

*모델 스케일링

어떠한 큰 물체나 건물을 작은 규모로 변형해 만드는 것.

*컴퓨팅 자원

: 그래피 저장 장치 (GPU)

EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks

Mingxing Tan¹ Quoc V. Le¹

CNN은 학습된 자원 카운트에서 발전 & 성능도↑(자원 사용)
모델 스케일링 & 깊이·폭·해상도 balance

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 ResNets.

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>.

→ Neural architecture Search 사용 + EFFICIENT Net
< 작은 모델 대신 대형 모델 사용 가능!>

1. Introduction

Scaling up ConvNets is widely used to achieve better accuracy. For example, ResNet (He et al., 2016) can be scaled up from ResNet-18 to ResNet-200 by using more layers. Recently, GPipe (Huang et al., 2018) achieved 84.3% ImageNet top-1 accuracy by scaling up a baseline model four

¹Google Research, Brain Team, Mountain View, CA. Correspondence to: Mingxing Tan <tanmingxing@google.com>.

Proceedings of the 36th International Conference on Machine Learning, Long Beach, California, PMLR 97, 2019.

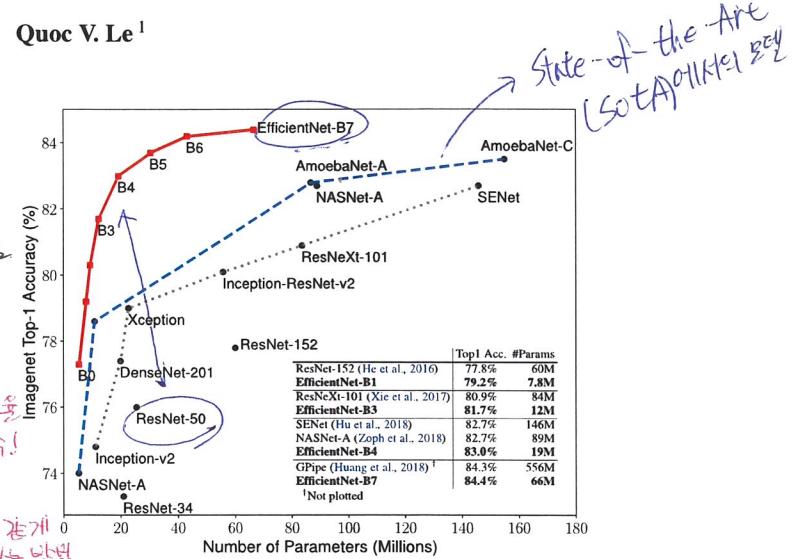


Figure 1. Model Size vs. ImageNet Accuracy. All numbers are for single-crop, single-model. Our EfficientNets significantly outperform other ConvNets. In particular, EfficientNet-B7 achieves new state-of-the-art 84.4% top-1 accuracy but being 8.4x smaller and 6.1x faster than GPipe. EfficientNet-B1 is 7.6x smaller and 5.7x faster than ResNet-152. Details are in Table 2 and 4.

time larger. However, the process of scaling up ConvNets has never been well understood and there are currently many ways to do it. The most common way is to scale up ConvNets by their depth (He et al., 2016) or width (Zagoruyko & Komodakis, 2016). Another less common, but increasingly popular, method is to scale up models by image resolution (Huang et al., 2018). In previous work, it is common to scale only one of the three dimensions (depth, width, and image size.) Though it is possible to scale two or three dimensions arbitrarily, arbitrary scaling requires tedious manual tuning and still often yields sub-optimal accuracy and efficiency.

In this paper, we want to study and rethink the process of scaling up ConvNets. In particular, we investigate the central question: is there a principled method to scale up ConvNets that can achieve better accuracy and efficiency? Our empirical study shows that it is critical to balance all dimensions of network width/depth/resolution, and surprisingly such balance can be achieved by simply scaling each of them with constant ratio. Based on this observation, we propose a simple yet effective compound scaling method. Unlike conventional practice that arbitrary scales these factors, our method uniformly scales network width, depth,

전통적인 방법!

Conv Net 스케일↑
1. 깊이·폭↑
2. 이미지 해상도↑

특정 이미지 크기↑
스케일링 계수↑
→ 자동화, manual tuning
정밀도, 효율성↑ sub-optimal scaling↑
문제점.

E.N.
특정 이미지 크기↑
배치 증가!
↑
상수 비로 초기
스케일링!

* FLOPS (Floating Point Operations Per Second)

컴퓨터의 성능을 수치로 나타낼 때 주로

사용되는 단위. 초당 부동 소수점 연산이라는 뜻.

(Conventional)

scaling

* Partitioning (데이터 베이스 분할)

: 데이터 베이스 여러 부분으로 분할.

중요한 튜닝 기법. 데이터 너무 커져서 처리 시간이
관리 불편, 성능, 가용성 향상 이유.

EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks

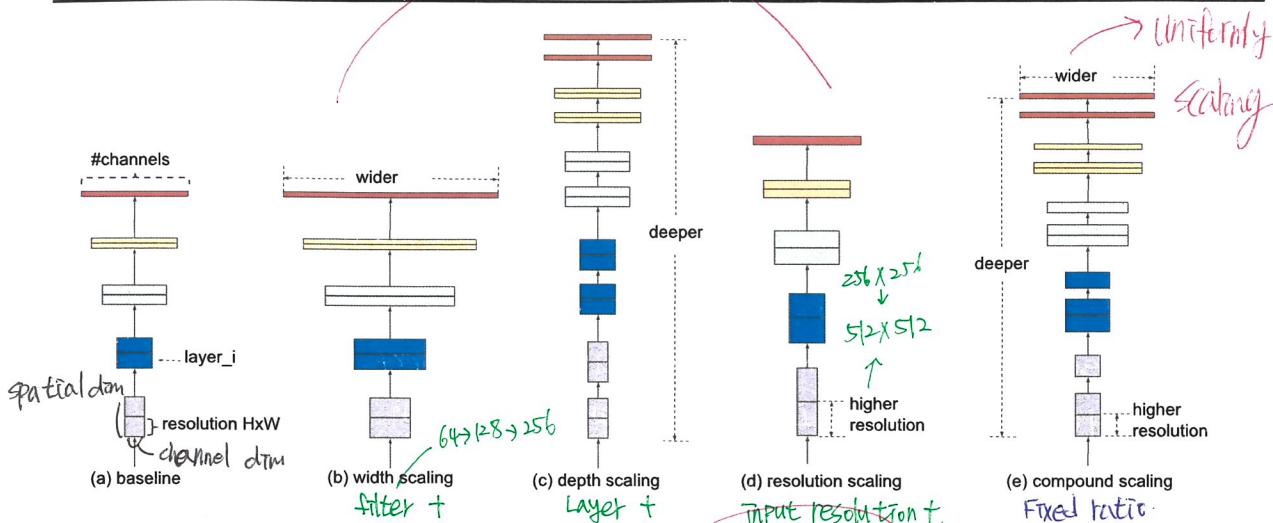


Figure 2. Model Scaling. (a) is a baseline network example; (b)-(d) are conventional scaling that only increases one dimension of network width, depth, or resolution. (e) is our proposed compound scaling method that uniformly scales all three dimensions with a fixed ratio.

고정된 스케일링 규칙
특정 계산 단계로 통일하기
스케일링!

and resolution/with a set of fixed scaling coefficients. For example, if we want to use 2^N times more computational resources, then we can simply increase the network depth by α^N , width by β^N , and image size by γ^N , where α, β, γ are constant coefficients determined by a small grid search on the original small model. Figure 2 illustrates the difference between our scaling method and conventional methods.

Intuitively, the compound scaling method makes sense because if the input image is bigger, then the network needs more layers (to increase the receptive field) and more channels to capture more fine-grained patterns on the bigger image. In fact, previous theoretical (Raghu et al., 2017; Lu et al., 2018) and empirical results (Zagoruyko & Komodakis, 2016) both show that there exists certain relationship between network width and depth, but to our best knowledge, we are the first to empirically quantify the relationship among all three dimensions of network width, depth, and resolution.

We demonstrate that our scaling method work well on existing MobileNets (Howard et al., 2017; Sandler et al., 2018) and ResNet (He et al., 2016). Notably, the effectiveness of model scaling heavily depends on the baseline network; to go even further, we use neural architecture search (Zoph et al., 2017; Le, 2017; Tan et al., 2019) to develop a new baseline network, and scale it up to obtain a family of models, called EfficientNets. Figure 1 summarizes the ImageNet performance, where our EfficientNets significantly outperform other ConvNets. In particular, our EfficientNet-B7 surpasses the best existing GPipe accuracy (Huang et al., 2018), but using 8.4x fewer parameters and running 6.1x faster on inference. Compared to the widely used ResNet-50 (He et al., 2016), our EfficientNet-B4 improves the top-1 accuracy from 76.3% to 83.0% (+6.7%) with similar FLOPS. Besides ImageNet, EfficientNets also transfer well and achieve state-of-the-art performance on various downstream tasks.

EFFICIENT-NET-B7 VS ResNet-50
정확도 83%

76.3%

* Model compression (모델 압축)

사전 학습된 모델에 SVD (특수값 분해) 적용.

사전 학습된 모델을 조금씩 차기내며 특수값 기반으로

내려가는 과정이다. 전부 0 으로 바꿔

sparse matrix 생성하는 network pruning 방식 등을.

of-the-art accuracy on 5 out of 8 widely used datasets, while reducing parameters by up to 21x than existing ConvNets.

→ 다른 모델과 비교했을 때 전략은 다른 패러미터로 다른 분야에서도 효과.

2. Related Work

ConvNet Accuracy: Since AlexNet (Krizhevsky et al., 2012) won the 2012 ImageNet competition, ConvNets have become increasingly more accurate by going bigger: while the 2014 ImageNet winner GoogleNet (Szegedy et al., 2015) achieves 74.8% top-1 accuracy with about 6.8M parameters, the 2017 ImageNet winner SENet (Hu et al., 2018) achieves 82.7% top-1 accuracy with 145M parameters. Recently, GPipe (Huang et al., 2018) further pushes the state-of-the-art ImageNet top-1 validation accuracy to 84.3% using 557M parameters; it is so big that it can only be trained with a specialized pipeline parallelism library by partitioning the network and spreading each part to a different accelerator. While these models are mainly designed for ImageNet, recent studies have shown better ImageNet models also perform better across a variety of transfer learning datasets (Kornblith et al., 2019), and other computer vision tasks such as object detection (He et al., 2016; Tan et al., 2019). Although higher accuracy is critical for many applications, we have already hit the hardware memory limit, and thus further accuracy gain needs better efficiency.

→ ImageNet以外의 데이터에서도 좋은 정확도 가져야 함. but 카드에서 메모리 문제.

ConvNet Efficiency: Deep ConvNets are often over-parameterized. Model compression (Han et al., 2016; He et al., 2018; Yang et al., 2018) is a common way to reduce model size by trading accuracy for efficiency. As mobile phones become ubiquitous, it is also common to hand-craft efficient mobile-size ConvNets, such as SqueezeNets (Iandola et al., 2016; Gholami et al., 2018), MobileNets (Howard et al., 2017; Sandler et al., 2018), and ShuffleNets (Zhang et al., 2018).

AlexNet 승리

크게 하면서 정확도↑
(Conv Net)

GoogleNet

SENet

G Pipe: 정확도↑
파라미터↓

Image Net 위한

최근 연구

다양한 데이터셋

수행: ImageNet 정확도↑

다른 분야

Object detection

최근 연구

다양한 데이터셋

수행: ImageNet 정확도↑

다른 분야

주제적인 정확도는 중요! 필요!

Model compression

효율성↑

정확도↑

↓

트레이닝 → 테스트(예측)

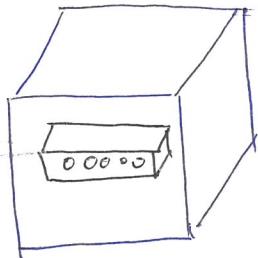
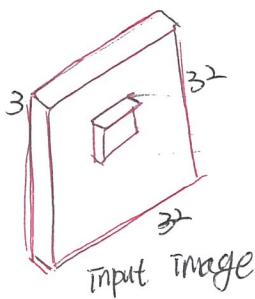
ex) squeeze Net

MobileNet

ShuffleNet

* receptive field (수용공간)

: 어떤 출력레이어의 뉴런 하나에 영향을 미치는 입력뉴런들의 공간의 크기 \Rightarrow 필터의 크기



입력이 $32 \times 32 \times 3$ 인 경우 필터 $5 \times 5 \times 3$ 일 때
receptive field는 $5 \times 5 \times 3$. (\Rightarrow 필터의 크기).

* fine grained

: Vision algorithm에서 coarse-to-fine-strategy 존재.

\rightarrow 간략화된 이미지 (lowered-resolution)에서 세밀화된 이미지 (the finest resolution)으로 변화시켜 KLT tracker 등 optimization 기반 알고리즘의 convergence basin 넓여주는 역할.

coarse-grained programming : 큰 덩어리로 나누는 것

fine-grained programming : 세밀한 단위로 프로그래밍 나눠서 고려

* GPipe

: 파이프라인 병렬화 및 체크 포인팅 학습.

tensorflow로 구현된 GPipe는 1개의 GPU로 학습 불가능한 거대한 모델을 가능한 효율적으로 훈련 위한 모델 학습 방법
GPipe 도입하면 메모리 사용량이 큰 모델을 여러 GPU에서 병렬처리 해 GPU를 효율적으로 사용하면서 빠른 학습 가능.

① 파이프라인 병렬화.

: 여러 부분으로 나눈 모델 각 GPU에 올리는 방식 (모델 병렬화와 같음). But 모델 훈련이 여러 GPU에서 일정시간 동안 동시에 연산될 수 있도록 미니 배치를 여러개로 조건 마이크로 배치를 학습에 활용.
 \rightarrow 중첩되어 처리되는 작업량 多. 같은 양의 데이터 더 짧은 시간 내 처리.

ex) 각 GPU가 100장 미니 배치 처리 10초. \rightarrow 2nd GPU는 10초 뒤 작업 \rightarrow 여러 GPU 최대
But, 10장 마이크로 배치 10초 훈련 \rightarrow 2nd GPU는 1초 뒤 작업 동시 동원.

마이크로 배치 크기 작아질수록 다음 GPU가 더 빨리 시작 \rightarrow 유체 시간 줄여짐. But 마이크로 배치 크기 무한정 줄일 수 X.

② 체크 포인팅.

: 은닉층 결과값 일부 버려 모델 차지하는 메모리 크기 ↓, 모델 나눈 부분에 체크 포인트 적으로 분리된 모델 서로 연결하는 은닉층만 메모리에 들림.
연결하는 은닉층만 모델 내부의 은닉층 결과값을 기억하지 않고 뇌. 이 같은 역전파 단계에서 입력층 결과값과 계산 로직을 가지고 재계산 가능.

\Rightarrow 모든 은닉층 결과값을 메모리에 올려놓고 없기 때문에 모델이 훈련에 사용하는 메모리 공간 줄일 수 O.

역전파 과정에서 순전파 계산을 한번 더 진행하기에 모델 학습 속도 약 25% 느려지지만,
마이크로 배치 수에 반비례하기 때문에 메모리 사용량 줄일 수 있어 속도 향상가치 在

* Overhead : 어떤 처리를 하기 위해 들어가는 간접적인 처리 시간. 메모리 ex) A' 처리 실행 10초

* expressive power : 컴퓨터 과학에서 언어의 표현력은 해당 언더로 표현하고 전달될 수 있는 아이디어의 폭. 표현력이 높을수록 표현하는데 사용할 수 있는 아이디어의 다양성과 양이 많아짐.

Neural architecture 3.1

EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks

최근에는 풀 커넥션
 디비전을 통한
 네트워크 폭, 크기
 및 디비전 크기
 조정하면서
 커널의 적용 범위
 확장↑
 특수 베이스↑
 모델 훈련 성능
 Model Scaling
 다른 지원 강화
 ↓
 ResNet : 18, 200
 (최신 고성능)
 WideResNet
 MobileNet
 ↳ 커널 수 증가
 ↓
 스케일 바벨 수
 Inv. Net
 inputting ↑

(Zhang et al., 2018; Ma et al., 2018). Recently, neural architecture search becomes increasingly popular in designing efficient mobile-size ConvNets (Tan et al., 2019; Cai et al., 2019), and achieves even better efficiency than hand-crafted mobile ConvNets by extensively tuning the network width, depth, convolution kernel types and sizes. However, it is unclear how to apply these techniques for larger models that have much larger design space and much more expensive tuning cost. In this paper, we aim to study model efficiency for super large ConvNets that surpass state-of-the-art accuracy. To achieve this goal, we resort to model scaling. ↳ 본 연구는 과정과 동일하고 있는 것 같은 연구인가?
Model Scaling: There are many ways to scale a ConvNet for different resource constraints. ResNet (He et al., 2016) can be scaled down (e.g., ResNet-18) or up (e.g., ResNet-200) by adjusting network depth (#layers), while WideResNet (Zagoruyko & Komodakis, 2016) and MobileNets (Howard et al., 2017) can be scaled by network width (#channels). It is also well-recognized that bigger input image size will help accuracy with the overhead of more FLOPS. Although prior studies (Raghu et al., 2017; Lin & Jegelka, 2018; Sharir & Shashua, 2018; Lu et al., 2018) 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. Our work systematically and empirically studies ConvNet scaling for all three dimensions of network width, depth, and resolutions.

i. Figure 2(a) illustrate a representative ConvNet, where the spatial dimension is gradually shrunk but the channel dimension is expanded over layers, for example, from initial input shape $(224, 224, 3)$ to final output shape $(7, 7, 512)$. ↳ ConvNet의 차원의 spatial dim↓, channel↑
 Unlike regular ConvNet designs that mostly focus on finding the best layer architecture F_i , model scaling tries to expand the network length (L_i), width (C_i), and/or resolution (H_i, W_i) without changing F_i predefined in the baseline network. By fixing F_i , model scaling simplifies the design problem for new resource constraints, but it still remains a large design space to explore different L_i, C_i, H_i, W_i for each layer. In order to further reduce the design space, we restrict that all layers must be scaled uniformly with constant ratio. Our target is to maximize the model accuracy for any given resource constraints, which can be formulated as an optimization problem:

$$\begin{aligned}
 & \text{target} \\
 & \max_{d, w, r} \text{Accuracy}(\mathcal{N}(d, w, r)) \quad (\text{자원 제약 조건}) \\
 & \text{s.t.} \quad \mathcal{N}(d, w, r) = \bigodot_{i=1 \dots s} \hat{\mathcal{F}}_i^{d, \hat{L}_i} (X_{(r \cdot \hat{H}_i, r \cdot \hat{W}_i, w \cdot \hat{C}_i)}) \\
 & \quad \text{Memory}(\mathcal{N}) \leq \text{target_memory} \\
 & \quad \text{FLOPS}(\mathcal{N}) \leq \text{target_flops}
 \end{aligned} \tag{2}$$

where w, d, r are coefficients for scaling network width, depth, and resolution; $\hat{\mathcal{F}}_i, \hat{L}_i, \hat{H}_i, \hat{W}_i, \hat{C}_i$ are predefined parameters in baseline network (see Table 1 as an example).

3. Compound Model Scaling

In this section, we will formulate the scaling problem, study different approaches, and propose our new scaling method.

3.1. Problem Formulation

A ConvNet Layer i can be defined as a function: $Y_i = \mathcal{F}_i(X_i)$, where \mathcal{F}_i is the operator, Y_i is output tensor, X_i is input tensor, with tensor shape $\langle H_i, W_i, C_i \rangle^1$, where H_i and W_i are spatial dimension and C_i is the channel dimension. A ConvNet \mathcal{N} can be represented by a list of composed layers: $\mathcal{N} = \mathcal{F}_k \odot \dots \odot \mathcal{F}_2 \odot \mathcal{F}_1(X_1) = \bigcirc_{j=1\dots k} \mathcal{F}_j(X_1)$. In practice, ConvNet layers are often partitioned into multiple stages and all layers in each stage share the same architecture; for example, ResNet (He et al., 2016) has five stages, and all layers in each stage has the same convolutional type except the first layer (performs down-sampling). Therefore, we can define a ConvNet as:

$$N = \bigodot_{i=1}^s F_i^{L_i}(X_{(H_i, W_i)C_i}) \quad (1)$$

where $F_i^{L_i}$ denotes layer F_i is repeated L_i times in stage i .
 $\langle H_i, W_i, C_i \rangle$ denotes the shape of input tensor X of layer i .

¹For the sake of simplicity, we omit batch dimension.

$$Y_i = F_i(X_i) \quad \xrightarrow{\text{convNet}} \quad N = F_0 \odot \dots \odot F_k \odot F_i(X_i)$$

output Input (tensor : H_i, W_i, C_i)

$i \in \text{ConvNet layer.}$

$F_i^{L_i}$: F_i 를 L_i 번 반복 (\neg 단계에서) \checkmark
 $\langle H_i, W_i, G_i \rangle$: 구조에서 인플루에서 x 의 크기

- $\text{ConvNet} = \bigcup_{T=1, \dots, S} F_T^{L_i}(x_{\langle H_i, W_i, C_i \rangle})$
 ConvNet은 여러 단계로 나누고
 각 단계의 층은 같은 구조 갖기 때문

* expressive power : 컴퓨터 과학에서 언어의 표현력은 해당 언어로 표현되고 전달될 수 있는 아이디어의 폭. 표현력이 높을수록 표현하는데 사용할 수 있는 아이디어의 다양성과 양이 많아짐.

i. Figure 2(a) illustrate a representative ConvNet, where the spatial dimension is gradually shrunk but the channel dimension is expanded over layers, for example, from initial input shape $(224, 224, 3)$ to final output shape $(7, 7, 512)$.
 Unlike regular ConvNet designs that mostly focus on finding the best layer architecture \mathcal{F}_i , model scaling tries to expand the network length (L_i), width (C_i), and/or resolution (H_i, W_i) without changing \mathcal{F}_i predefined in the baseline network. By fixing \mathcal{F}_i , model scaling simplifies the design problem for new resource constraints, but it still remains a large design space to explore different (L_i, C_i, H_i, W_i) for each layer. In order to further reduce the design space, we restrict that all layers must be scaled uniformly with constant ratio. Our target is to maximize the model accuracy for any given resource constraints, which can be formulated as an optimization problem:

$$\begin{aligned} & \max_{d,w,r} \text{Accuracy}(\mathcal{N}(d, w, r)) \quad \langle \text{자연 제한 } \text{X} \rangle \\ \text{s.t.} \quad & \mathcal{N}(d, w, r) = \bigodot_{i=1 \dots s} \hat{\mathcal{F}}_i^{d \cdot \hat{L}_i} (X_{\langle r \cdot \hat{H}_i, r \cdot \hat{W}_i, w \cdot \hat{C}_i \rangle}) \\ & \text{Memory}(\mathcal{N}) \leq \text{target_memory} \\ & \text{FLOPS}(\mathcal{N}) \leq \text{target_flops} \end{aligned}$$

where w, d, r are coefficients for sealing network width, depth, and resolution; $\hat{F}_i, \hat{L}_i, \hat{H}_i, \hat{W}_i, \hat{C}_i$ are predefined parameters in baseline network (see Table 1 as an example).

3.2 Scaling Dimension

The main difficulty of problem 2 is that the optimal d, w, r depend on each other and the values change under different resource constraints. Due to this difficulty, conventional methods mostly scale ConvNets in one of these dimensions:

Depth (d): Scaling network depth is the most common way used by many ConvNets (He et al., 2016; Huang et al., 2017; Szegedy et al., 2015; 2016). The intuition is that deeper ConvNet can capture richer and more complex features, and generalize well on new tasks. However, deeper networks are also more difficult to train due to the vanishing gradient problem (Zagoruyko & Komodakis, 2016). Although several techniques, such as skip connections (He et al., 2016) and batch normalization (Ioffe & Szegedy, 2015), alleviate the training problem, the accuracy gain of very deep network diminishes: for example, ResNet-1000 has similar accuracy as ResNet-101 even though it has much more layers. Figure 3 (middle) shows our empirical study on scaling a baseline model with different depth coefficient d , further suggesting the diminishing accuracy return for very deep ConvNets. ↗

Width (w): Scaling network width is commonly used for small size models (Howard et al., 2017; Sandler et al., 2018;

Scaling Network

↳ 작은 크기의 모데 사용

*high level feature

: 고차원의 특징

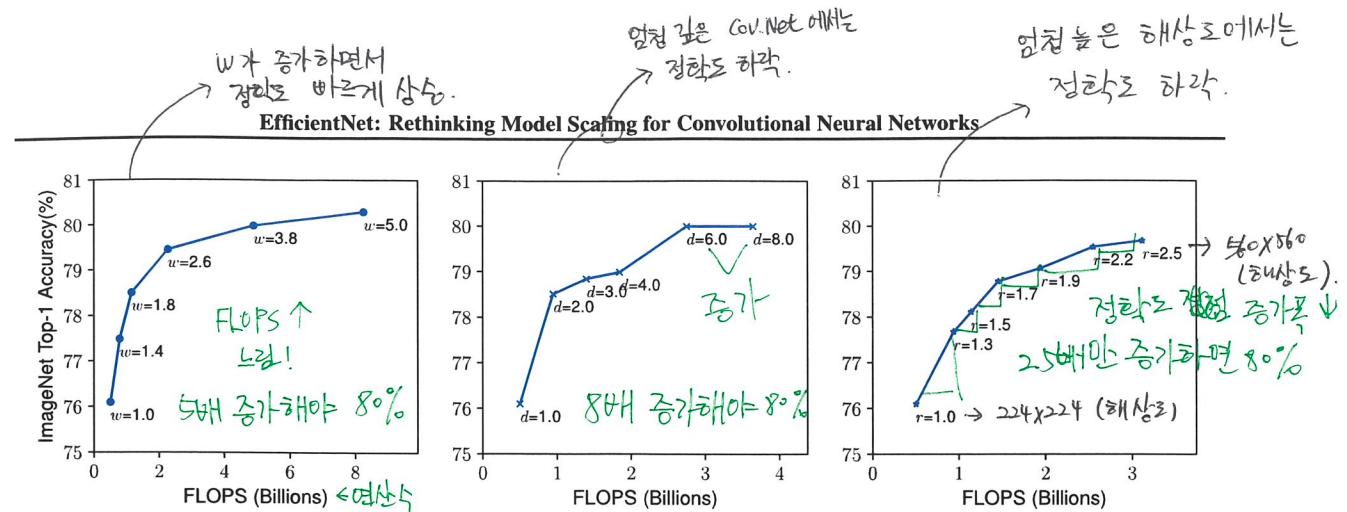


Figure 3. Scaling Up a Baseline Model with Different Network Width (w), Depth (d), and Resolution (r) Coefficients. Bigger networks with larger width, depth, or resolution tend to achieve higher accuracy, but the accuracy gain quickly saturates after reaching 80%, demonstrating the limitation of single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

Width

통령
세대별 feature 차이

• 출전 식물
• 전반적
• 매우 넓고 얕은
Network
• 높은 수준의 퍼포먼스
• 증기 어렵

Resolution
• img 해상도 ↑
• fine grained pattern 잘 잡아냄
• GConvNet
• 224x224 48Flops
• 299x299 134Flops
• 331x331 148Flops
• 480x480 480Flops
• 600x600 600Flops
• 높은 해상도로
• object detection

• 224x224 48Flops
• 299x299 134Flops
• 331x331 148Flops
• 480x480 480Flops
• 600x600 600Flops
• 높은 해상도로
• object detection

• GConvNet 48Flops

Observation 1 – Scaling up any dimension of network width, depth, or resolution improves accuracy, but the accuracy gain diminishes for bigger models.

→ width, depth, resolution 크게 카면 정학도 ↑

3.3. Compound Scaling

but 더 큰 모델 때에는 정학도 ↓

We empirically observe that different scaling dimensions are not independent. Intuitively, for higher resolution images, we should increase network depth, such that the larger receptive fields can help capture similar features that include more pixels in bigger images. Correspondingly, we should also increase network width when resolution is higher, in

²In some literature, scaling number of channels is called “depth multiplier”, which means the same as our width coefficient w .

해상도 ↑ → 깊이 ↑

큰 이미지에서 더 많은 픽셀 포함하는

더 큰 수용영역(필터) 가 배수한 특징

포착 가능

해상도 ↑ → 폭 ↑.

더 많은 픽셀에서 세밀한 패턴

발견 가능

엄청 같은 ConvNet에서는 정학도 하락.

엄청 높은 해상도에서는 정학도 하락.

EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks

→ width scaling under different network depths and resolution

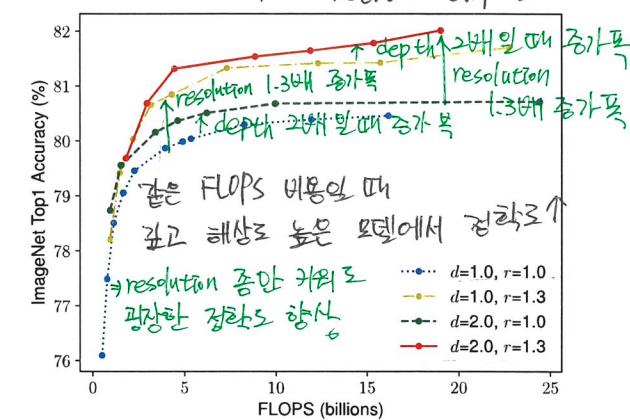


Figure 4. Scaling Network Width for Different Baseline Networks. Each dot in a line denotes a model with different width coefficient (w). All baseline networks are from Table 1. The first baseline network ($d=1.0, r=1.0$) has 18 convolutional layers with resolution 224x224, while the last baseline ($d=2.0, r=1.3$) has 36 layers with resolution 299x299.

order to capture more fine-grained patterns with more pixels in high resolution images. These intuitions suggest that we need to coordinate and balance different scaling dimensions rather than conventional single-dimension scaling. → 3 dim 조화 중요

To validate our intuitions, we compare width scaling under different network depths and resolutions, as shown in Figure 4. If we only scale network width (w) without changing depth ($d=1.0$) and resolution ($r=1.0$), the accuracy saturates quickly. With deeper ($d=2.0$) and higher resolution ($r=2.0$), width scaling achieves much better accuracy under the same FLOPs cost. These results lead us to the second observation:

Observation 2 – In order to pursue better accuracy and efficiency, it is critical to balance all dimensions of network width, depth, and resolution during ConvNet scaling.

→ 정학도, 효율성 더 낫게 하려면 3 dim 조화 중요.

스케일링

적당히

적당히

적당히

적당히

적당히

적당히

* multi-objective

여러 개의 객체 검출.

* Squeeze and excitation

: SE Net

* Building Block

정보 통신 분야에서 컴퓨터를 구성하는 여러 종류의 기본 회로 의미.

$$\text{equation} \Rightarrow \max_{d, w, r} \text{Accuracy}(N(d, w, r)) \rightarrow N(d, w, r) = \bigcup_{T=1}^L (x_{r, A_T, T}, f_{r, A_T, T}, d, w, r)$$

$$\text{Memory}(N) \leq \text{target_memory}, \text{FLOPS}(N) \leq \text{target_flops}$$

EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks

이전에도 이런 실험 존재.
ut, 예상된 결과와 In fact, a few prior work (Zoph et al., 2018; Real et al., 2019) have already tried to arbitrarily balance network width and depth, but they all require tedious manual tuning.

Compound scaling method In this paper, we propose a new **compound scaling method**, which uses a compound coefficient ϕ to uniformly scale network width, depth, and resolution (in a principled way):

원칙적인 방식으로!

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

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

$$\text{resolution: } r = \gamma^\phi \quad (3)$$

$$\frac{d}{w} \cdot \frac{w}{r} \cdot \frac{r}{d} \approx 2 \quad \text{s.t. } \alpha \cdot \beta^2 \cdot \gamma^2 \approx 2$$

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

α, β, γ : small such grid 아니 결정된 상수

where α, β, γ are constants that can be determined by a small grid search. Intuitively, ϕ is a user-specified coefficient that controls how many more resources are available for model scaling while α, β, γ specify how to assign these extra resources to network width, depth, and resolution respectively. Notably, the FLOPS of a regular convolution op is proportional to d, w^2, r^2 , i.e., doubling network depth will double FLOPS, but doubling network width or resolution will increase FLOPS by four times. Since convolution ops usually dominate the computation cost in ConvNets, scaling a ConvNet (with equation 3) will approximately increase total FLOPS by $(\alpha \cdot \beta^2 \cdot \gamma^2)^\phi$. In this paper, we constraint $\alpha \cdot \beta^2 \cdot \gamma^2 \approx 2$ such that for any new ϕ , the total FLOPS will approximately³ increase by 2^ϕ .

$\rightarrow d \times 2 \rightarrow \text{FLOPS} \times 2 \rightarrow w \times 2 \rightarrow \text{FLOPS} \times 4 \rightarrow r \times 2 \rightarrow \text{FLOPS} \times 16 \rightarrow \text{FLOPS} = 2^{\phi} \text{까지 증가.}$

4. EfficientNet Architecture

\rightarrow MNAS Net.

Since model scaling does not change layer operators \hat{F}_i (in baseline network) having a good baseline network is also critical. We will evaluate our scaling method using existing ConvNets, but in order to better demonstrate the effectiveness of our scaling method, we have also developed a new mobile-size baseline, called EfficientNet.

Inspired by (Tan et al., 2019), we develop our baseline network by leveraging a **multi-objective** neural architecture search that optimizes both accuracy and FLOPS. Specifically, we use the same search space as (Tan et al., 2019), and use $ACC(m) \times [FLOPS(m)/T]^{\psi}$ as the optimization goal, where $ACC(m)$ and $FLOPS(m)$ denote the accuracy and FLOPS of model m , T is the target FLOPS and $\psi = -0.07$ is a hyperparameter for controlling the trade-off between accuracy and FLOPS. Unlike (Tan et al., 2019; Cai et al., 2019), here we optimize FLOPS rather than latency since we are not targeting any specific hardware device. Our search produces an efficient network, which we name EfficientNet-B0. Since we use the same search space as (Tan et al., 2019), the architecture is similar to Mnas-

³FLOPS may differ from theoretical value due to rounding.

최적화 모드 사용

Model에 대한 정량화된 성능

하이퍼파라미터 (정확도 & FLOPS 모두 함께 조절)

→ 차이/간 모다 FLOPS 타겟으로 쓰는 이유?

특별한 하드웨어 타겟팅 X.

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. → 가중!

Stage i	Operator \hat{F}_i	Resolution $\hat{H}_i \times \hat{W}_i$	#Channels \hat{C}_i	#Layers \hat{L}_i
1	Conv3x3	224 × 224	32	1
2	MBCConv1, k3x3	112 × 112	16	1
3	MBCConv, k3x3	112 × 112	24	2
4	MBCConv6, k5x5	56 × 56	40	2
5	MBCConv6, k3x3	28 × 28	80	3
6	MBCConv6, k5x5	14 × 14	112	3
7	MBCConv6, k5x5	14 × 14	192	4
8	MBCConv6, k3x3	7 × 7	320	1
9	Conv1x1 & Pooling & FC	7 × 7	1280	1

Baseline

EfficientNet V0의
- MB conv
: mobile inverted
bottleneck

- SENet

↓

Compound scaling

EfficientNet 구조 설계 (baseline)
성능 적용

Net, except our EfficientNet-B0 is slightly bigger due to the larger FLOPS target (our FLOPS target is 400M). Table 1 shows the architecture of EfficientNet-B0. Its main building block is mobile inverted bottleneck MBCConv (Sandler et al., 2018; Tan et al., 2019), to which we also add squeeze-and-excitation optimization (Hu et al., 2018).

Starting from the baseline EfficientNet-B0, we apply our compound scaling method to scale it up with two steps:

- STEP 1: we first fix $\phi = 1$, assuming twice more resources available, and do a small grid search of α, β, γ based on Equation 2 and ③. In particular, we find the best values for EfficientNet-B0 are $\alpha = 1.2, \beta = 1.1, \gamma = 1.15$, under constraint of $\alpha \cdot \beta^2 \cdot \gamma^2 \approx 2$.
- STEP 2: we then fix α, β, γ as constants and scale up baseline network with different ϕ using Equation 3, to obtain EfficientNet-B1 to B7 (Details in Table 2).

→ EfficientNet baseline 기반 스케일링!

Notably, it is possible to achieve even better performance by searching for α, β, γ directly around a large model, but the search cost becomes prohibitively expensive on larger models. Our method solves this issue by only doing search once on the small baseline network (step 1), and then use the same scaling coefficients for all other models (step 2).

→ search cost 해결법.

비교 크로스의 α, β, γ

찾으며 성능.

하지만, search 하는

회사 ↑ 홍보.

Step 1에서 작은

1번시작 — 베이스인 네트워크 시작해.

Step 2에서 같은 계수를

다른 모델에 대해서도

사용한다!

5. Experiments

In this section, we will first evaluate our scaling method on existing ConvNets and the new proposed EfficientNets.

5.1. Scaling Up MobileNets and ResNets

As a proof of concept, we first apply our scaling method to the widely-used MobileNets (Howard et al., 2017; Sandler et al., 2018) and ResNet (He et al., 2016). Table 3 shows the ImageNet results of scaling them in different ways. Compared to other single-dimension scaling methods, our compound scaling method improves the accuracy on all these models, suggesting the effectiveness of our proposed scaling method for general existing ConvNets.

→ Compound scaling 정확도 높고

effectiveness도 좋다!

MnasNet과 비슷

but, 더 큰 FLOPs 타겟이니까

EfficientNet-B0가 약간 더 큼

* MnasNet

Neural Architecture Search for Mobile

모바일 위한 최적화된 모델 찾아주는 모델 X

→ 두에 설명!

* Mobile Inverted Bottle Neck.

→ MobileNet V2 //

(DSC)

• 전의 MobileNet : 일반적 conv 무겁 → depthwise separable conv 사용

Mobile Net V2 기법

① depthwise separable conv. ← Mobile Net 참고하세요!

② Inverted Residuals

Linear Bottlenecks

- 방식 : Residual Learning 방식을 통해 네트워크를 깊게 구성할 경우, 연산량 ↑

1 by 1 conv 활용해 채널 수 감소. cf. Residual Learning? 학습을 위한 weight layer 거치고 초기 input 값을 더해주는 형태

- 일반적 : 3x3 conv 실시

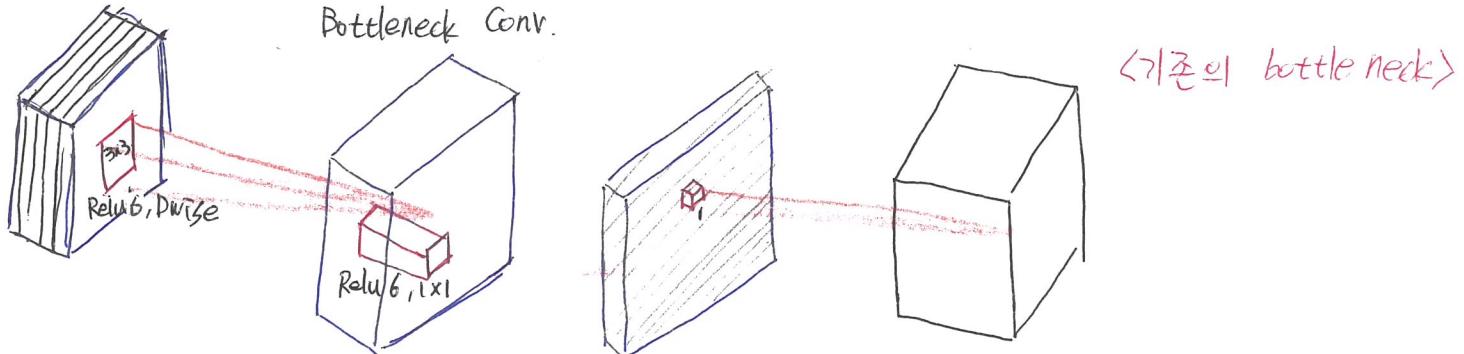
- Bottle Neck 방식 : 1x1 conv로 채널수 ↓ → 3x3 conv 실시 → 다시 1x1 conv로 채널수 복구

- 왜 Inverted Residuals인가?

: 기존의 bottleNeck은 채널감소 → 학습(DSC) → 채널복구.

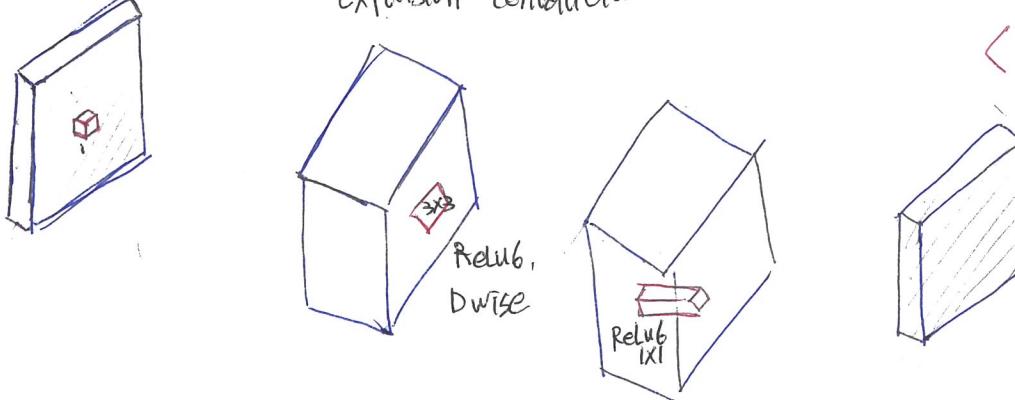
Inverted Residual은 채널 증가 → 학습(DSC) → 채널감소

<Inverted 됨> ⇒ input data와 output data의 크기 작아 메모리 효율적.



Expansion Convolution Block

<Inverted Residual>



* MnasNet

: Neural Architecture Search for Mobile.

Nas → Neural Architecture search
구글에서 딥러닝 모델 아키텍처 탐색
보다는 딥러닝

MnasNet

: 모바일에 최적화된 네트워크
찾아주는 모델

EfficientNet 모델은

① 더 적은 파라미터 수 위해 사용.

② 비슷한 정확도 가진 다른 ConvNet 보다

더 적은 FLOPS

모든 EfficientNet 모델의 성능 : scale

EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks

Table 2. EfficientNet Performance Results on ImageNet (Russakovsky et al., 2015). All EfficientNet models are scaled from our baseline EfficientNet-B0 using different compound coefficient ϕ in Equation 3. ConvNets with similar top-1/top-5 accuracy are grouped together for efficiency comparison. Our scaled EfficientNet models consistently reduce parameters and FLOPS by an order of magnitude (up to 8.4x parameter reduction and up to 16x FLOPS reduction) than existing ConvNets.

Model	Top-1 Acc.	Top-5 Acc.	#Params	Ratio-to-EfficientNet	#FLOPS	Ratio-to-EfficientNet
EfficientNet-B0 ← 베이스 모델	77.3%	93.5%	5.3M	1x	0.39B	1x
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
EfficientNet-B1	79.2%	94.5%	7.8M	1x	0.70B	1x
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
EfficientNet-B2	80.3%	95.0%	9.2M	1x	1.0B	1x
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
EfficientNet-B3	81.7%	95.6%	12M	1x	1.8B	1x
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
EfficientNet-B4	83.0%	96.3%	19M	1x	4.2B	1x
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
EfficientNet-B5	83.7%	96.7%	30M	1x	9.9B	1x
AmoebaNet-C (Cubuk et al., 2019)	83.5%	96.5%	155M	5.2x	41B	4.1x
EfficientNet-B6	84.2%	96.8%	43M	1x	19B	1x
EfficientNet-B7 ← 우거운 모델	84.4% ↑	97.1% ↑	66M ↓	1x	37B	1x
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).

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%
compound scale ($d=1.4$, $w=1.2$, $r=1.3$)	2.3B	75.6%
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%
MobileNetV2 compound scale	1.3B	77.4%
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%
ResNet-50 compound scale	16.7B	78.8%

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
Speedup	5.7x	Speedup	6.1x

더 빠름!

더 빠름!

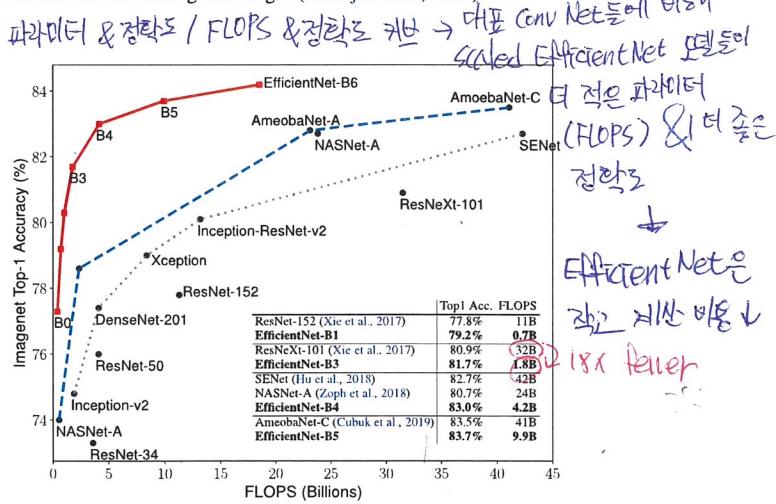


Figure 5. FLOPS vs. ImageNet Accuracy – Similar to Figure 1 except it compares FLOPS rather than model size.

5.2. ImageNet Results for EfficientNet

We train our EfficientNet models on ImageNet using similar settings as (Tan et al., 2019): RMSProp optimizer with decay 0.9 and momentum 0.9; batch norm momentum 0.99;

* RMSProp

: AdaGrad는 최소값에 도달하기 전에 학습률을 0에 수렴하게 만들 수도 있다.
AdaGrad가 간단한 convex function에서는 잘 동작 but 복잡한 다차원 곡면함수를
잘 탐색하도록 설계되지 X. 기울기의 단순한 누적만으로 충분X.
RMSProp은 기울기 단순 누적X. 지수 가중 이동 평균 (Exponentially weighted moving average)
사용해 최근 기울기들이 더 크게 반영되도록 함.

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

이전 AdaGrad의 h 에 새로운 hyper parameter ρ 추가

$\Rightarrow h$ 가 무한히 커지지 않으면서 ρ 가 작을 수록 가장 최근의 기울기를 더 크게 반영.

Inference

학습된 model을 사용해 입력에 대한 답 추출

EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks

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.

Model	Comparison to best public-available results				Model	Comparison to best reported results				Model	Acc.	#Param
	Acc.	#Param	Our Model	Acc.		Model	Acc.	#Param	Our Model			
CIFAR-10	NASNet-A	98.0%	85M	EfficientNet-B0	98.1%	4M (21x)	↑GPipe	99.0%	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	91.7%	64M (8.7x)
Birdsnap	Inception-v4	81.8%	41M	EfficientNet-B5	82.0%	28M (1.5x)	↑GPipe	83.6%	556M	EfficientNet-B7	84.3%	64M (8.7x)
Stanford Cars	Inception-v4	93.4%	41M	EfficientNet-B3	93.6%	10M (4.1x)	↑DAT	94.8%	-	EfficientNet-B7	94.7%	-
Flowers	Inception-v4	98.5%	41M	EfficientNet-B5	98.5%	28M (1.5x)	DAT	97.7%	-	EfficientNet-B7	98.8%	-
FGVC Aircraft	Inception-v4	90.9%	41M	EfficientNet-B4	90.7%	10M (4.1x)	DAT	92.9%	-	EfficientNet-B7	92.9%	-
Oxford-IIIT Pets	ResNet-152	94.5%	58M	EfficientNet-B4	94.8%	17M (5.6x)	↑GPipe	95.9%	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	93.0%	64M (8.7x)
Geo-Mean												(9.6x)

[↑]GPipe (Huang et al., 2018) trains giant models with specialized pipeline parallelism library.

[↑]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).

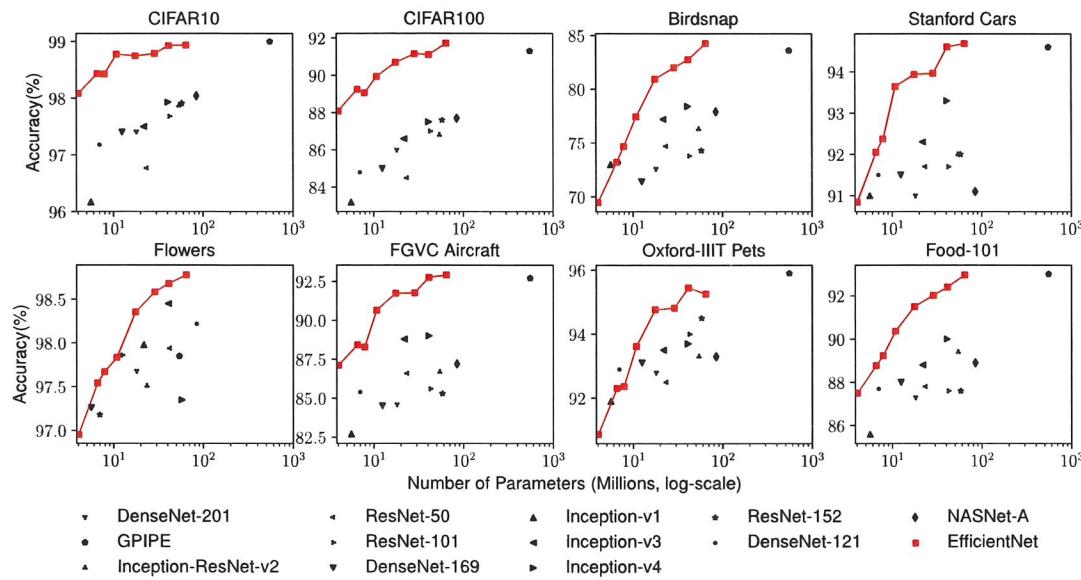


Figure 6. Model Parameters vs. Transfer Learning Accuracy – All models are pretrained on ImageNet and finetuned on new datasets.

weight decay 1e-5; initial learning rate 0.256 that decays by 0.97 every 2.4 epochs. We also use swish activation (Ramachandran et al., 2018; Elfwing et al., 2018), fixed AutoAugment policy (Cubuk et al., 2019), and stochastic depth (Huang et al., 2016) with survival probability 0.8. As commonly known that bigger models need more regularization, we linearly increase dropout (Srivastava et al., 2014) ratio from 0.2 for EfficientNet-B0 to 0.5 for EfficientNet-B7.

Table 2 shows the performance of all EfficientNet models that are scaled from the same baseline EfficientNet-B0. Our EfficientNet models generally use an order of magnitude fewer parameters and FLOPS than other ConvNets with similar accuracy. In particular, our EfficientNet-B7 achieves 84.4% top 1 / 97.1% top-5 accuracy with 66M parameters and 37B FLOPS, being more accurate but 8.4x smaller than the previous best GPipe (Huang et al., 2018).

→ Table 2 설명

EfficientNet 모델은 일반적으로 같은 규모의
파라미터, 작은 FLOPS 위해 사용
(비슷한 정확도 갖는 다른 ConvNet들 비해)

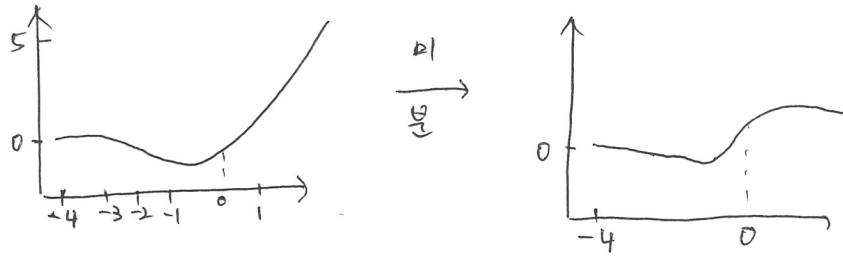
Figure 1 and Figure 5 illustrates the parameters-accuracy and FLOPS-accuracy curve for representative ConvNets, where our scaled EfficientNet models achieve better accuracy with much fewer parameters and FLOPS than other ConvNets. Notably, our EfficientNet models are not only small, but also computational cheaper. For example, our EfficientNet-B3 achieves higher accuracy than ResNet-101 (Xie et al., 2017) using 18x fewer FLOPS.

To validate the computational cost, we have also measured the inference latency for a few representative ConvNets on a real CPU as shown in Table 4, where we report average latency over 20 runs. Our EfficientNet-B1 runs 5.7x faster than the widely used ResNet-152 (He et al., 2016), while EfficientNet-B7 runs about 6.1x faster than GPipe (Huang et al., 2018), suggesting our EfficientNets are indeed fast on real hardware.

→ Table 4 설명 (지연시간 비교)

CPU에 대해 다른 ConvNet보다 훨씬 지연시간이
EfficientNet 더 적음

* Sigmoid Linear Unit, SiLU



$$f(x) = \frac{x}{1+e^{-x}} = x \cdot f(x)$$

$$f'(x) = f(x) + f(x)(1-f(x))$$

- 2차원에서 확인 했을 때, Linear 또는 Relu보다 훨씬 부드러운 형태.
- Relu 및 다른 학성화함수를 대체하기 위해 만든 함수.
- 논문에서 CIFAR 예시에서 실험 결과, ReLU 및 다른 학성화 함수보다 좋은 성능

* Auto Augment

- Augmentation

- : 딥러닝 모델을 충분히 훈련하는데 필요한 데이터를 확보하는 방법 중 하나.
적은 양의 훈련 데이터에 인위적 변화를 가해 새로운 훈련 데이터를 대량 확보하는 방법.
- ex) 이미지 상하좌우로 뒤집기 (Flipping), 자르기 (cropping) → 새로운 데이터 확보
- 현실 세계에서도 실제로 존재할 법한 데이터 생성하여 좀 더 일반화된 모델 얻는 것 목표.

문제점

- : 데이터에 최적화된 Augmentation 방법인 탐색 어렵. 현실과 동떨어지거나 기존 특징 태그 할 수도 있는 기법은 오히려 학습 나이도와 성능에 안 좋은 영향.
- 데이터를 취급하는 방식이 다른 점도 Augmentation 기법에 영향.
- Ex) 이미지 확장 때 대상 물체 크기 달라짐. 이미지 여러 크기로 자동 조절하는 기법이 도움.
의료영상 - 장기의 크기는 거의 바뀌지 않아도 데이터 확장 방식이나 환자의 상태에 따라 장기에 변이 경우.
: 노이즈 적절하게 생성하는 기법 필요.

- Auto Augment

- : ML로 데이터 수 늘리는 자동화 방법 찾는 연구. 강화학습을 통해 이미지 데이터셋에 가장 적합한 Augmentation 정책 (Policy)를 자동으로 찾았다는 알고리즘. 탐색 공간 (Search Space)는 데이터에 적합한 Augmentation 기법을 다양하게 찾으려는 목적에서 하나의 정책을 5개의 하위정책으로 구성.
- 하위 정책은 그개의 이미지 오퍼레이션 (operation)으로 구성. 각 오퍼레이션은 augmentation 기법
- ($T = \{\text{반전, 대비, 밝기, } \dots, \text{침여도}\}, n(T)=6$)을 어떤 확률 ($P=\{0, 0.1, \dots, 1\}, n(P)=17$) 어느 수준의 강도 ($M = \{0, 1, \dots, 9\}, n(M)=10$)로 실행할지 결정. 한 하위 정책이 생성하는 augmentation 기법은 총 $(16 \times 11 \times 10)^5$ 개 ($M=10, T=17$)가 결과적으로 AutoAugmentation에서 정의된 이미지 Augmentation에는 총 2.9×10^{32} 개 후보군.
- 결과적으로 AutoAugmentation에서 정의된 이미지 Augmentation에는 총 2.9×10^{32} 개 후보군.
- 일일이 학인 불가능, 탐색 깊이, 네비 줄이는 알고리즘 필요! AutoAugmentation은 임의로 선택한 Augmentation 기법의 일일이 학인 불가능, 탐색 깊이, 네비 줄이는 알고리즘 필요! AutoAugmentation은 임의로 선택한 Augmentation 기법의 일일이 학인 불가능, 탐색 깊이, 네비 줄이는 알고리즘 필요! AutoAugmentation은 임의로 선택한 Augmentation 기법의 일일이 학인 불가능, 탐색 깊이, 네비 줄이는 알고리즘 필요! AutoAugmentation은 임의로 선택한 Augmentation 기법의 일일이 학인 불가능, 탐색 깊이, 네비 줄이는 알고리즘 필요!
- 훈련 데이터를 증강, 네트워크 훈련 초기는 성능 Good, bad 모두 평균 향상. 이후 모델 성능 높이는 방향으로 보상
- 최득해가며 AutoAugmentation은 점점 더 높은 성능을 내는 기법 찾아감. ⇒ AutoAugmentation은 동일한 모델 구조
- 최적화 기술 활용할 때 학습에 적합한 데이터 증강 기법 적용하는 것만으로도 성능 개선 의의.
- AutoAugmentation은 데이터셋 규모를 축소했을 때 더 큰 성능 개선 효과. 비라벨링 데이터를 사용한 기존의
- 증강 학습 (semi-supervised learning)과 달리 모여 데이터 규모 축소 상태에서 AutoAugmentation이 찾아낸 방법들은
- 성능 훨씬 도드라짐. But, 계산 복잡도 ↑, 일반적인 연구환경에 적용 어려움.

*transfer learning

기존의 만들어진 모델을 사용해 새로운 모델 만들 때 학습 빠르게 하며 예측을 더 높이는 방법.

→ transfer learning 데 사용?

- ① 실질적으로 Conv. Net 처음부터 학습시킬 필요 없지 X. 대부분의 문제는 이미 학습된 모델 사용해 문제 해결 가능.

EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks

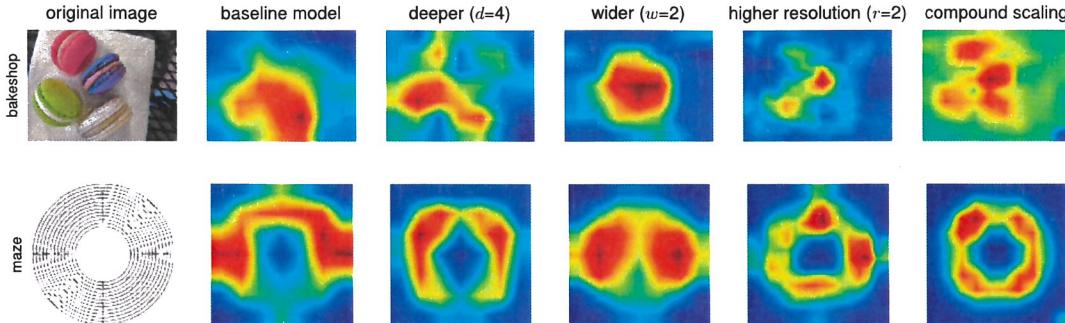


Figure 7. Class Activation Map (CAM) (Zhou et al., 2016) for Models with different scaling methods- Our compound scaling method allows the scaled model (last column) to focus on more relevant regions with more object details. Model details are in Table 7.

ImageNet 데이터 Table 6. Transfer Learning Datasets.

Dataset	Train Size	Test Size	#Classes
CIFAR-10 (Krizhevsky & Hinton, 2009)	50,000	10,000	10
CIFAR-100 (Krizhevsky & Hinton, 2009)	50,000	10,000	100
Birdsnap (Berg et al., 2014)	47,386	2,443	500
Stanford Cars (Krause et al., 2013)	8,144	8,041	196
Flowers (Nilsback & Zisserman, 2008)	2,040	6,149	102
FGVC Aircraft (Maji et al., 2013)	6,667	3,333	100
Oxford-IIIT Pets (Parkhi et al., 2012)	3,680	3,369	37
Food-101 (Bossard et al., 2014)	75,750	25,250	101

5.3 Transfer Learning Results for EfficientNet

We have also evaluated our EfficientNet on a list of commonly used transfer learning datasets, as shown in Table 6. We borrow the same training settings from (Kornblith et al., 2019) and (Huang et al., 2018), which take ImageNet pretrained checkpoints and finetune on new datasets.

Table 6 shows the transfer learning performance:

(1) Compared to public available models, such as NASNet-A (Zoph et al., 2018) and Inception-v4 (Szegedy et al., 2017), our EfficientNet models achieve better accuracy with 4.7x average (up to 21x) parameter reduction. (2) Compared to state-of-the-art models, including DAT (Ngiam et al., 2018) that dynamically synthesizes training data and GPipe (Huang et al., 2018) that is trained with specialized pipeline parallelism, our EfficientNet models still surpass their accuracy in 5 out of 8 datasets, but using 9.6x fewer parameters.

Figure 6 compares the accuracy-parameters curve for a variety of models. In general, our EfficientNets consistently achieve better accuracy with an order of magnitude fewer parameters than existing models, including ResNet (He et al., 2016), DenseNet (Huang et al., 2017), Inception (Szegedy et al., 2017), and NASNet (Zoph et al., 2018).

Figure 6 shows that our EfficientNets are able to achieve similar accuracy to other models but with significantly fewer parameters.

6. Discussion

To disentangle the contribution of our proposed scaling method from the EfficientNet architecture, Figure 8 compares the ImageNet performance of different scaling methods.

*pretrained

: 이미 학습이 끝난 모델 (모델 구성요소 (ex. 임베딩)).

선행 학습된 임베딩을 신경망에 입력하는 경우

모델에서 선행 학습된 임베딩에 의존하지 않고

임베딩 자체를 학습하는 경우

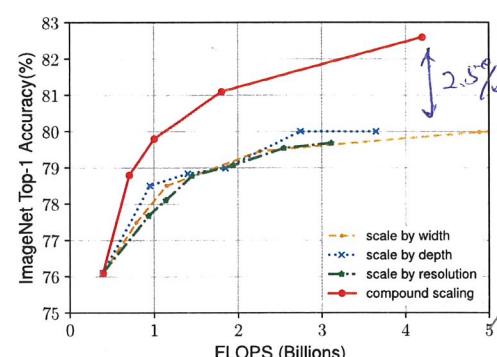


Figure 8. Scaling Up EfficientNet-B0 with Different Methods.

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%

ods for the same EfficientNet-B0 baseline network. In general, all scaling methods improve accuracy with the cost of more FLOPS, but our compound scaling method can further improve accuracy, by up to 2.5%, than other single-dimension scaling methods, suggesting the importance of our proposed compound scaling.

In order to further understand why our compound scaling method is better than others, Figure 7 compares the class activation map (Zhou et al., 2016) for a few representative models with different scaling methods. All these models are scaled from the same baseline, and their statistics are shown in Table 7. Images are randomly picked from ImageNet validation set. As shown in the figure, the model with compound scaling tends to focus on more relevant regions with more object details, while other models are either lack of object details or unable to capture all objects in the images.

*checkpoint

: 특정 시점에 모델 변수의 상태를 저장한 데이터. 체크포인트를 통해 모델 가중치를 내보내고 여러 세션 (session)을 넘나들며 학습 수행할 수 있음. 또한 체크포인트로부터 과거의 오류 (ex. 작업 설정)를 이어받아 학습할 수 있음. 단, 그때 자체는 체크포인트에 포함X.

* transfer learning.

: 기존의 만들어진 모델을 이용하여 새로운 모델 만들 시 빠르게 학습하며, 예측 더 높임.

왜 사용?

① 실질적으로 처음부터 CNN 학습하는 일 많지 X. 대부분의 문제는 이미 학습된 모델을 사용해 문제 해결 가능.

② 복잡한 모델의 수록 학습 시키기 어렵. 어떤 모델은 그 이상, 비싼 GPU 여려대 사용.

③ layer의 개수, activation, hyperparameter 등 고려사항 多.

실질적으로 처음부터 학습시키려면 많은 시간 필요.

→ 이미 잘 훈련된 모델이 있고, 해당 Model과 유사한 문제 해결 시 transfer learning 사용

실질적 사용

- 새로 훈련할 데이터 적지만, original 데이터와 유사한 경우.

: 데이터 양이 적어 fine-tune (전체 모델에 대해 backpropagation 진행) 은 over-fitting 위험이 있어 하지 X. 새로 학습할 데이터는 original 데이터와 유사하기 때문에 이 경우 최종 linear classifier 레이어만 학습.

- 새로 훈련할 데이터 매우 많으며 original 데이터와 유사한 경우.

: 새로 학습할 데이터의 양이 많다는 것은 over-fitting의 위험이 낫다는 뜻으로, 전체 레이어에 대해 fine-tune 함.

- 새로 훈련할 데이터가 적으며 original 데이터와 다른 경우.

: 데이터의 양이 적기 때문에 최종 단계의 linear classifier 레이어 학습 좋을 것.

반면, 데이터가 서로 다르기 때문에 거의 마지막 부분 (the top of the network) 만 학습 좋지 X. 서로 상충되는 내용이지만, 이 경우 네트워크 초기부분 어딘가 activation 이후에 특정 레이어 학습 시키는게 좋음.

- 새로 훈련할 데이터가 많지만 original 데이터와 다른 경우.

: 데이터 많기 때문에 아예 새로운 ConvNet 만들 수 있지만, 실질적으로 transfer learning이 더 효율 좋음. 전체 네트워크에 대해 fine-tune 해도 됨.

* Class Activation Map

: 기존 CNN의 모델은 classification 을 이미지 내에서 어떤 특징 (feature) 을 보고 판단하는지 알 수X 하지만 출격의 결과를 시각화 하여 실제 모델이 이미지의 어떤 부분을 보고 분류하는지 알 수 있게 하는 것이 Class Activation Map (CAM)

Compound scaling 방법으로
모바일 크기의 EfficientNet 모델이 효과적으로 스케일 ↑ 가능.

→ 정확도 ↑, 파라미터 ↓, FLOPS ↓

(ImageNet과 5개의 다른 사용되는 transfer learning dataset)

EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks

7 Conclusion

In this paper, we systematically study ConvNet scaling and identify that carefully balancing network width, depth, and resolution is an important but missing piece, preventing us from better accuracy and efficiency. To address this issue, we propose a simple and highly effective compound scaling method, which enables us to easily scale up a baseline ConvNet to any target resource constraints in a more principled way, while maintaining model efficiency. Powered by this compound scaling method, we demonstrate that a mobile-size EfficientNet model can be scaled up very effectively, surpassing state-of-the-art accuracy with an order of magnitude fewer parameters and FLOPS, on both ImageNet and five commonly used transfer learning datasets.

Acknowledgements

We thank Ruoming Pang, Vijay Vasudevan, Alok Aggarwal, Barret Zoph, Hongkun Yu, Xiaodan Song, Samy Bengio, Jeff Dean, and Google Brain team for their help.

References

- Berg, T., Liu, J., Woo Lee, S., Alexander, M. L., Jacobs, D. W., and Belhumeur, P. N. Birdsnap: Large-scale fine-grained visual categorization of birds. *CVPR*, pp. 2011–2018, 2014.
- Bossard, L., Guillaumin, M., and Van Gool, L. Food-101-mining discriminative components with random forests. *ECCV*, pp. 446–461, 2014.
- Cai, H., Zhu, L., and Han, S. Proxylessnas: Direct neural architecture search on target task and hardware. *ICLR*, 2019.
- Chollet, F. Xception: Deep learning with depthwise separable convolutions. *CVPR*, pp. 1610–02357, 2017.
- Cubuk, E. D., Zoph, B., Mane, D., Vasudevan, V., and Le, Q. V. Autoaugment: Learning augmentation policies from data. *CVPR*, 2019.
- Elfwing, S., Uchibe, E., and Doya, K. Sigmoid-weighted linear units for neural network function approximation in reinforcement learning. *Neural Networks*, 107:3–11, 2018.
- Gholami, A., Kwon, K., Wu, B., Tai, Z., Yue, X., Jin, P., Zhao, S., and Keutzer, K. SqueezeNext: Hardware-aware neural network design. *ECV Workshop at CVPR'18*, 2018.
- Han, S., Mao, H., and Dally, W. J. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. *ICLR*, 2016.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. *CVPR*, pp. 770–778, 2016.
- He, K., Gkioxari, G., Dollár, P., and Girshick, R. Mask r-cnn. *ICCV*, pp. 2980–2988, 2017.
- He, Y., Lin, J., Liu, Z., Wang, H., Li, L.-J., and Han, S. Amc: Automl for model compression and acceleration on mobile devices. *ECCV*, 2018.
- Howard, A. G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., Andreetto, M., and Adam, H. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.
- Hu, J., Shen, L., and Sun, G. Squeeze-and-excitation networks. *CVPR*, 2018.
- Huang, G., Sun, Y., Liu, Z., Sedra, D., and Weinberger, K. Q. Deep networks with stochastic depth. *ECCV*, pp. 646–661, 2016.
- Huang, G., Liu, Z., Van Der Maaten, L., and Weinberger, K. Q. Densely connected convolutional networks. *CVPR*, 2017.
- Huang, Y., Cheng, Y., Chen, D., Lee, H., Ngiam, J., Le, Q. V., and Chen, Z. Gpipe: Efficient training of giant neural networks using pipeline parallelism. *arXiv preprint arXiv:1808.07233*, 2018.
- Iandola, F. N., Han, S., Moskewicz, M. W., Ashraf, K., Dally, W. J., and Keutzer, K. SqueezeNext: Alexnet-level accuracy with 50x fewer parameters and <0.5 mb model size. *arXiv preprint arXiv:1602.07360*, 2016.
- Ioffe, S. and Szegedy, C. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *ICML*, pp. 448–456, 2015.
- Kornblith, S., Shlens, J., and Le, Q. V. Do better imagenet models transfer better? *CVPR*, 2019.
- Krause, J., Deng, J., Stark, M., and Fei-Fei, L. Collecting a large-scale dataset of fine-grained cars. *Second Workshop on Fine-Grained Visual Categorizatio*, 2013.
- Krizhevsky, A. and Hinton, G. Learning multiple layers of features from tiny images. *Technical Report*, 2009.
- Krizhevsky, A., Sutskever, I., and Hinton, G. E. Imagenet classification with deep convolutional neural networks. In *NIPS*, pp. 1097–1105, 2012.
- Lin, H. and Jegelka, S. Resnet with one-neuron hidden layers is a universal approximator. *NeurIPS*, pp. 6172–6181, 2018.

- Lin, T.-Y., Dollár, P., Girshick, R., He, K., Hariharan, B., and Belongie, S. Feature pyramid networks for object detection. *CVPR*, 2017.
- Liu, C., Zoph, B., Shlens, J., Hua, W., Li, L.-J., Fei-Fei, L., Yuille, A., Huang, J., and Murphy, K. Progressive neural architecture search. *ECCV*, 2018.
- Lu, Z., Pu, H., Wang, F., Hu, Z., and Wang, L. The expressive power of neural networks: A view from the width. *NeurIPS*, 2018.
- Ma, N., Zhang, X., Zheng, H.-T., and Sun, J. Shufflenet v2: Practical guidelines for efficient cnn architecture design. *ECCV*, 2018.
- Mahajan, D., Girshick, R., Ramanathan, V., He, K., Paluri, M., Li, Y., Bharambe, A., and van der Maaten, L. Exploring the limits of weakly supervised pretraining. *arXiv preprint arXiv:1805.00932*, 2018.
- Maji, S., Rahtu, E., Kannala, J., Blaschko, M., and Vedaldi, A. Fine-grained visual classification of aircraft. *arXiv preprint arXiv:1306.5151*, 2013.
- Ngiam, J., Peng, D., Vasudevan, V., Kornblith, S., Le, Q. V., and Pang, R. Domain adaptive transfer learning with specialist models. *arXiv preprint arXiv:1811.07056*, 2018.
- Nilsback, M.-E. and Zisserman, A. Automated flower classification over a large number of classes. *ICVGIP*, pp. 722–729, 2008.
- Parkhi, O. M., Vedaldi, A., Zisserman, A., and Jawahar, C. Cats and dogs. *CVPR*, pp. 3498–3505, 2012.
- Raghu, M., Poole, B., Kleinberg, J., Ganguli, S., and Sohl-Dickstein, J. On the expressive power of deep neural networks. *ICML*, 2017.
- Ramachandran, P., Zoph, B., and Le, Q. V. Searching for activation functions. *arXiv preprint arXiv:1710.05941*, 2018.
- Real, E., Aggarwal, A., Huang, Y., and Le, Q. V. Regularized evolution for image classifier architecture search. *AAAI*, 2019.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al. Imagenet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3): 211–252, 2015.
- Sandler, M., Howard, A., Zhu, M., Zhmoginov, A., and Chen, L.-C. Mobilenetv2: Inverted residuals and linear bottlenecks. *CVPR*, 2018.
- Sharir, O. and Shashua, A. On the expressive power of overlapping architectures of deep learning. *ICLR*, 2018.
- Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. Dropout: a simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958, 2014.
- Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., and Rabinovich, A. Going deeper with convolutions. *CVPR*, pp. 1–9, 2015.
- Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., and Wojna, Z. Rethinking the inception architecture for computer vision. *CVPR*, pp. 2818–2826, 2016.
- Szegedy, C., Ioffe, S., Vanhoucke, V., and Alemi, A. A. Inception-v4, inception-resnet and the impact of residual connections on learning. *AAAI*, 4:12, 2017.
- Tan, M., Chen, B., Pang, R., Vasudevan, V., Sandler, M., Howard, A., and Le, Q. V. MnasNet: Platform-aware neural architecture search for mobile. *CVPR*, 2019.
- Xie, S., Girshick, R., Dollár, P., Tu, Z., and He, K. Aggregated residual transformations for deep neural networks. *CVPR*, pp. 5987–5995, 2017.
- Yang, T.-J., Howard, A., Chen, B., Zhang, X., Go, A., Sze, V., and Adam, H. Netadapt: Platform-aware neural network adaptation for mobile applications. *ECCV*, 2018.
- Zagoruyko, S. and Komodakis, N. Wide residual networks. *BMVC*, 2016.
- Zhang, X., Li, Z., Loy, C. C., and Lin, D. Polynet: A pursuit of structural diversity in very deep networks. *CVPR*, pp. 3900–3908, 2017.
- Zhang, X., Zhou, X., Lin, M., and Sun, J. Shufflenet: An extremely efficient convolutional neural network for mobile devices. *CVPR*, 2018.
- Zhou, B., Khosla, A., Lapedriza, A., Oliva, A., and Torralba, A. Learning deep features for discriminative localization. *CVPR*, pp. 2921–2929, 2016.
- Zoph, B. and Le, Q. V. Neural architecture search with reinforcement learning. *ICLR*, 2017.
- Zoph, B., Vasudevan, V., Shlens, J., and Le, Q. V. Learning transferable architectures for scalable image recognition. *CVPR*, 2018.