

**CVPR 2020**

**EfficientNet: Rethinking Model Scaling for  
Convolutional Neural Networks**

2022.08.01

논문 리뷰

배성훈

# EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks (CVPR 2020)

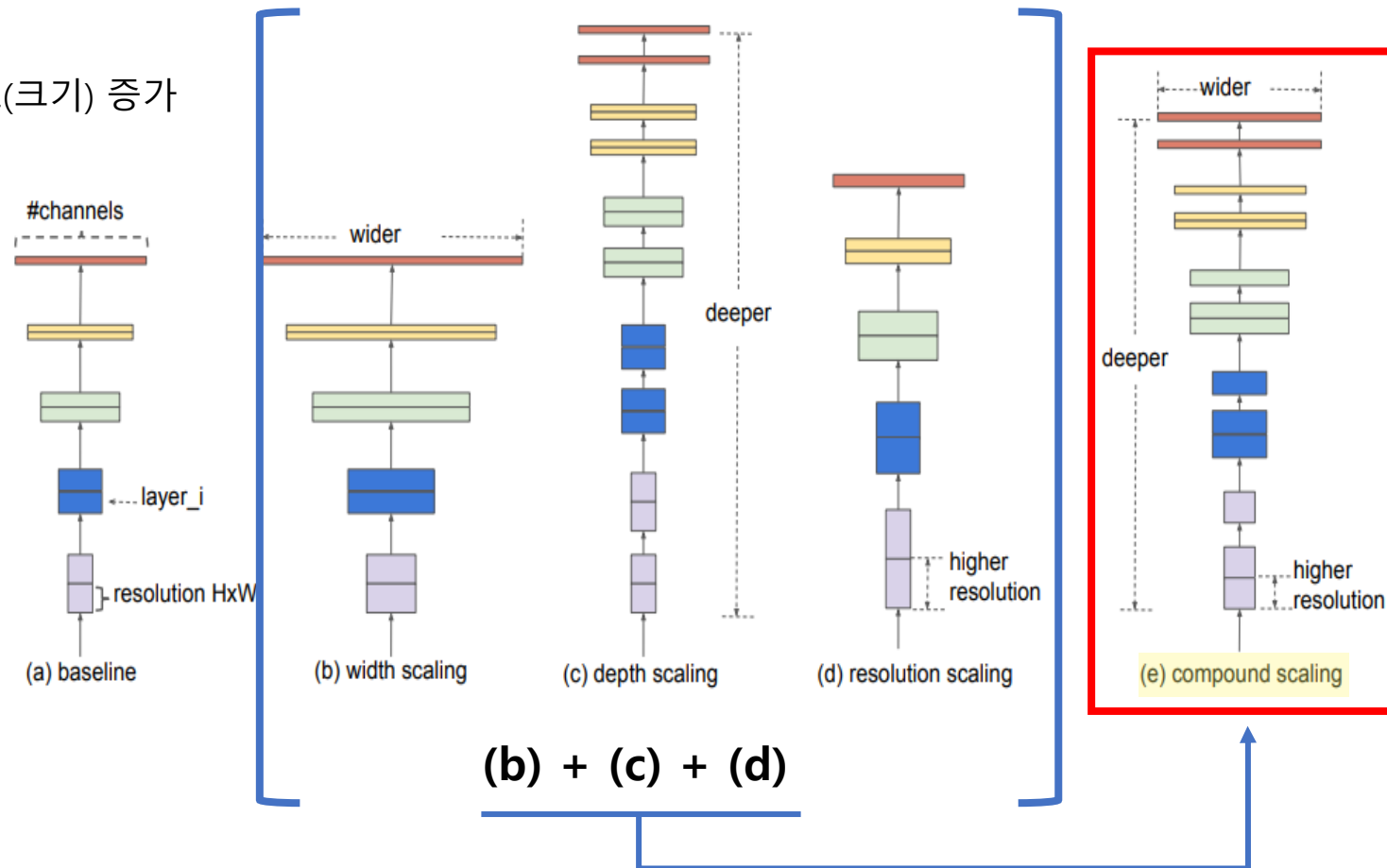
- **Research Background:**

- 기존의 연구는 CNN 정확도를 높이기 위해 Depth, Width, Resolution 중 한 가지 방법만을 선택해 성능을 향상
- 이에 따라 저자는 Convolution Network의 **정확성과 효율성**을 향상시키는 원칙에 따른 **Scale up 방법**에 대한 의문
- **Compound Coefficient**를 통해 **Depth, Width, Resolution**을 **균형 있게 Scaling**해 성능 향상

\***Depth:** Layer 수 증가

\***Width:** Filter(channel 수) 증가

\***Resolution:** Input image의 해상도(크기) 증가



# EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks (CVPR 2020)

- **Method:**

- Compound Scaling: Small Grid Search에 의해 결정되는 Width, Depth, Resolution scaling을 위한 **Constant coefficients**( $\alpha, \beta, \gamma$ )
- 각 Layer에서 수행하는 연산 F는 고정해 Architecture 설계를 단순화
- 제한된 Resource 환경에서 모델의 정확도는 최대한 높이면서, 연산량은 최대한 줄임

(a) Standard Convolution

$$N = \odot_{i=1 \dots s} F_i^{L_i}(X_{\langle H_i, W_i, C_i \rangle})$$



(e) Compound Scaling

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

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

$\underline{F}$  고정

$$\left. \begin{array}{l} \text{Memory}(\mathcal{N}) \leq \text{target\_memory} \\ \text{FLOPS}(\mathcal{N}) \leq \text{target\_flops} \end{array} \right\} \begin{array}{l} \text{제한된} \\ \text{Resource} \end{array}$$

$d, w, r$  은 Network Depth, Width, Resolution scaling을 위한 coefficients

# EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks (CVPR 2020)

- Method:

- <Compound Scaling 단계>

- Step 1.  $\phi = 1$ 로 고정한 상태로, 앞서 정의한 수식에  $\alpha, \beta, \gamma$ 을 순차적으로 입력해 최적 값을 구함

- Step 2. step 1에서만 진행해 얻어진  $\alpha, \beta, \gamma$  고정하고  $\phi$ 를 임의적으로 입력해 비교

- 제한된 범위에서  $\phi$ 를 사용해  $\alpha, \beta, \gamma$  Scaling

## Compound Scaling

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

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

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

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

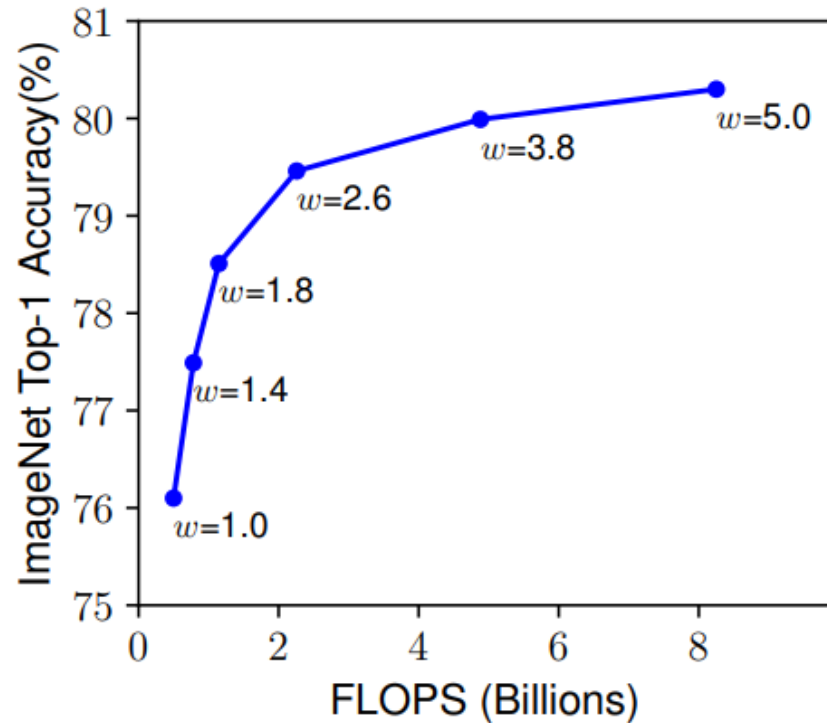
$$\alpha \geq 1, \beta \geq 1, \gamma \geq 1$$

제한 범위

# EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks (CVPR 2020)

- **Method:**

Scaling Up a Baseline Model with Different Network **Width (w)**, Depth (d), and Resolution (r) Coefficients



$w \uparrow$

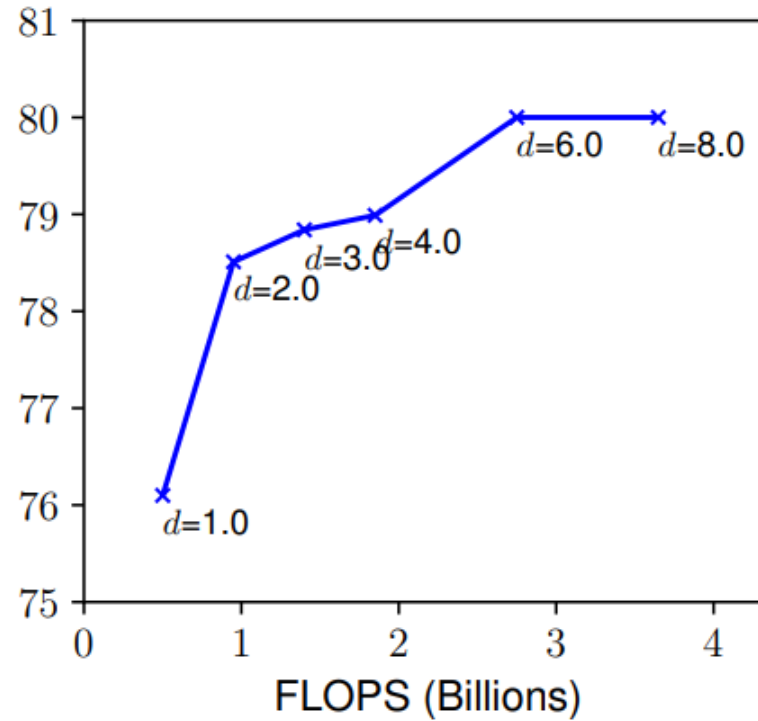
Fine grained feature 잘 포착, 학습이 쉬움

Higher level feature 포착이 어려움

# EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks (CVPR 2020)

- Method:

Scaling Up a Baseline Model with Different Network Width (w), **Depth (d)**, and Resolution (r) Coefficients



$d \uparrow$

Rich, complex feature 잘 포착

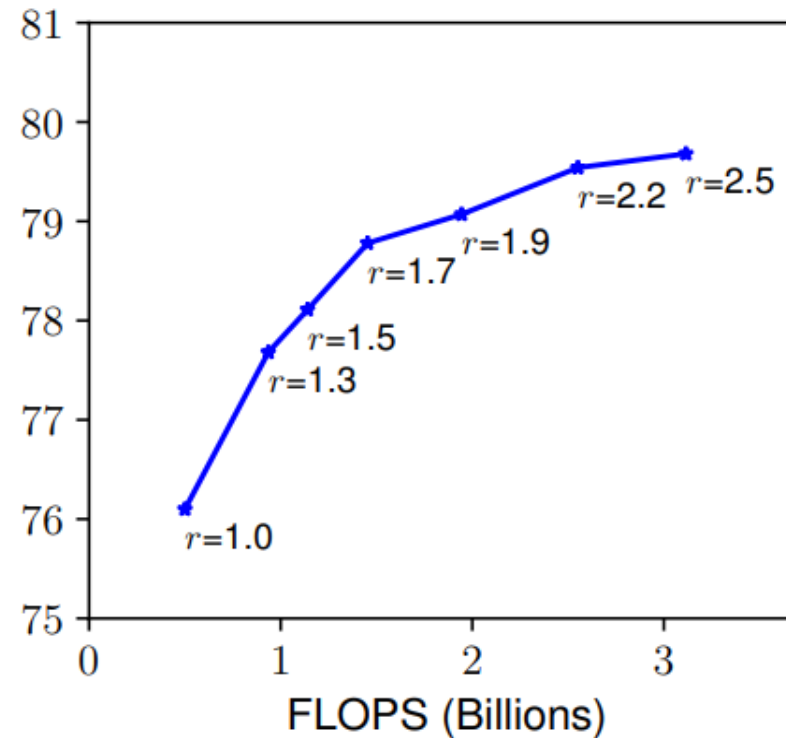
Vanishing gradient 증가

Accuracy gain 감소

# EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks (CVPR 2020)

- Method:

Scaling Up a Baseline Model with Different Network Width (w), Depth (d), and **Resolution (r)** Coefficients



$r \uparrow$   
Accuracy gain 감소

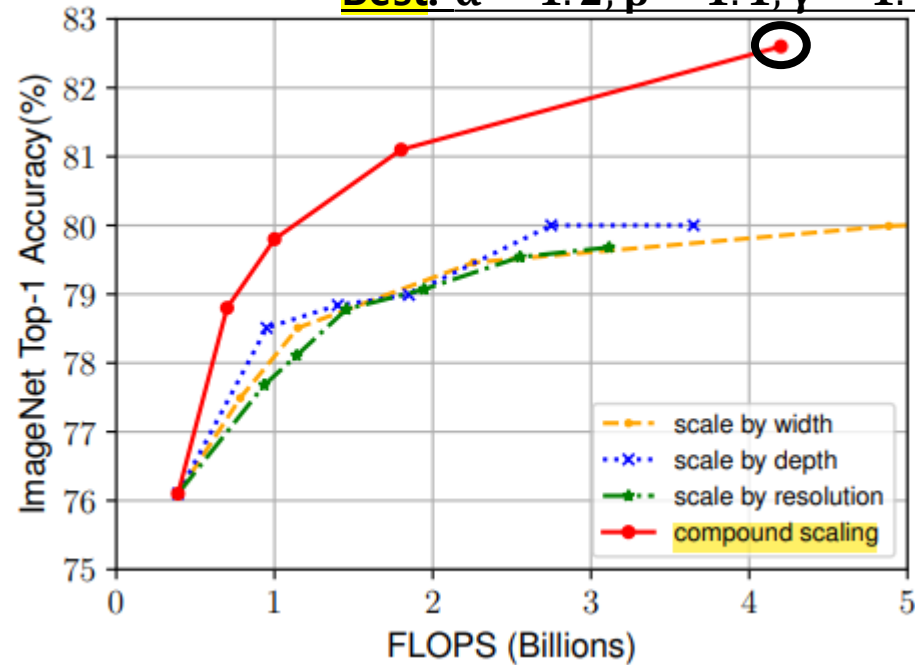
# EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks (CVPR 2020)

- **Method:**

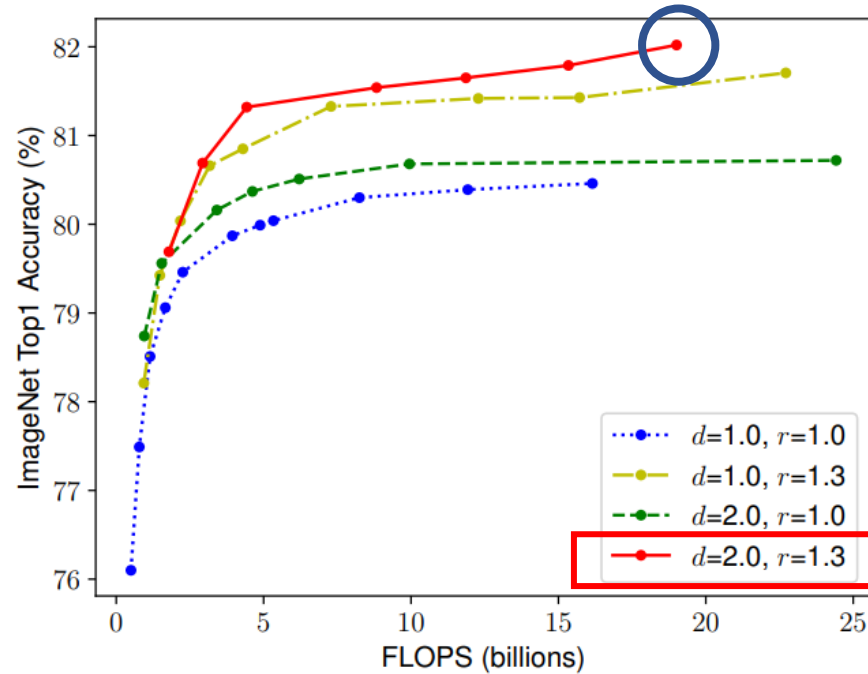
- 3개의 dimension은 독립적이지 않음 => **w, d, r** 각각에 다른 값을 scaling해 균형을 맞추는 필요가 있음
- 네트워크가 깊어지고, 해상도가 높아질수록 **정확도 향상**

## Scaling Network Factor vs Compound Scaling

**Best:**  $\alpha = 1.2$ ;  $\beta = 1.1$ ;  $\gamma = 1.15$



## Scaling Network Width for Different Baseline Networks



정확도 향상

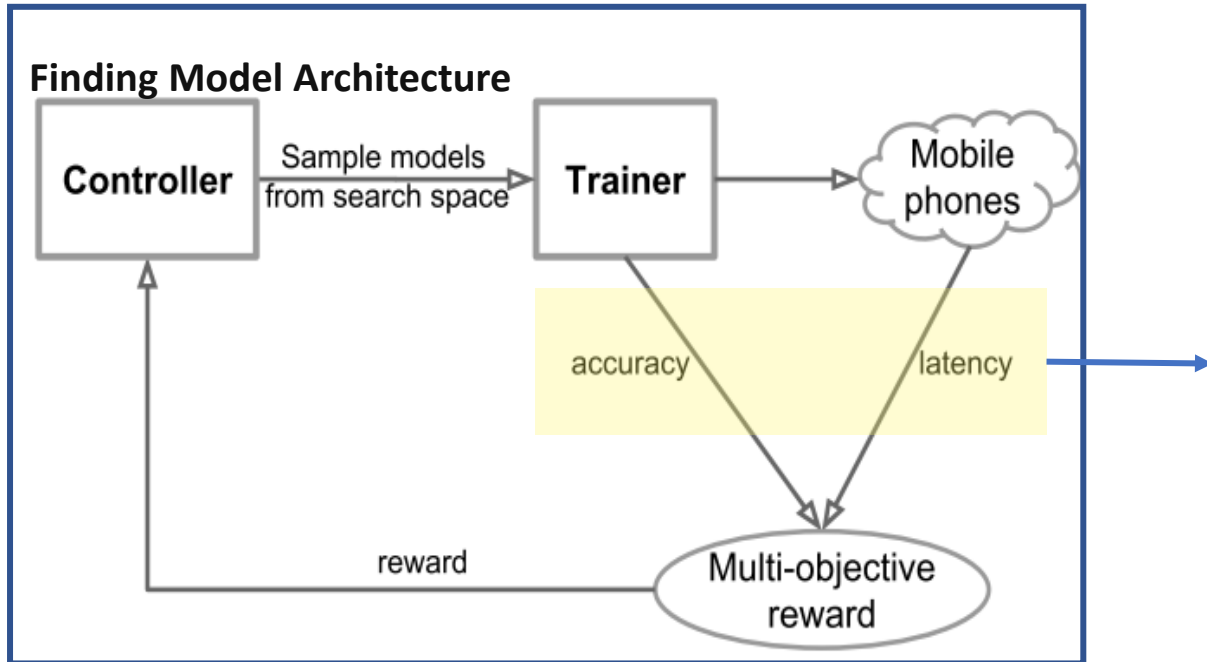


# EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks (CVPR 2020)

- **Method:**

- MNasNet와 유사한 **Neural Architecture Search** 접근 방식을 사용
- **Accuracy**와 **FLOPS**를 모두 최적화하는 multi-objective search를 활용

## MNasNet Approach



Model accuracy (on ImageNet)와 latency을 모델 목표로 사용해 **Best architecture**를 찾음

## Model Architecture :

ImageNet에서 학습하는 데 사용  
Accuracy and Latency 계산

$$\underset{m}{\text{maximize}} \quad ACC(m) \times \left[ \frac{LAT(m)}{T} \right]^w$$

where  $w$  is the weight factor defined as:

$$w = \begin{cases} \alpha, & \text{if } LAT(m) \leq T \\ \beta, & \text{otherwise} \end{cases}$$

Latency가 특정 지정된 값보다 낮을 때 Accuracy가 최대가 되도록 최적의 아키텍처가 달성될 때까지 반복.

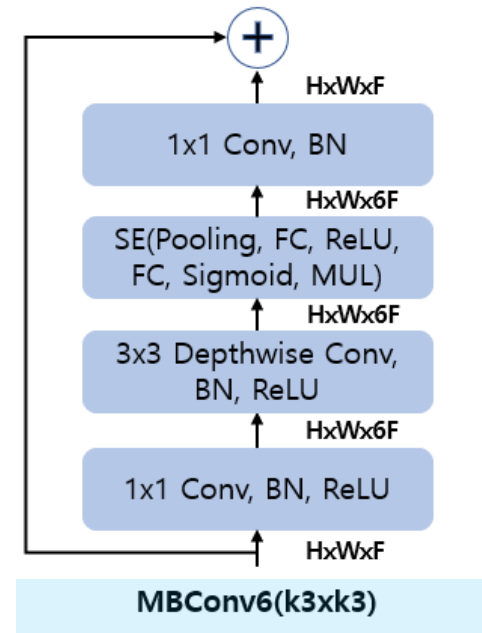
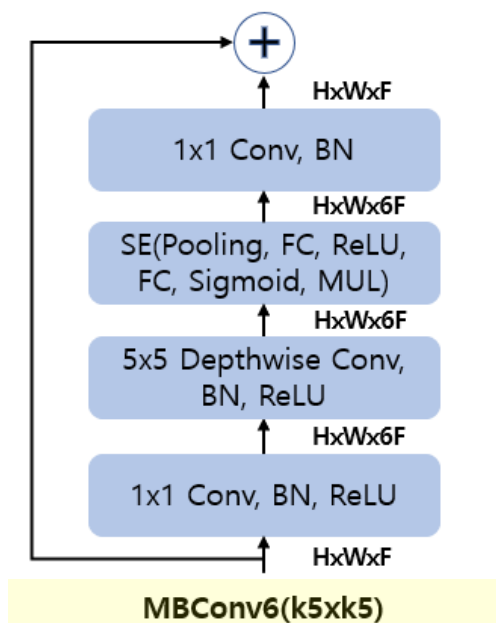
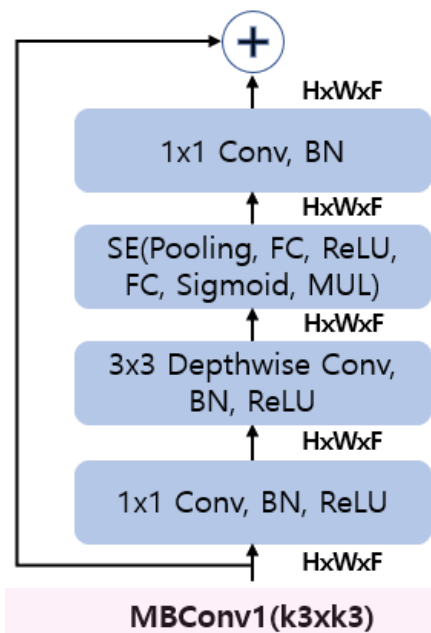
# EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks (CVPR 2020)

- **Method:** NAS에서 정확도 및 FLOPS를 모두 챙길 수 있도록 최적화한 **EfficientNet**

## EfficientNet-B0 Baseline Network

( $\emptyset$ 의 값에 따라 EfficientNet-B0, 1, ..., 7로 분류)

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



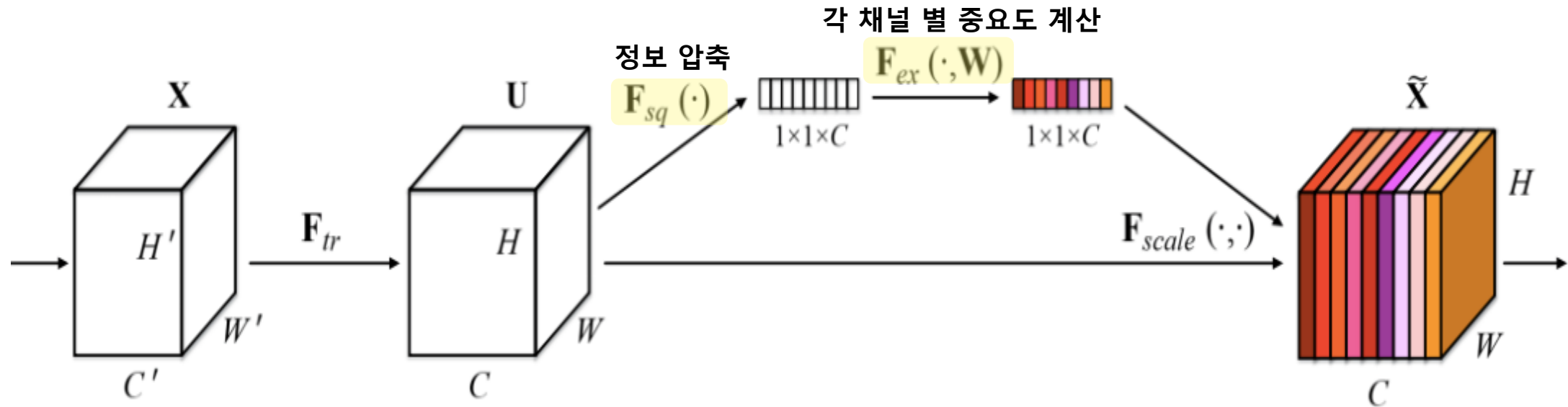
# EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks (CVPR 2020)

- Method:

- Squeeze-and-excitation optimization (SE) 추가

\***Excitation:** 압축된 정보를 weighted layer와 non-linear activation function으로 각 채널 별 중요도를 계산해 기존 input에 곱을 해주는 방식

## Squeeze-and-excitation optimization (SE)



# EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks (CVPR 2020)

- **Experiment:**

- EfficientNet-B0 ~ B7 까지의 성능을 SOTA 모델과 비교
- 같은 성능 대비 Parameters와 FLOPS가 EfficientNet에서 더 낮음 (파라미터 수 약 8.4배 적음)

**EfficientNet Performance Results on ImageNet**

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	-	-

# EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks (CVPR 2020)

- Experiment:

- 기존에 연구된 모델에 compound scaling을 적용한 결과, single-dimension scaling보다 더 좋은 성능을 보임

## Scaling Up MobileNets and ResNet

Single-dimension Scaling vs **Compound Scaling** Better Performance

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>

# EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks (CVPR 2020)

- **Experiment:** 기존의 Top level 모델보다 적은 Parameters와 FLOPS로도 좋은 성능을 보임

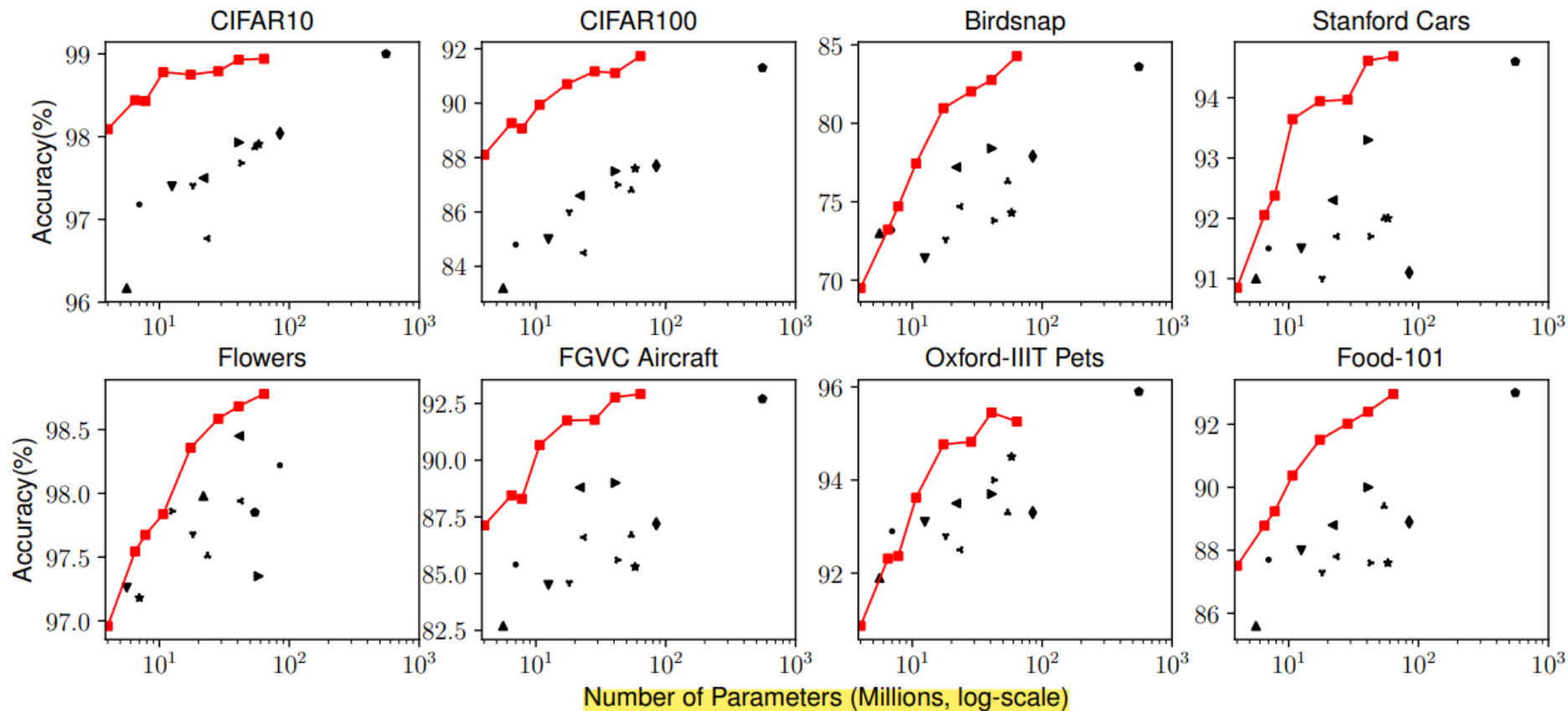
EfficientNet Performance Results on Transfer Learning Datasets

	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	<b>EfficientNet-B7</b>	<b>91.7%</b>	<b>64M</b> (8.7x)
Birdsnap	Inception-v4	81.8%	41M	EfficientNet-B5	82.0%	28M (1.5x)	GPipe	83.6%	556M	<b>EfficientNet-B7</b>	<b>84.3%</b>	<b>64M</b> (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%	-	<b>EfficientNet-B7</b>	<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	<b>EfficientNet-B7</b>	<b>93.0%</b>	<b>64M</b> (8.7x)
Geo-Mean				<b>(4.7x)</b>			<b>(9.6x)</b>					

# EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks (CVPR 2020)

- **Experiment:** 기존의 Top level 모델보다 적은 Parameters와 FLOPS로도 좋은 성능을 보임

Model Parameters vs Transfer Learning Accuracy



▼ DenseNet-201

● GPIPE

▲ Inception-ResNet-v2

◄ ResNet-50

◄ ResNet-101

▼ DenseNet-169

▲ Inception-v1

◄ Inception-v3

► Inception-v4

★ ResNet-152

• DenseNet-121

◆ NASNet-A

■ EfficientNet

Best Performance

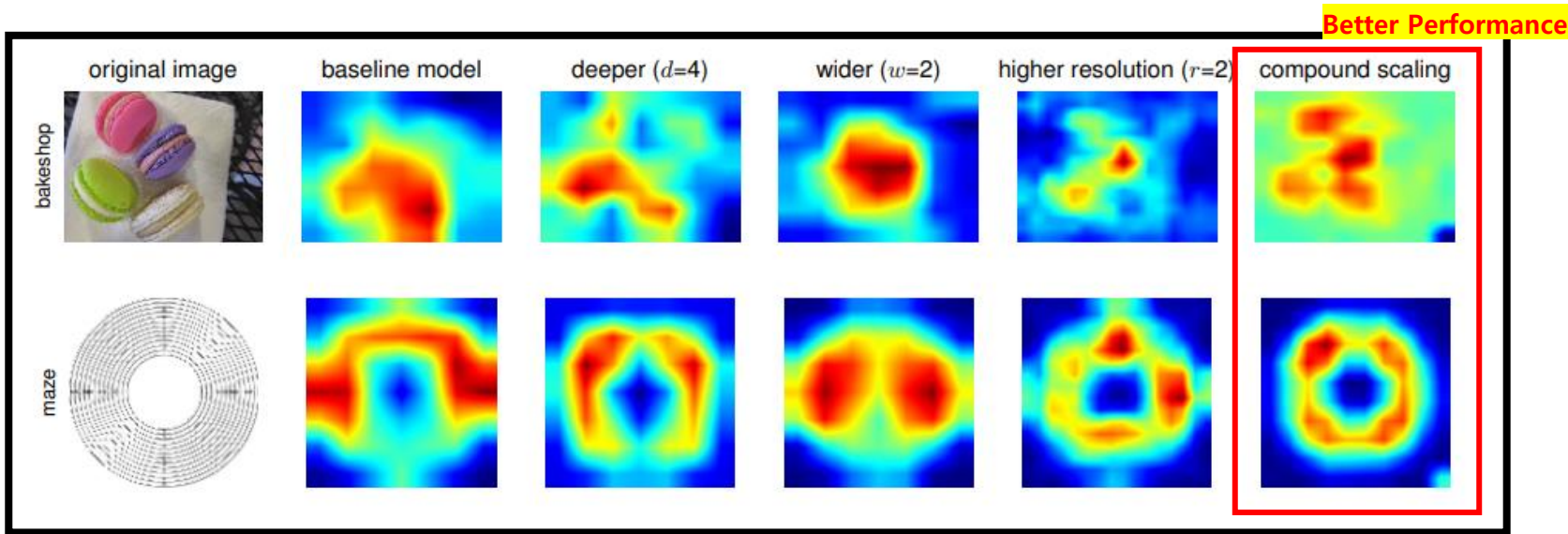


# EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks (CVPR 2020)

- **Experiment:**

- Class Activation Map을 시각화 했을 때 Depth, Width, Resolution을 각각 따로 Scaling up한 것 보다 **Compound Scaling** 했을 때 좀 더 객체들을 잘 담고 있고 정확함

## EfficientNet Performance Results on Transfer Learning Datasets



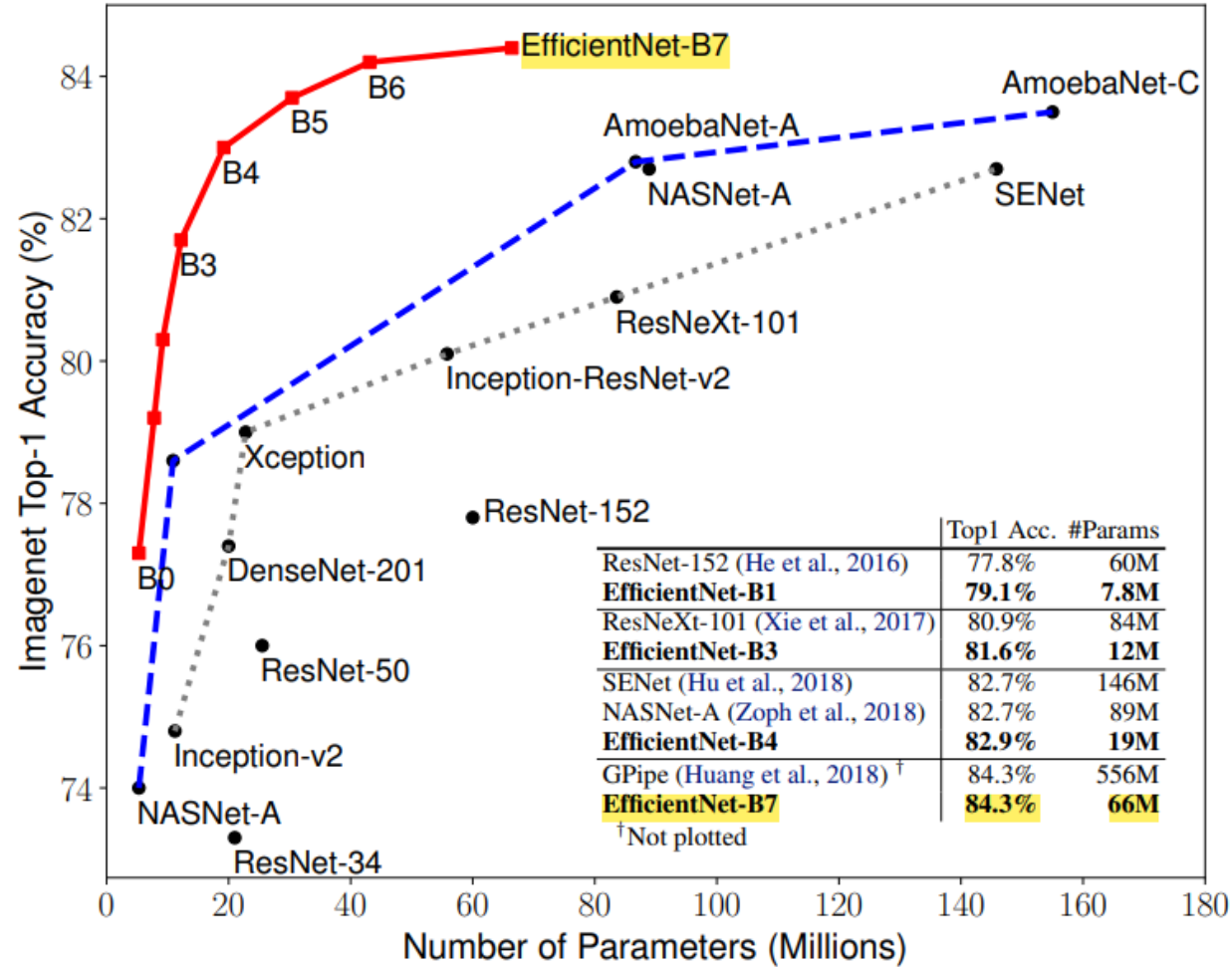
**Scaled Models**

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>



# EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks (CVPR 2020)

Model Size vs ImageNet Accuracy



**한줄평:** 기존의 CNN Scaling 방식을 효과적인 Scaling 방법인 Compound Scaling을 제안해 성능을 효율적으로 향상했다. 단순하면서 직관적인 방법으로 좋은 성능을 달성한 만큼 인상 깊은 논문이다.