



# EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks

Mingxing Tan, Quoc V. Le



Youngwoo Kim



# Index

---

- Introduction
- Related Work
- Compound Model Scaling
- Architecture
- Experiments
- Conclusion

# Introduction

---

*It is common to scale only one of the three dimensions - depth, width, and image size.*

- depth - *Deep Residual Learning for Image Recognition (He et al., 2016)*
- width - *Wide Residual Networks (Zagoruyko & Komodakis, 2016)*
- resolution - *GPipe: Efficient training of giant neural networks using pipeline parallelism (Huang et al., 2018)*

# Introduction

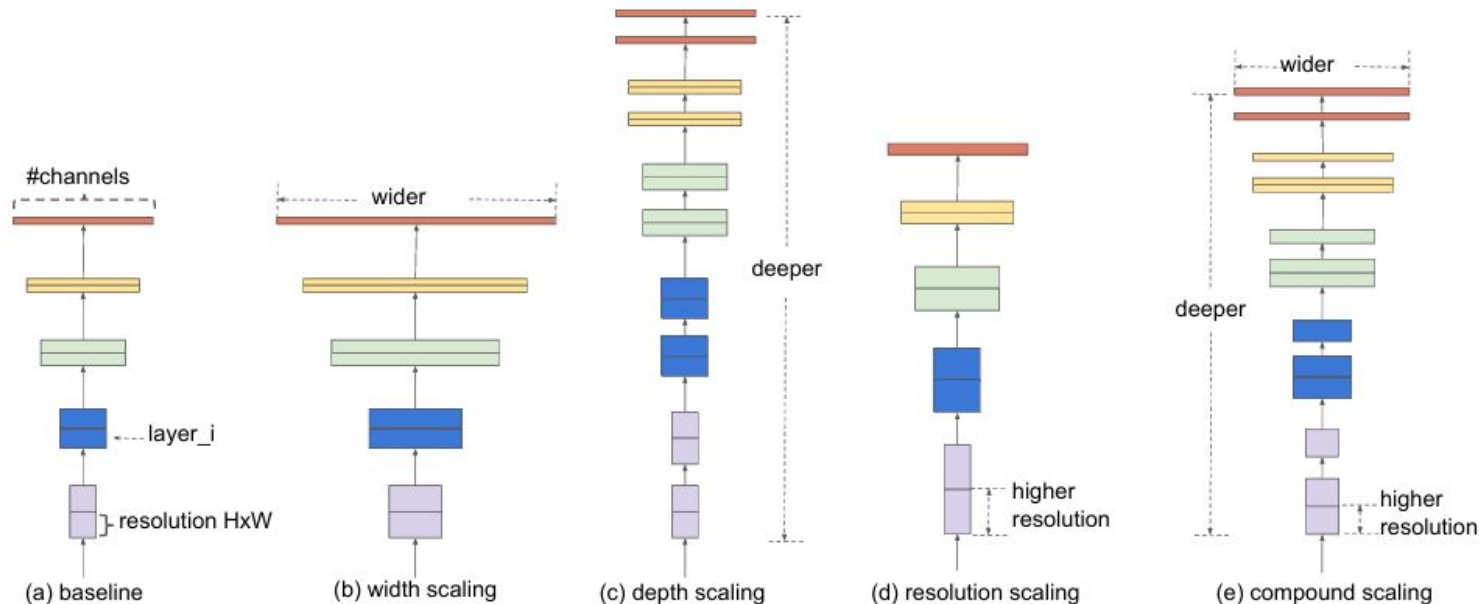


Figure 2. **Model Scaling.** (a) is a baseline network example; (b)-(d) are conventional scaling that only increases one dimension of network width, depth, or resolution. (e) is our proposed compound scaling method that uniformly scales all three dimensions with a fixed ratio.

# Introduction

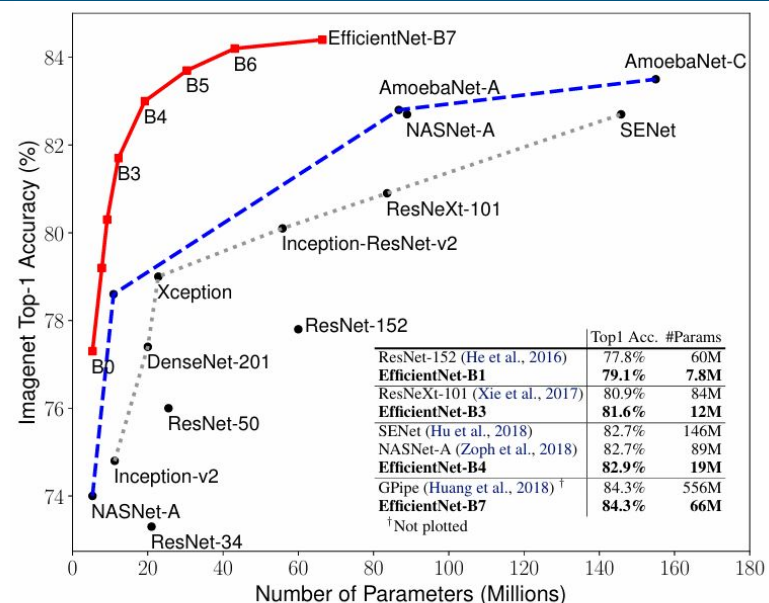
---

Is there a principled method to scale up ConvNets that can achieve better accuracy and efficiency?

- Compound Scaling: Uniformly scales network width, depth, and resolution with a set of fixed scaling coefficients
- For  $2^N$  times more computational resources, increase network depth by  $\alpha^N$ , width by  $\beta^N$ , and image size by  $\gamma^N$

# Introduction

- EfficientNets outperform other ConvNets
- Surpasses the best existing GPipe accuracy but using 8.4x fewer parameters and running 6.1x faster on inference



**Figure 1. Model Size vs. ImageNet Accuracy.** All numbers are for single-crop, single-model. Our EfficientNets significantly outperform other ConvNets. In particular, **EfficientNet-B7 achieves new state-of-the-art 84.3% top-1 accuracy but being 8.4x smaller and 6.1x faster than GPipe.** EfficientNet-B1 is 7.6x smaller and 5.7x faster than ResNet-152. Details are in Table 2 and 4.

# Related Work - ConvNet Accuracy

---

GPipe pushes the state-of-the-art ImageNet top-1 accuracy to 84.3% using 557M parameters

- it is so big that it can only be trained with a specialized pipeline parallelism library

# Related Work - ConvNet Efficiency

---

Neural Architecture Search becomes increasingly popular in designing efficient mobile-size ConvNets

- But it is unclear how to apply these techniques for larger models that have much larger design space and much more expensive tuning cost.



# Related Work - Model Scaling

---

Prior studies such as WideResNet and scaled ResNet have shown that network depth and width are both important for ConvNets' expressive power

- it still remains an open question of how to effectively scale a ConvNet to achieve better efficiency and accuracy

# Compound Model Scaling - Problem Formulation

$$\mathcal{N} = \bigodot_{i=1 \dots s} \mathcal{F}_i^{L_i} (X_{\langle H_i, W_i, C_i \rangle})$$

$$\max_{d, w, r} \text{Accuracy}(\mathcal{N}(d, w, r))$$

$$s.t. \quad \mathcal{N}(d, w, r) = \bigodot_{i=1 \dots s} \hat{\mathcal{F}}_i^{d \cdot \hat{L}_i} (X_{\langle r \cdot \hat{H}_i, r \cdot \hat{W}_i, w \cdot \hat{C}_i \rangle})$$

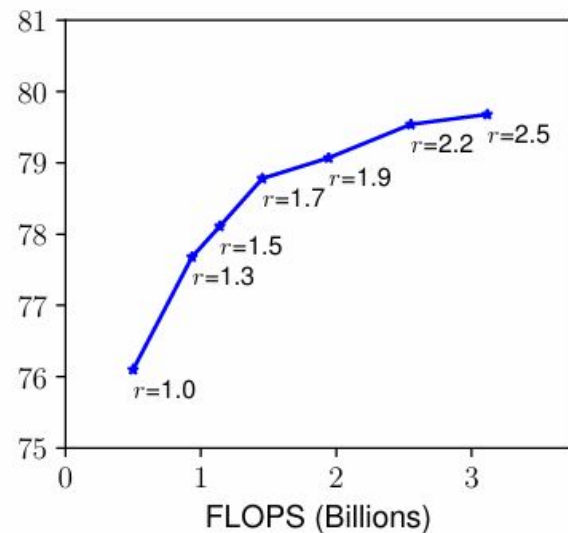
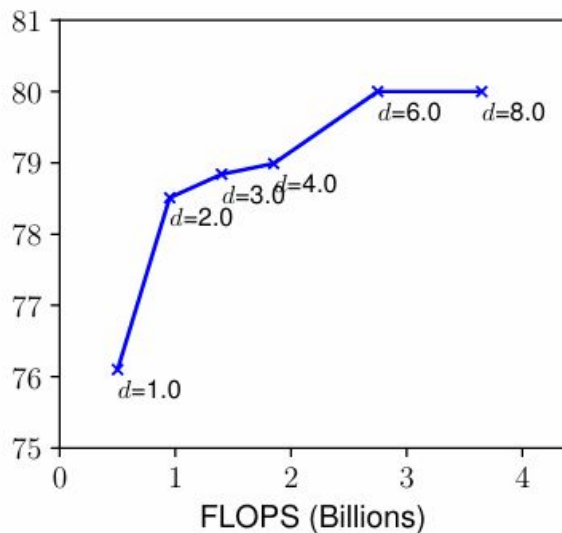
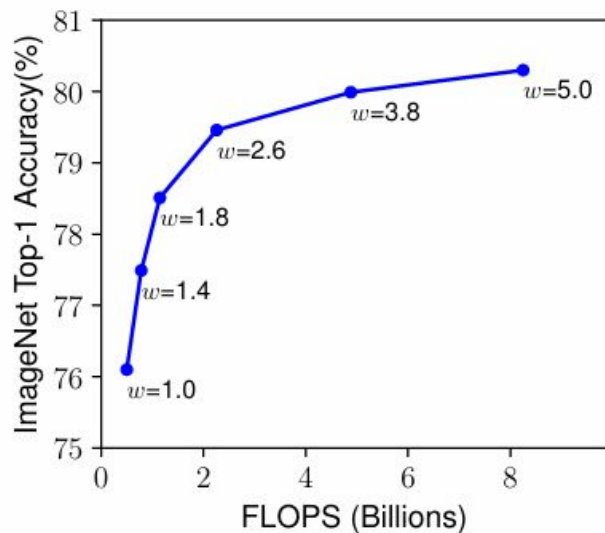
$$\text{Memory}(\mathcal{N}) \leq \text{target\_memory}$$

$$\text{FLOPS}(\mathcal{N}) \leq \text{target\_flops}$$

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

Stage $i$	Operator $\hat{\mathcal{F}}_i$	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

# Compound Model Scaling - Scaling Dimensions

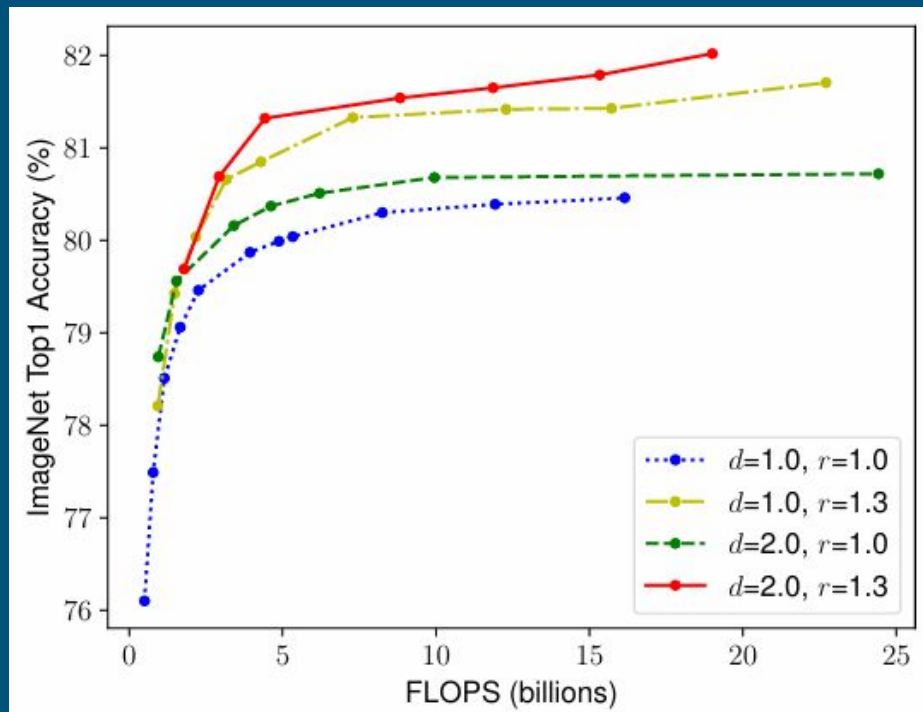


Scaling up any dimension of network width, depth, or resolution improves accuracy, but the accuracy gain diminishes for bigger models

# Compound Model Scaling - Compound Scaling

In order to pursue better accuracy and efficiency, it is critical to balance all dimensions of network width, depth, and resolution during ConvNet scaling

$d=1.0, r=1.0 \Rightarrow$  18 layers with  $224 \times 224$



# Compound Model Scaling

EfficientNet aims to double the FLOPs with each scaling step

$\alpha, \beta, \gamma$  are constants that can be determined by a small grid search

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

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

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

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

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

$$FLOPs \propto d \cdot w^2 \cdot r^2$$

# Architecture

---

Since model scaling does not change layer operators in baseline network, having a good baseline network is also critical.

- Evaluated scaling method using existing ConvNets
- also developed a new mobile-size baseline for better demonstration of effectiveness (EfficientNet)

# Architecture

---

Developed baseline network with Multi-Objective Neural Architecture Search (MO-NAS)

- Optimizes both accuracy and FLOPs
- $ACC(m) \times [FLOPs(m)/T]^w$
- $ACC(m)$ ,  $FLOPs(m)$  denote the accuracy and FLOPs of model  $m$
- $T$  is the target FLOPs
- $w=-0.07$  is a hyperparameter for controlling the trade-off between accuracy and FLOPs
- Author optimized FLOPs rather than latency
- Architecture is similar to MnasNET, except the size differs due to the larger FLOPs target

# Architecture

---

STEP 1: fix  $\Phi = 1$ , assuming twice more resources available. Do a small grid search of  $\alpha, \beta, \gamma$  based on previous equations. For EfficientNet-B0,  $\alpha=1.2$ ,  $\beta=1.1$ ,  $\gamma=1.15$

STEP 2: fix  $\alpha, \beta, \gamma$  as constants and scale up baseline network with different  $\Phi$



# Experiments

Model	Top-1 Acc.	Top-5 Acc.	#Params	Ratio-to-EfficientNet	#FLOPs	Ratio-to-EfficientNet
<b>EfficientNet-B0</b>	<b>77.1%</b>	<b>93.3%</b>	<b>5.3M</b>	<b>1x</b>	<b>0.39B</b>	<b>1x</b>
ResNet-50 (He et al., 2016)	76.0%	93.0%	26M	4.9x	4.1B	11x
DenseNet-169 (Huang et al., 2017)	76.2%	93.2%	14M	2.6x	3.5B	8.9x
<b>EfficientNet-B1</b>	<b>79.1%</b>	<b>94.4%</b>	<b>7.8M</b>	<b>1x</b>	<b>0.70B</b>	<b>1x</b>
ResNet-152 (He et al., 2016)	77.8%	93.8%	60M	7.6x	11B	16x
DenseNet-264 (Huang et al., 2017)	77.9%	93.9%	34M	4.3x	6.0B	8.6x
Inception-v3 (Szegedy et al., 2016)	78.8%	94.4%	24M	3.0x	5.7B	8.1x
Xception (Chollet, 2017)	79.0%	94.5%	23M	3.0x	8.4B	12x
<b>EfficientNet-B2</b>	<b>80.1%</b>	<b>94.9%</b>	<b>9.2M</b>	<b>1x</b>	<b>1.0B</b>	<b>1x</b>
Inception-v4 (Szegedy et al., 2017)	80.0%	95.0%	48M	5.2x	13B	13x
Inception-resnet-v2 (Szegedy et al., 2017)	80.1%	95.1%	56M	6.1x	13B	13x
<b>EfficientNet-B3</b>	<b>81.6%</b>	<b>95.7%</b>	<b>12M</b>	<b>1x</b>	<b>1.8B</b>	<b>1x</b>
ResNeXt-101 (Xie et al., 2017)	80.9%	95.6%	84M	7.0x	32B	18x
PolyNet (Zhang et al., 2017)	81.3%	95.8%	92M	7.7x	35B	19x
<b>EfficientNet-B4</b>	<b>82.9%</b>	<b>96.4%</b>	<b>19M</b>	<b>1x</b>	<b>4.2B</b>	<b>1x</b>
SENet (Hu et al., 2018)	82.7%	96.2%	146M	7.7x	42B	10x
NASNet-A (Zoph et al., 2018)	82.7%	96.2%	89M	4.7x	24B	5.7x
AmoebaNet-A (Real et al., 2019)	82.8%	96.1%	87M	4.6x	23B	5.5x
PNASNet (Liu et al., 2018)	82.9%	96.2%	86M	4.5x	23B	6.0x
<b>EfficientNet-B5</b>	<b>83.6%</b>	<b>96.7%</b>	<b>30M</b>	<b>1x</b>	<b>9.9B</b>	<b>1x</b>
AmoebaNet-C (Cubuk et al., 2019)	83.5%	96.5%	155M	5.2x	41B	4.1x
<b>EfficientNet-B6</b>	<b>84.0%</b>	<b>96.8%</b>	<b>43M</b>	<b>1x</b>	<b>19B</b>	<b>1x</b>
<b>EfficientNet-B7</b>	<b>84.3%</b>	<b>97.0%</b>	<b>66M</b>	<b>1x</b>	<b>37B</b>	<b>1x</b>
GPipe (Huang et al., 2018)	84.3%	97.0%	557M	8.4x	-	-

We omit ensemble and multi-crop models (Hu et al., 2018), or models pretrained on 3.5B Instagram images (Mahajan et al., 2018).

# Experiments

Apply scaling method to  
widely-used  
architectures ( ResNet,  
MobileNets)

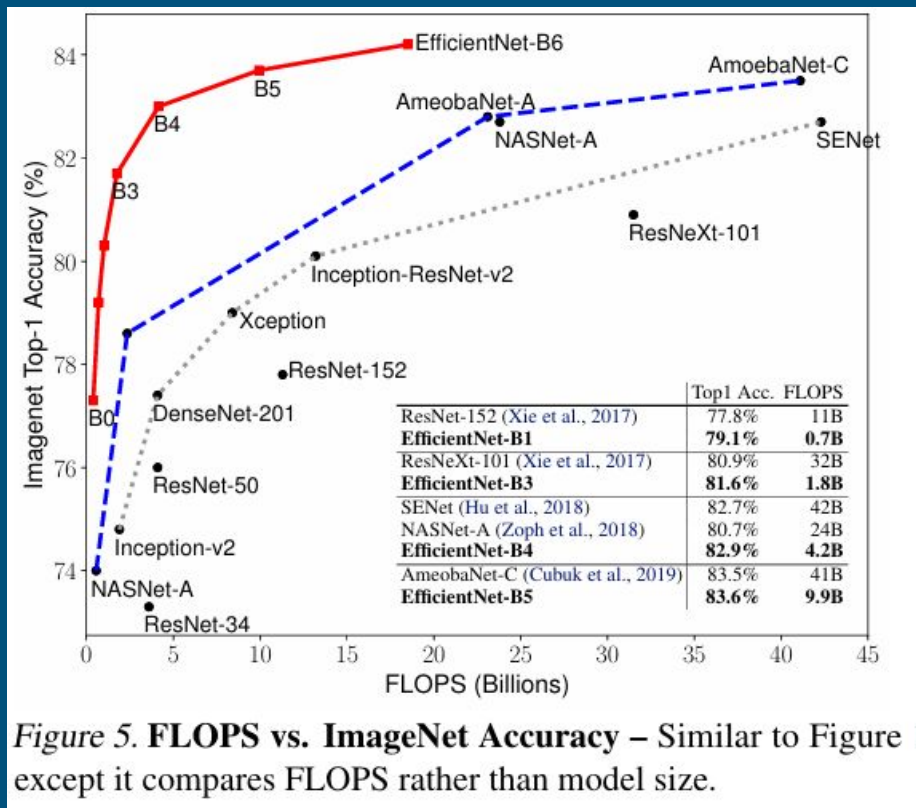
*Table 3. Scaling Up MobileNets and ResNet.*

Model	FLOPS	Top-1 Acc.
Baseline MobileNetV1 (Howard et al., 2017)	0.6B	70.6%
Scale MobileNetV1 by width ( $w=2$ )	2.2B	74.2%
Scale MobileNetV1 by resolution ( $r=2$ )	2.2B	72.7%
<b>compound scale (<math>d=1.4, w=1.2, r=1.3</math>)</b>	<b>2.3B</b>	<b>75.6%</b>
Baseline MobileNetV2 (Sandler et al., 2018)	0.3B	72.0%
Scale MobileNetV2 by depth ( $d=4$ )	1.2B	76.8%
Scale MobileNetV2 by width ( $w=2$ )	1.1B	76.4%
Scale MobileNetV2 by resolution ( $r=2$ )	1.2B	74.8%
<b>MobileNetV2 compound scale</b>	<b>1.3B</b>	<b>77.4%</b>
Baseline ResNet-50 (He et al., 2016)	4.1B	76.0%
Scale ResNet-50 by depth ( $d=4$ )	16.2B	78.1%
Scale ResNet-50 by width ( $w=2$ )	14.7B	77.7%
Scale ResNet-50 by resolution ( $r=2$ )	16.4B	77.5%
<b>ResNet-50 compound scale</b>	<b>16.7B</b>	<b>78.8%</b>

# Experiments

parameters-accuracy and  
FLOPs-accuracy curve

EfficientNet-B3 achieves higher  
accuracy than ResNeXt-101  
using 18x fewer FLOPs



# Experiments

**Table 4. Inference Latency Comparison** – Latency is measured with batch size 1 on a single core of Intel Xeon CPU E5-2690.

Acc. @ Latency		Acc. @ Latency	
ResNet-152	77.8% @ 0.554s	GPipe	84.3% @ 19.0s
EfficientNet-B1	78.8% @ 0.098s	EfficientNet-B7	84.4% @ 3.1s
<b>Speedup</b>	<b>5.7x</b>	<b>Speedup</b>	<b>6.1x</b>

# Experiments

**Table 5. EfficientNet Performance Results on Transfer Learning Datasets.** Our scaled EfficientNet models achieve new state-of-the-art accuracy for 5 out of 8 datasets, with 9.6x fewer parameters on average.

	Comparison to best public-available results						Comparison to best reported results					
	Model	Acc.	#Param	Our Model	Acc.	#Param(ratio)	Model	Acc.	#Param	Our Model	Acc.	#Param(ratio)
CIFAR-10	NASNet-A	98.0%	85M	EfficientNet-B0	98.1%	4M (21x)	<sup>†</sup> Gpipe	<b>99.0%</b>	556M	EfficientNet-B7	98.9%	64M (8.7x)
CIFAR-100	NASNet-A	87.5%	85M	EfficientNet-B0	88.1%	4M (21x)	Gpipe	91.3%	556M	EfficientNet-B7	<b>91.7%</b>	64M (8.7x)
Birdsnap	Inception-v4	81.8%	41M	EfficientNet-B5	82.0%	28M (1.5x)	Gpipe	83.6%	556M	EfficientNet-B7	<b>84.3%</b>	64M (8.7x)
Stanford Cars	Inception-v4	93.4%	41M	EfficientNet-B3	93.6%	10M (4.1x)	<sup>‡</sup> DAT	<b>94.8%</b>	-	EfficientNet-B7	94.7%	-
Flowers	Inception-v4	98.5%	41M	EfficientNet-B5	98.5%	28M (1.5x)	DAT	97.7%	-	EfficientNet-B7	<b>98.8%</b>	-
FGVC Aircraft	Inception-v4	90.9%	41M	EfficientNet-B3	90.7%	10M (4.1x)	DAT	92.9%	-	EfficientNet-B7	<b>92.9%</b>	-
Oxford-IIIT Pets	ResNet-152	94.5%	58M	EfficientNet-B4	94.8%	17M (5.6x)	Gpipe	<b>95.9%</b>	556M	EfficientNet-B6	95.4%	41M (14x)
Food-101	Inception-v4	90.8%	41M	EfficientNet-B4	91.5%	17M (2.4x)	Gpipe	93.0%	556M	EfficientNet-B7	<b>93.0%</b>	64M (8.7x)
Geo-Mean	<b>(4.7x)</b>						<b>(9.6x)</b>					

<sup>†</sup>Gpipe (Huang et al., 2018) trains giant models with specialized pipeline parallelism library.

<sup>‡</sup>DAT denotes domain adaptive transfer learning (Ngiam et al., 2018). Here we only compare ImageNet-based transfer learning results.

Transfer accuracy and #params for NASNet (Zoph et al., 2018), Inception-v4 (Szegedy et al., 2017), ResNet-152 (He et al., 2016) are from (Kornblith et al., 2019).

# Experiments

ImageNet performance of different scaling methods for the same EfficientNet-B0 baseline network

- compound scaling can further improve accuracy by up to 2.5% than other single-dimension scaling methods

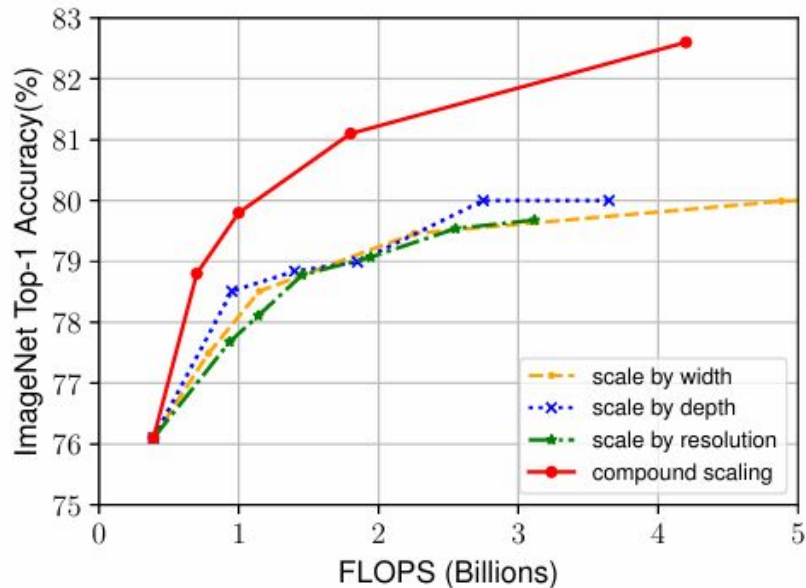


Figure 8. Scaling Up EfficientNet-B0 with Different Methods.



# Experiments

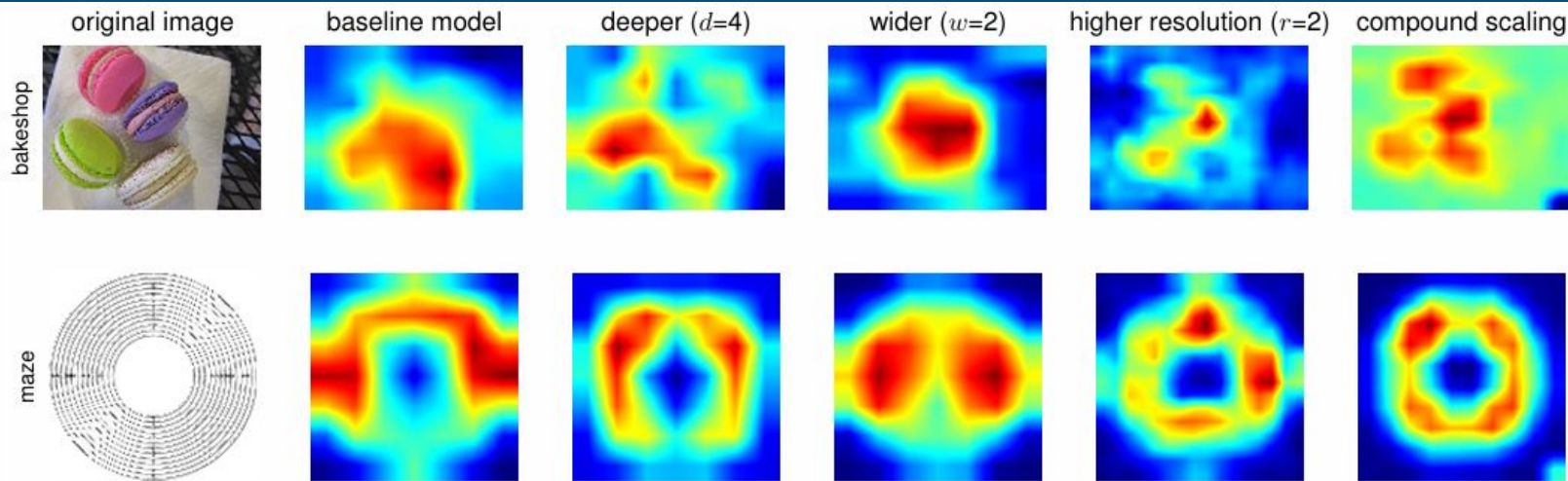


Figure 7. **Class Activation Map (CAM) (Zhou et al., 2016) for Models with different scaling methods-** Our compound scaling method allows the scaled model (last column) to focus on more relevant regions with more object details. Model details are in Table 7.

# Experiments

*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%
<b>Compound Scale (<math>d=1.4, w=1.2, r=1.3</math>)</b>	<b>1.8B</b>	<b>81.1%</b>



# Conclusion

---

- Balanced scaling of width, depth, and resolution is crucial for optimizing accuracy and efficiency in ConvNets.
- Prior ConvNet models lacked a systematic scaling method, leading to suboptimal performance.
- The proposed Compound Scaling method efficiently scales models while maintaining high performance.
- EfficientNet significantly reduces FLOPs and parameters while achieving SOTA accuracy.
- EfficientNet generalizes well across various tasks, excelling in both ImageNet and transfer learning datasets.