

EfficientNet

Rethinking Model Scaling for Convolutional Neural Networks

인공지능 논문 리딩/리뷰 스터디

Abstract

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 84.3% top-1 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>.

Abstract

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.

CNN 모델의 크기를 안정적으로 늘리는 방법 연구함

To go beyond this, 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 84.3% top-1 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>.

Abstract

Abstract

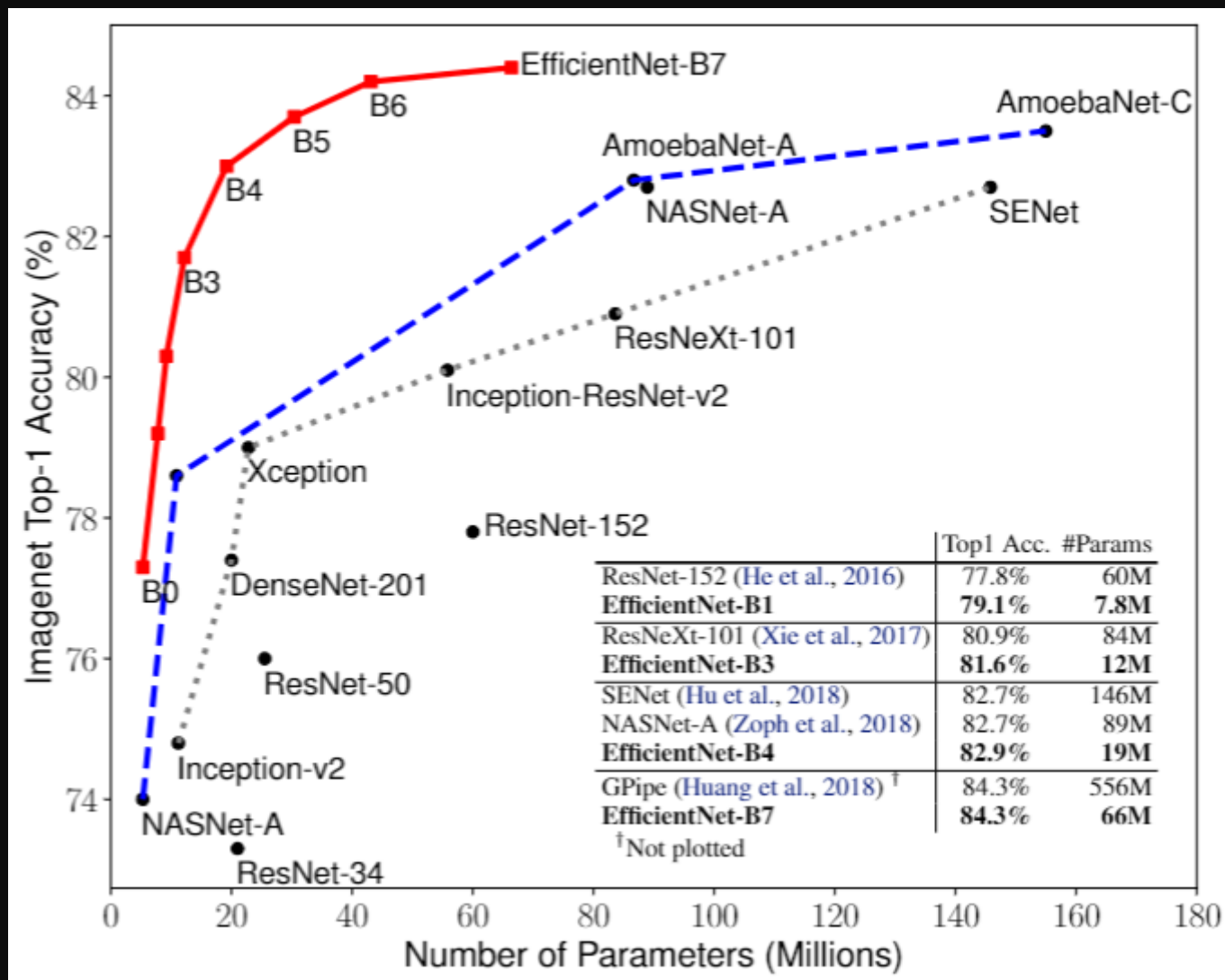
Convolutional Neural Networks (ConvNets) are commonly developed at a fixed resource budget, and then scaled up for better accuracy if more resources are available. In this paper, we systematically study model scaling and identify that carefully balancing network depth, width, and resolution can lead 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 our method on scaling up MobileNets and ResNet.

To go even further, we use neural architecture search to design a new baseline network and use it as a basis of training models called EfficientNets, which achieve much better accuracy and efficiency than previous ConvNets. In particular, our EfficientNet-B7 achieves state-of-the-art 84.3% top-1 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>.

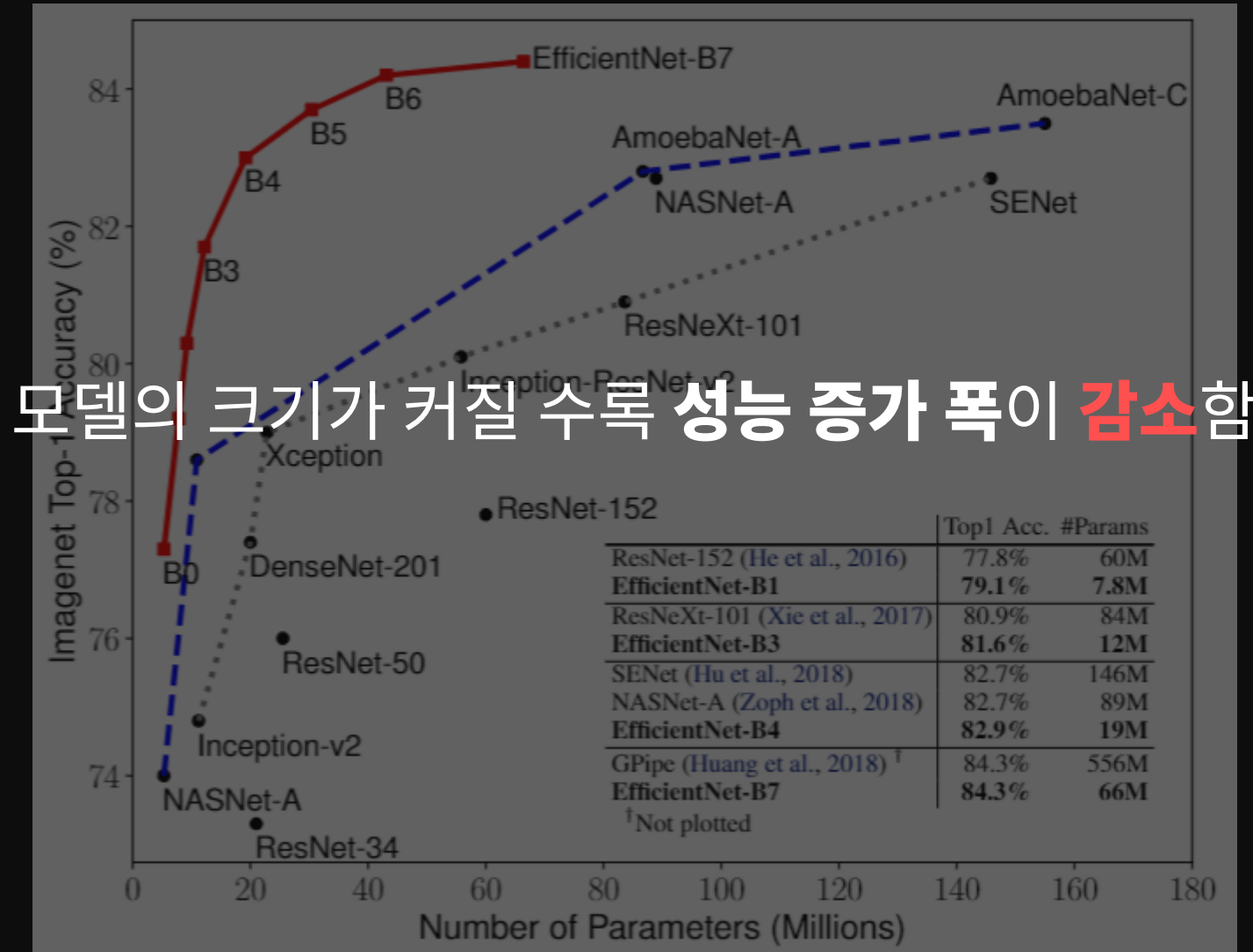
CNN 모델의 크기를 안정적으로 늘리는 방법 연구함

그 과정에서 **EfficientNet** 모델 개발

Introduction



Introduction



Introduction

Width/Depth/Resolution

위 세가지 수치를 균형 있게 키우자

Introduction

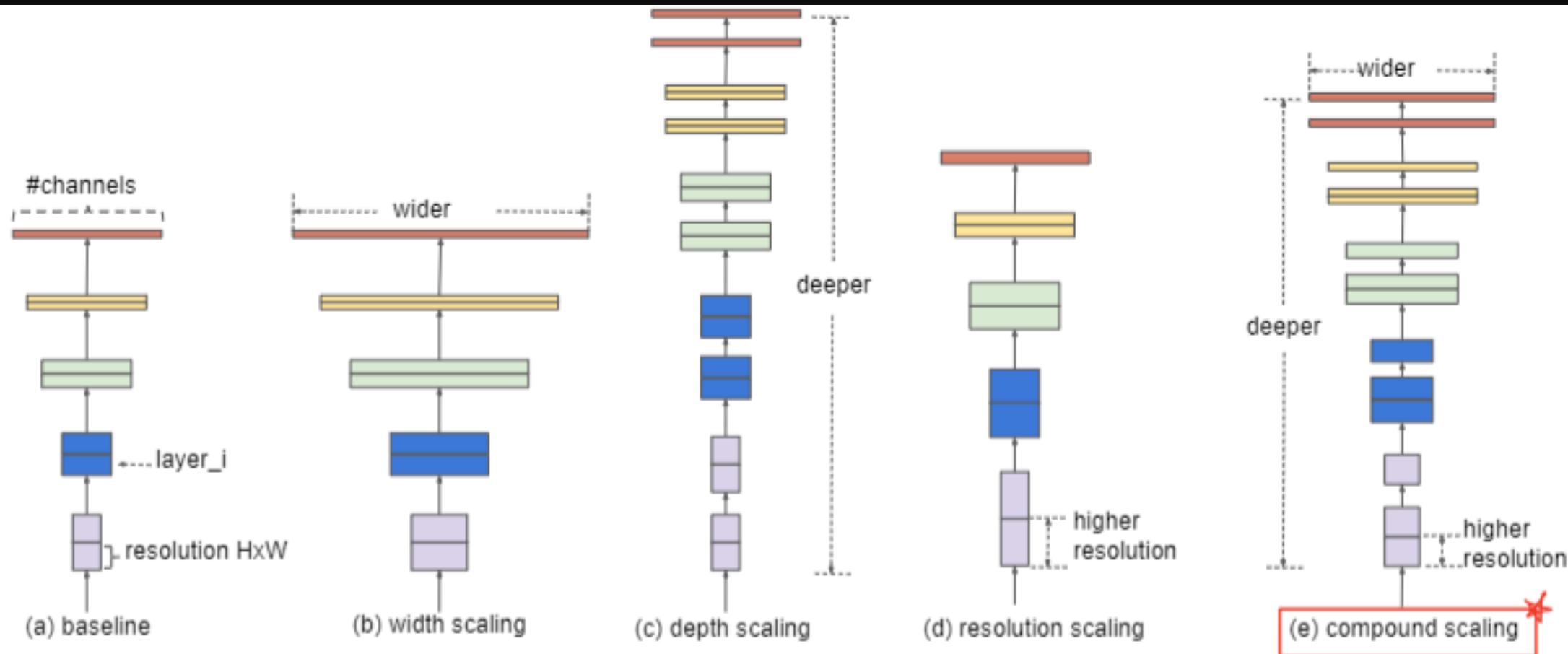


Figure 2. Model Scaling. (a) is a baseline network example; (b)-(d) are conventional scaling that only increases one dimension of network width, depth, or resolution. (e) is our proposed compound scaling method that uniformly scales all three dimensions with a fixed ratio.

Compound Model Scaling

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

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

$$\text{Memory}(N) \leq \text{target_memory}$$

$$\text{FLOPS}(N) \leq \text{target_flops}$$

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

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

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

$$s.t. \quad \alpha \cdot \beta^2 \cdot \gamma^2 \approx 2$$

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

1. grid search

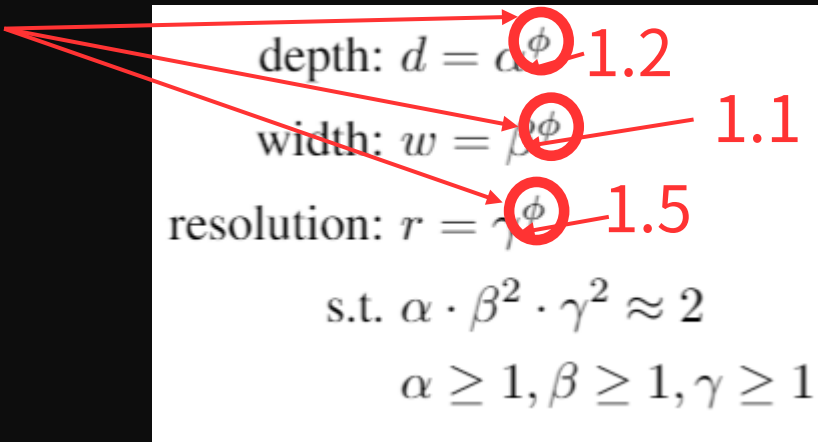
Compound Model Scaling

$$\begin{aligned} \text{depth: } d &= \alpha^{\phi} \quad 1.2 \\ \text{width: } w &= \beta^{\phi} \quad 1.1 \\ \text{resolution: } r &= \gamma^{\phi} \quad 1.5 \\ \text{s.t. } \alpha \cdot \beta^2 \cdot \gamma^2 &\approx 2 \\ \alpha \geq 1, \beta \geq 1, \gamma &\geq 1 \end{aligned}$$

1. grid search

Compound Model Scaling

Phi 값을 증가시켜
Scaling

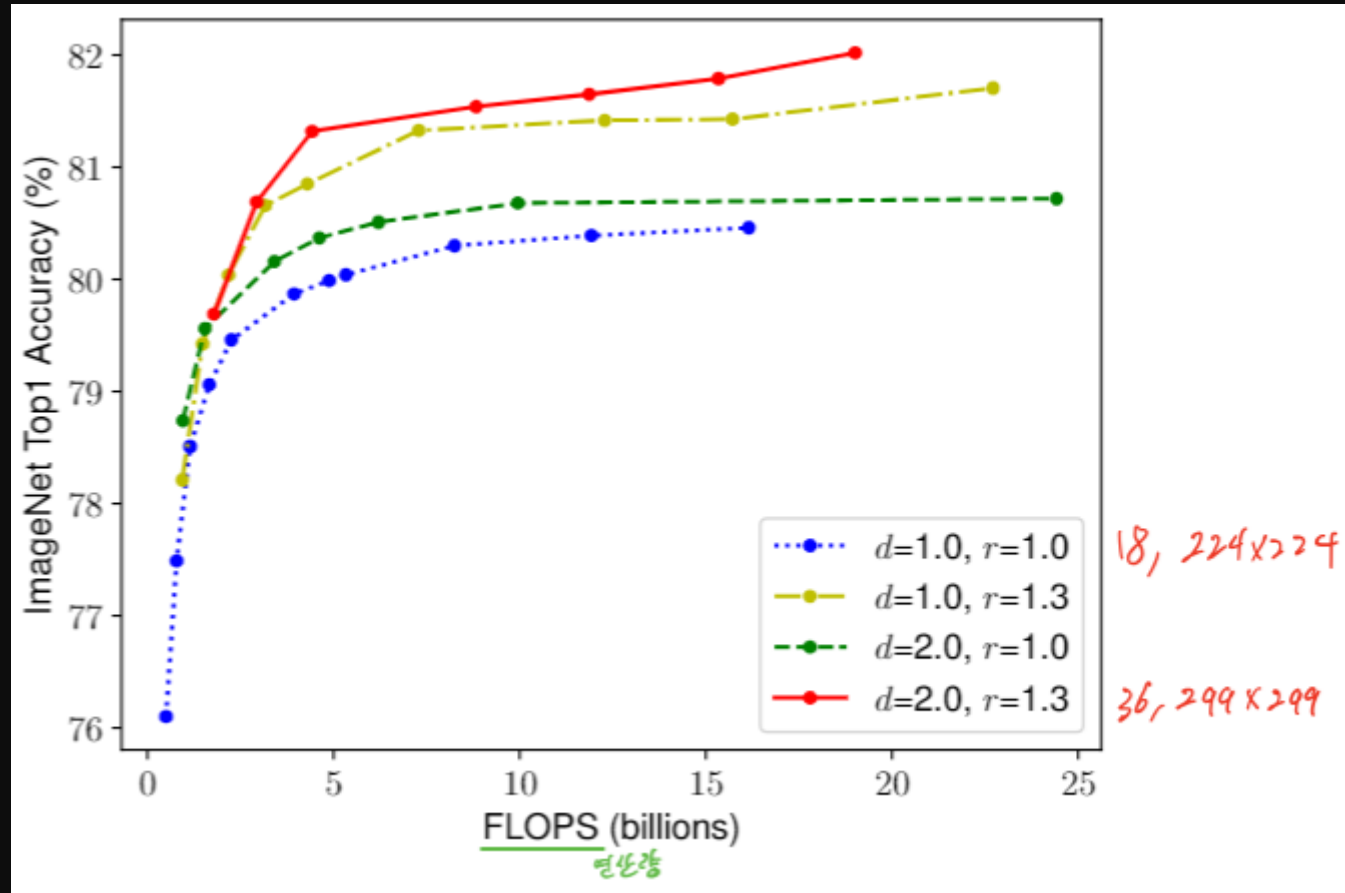


The diagram shows three scaling equations with red circles around the ϕ term and red arrows pointing to them from the text 'Phi 값을 증가시켜 Scaling'. The scaling factors are 1.2 for depth, 1.1 for width, and 1.5 for resolution.

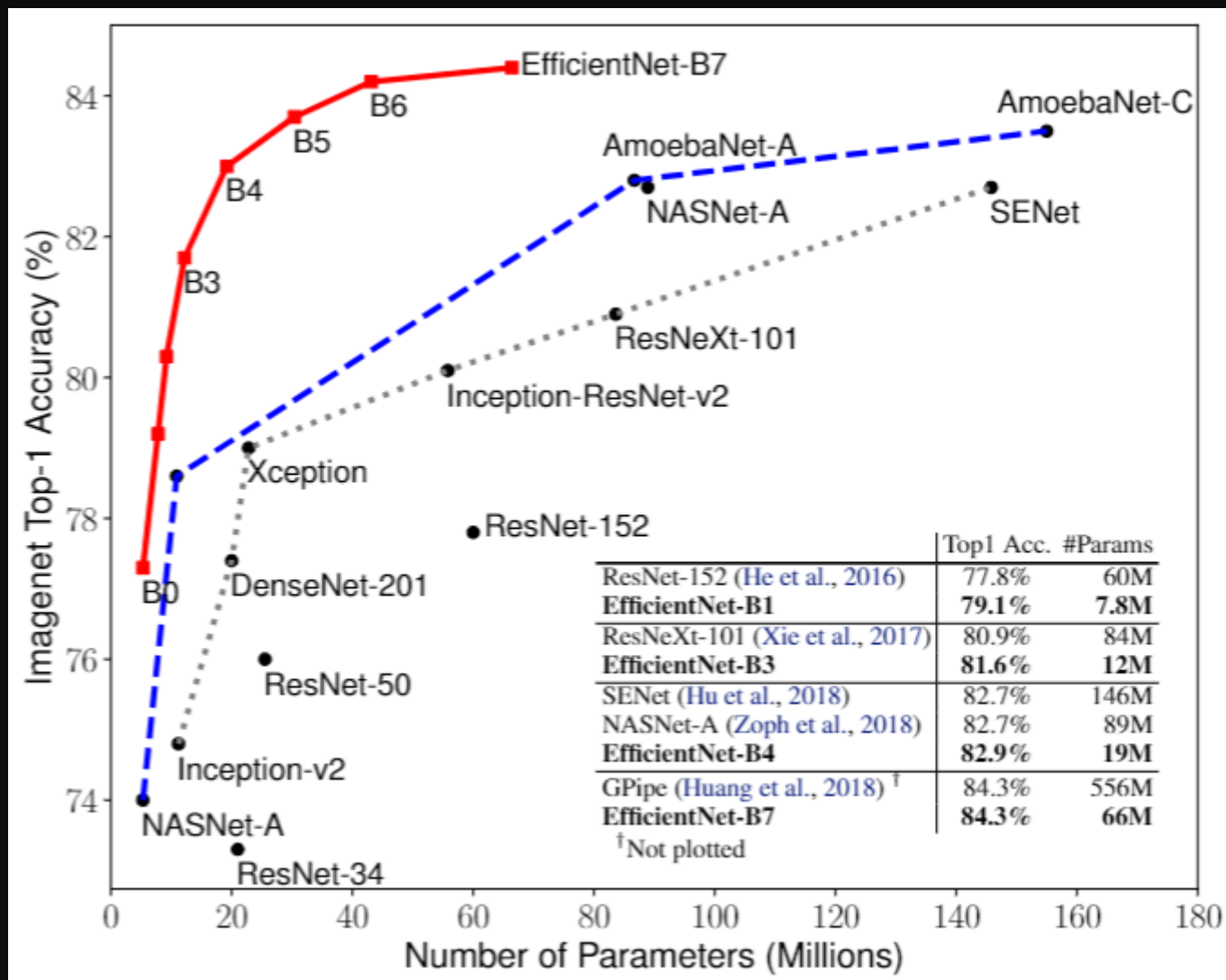
$$\begin{aligned}\text{depth: } d &= \alpha \phi^{1.2} \\ \text{width: } w &= \beta \phi^{1.1} \\ \text{resolution: } r &= \gamma \phi^{1.5}\end{aligned}$$
$$\begin{aligned}\text{s.t. } \alpha \cdot \beta^2 \cdot \gamma^2 &\approx 2 \\ \alpha \geq 1, \beta \geq 1, \gamma &\geq 1\end{aligned}$$

2. Scaling

Compound Model Scaling

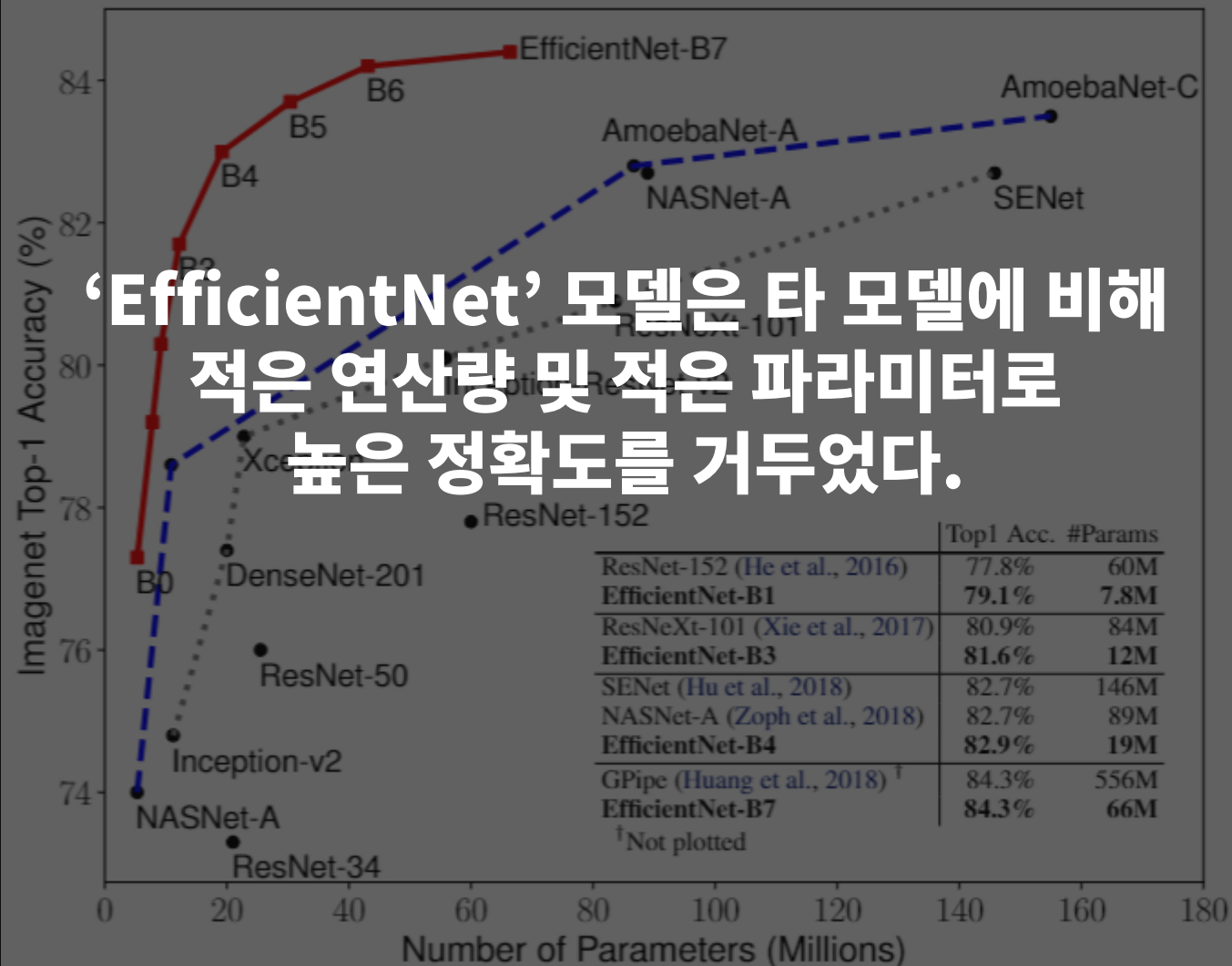


결론



결론

‘EfficientNet’ 모델은 타 모델에 비해 적은 연산량 및 적은 파라미터로 높은 정확도를 거두었다.



결론

‘EfficientNet’ 모델은 타 모델에 비해 적은 연산량 및 적은 파라미터로 높은 정확도를 거두었다.

= 효율적이다

