



# Article EvolveNet: Evolving Networks by Learning Scale of Depth and Width

Athul Shibu and Dong-Gyu Lee \*D

Department of Artificial Intelligence, Kyungpook National University, Daegu 41566, Republic of Korea; athulshibu@knu.ac.kr

\* Correspondence: dglee@knu.ac.kr

Abstract: Convolutional neural networks (CNNs) have shown decent performance in a variety of computer vision tasks. However, these network configurations are largely hand-crafted, which leads to inefficiency in the constructed network. Various other algorithms have been proposed to address this issue, but the inefficiencies resulting from human intervention have not been addressed. Our proposed EvolveNet algorithm is a task-agnostic evolutionary search algorithm that can find optimal depth and width scales automatically in an efficient way. The optimal configurations are not found using grid search, and are instead evolved from an existing network. This eliminates inefficiencies that emanate from hand-crafting, thus reducing the drop in accuracy. The proposed algorithm is a framework to search through a large search space of subnetworks until a suitable configuration is found. Extensive experiments on the ImageNet dataset demonstrate the superiority of the proposed method by outperforming the state-of-the-art methods.

Keywords: convolutional neural network; network scaling; evolutionary computation

MSC: 68T07

## 1. Introduction

The convolutional neural network (CNN) is one of the most popular networks in the field of deep learning, showing decent performance in various computer vision tasks, including classification [1,2], semantic segmentation [3,4], and action recognition [5,6]. They were designed to extract two-dimensional features by taking structured data such as images as input and then processing them using convolutional operators [7,8]. Studies [1,9,10] have shown that a larger number of layers results in increased receptive fields and, therefore, captures more detail of the image. Recent networks have achieved higher accuracy by increasing the filter bank [11–13]. There have also been cases of improved block architecture that yielded higher accuracy without significantly increasing the size of the networks [10,14].

Scaling is a widely used technique to achieve better accuracy and numerous methods have been utilized to scale networks. Upscaling depth is the most prevalent method, al-though scaling models by image resolution is also becoming increasingly popular. Figure 1 represents the depth and width of a network. EfficientNet [15] was created by compound scaling MobileNets [14,16] and ResNet [9] networks, i.e., scaling the network width, depth, and resolution by fixed coefficients. Network width refers to the number of filter banks in any layer of a network, depth refers to the number of layers in a network, and the resolution represents the number of pixels in the input image [15]. Scaling is the increasing or decreasing of any of these factors, resulting in a change in size and accuracy. Upscaling, however, could result in the configurations of the upscaled networks being ill-suited to their tasks because hand-crafting networks lead to human errors and consequent inaccuracies, resulting in an inefficient network. Filter pruning has long been considered a good alternative to accelerate deep neural networks, but this does not solve the core inefficiencies in the



**Citation:** Shibu, A.; Lee, D.-G. EvolveNet: Evolving Networks by Learning Scale of Depth and Width. *Mathematics* **2023**, *11*, 3611. https:// doi.org/10.3390/math11163611

Academic Editors: Vitaly Kober and Tae Sun Choi

Received: 24 July 2023 Revised: 18 August 2023 Accepted: 18 August 2023 Published: 21 August 2023



**Copyright:** © 2023 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (https:// creativecommons.org/licenses/by/ 4.0/).



construction of the network. In other words, the accuracy gain diminishes as networks get larger because all parameters in a network have different sensitivity to its accuracy.

**Figure 1.** Configuration of a network. The depth represents the number of layers in a network, whereas the width of a layer represents the number of filters present in it.

Identifying these superfluous parameters in a network is crucial to the optimization process. The Lottery Ticket Hypothesis [17] states that any trained dense neural network contains a subnetwork (called a winning ticket), which, when trained in isolation for at most the same number of epochs, can match the test accuracy of the original network. Automated machine learning (AutoML) [18] automated steps in the machine learning pipeline. This concept has been applied for neural architecture search (NAS) [19–21] to optimize the search for winning tickets within large networks. However, applying AutoML concepts for NAS has wider consequences that can be aggravated during the search stage. The No Free Lunch Theorem [22] posits that no universal optimization algorithm consistently outperforms all other algorithms across every optimization problem. Therefore, it is necessary to optimize networks with respect to the tasks at hand.

EfficientNet [15] showed that networks can be optimized by removing all channel connections in the depth-wise layer, and instead increasing the number of channels to boost capacity. This reduces the number of parameters, but significantly increases data movement, resulting in poor performance on hardware accelerators. The final scaling coefficients are determined by grid search. Grid search optimizes scaled networks by conducting a complete search over a given subset of the search space [23]. However, the search cost of using grid search is very expensive when scaling larger networks. Therefore, scaling algorithms are primarily applied to small networks, and large networks are created by massively upscaling small networks. The proposed algorithm instead downscales large networks to generate slightly smaller networks. This is because larger networks are downscaled from models designed for larger channel connections, which would not increase the number of channels. Thus, the generated network is similar in both size and accuracy to the large networks, but is much closer in size to its original network.

Evolutionary algorithms have been found to significantly outperform random and systemic search methods when searching in large search spaces [24]. Over the years, various multi-objective evolutionary algorithms have been proposed to varying degrees of success [25,26]. However, they tend to suffer from a weak global search ability in low inter-task relevance problems [27] because the cross-over operator is unable to distinguish between information and noise. This problem can be addressed by introducing multiple search strategies into the objective function to evolve an efficient network with high accuracy. In this paper, we present a framework to automatically generate the optimal number of layers and channels for a network without manual interference. We use evolutionary algorithms to search through the large sample space of possible subnetworks to address this issue. The proposed method uses evolutionary search to find an optimized subnetwork that keeps the number of parameters low without compromising accuracy. Instead of upscaling networks, a collection of layers and channels are downscaled to find the optimal configuration. The evolved network is built to counteract the lack of expressiveness and effectiveness that is inherent to hand-crafted and grid-searched networks. Generated networks counteract these drawbacks by integrating pruning concepts into the creation of new networks, resulting in more efficient networks.

Our contribution lies in three folds:

- We proposed an algorithm to counter inefficiencies in subnetworks by evolving taskagnostic networks of ideal depth and width for a given architecture.
- We created a framework to efficiently search through a large sample space of subnetworks to identify smaller networks without a major loss in accuracy.
- We experimentally showed the superiority of the network generated by the proposed method on publicly available pre-trained CNNs.

The remainder of this study is organized as follows: Section 3 briefly discusses the algorithm and its working concept, and Section 4 describes the result of experiments conducted on networks evolved using the EvolveNet algorithm. We also discuss the advantages of this algorithm in Section 4.4, and conclude this study in Section 5.

#### 2. Related Works

#### 2.1. Convolutional Neural Networks

CNN is the popular design choice for visual recognition tasks and has gone through many upgrades over the years. Simonyan et al., 2014 [28] introduced VGG, which used very small convolution filters and increased the depth of the model to achieve high accuracy. Szegedy et al., 2015 [29] introduced GoogLeNet, which proposed the idea of an Inception module to find the optimal local sparse structure in a CNN that can be approximated by its dense components. He et al., 2016 [9] introduced ResNet stacked layers to fit a residual mapping using skip connections, which further improved the accuracy. Howard et al., 2017 [16] introduced depth-wise separable convolution layers, which separated the filtering and combining operations of the convolution operation, which, in turn, reduced the computational complexity and model size in MobileNet. Huang et al., 2017 [10] introduced DenseNet, which connected all layers to combine the feature maps together instead of just the feature summations, resulting in larger models with higher accuracy. Liu et al., 2022 [1] introduced ConvNeXt, which improved a standard ResNet by gradually modernizing the architecture to construct a hierarchical vision transformer using Swin-T [30]. These models have one thing in common; beyond just the architecture or convolution layers, the depth of the networks and width of each layer were handcrafted by humans.

### 2.2. Neural Architecture Search

Searching for optimal network structures has been studied using reinforcement learning [12,31], gradient-based approaches [32], parameter sharing [33], weight prediction [34], and genetic algorithms [35,36]. Zoph et al., 2016 [12] uses reinforcement learning to optimize the networks generated from the model descriptions given by a recurrent network. Baker et al., 2016 [31] introduced MetaQNN, which uses reinforcement learning with a greedy exploration strategy to generate high-performance CNN architectures. Gradientbased learning allows a network to efficiently optimize new instances of a task [37]. Liu et al., 2018 [38] introduced DARTS, which formulate tasks in a differentiable manner to address the scalability of the architecture search. Elsken et al., 2018 [39] introduced LEMON-ADE, which penalizes excessive resource consumption by approximating the network morphism operators while generating subnetworks. Wu et al., 2019 [32] introduced FBNets (Facebook-Berkeley-Nets), which uses a gradient-based method to optimize CNNs created by a differentiable neural architecture search framework. Pham et al., 2018 [40] introduced ENAS, which constructs a large computational graph, where each subgraph represents a neural network architecture; thus, all subgraphs share parameters, delivering strong empirical performances while using a lower amount of resources. Chen et al., 2019 [41] introduced RENAS, which integrates reinforced mutation into the evolutionary search algorithm for architectural exploration, which efficiently evolves the model. Chen et al., 2021 [42] later introduced TE-NAS, which ranks the architectures by analyzing the spectrum of the neural tangent kernel and the number of linear regions in the input space, which imply the trainability and expressivity of the networks.

Evolutionary algorithms are widely used to deal with complex and non-linear optimization problems [43–45]. It is a common solution for difficult real-world problems where the sample space is large. Genetic algorithms are used to solve search and optimization problems using bio-inspired operators [35]. Real et al. [36] uses genetic algorithms to discover neural network architectures, minimizing the role of humans in the design.

### 2.3. Network Scaling

Network scaling is necessary to keep up with the growing datasets that have very large samples and large memory requirements. Network scaling is usually implemented on network architectures after they are constructed. Residual networks, introduced by He et al., 2016 [9], were scaled up to ResNet-200, and down to ResNet-18. Zagoruyko et al., 2016 [46] introduced WideResNet, which is a width-scaled interpretation of the original residual network, with the number of channels in various layers increased to increase the resolution of feature maps. Modern CNNs have also been shown to use higher-resolution input images. Higher accuracies were obtained when using higher-resolution images at  $299 \times 299$  by Szegedy et al., 2016 [47], and  $331 \times 331$  by Zoph et al., 2018 [13]. Dryden et al., 2019 [48] exploits parallelism in convolutional layers beyond data parallelism to tackle scaling and memory issues. However, this does not focus on batch normalization, ReLUs, pooling, and fully-connected layers, which are also present in conventional networks.

## 3. EvolveNet

EvolveNet algorithm attempts to build new networks from scratch by evolving an ideal configuration of layers for a pre-defined architecture. There are four major steps to EvolveNet: (1) filter training to strengthen the individuality of layers, (2) depth evolution to find the ideal number of layers, (3) width evolution to compute the ideal width for each layer, and (4) retraining to fine-tune the evolved network.

Pre-built networks used bottleneck blocks used by EfficientNet. The bottleneck blocks allow the network to reduce the number of parameters and, consequently, the number of floating-point operations. This makes the network more compact and efficient. The bottleneck operation consists of three operators: a linear transformation followed by a non-linear transformation, and then another linear transformation. Each bottleneck first expands a low-dimensional feature map into a high-dimensional feature map using a point-wise convolution. A depth-wise convolution then performs spatial filtering on the high-dimensional tensor. Finally, another point-wise convolution projects the spatially-filtered map back down into a low-dimensional tensor. We try to change as little about the original network as possible, so that much of the changes made have been made by evolution and not by human intervention.

#### 3.1. Filter Training

An initial network N is constructed as a set of layers and filters whose configurations are ideal as found from existing network architectures as follows:

$$\mathcal{N} = \{ L_j : j \in (0, N) \}, \tag{1}$$

where  $L_j$  represents each layer with channel and kernel sizes. *N* is the total number of layers, with each group of a layer consisting of  $N^*$  layers. Using the weights  $\theta$ , the new

collection of layers is derived into a pseudo network  $\mathcal{N}_{\theta}^*$ , which is a parameterized subset of the layers in  $\mathcal{N}$ , i.e.,  $\mathcal{N}_{\theta}^* \subseteq \theta(\mathcal{N})$ . During each epoch of the training stage, random layers and filters are chosen to be trained which omits certain layers and filters from  $\mathcal{N}$ , as shown in Figure 2.  $\mathcal{N}_{\theta}^*$  is trained and its cross-entropy loss is computed as follows:

$$\mathcal{L}_{D_t}(\mathcal{N}^*_{\theta}) = -\sum_{i=1}^n t_i \log p_i,$$
(2)

where  $\mathcal{L}_{D_t}$  represents the cross entropy loss w.r.t. training data  $D_t$ , and  $t_i$  and  $p_i$  are the truth label and softmax probability of class *i* w.r.t.  $\mathcal{N}_{\theta}^*$ . This loss is smoothed as a consequence of cross-entropy and then backpropagated through every layer of  $\mathcal{N}$  including the omitted layers as follows:

$$\theta \leftarrow \theta - \beta \frac{\partial \mathcal{L}(\mathcal{N}_{\theta})}{\partial \theta},$$
(3)

where  $\beta$  is the learning rate, adjusted by the Lambda scheduler [49] to converge quickly and optimally. This trains the larger network, which will have layers that can efficiently be recalibrated into smaller networks composed of only a few of its layers without impacting the overall accuracy. Each layer contributes to the feature map without taking away from the feature map of the larger network, which is made up of other layers trained in a similar fashion.





#### 3.2. Evolving Depth

The network  $\mathcal{N}_{\theta}$ , which is now a collection of layers that have been trained to work independently of the network, is used to evolve a recalibrated network with ideal depth configurations. Configurations of the architecture of each block, such as the number of out-channels, kernel sizes, and strides of each layer, remain the same as before. The depth of the recalibrated network, i.e., the layers chosen to be trained, are chosen by depth encoding vectors (DEVs). DEV is a vector generated using genetic algorithms that has the depths of each layer of the recalibrated network encoded into it as the presence or absence of a layer in the network. These DEVs generate networks of parameter sizes within preset constraints, and the computed reward is assigned as the reward of each DEV. Since  $\mathcal{N}_{\theta}$ is only a collection of layers strung together to build a network, the recalibrated layer, the layers of which are chosen by the DEV, shows the various layers that will be present in the network to be built.

The reward function computes the pre-train accuracy of the generated network. The networks are then evolved as chromosomes under the same preset constraints to maximize their rewards, as seen in Algorithm 1. The best chromosomes are chosen, then mutated, and crossed over to be propagated through to the next chromosome pool. Muta-

tion is the genetic operation of flipping arbitrary genes in a parent chromosome to generate offspring. Cross-over is the genetic operation of combining the genetic information of the two parent genes to generate offspring. In EvolveNet, the mutation is implemented by giving every gene on a random chromosome a 10% chance to be switched to a random new gene. Cross-over is implemented by choosing two random chromosomes and then selecting a gene from either chromosome with a 50% chance of being selected from either parent. We assume that every genetically modified offspring will not be better than every parent; hence, the subsequent chromosome pools are selected from the overall pool and not just the genetically modified pool.

After multiple chromosome pools are generated, the best chromosome is chosen and the network generated from it has the ideal depth configuration for a network block architecture for the given task.

### Algorithm 1 Algorithm for evolving depth

**Hyperparameters:** Number of searching epochs  $N_s$ , number of fine-tuning epochs  $N_f$ , number of chromosomes in C n, number of chromosomes selected for evolution k **Input:**  $D_{train}$ : training images that can be split into *batches*,  $X_f$ : filter trained network **Functions:** reward(dev) computes reward of the network created using *dev*, *mutation* and *crossover* are evolutionary operations performed on a list of chromosomes, layers(dev, network) converts n *dev* into a list of layers using trained network, *create*(*layers*) creates a network from a list of layers, f(model, data) trains *model* using given data,  $\nabla$  computes gradient of loss of trained network

**Output:** Depth-evolved network *x* 

C = List of n random devs $C_k^{top} = \{\}$ 

```
for i = 0 to N_s do
   R = [r_1, r_2, ..., r_n]
   for j = 0 to n do
       r_i = \text{reward}(dev_i)
   end for
   C_k^{iop}.append(C)
   sort C_k^{top} in descending order of R
   C_{mu} = \text{mutation}(C_k^{top}[:10])

C_{co} = \text{crossover}(C_k^{top}[:10])
   C = C_{mu} + C_{co}
   C_k^{top} = C_k^{top}[:k]
end for
[l_1, l_2, ..., l_n] = \text{layers}(C_k^{top}[0], X_f)
x = create([l_1, l_2, ..., l_n])
for i = 0 to N_f do
   \mathbf{x} += \nabla \mathbf{f}(x, \hat{D}_{train})
end for
```

## 3.3. Evolving Width

The best DEV after depth evolution is used to build a recalibrated network for width evolution, as explained in Algorithm 1. The recalibrated network is then used to compute rewards for the networks derived from it using the width encoding vectors (WEVs). Unlike DEVs, each WEV is injectively mapped to a network of specific depth and layers of outchannels, i.e., the networks created by every WEV are unique to each element in it and also to the DEV it is evolved from. As shown in Figure 3, 50 WEVs are built from the recalibrated network, and their rewards are computed. Each gene of a chromosome, represented by the WEV, is encoded with the ratio of channels of each layer compared to the original. The individual genes chosen do not matter, as this evolution is used to resolve the size of the final network, not the specific configuration. The size of each chromosome is dependent on the size of the recalibrated network. At every step before evolution, random chromosomes are used to generate sub-networks for the recalibrated network, and their reward is computed agnostic to the training dataset. The chromosomes with the highest rewards are mutated and crossed over, and the subsequent chromosome pool is created by selecting the best chromosomes from a pool of the parents and their offspring.



**Figure 3.** Overview of width evolution. White layers denote the layers omitted by DEVs, whereas the width of green and red channels denotes the channels selected and omitted by WEV.

The chromosome with the highest reward from the final chromosome pool is used to derive the final network from the recalibrated network. The final derived network has had its depth and width evolved to be ideal for the given task and architecture.

## 3.4. Retraining

Once the final configurations of the network have been evolved, the derived network is retrained to achieve competitive accuracy with minimum parameters. The final network is a subnetwork; hence, according to Frankle et al., it has to be trained to achieve similar accuracy. The number of retraining epochs is determined by the network whose architecture is used to build  $N_{\theta}$ , i.e., EfficientNet [15]. The number of layers and their configurations have been determined by the DEV, whereas the number of channels in each specific layer is determined by the WEV. Each unique value of DEV and WEV thus generate a unique network. Thus, unlike in EfficientNet, where the network configuration was found using grid search, the channel configurations of the final network here are selected by evolutionary computation without manual influence.

## 4. Experiments

In this section, we demonstrate the superiority of networks generated by EvolveNet. We compare the results obtained with other major state-of-the-art models. Lastly, we discuss the impact of various hyperparameters to understand their impact on the proposed method.

## 4.1. Experimental Settings

The generated networks were trained on ImageNet dataset [50] for image classification tasks. The kernel sizes were selected after conducting ablation studies and the best results were found for kernels of sizes  $1 \times 1$  and  $3 \times 3$ . In other words, the ideal resolution was hand-crafted. The data augmentation involves random cropping after randomly resizing from  $256 \times 256$  to  $224 \times 224$ . The images are also randomly rotated and flipped to further augment the data. Experiments were conducted on four NVIDIA RTX A6000 GPUs with 40 workers and a batch size of 512.

### 4.2. Evaluation Protocol

We measured the Top-1 and Top-5 accuracies, as well as the number of parameters, to evaluate the evolved networks, and compared them to existing state-of-the-art networks. The primary aim of this experiment is to showcase an improvement in the performance and efficiency of a network before and after the depth and width have evolved. Accuracy is the proportion of images that have been labeled correctly. For each image, the network computes the probability of them being classified into each label. Top-1 accuracy is the proportion of images in which the predicted label is the same as the actual label. Top-5 accuracy is the proportion of images where the actual label is present as at least one of the top five predictions. A fewer number of parameters result in a more streamlined and efficient network.

#### 4.3. Experimental Results

We present the performance of different networks created using EvolveNet. The networks have been created by setting the parameter size as a constraint and generating networks with an ideal network configuration. The generated networks have been called EvolveNet-XS (Extra small EvolveNet), EvolveNet-S (Small EvolveNet), EvolveNet-M (Medium EvolveNet), and EvolveNet-L (Large EvolveNet). These networks are then compared with state-of-the-art methods of similar size.

## 4.3.1. Performance against Very Small Networks

A set of very small state-of-the-art networks, including DenseNet121 [10], HRFormer-T [51], EfficientNetB1 [15], and EfficientNetV2B1 [52], are selected for comparing to EvolveNet-XS. Table 1 presents the result of the comparison with those networks. EvolveNet-XS shows Top-1 and Top-5 accuracy of 80.4% and 95.1%, respectively, with 7.8 M parameters. Overall, it shows competitive performance compared to other methods with the least number of parameters. It outperforms DenseNet121 and HRFormer-T by 5.4% and 1.9% in Top-1 accuracy. Additionally, it has 0.3 M fewer parameters compared to DenseNet121. EvolveNet-XS shows comparable performance to the EfficientNetB1 and EfficientNetV2B1. The difference in number of parameters between EfficientNetB1 and EvolveNet-XS is only 0.1 M. However, EvolveNet-XS shows a gain of 1.3% and 0.7% in Top-1 and Top-5 accuracies, respectively. It outperforms EfficientNetV2B1 by small margins of 0.6% and 0.1% in Top-1 and Top-5 accuracies, respectively, in spite of having 0.4 M fewer parameters.

Table 1. Performance Against Very Small Networks on ImageNet Dataset

Model	Top-1 Accuracy	Top-5 Accuracy	#Parameters
EfficientNetB1 [15]	79.1%	94.4%	7.9 M
HRFormer-T [51]	78.5%	-	8.0 M
DenseNet121 [10]	75.0%	92.3%	8.1 M
EfficientNetV2B1 [52]	79.8%	95.0%	8.2 M
EvolveNet-XS	80.4%	95.1%	7.8 M

4.3.2. Performance against Small Networks

Table 2 presents the results of the EvolveNet-S network compared to small networks, including EfficientNetB2 [15], and EfficientNetV2B1 [52]. EvolveNet-S shows Top-1 and Top-5 accuracy of 81.1% and 95.6%, respectively, with 8.6 M parameters. It shows competitive performance against those networks with the least number of parameters, outperforming LeViT-128 and ConViT-Ti+ by 1.5% and 4.4%, respectively, in Top-1 accuracy while using 0.2 M and 1.4 M parameters lower. EfficientNetV2B1 has 10.2 M parameters with Top-1 and Top-5 accuracy of 80.5% and 95.1%, respectively. Although EvolveNet-S outperforms it by a small margin of 0.6% and 0.5% in Top-1 and Top-5 accuracy, respectively, it has 1.6 M fewer parameters than EfficientNetV2B1. Similarly, EvolveNet-S outperforms EfficientNetB2 by a small margin of 0.7% in Top-5 accuracy. It also shows a decent gain of 1.0% in Top-1 accuracy with 0.6 M fewer parameters compared to EfficientNetB2.

Model	Top-1 Accuracy	<b>Top-5 Accuracy</b>	<b>#Parameters</b>
LeViT-128 [53]	79.6%	-	8.8 M
EfficientNetB2 [15]	80.1%	94.9%	9.2 M
ConViT-Ti+ [54]	76.7%	-	10.0 M
EfficientNetV2B2 [52]	80.5%	95.1%	10.2 M
RevBiFPN [55]	79.0%	-	10.6 M
EvolveNet-S	81.1%	95.6%	8.6 M

Table 2. Performance Against Small Networks on ImageNet Dataset

4.3.3. Performance against Medium-Sized Networks

A set of medium-sized networks, including DenseNet169 [10], TinyNet [56], Efficient-NetB3 [15], and EfficientNetV2B3 [52], are selected for comparing EvolveNet-M network. Table 3 shows the results for EvolveNet-M compared to these networks. It shows top-1 and top-5 accuracy of 82.8% and 96.3%, respectively. Overall, it outperforms other networks by a decent margin, with fewer parameters. It outperforms SAMix ResNet-18, which has 0.4 M more parameters, by 10.5%. DenseNet169, with 14.3 M parameters, and shows 76.2% and 93.2% of Top-1 and Top-5 accuracy, respectively. However, EvolveNet-M, with 3 M fewer parameters, outperforms it by a significant margin of 6.6% and 3.1% in Top-1 and Top-5 accuracy, respectively, despite having 0.6 M fewer parameters. EvolveNet-M shows comparable performance with EfficientNetB3 and EfficientNetV2B3 in Top-5 accuracy. It outperforms them by a small margin of 0.6% and 0.5%. However, EvolveNet-M has 1 M and 3.2 M fewer parameters, respectively. Additionally, the Top-1 accuracy of EvolveNet-M is 1.2% and 0.8% higher than EfficientNetB3 and EfficientNetV2B3, respectively.

Model	<b>Top-1</b> Accuracy	<b>Top-5</b> Accuracy	<b>#Parameters</b>
SAMix ResNet-18 [57]	72.33%	91.8%	11.7 M
Densenet169 [10]	76.2%	93.2%	14.3 M
TinyNet [56]	79.4%	94.5%	11.9 M
EfficientNetB3 [15]	81.6%	95.7%	12.3 M
EfficientNetV2B3 [52]	82.0%	95.8%	14.5 M
EvolveNet-M	82.8%	96.3%	11.3 M

Table 3. Performance Against Medium-sized Networks on ImageNet Dataset

## 4.3.4. Performance against Large Networks

For the comparison of EvolveNet-L, a collection of larger networks with a significantly large number of parameters is selected. These include Xception [58], ConNeXtTiny [1], ConvNeXtSmall [1], NASNETLarge [13], and EfficientNetB4 [15]. The results of EvolveNet-L compared to the above networks are presented in Table 4. EvolveNet-L outperforms other networks with a decent margin in Top-1 accuracy, with the least number of parameters. It shows Top-1 and Top-5 accuracy of 83.2% and 96.5%, respectively, with 17.6 M parameters. The performance of EfficientNetB4, with Top-1 and Top-5 accuracy of 82.9% and 96.4%, respectively, is comparable to EvolveNet-L. However, EvolveNet-L has significantly fewer parameters compared to EfficientNetB4. With 1.9 M fewer parameters, EvolveNet-L outperforms it by 0.3% and 0.1% in Top-1 and Top-5 accuracy, respectively. EvolveNet-L has significantly fewer parameters than NASNETLarge. In spite of having 71.3 M fewer parameters, it outperforms NASNETLarge by 0.7% and 0.5% in Top-1 and Top-5 accuracy, respectively. Similarly, with 32.6 M fewer parameters, it outperforms ConvNeXtSmall by 0.9% in Top-1 accuracy. EvolveNet-L outperforms ConvNeXtSmall by 0.9% in Top-1 accuracy despite having 11 M fewer parameters. Xception has 22.9 M pa-

10 of 14

rameters and shows Top-1 and Top-5 accuracy of 79.0% and 94.5%, respectively. However, EvolveNet-L, with 5.3 M fewer parameters, outperforms it by a decent margin of 4.2% and 2.0% in Top-1 and Top-5 accuracy, respectively.

Model	Top-1 Accuracy	Top-5 Accuracy	#Parameters
Xception [58]	79.0%	94.5%	22.9 M
ConvNeXtTiny [1]	81.3%	-	28.6 M
ConvNeXtSmall [1]	82.3%	-	50.2 M
NasNetLarge [13]	82.5%	96.0%	88.9 M
EfficientNetB4 [15]	82.9%	96.4%	19.5 M
EvolveNet-L	83.2%	96.5%	17.6 M

Table 4. Performance Against Large Networks on ImageNet Dataset

#### 4.4. Discussion

We have experimentally shown that the networks generated by the depth and width encoding vectors evolved using the EvolveNet method consistently show better performance when compared to EfficientNet, while maintaining their architecture. The improvement is significant, and the generated network can still be pruned using the same methods that are used on EfficientNet and other similar CNNs. Hence, it can be inferred that the architecture computed by evolution outperforms the architectures computed using the grid-search method. MobileNetV2 introduced inverted residuals and bottlenecks and improved the accuracy of MobileNetV1 using the new architecture. EfficientNet was an improvement on the MobileNetV2 architecture, where the accuracy of the network was improved by scaling the width, depth, and resolution of MobileNetV2 using grid-search. By evolving networks with higher accuracy and efficiency, EvolveNet has experimentally proven that hand-crafting and grid search are not ideal methods to build networks. Pruning algorithms have shown that a randomly initialized dense network contains multiple sub-networks with fewer parameters and comparable accuracies. However, most pruning algorithms limit themselves by trying to reduce the number of parameters. Since the proposed algorithm evolves networks to emphasize ideal configurations while maximizing rewards, the focus is placed on accuracy, and efficiency is taken care of as a consequence of it. This allows for high accuracies on a relatively smaller network.

The generated network is also independent of the original network, but given the original network and the depth and width encoding vectors, the network can be regenerated. The generated network is injectively mapped to each DEV and WEV and the number of layers in the larger network. Therefore, changing any of these encoding vectors will significantly change the final network. The number of randomly generated layers, from which recalibrated networks are generated, is a hyperparameter used to control the size of the final network, but it has no other bearing on the evolution of the final network. This can be seen in Table 5. Before evolving the width of the final network, the number of out-channels in each layer is equal to the number of out-channels in MobileNetV2. The structure is the same as that of each block in EfficientNet and MobileNetV2. There is one fixed-out channel for each block layer, but the number of additional blocks is determined by the network encoding vector.

Input	Operation	<b>Out-Channels</b>	#Layers
$224^2 \times 3$	Conv2D	32	1
$112^{2} \times 16$	Bottleneck	24	DEV <sub>0</sub>
$112^{2} \times 16$	Bottleneck	24	1
$56^2 \times 24$	Bottleneck	32	$DEV_1$
$56^2 \times 24$	Bottleneck	32	1
$28^2 \times 32$	Bottleneck	64	DEV <sub>2</sub>
$28^2 \times 32$	Bottleneck	64	1
$14^2 \times 64$	Bottleneck	96	DEV <sub>3</sub>
$14^2  imes 64$	Bottleneck	96	1
$14^2 \times 96$	Bottleneck	160	$DEV_4$
$14^2 \times 96$	Bottleneck	160	1
$7^{2} \times 160$	Bottleneck	320	$DEV_5$
$7^{2} \times 160$	Bottleneck	320	1
$7^{2} \times 320$	Conv2D	1280	1
$7^{2} \times 1280$	AvgPool	-	1
$112^2 \times 32$	Conv2D	п	1

**Table 5.** Each layer used in the evolved model.  $DEV_i$  represents the *i*-th element of the depth encoding vector, which indicates the number of layers, and *n* represents the number of classes.

## 5. Conclusions

In this study, we described a simple evolutionary algorithm to evolve ideal depths and widths that are task-agnostic for a given architecture. We also experimentally proved its superiority over networks where the architecture was hand-crafted or grid-searched. EvolveNet is dependent on the architectures of individual blocks in a network, but not the network as a whole. The architecture used can be replaced with better-performing architectures, but the hyperparameters of the algorithm, such as the initial number of layers and range of number of parameters for the output network, should be adjusted appropriately. In future, further research can be conducted on improving the algorithm for the selection of chromosomes. This algorithm can also be applied to networks that process other datasets such as those using speech patterns, human-object interactions, actions etc. Some constraints are more important to the overall schema than others, so treating them as the same would cause inefficiencies in the recalibrated network. EvolveNet can be made more robust with more experiments to understand the importance of the various hyperparameters that make it function.

**Author Contributions:** This research was conceptualized and the methodology and software were written by A.S. It was then validated and formally analyzed by D.-G.L. The writing was drafted, reviewed, and edited by A.S. under the supervision of D.-G.L. The project was administered and funding was secured by D.-G.L. All authors have read and agreed to the published version of the manuscript.

**Funding:** This work was supported by the National Research Foundation of Korea (NRF) grant funded by the Korean Government (MSIT) (No. 2021R1C1C1012590), (No. 2022R1A4A1023248), Project BK21 FOUR, and the Information Technology Research Center (ITRC) support program supervised by the Institute of Information Communications & Technology Planning & Evaluation (IITP) grant funded by the Korean Government (MSIT) (IITP-2023-2020-01808).

Institutional Review Board Statement: Not applicable.

Informed Consent Statement: Not applicable.

**Data Availability Statement:** The Imagenet dataset can be found at: https://www.image-net.org/ download.php, accessed on 22 July 2023.

**Conflicts of Interest:** The authors declare no conflict of interest. The funders had no role in the design of the study; in the collection, analyses, or interpretation of data; in the writing of the manuscript; or in the decision to publish the results.

## References

- 1. Liu, Z.; Mao, H.; Wu, C.Y.; Feichtenhofer, C.; Darrell, T.; Xie, S. A convnet for the 2020s. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, New Orleans, LA, USA, 18–24 June 2022; pp. 11976–11986.
- Krizhevsky, A.; Sutskever, I.; Hinton, G.E. Imagenet classification with deep convolutional neural networks. *Commun. ACM* 2017, 60, 84–90. [CrossRef]
- 3. Sultana, F.; Sufian, A.; Dutta, P. Evolution of image segmentation using deep convolutional neural network: A survey. *Knowl.-Based Syst.* **2020**, *201*, 106062. [CrossRef]
- Kumar, S.; Kumar, A.; Lee, D.G. Semantic Segmentation of UAV Images Based on Transformer Framework with Context Information. *Mathematics* 2022, 10, 4735. [CrossRef]
- Duan, H.; Zhao, Y.; Chen, K.; Lin, D.; Dai, B. Revisiting skeleton-based action recognition. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, New Orleans, LA, USA, 18–24 June 2022; pp. 2969–2978.
- Lee, D.G.; Lee, S.W. Human activity prediction based on sub-volume relationship descriptor. In Proceedings of the 2016 23rd International Conference on Pattern Recognition (ICPR), Cancun, Mexico, 4–8 December 2016; IEEE: Piscataway, NJ, USA, 2016; pp. 2060–2065.
- 7. Zhang, Q.; Zhang, M.; Chen, T.; Sun, Z.; Ma, Y.; Yu, B. Recent advances in convolutional neural network acceleration. *Neurocomputing* **2019**, *323*, 37–51. [CrossRef]
- 8. Liu, Y.; Pu, H.; Sun, D.W. Efficient extraction of deep image features using convolutional neural network (CNN) for applications in detecting and analysing complex food matrices. *Trends Food Sci. Technol.* **2021**, *113*, 193–204. [CrossRef]
- 9. He, K.; Zhang, X.; Ren, S.; Sun, J. Deep residual learning for image recognition. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Las Vegas, NV, USA, 27–30 June 2016; pp. 770–778.
- 10. Huang, G.; Liu, Z.; Van Der Maaten, L.; Weinberger, K.Q. Densely connected convolutional networks. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Honolulu, HI, USA, 21–26 July 2017; pp. 4700–4708.
- Xie, S.; Girshick, R.; Dollár, P.; Tu, Z.; He, K. Aggregated residual transformations for deep neural networks. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Honolulu, HI, USA, 21–26 July 2017; pp. 1492–1500.
- 12. Zoph, B.; Le, Q.V. Neural architecture search with reinforcement learning. *arXiv* **2016**, arXiv:1611.01578.
- Zoph, B.; Vasudevan, V.; Shlens, J.; Le, Q.V. Learning transferable architectures for scalable image recognition. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. CVPR, Salt Lake City, UT, USA, 18–23 June 2018; pp. 8697–8710.
- Sandler, M.; Howard, A.; Zhu, M.; Zhmoginov, A.; Chen, L.C. Mobilenetv2: Inverted residuals and linear bottlenecks. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Salt Lake City, UT, USA, 18–23 June 2018; pp. 4510–4520.
- 15. Tan, M.; Le, Q. Efficientnet: Rethinking model scaling for convolutional neural networks. In Proceedings of the International Conference on Machine Learning. PMLR, Long Beach, CA, USA, 9–15 June 2019; pp. 6105–6114.
- 16. Howard, A.G.; Zhu, M.; Chen, B.; Kalenichenko, D.; Wang, W.; Weyand, T.; Andreetto, M.; Adam, H. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv* **2017**, arXiv:1704.04861.
- 17. Frankle, J.; Carbin, M. The lottery ticket hypothesis: Finding sparse, trainable neural networks. arXiv 2018, arXiv:1803.03635.
- 18. Hutter, F.; Kotthoff, L.; Vanschoren, J. *Automated Machine Learning: METHODS, Systems, Challenges*; Springer Nature: Cham, Switzerland, 2019.
- 19. Elsken, T.; Metzen, J.H.; Hutter, F. Neural architecture search: A survey. J. Mach. Learn. Res. 2019, 20, 1997–2017.
- 20. Yu, K.; Sciuto, C.; Jaggi, M.; Musat, C.; Salzmann, M. Evaluating the search phase of neural architecture search. *arXiv* 2019, arXiv:1902.08142.
- 21. Mellor, J.; Turner, J.; Storkey, A.; Crowley, E.J. Neural architecture search without training. In Proceedings of the International Conference on Machine Learning. PMLR, Virtual, 18–24 July 2021; pp. 7588–7598.
- Ho, Y.C.; Pepyne, D.L. Simple explanation of the no free lunch theorem of optimization. In Proceedings of the 40th IEEE Conference on Decision and Control (Cat. No. 01CH37228), Orlando, FL, USA, 4–7 December 2001; IEEE: Piscataway, NJ, USA, 2001; Volume 5, pp. 4409–4414.
- 23. Liashchynskyi, P.; Liashchynskyi, P. Grid search, random search, genetic algorithm: A big comparison for NAS. *arXiv* 2019, arXiv:1912.06059.
- Godefroid, P.; Khurshid, S. Exploring very large state spaces using genetic algorithms. In Proceedings of the Tools and Algorithms for the Construction and Analysis of Systems: 8th International Conference, TACAS 2002 Held as Part of the Joint European Conferences on Theory and Practice of Software, ETAPS 2002, Grenoble, France, 8–12 April 2002; Proceedings 8; Springer: Berlin/Heidelberg, Germany, 2002; pp. 266–280.
- 25. Zhang, T.; Qi, W.; Zhao, X.; Yan, Y.; Cao, Y. A local dimming method based on improved multi-objective evolutionary algorithm. *Expert Syst. Appl.* **2022**, 204, 117468. [CrossRef]
- 26. Zheng, W.; Sun, J. Two-stage hybrid learning-based multi-objective evolutionary algorithm based on objective space decomposition. *Inf. Sci.* 2022, *610*, 1163–1186. [CrossRef]
- 27. Chen, Q.; Ma, X.; Yu, Y.; Sun, Y.; Zhu, Z. Multi-objective evolutionary multi-tasking algorithm using cross-dimensional and prediction-based knowledge transfer. *Inf. Sci.* 2022, *586*, 540–562. [CrossRef]
- 28. Simonyan, K.; Zisserman, A. Very deep convolutional networks for large-scale image recognition. arXiv 2014, arXiv:1409.1556.

- Szegedy, C.; Liu, W.; Jia, Y.; Sermanet, P.; Reed, S.; Anguelov, D.; Erhan, D.; Vanhoucke, V.; Rabinovich, A. Going deeper with convolutions. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Boston, MA, USA, 7–12 June 2015; pp. 1–9.
- Liu, Z.; Lin, Y.; Cao, Y.; Hu, H.; Wei, Y.; Zhang, Z.; Lin, S.; Guo, B. Swin transformer: Hierarchical vision transformer using shifted windows. In Proceedings of the IEEE/CVF International Conference on Computer Vision, Montreal, QC, Canada, 10–17 October 2021; pp. 10012–10022.
- 31. Baker, B.; Gupta, O.; Naik, N.; Raskar, R. Designing neural network architectures using reinforcement learning. *arXiv* 2016, arXiv:1611.02167.
- Wu, B.; Dai, X.; Zhang, P.; Wang, Y.; Sun, F.; Wu, Y.; Tian, Y.; Vajda, P.; Jia, Y.; Keutzer, K. Fbnet: Hardware-aware efficient convnet design via differentiable neural architecture search. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, Long Beach, CA, USA, 15–20 June 2019; pp. 10734–10742.
- 33. Cai, H.; Zhu, L.; Han, S. Proxylessnas: Direct neural architecture search on target task and hardware. arXiv 2018, arXiv:1812.00332.
- 34. Brock, A.; Lim, T.; Ritchie, J.M.; Weston, N. Smash: One-shot model architecture search through hypernetworks. *arXiv* 2017, arXiv:1708.05344.
- Xie, L.; Yuille, A. Genetic cnn. In Proceedings of the IEEE International Conference on Computer Vision, Venice, Italy, 22–29 October 2017; pp. 1379–1388.
- Real, E.; Moore, S.; Selle, A.; Saxena, S.; Suematsu, Y.L.; Tan, J.; Le, Q.V.; Kurakin, A. Large-scale evolution of image classifiers. In Proceedings of the International Conference on Machine Learning. PMLR, Sydney, NSW, Australia, 6–11 August 2017; pp. 2902–2911.
- Tancik, M.; Mildenhall, B.; Wang, T.; Schmidt, D.; Srinivasan, P.P.; Barron, J.T.; Ng, R. Learned initializations for optimizing coordinate-based neural representations. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, Nashville, TN, USA, 20–25 June 2021; pp. 2846–2855.
- 38. Liu, H.; Simonyan, K.; Yang, Y. Darts: Differentiable architecture search. arXiv 2018, arXiv:1806.09055.
- 39. Elsken, T.; Metzen, J.H.; Hutter, F. Efficient multi-objective neural architecture search via lamarckian evolution. *arXiv* 2018, arXiv:1804.09081.
- 40. Pham, H.; Guan, M.; Zoph, B.; Le, Q.; Dean, J. Efficient neural architecture search via parameters sharing. In Proceedings of the International Conference on Machine Learning. PMLR, Stockholm, Sweden, 10–15 July 2018; pp. 4095–4104.
- Chen, Y.; Meng, G.; Zhang, Q.; Xiang, S.; Huang, C.; Mu, L.; Wang, X. Renas: Reinforced evolutionary neural architecture search. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, Long Beach, CA, USA, 15–20 June 2019; pp. 4787–4796.
- 42. Chen, W.; Gong, X.; Wang, Z. Neural architecture search on imagenet in four gpu hours: A theoretically inspired perspective. *arXiv* **2021**, arXiv:2102.11535.
- 43. Mallipeddi, R.; Suganthan, P.N.; Pan, Q.K.; Tasgetiren, M.F. Differential evolution algorithm with ensemble of parameters and mutation strategies. *Appl. Soft Comput.* **2011**, *11*, 1679–1696. [CrossRef]
- 44. Nguyen, B.M.; Thi Thanh Binh, H.; The Anh, T.; Bao Son, D. Evolutionary algorithms to optimize task scheduling problem for the IoT based bag-of-tasks application in cloud–fog computing environment. *Appl. Sci.* **2019**, *9*, 1730. [CrossRef]
- 45. Bäck, T.; Schwefel, H.P. An overview of evolutionary algorithms for parameter optimization. *Evol. Comput.* **1993**, *1*, 1–23. [CrossRef]
- 46. Zagoruyko, S.; Komodakis, N. Wide residual networks. arXiv 2016, arXiv:1605.07146.
- Szegedy, C.; Vanhoucke, V.; Ioffe, S.; Shlens, J.; Wojna, Z. Rethinking the inception architecture for computer vision. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. CVPR, Las Vegas, NV, USA, 27–30 June 2016; pp. 2818–2826.
- Dryden, N.; Maruyama, N.; Benson, T.; Moon, T.; Snir, M.; Van Essen, B. Improving strong-scaling of CNN training by exploiting finer-grained parallelism. In Proceedings of the 2019 IEEE International Parallel and Distributed Processing Symposium (IPDPS), Rio de Janeiro, Brazil, 20–24 May 2019; IEEE: Piscataway, NJ, USA, 2019; pp. 210–220.
- 49. Lewkowycz, A. How to decay your learning rate. *arXiv* 2021, arXiv:2103.12682.
- Deng, J.; Dong, W.; Socher, R.; Li, L.J.; Li, K.; Fei-Fei, L. Imagenet: A large-scale hierarchical image database. In Proceedings of the 2009 IEEE Conference on Computer Vision and Pattern Recognition, Miami, FL, USA, 20–25 June 2009; IEEE: Piscataway, NJ, USA, 2009; pp. 248–255.
- 51. Yuan, Y.; Fu, R.; Huang, L.; Lin, W.; Zhang, C.; Chen, X.; Wang, J. Hrformer: High-resolution transformer for dense prediction. *arXiv* 2021, arXiv:2110.09408.
- 52. Tan, M.; Le, Q. Efficientnetv2: Smaller models and faster training. In Proceedings of the International Conference on Machine Learning. PMLR, Virtual, 18–24 July 2021; pp. 10096–10106.
- 53. Graham, B.; El-Nouby, A.; Touvron, H.; Stock, P.; Joulin, A.; Jégou, H.; Douze, M. Levit: A vision transformer in convnet's clothing for faster inference. In Proceedings of the IEEE/CVF International Conference on Computer Vision, Montreal, QC, Canada, 10–17 October 2021; pp. 12259–12269.
- d'Ascoli, S.; Touvron, H.; Leavitt, M.L.; Morcos, A.S.; Biroli, G.; Sagun, L. Convit: Improving vision transformers with soft convolutional inductive biases. In Proceedings of the International Conference on Machine Learning. PMLR, Virtual, 18–24 July 2021; pp. 2286–2296.

- 55. Chiley, V.; Thangarasa, V.; Gupta, A.; Samar, A.; Hestness, J.; DeCoste, D. RevBiFPN: The Fully Reversible Bidirectional Feature Pyramid Network. *arXiv* **2022**, arXiv:2206.14098.
- 56. Han, K.; Wang, Y.; Zhang, Q.; Zhang, W.; Xu, C.; Zhang, T. Model rubik's cube: Twisting resolution, depth and width for tinynets. *Adv. Neural Inf. Process. Syst.* **2020**, *33*, 19353–19364.
- 57. Li, S.; Liu, Z.; Wu, D.; Liu, Z.; Li, S.Z. Boosting discriminative visual representation learning with scenario-agnostic mixup. *arXiv* **2021**, arXiv:2111.15454.
- Chollet, F. Xception: Deep learning with depthwise separable convolutions. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, Honolulu, HI, USA, 21–26 July 2017; pp. 1251–1258.

**Disclaimer/Publisher's Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.