaling method that uniformly scales all dimensions of depth/width/resolution using a simple yet highly effective *compound coefficient*. We demonstrate the effectiveness of this method on scaling up MobileNets and ResNet.

To go even further, we use neural architecture search to design a new baseline network and scale it up to obtain a family of models, called *EfficientNets*, which achieve much better accuracy and efficiency than previous ConvNets. In particular, our EfficientNet-B7 achieves state-of-the-art 84.4% top-1 / 97.1% top-5 accuracy on ImageNet, while being 8.4x smaller and 6.1x faster on inference than the best existing ConvNet. Our EfficientNets also transfer well and achieve state-of-the-art accuracy on CIFAR-100 (91.7%), Flowers (98.8%), and 3 other transfer learning datasets, with an order of magnitude fewer parameters. Source code is at <https://github.com/tensorflow/tpu/tree/master/models/official/efficientnet>.

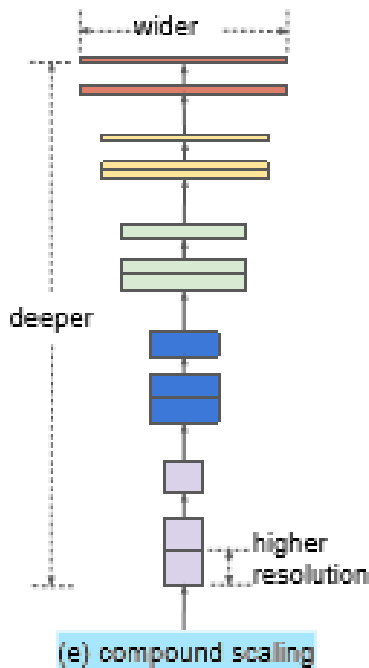


Table 1. EfficientNet-B0 baseline network – Each row describes a stage  $i$  with  $\hat{L}_i$  layers, with input resolution  $(\hat{H}_i, \hat{W}_i)$  and output channels  $\hat{C}_i$ . Notations are adopted from equation 2.

Stage $i$	Operator $\mathcal{F}_i$	Resolution $\hat{H}_i \times \hat{W}_i$	#Channels $\hat{C}_i$	#Layers $\hat{L}_i$
1	Conv3x3	$224 \times 224$	32	1
2	MBConv1, k3x3	$112 \times 112$	16	1
3	MBConv6, k3x3	$112 \times 112$	24	2
4	MBConv6, k5x5	$56 \times 56$	40	2
5	MBConv6, k3x3	$28 \times 28$	80	3
6	MBConv6, k5x5	$14 \times 14$	112	3
7	MBConv6, k5x5	$14 \times 14$	192	4
8	MBConv6, k3x3	$7 \times 7$	320	1
9	Conv1x1 & Pooling & FC	$7 \times 7$	1280	1

## ■ CNN Architecture

**Purpose:** Improving MobileNet, ResNet.

**Method:** proposing a new compound scaling method

$$\text{depth: } d = \alpha^\phi$$

$$\text{width: } w = \beta^\phi$$

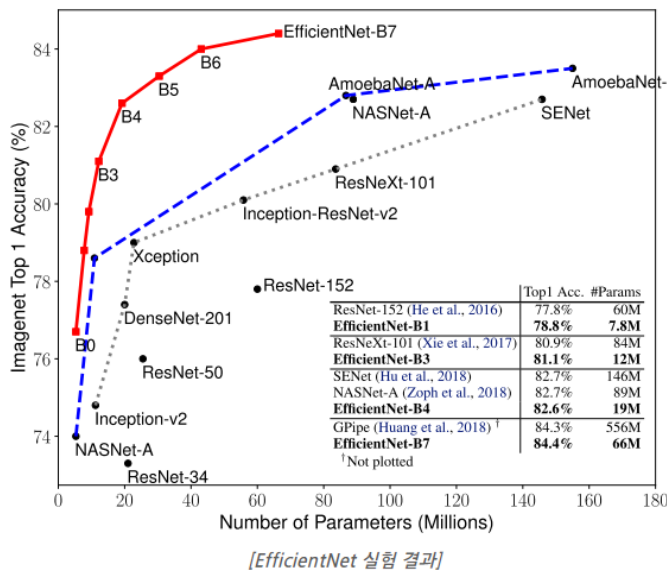
$$\text{resolution: } r = \gamma^\phi$$

$$\text{s.t. } \alpha \cdot \beta^2 \cdot \gamma^2 \approx 2$$

$$\alpha \geq 1, \beta \geq 1, \gamma \geq 1$$

**Problem:** it is common to scale only one of the three dimensions—depth, width, image size. But 세가지 방법을 조합해서 모델을 확장하는 경우나 적절한 조합을 찾는 방법에 대해서 명확하지 않음.

**Suggestion:** proposing a new scaling method that uniformly scales all dimensions of depth/width/resolution using a simple yet highly effective compound coefficient



## Model Scaling

1. Filter개수(Channel 개수) 늘리는 **width scaling**

(eg. MobileNet-224 1.0, MobileNet-224 0.5, ShuffleNet)

2. Layer 개수를 늘리는 **depth scaling**

(ex. ResNet-50, ResNet-101)

3. Input image 해상도를 높이는 **resolution scaling**

- **EfficientNet**: 3가지 scaling 방법을 동시에 고려

- 3가지 scaling 기법에 대해 각 scaling 기법마다 나머지는 고정해두고 1개의 scaling factor만 키워가며 정확도의 변화를 측정

- ◆ 실험적으로, 1-2가지 scaling factor만 키워주는 것 보다, 3 가지 scaling factor를 동시에 키워주는 것이 가장 성능이 좋을 보임

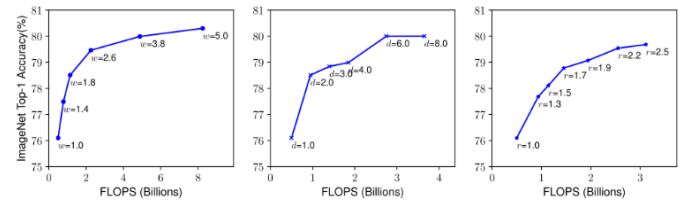


Figure 3. Scaling Up a Baseline Model with Different Network Width ( $w$ ), Depth ( $d$ ), and Resolution ( $r$ ) Coefficients. Bigger networks with larger width, depth, or resolution tend to achieve higher accuracy, but the accuracy gain quickly saturate after reaching 80%, demonstrating the limitation of single dimension scaling. Baseline network is described in Table 1.

[Single Dimension Model Scaling 실험 결과]

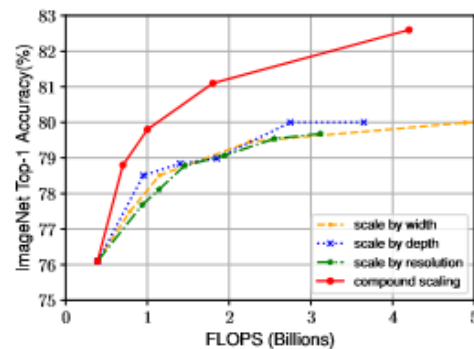


Figure 8. Scaling Up EfficientNet-B0 with Different Methods.

Table 7. Scaled Models Used in Figure 7.

Model	FLOPS	Top-1 Acc.
Baseline model (EfficientNet-B0)	0.4B	77.3%
Scale model by depth ( $d=4$ )	1.8B	79.0%
Scale model by width ( $w=2$ )	1.8B	78.9%
Scale model by resolution ( $r=2$ )	1.9B	79.1%
Compound Scale ( $d=1.4, w=1.2, r=1.3$ )	1.8B	81.1%