

# EfficientNet: Understanding the Compound Scaling Method

Siqi Zhu, Ethan Qiu, Yuexi Shen, Charlie Cheng

## 1 Introduction

*EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks*<sup>1</sup> introduces a novel model architecture scaling approach, optimizing for FLOPs and parameter efficiency. As a reminder, FLOPs stand for floating point operations per second, and parameters are a general measure for the complexity of a model. This architecture aims to capture the intuition that architecture dimensions (depth, width, resolution) should not be independent of each other, and indeed, should be balanced with each other. To this end, Efficient introduces compound model scaling, leveraging training-aware neural architecture search to develop a strong baseline model, which can then be scaled up to a variety of models.

## 2 Intuition of the Compound Scaling Method

In this problem, we will quickly walk through the Compound Scaling Method, and attempt to communicate the intuition behind the approach. For ConvNet layers, generally, there are three essential parameters that could be scaled: depth, width, and resolution:

**Depth:** Generally, deeper networks are able to capture more complex features. However, there is also a trade off considering the deeper networks take more resources (time, compute) to train. Notably, there appears to be diminishing returns with regards to continually increasing the depth of a network, such as ResNet-1000 having similar performance to ResNet-101.

**Width:** Wider networks are generally more capable of capturing complex features, they also tend to be easier to train. Unfortunately, wide and shallow networks tend to have difficulties in capturing high level features.

**Resolution:** Higher resolution images allow the capture of more complex features. But, again, there are resource constraints (time, compute), and diminishing returns.

**Problem 1.1.** Intuitively, should the different dimensions (depth, width, resolution) be independent? Let's say that you have some extra compute to improve a model, what dimension(s) will you tune?

Let us pose the construction of a ConvNet in the form of an optimization problem. Rather obviously, we are trying to maximize the performance of our model, thus accuracy. For sake of simplicity, we will be assuming that every ConvNet layer has the same architecture.

**Problem 1.2.** Let us formulate an ConvNet stage as a function  $f$ . Please fill in the blanks regarding the domain and codomain of the function. Let us have input tensor (to the layer) be  $X$  with shape  $H, W, C$  representing height, width, channels.  $L$  being the layer repetition on this stage (We ignore batching for simplicity).

$$\mathcal{F} : \text{Input Tensor}_{(H, W, C)} \rightarrow \text{Output Tensor}_{(H', W', C')}$$

<sup>1</sup>Mingxing Tan, Quoc V. Le *EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks*. <https://arxiv.org/abs/1905.11946>

**Problem 1.3.** Given our ConvNet stage, please derive an expression for an entire ConvNet  $\mathcal{N}$  composed of  $n$  stages. (Remember  $H$ ,  $W$ ,  $C$  and  $L$  could vary from stage to stage!) We use  $\odot$  to denote composition of stages.

$$\mathcal{N} = \underline{\hspace{1cm}}$$

**Problem 1.4.** Suppose that we have defined parameters  $\hat{H}_i, \hat{W}_i, e\hat{C}_i, \hat{L}_i$ , and tune able parameters  $d, w, r$  representing depth, width, and resolution respectively. And we are given functions *Accuracy*, *Memory* and *FLOPS* : *Network*  $\rightarrow \mathbb{Z}$  that takes in a network, and output respective values of the network, as well as hard constraints *targetmemory* and *targetFLOPS*, please fill in the blanks for the following optimization problem:

(Reusing your work is suggested).

$$\begin{aligned} & \max_{d,w,r} \text{Accuracy}(\text{---}) \\ & \text{s.t. } \mathcal{N}(d, w, r) = \text{---} \\ & \quad \text{---} \leq \text{targetmemory} \\ & \quad \text{---} \leq \text{targetFLOPS} \end{aligned}$$

As you have hopefully observed, changing  $d, w, r$  is a complex problem. And hence, the main contribution of the paper:

### 3 The Compound Scaling Method

Tan, Quoc and Le introduces the Compound Scaling Method, which they have formulated as an optimization problem below.

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

Looking at the formulation, we note that  $\alpha, \beta, \gamma$  are parameters obtained through a grid search, thus in application, the only variable to be varied is  $\phi$ , the compound coefficient. Thus, the compound scaling method produces a way for one to determine the dimensions of a network by only varying one paramter, the compound coefficient.

**Problem 2.1.** Looking at the optimization problem, what will happen to the depth, width, and resolution as you increase or decrease  $\phi$ ? What advantages could this approach bring compared to arbitrarily scaling these dimensions?

**Problem 2.2.** For an arbitrary ConvNet (thinking about our optimization problem formulation above, and the basic understanding of what a ConvNet is), fixing all other variables:

- What will doubling the depth do to the total FLOPS?
- What will doubling the width do to the total FLOPS?
- What will doubling the resolution do to the total FLOPS?

You may think of FLOPS just as operations that need to be performed. Thinking in quasi-big-o concepts might be useful.

**Problem 2.3.** Considering the constraint  $\alpha \cdot \beta^2 \cdot \gamma^2 \approx 2$ , how much will the total FLOPS be for some  $\phi$ ? Again, thinking in quasi-big-o concepts might be useful.

## 4 Implementing EfficientNet

The rest of this problem will be in jupyter notebook form. We also have the conceptual questions on latex. The coding questions are in the notebook!

**Problem 3.1.** Play around with this function, at what point do you start to see new depth layers emerging? (You might also want to use this opportunity to check your intuition about the previous conceptual questions).

**Problem 3.2.** What's the shape of data in train loader for a single batch? (in terms of  $[N, C, H, W]$  )

**Problem 3.3.** What is the best accuracy you can get? What is the best accuracy you can get with the same number of parameters as the EfficientNet-B2 model? Feel free to use different models and find the one with the best validation accuracy.

**Problem 3.4.** Looking at the accuracy of EfficientNet compared to other state-of-the-art model architectures, how does the validation accuracy compare? Anything else interesting you've noticed about these plots?

**Problem 3.4.** Examine the number of parameters of the different models, does anything stand out to you? If we take this as context, is there anything you want to comment about the validation performance of the models?