

# Convolutional Neural Networks with Dynamic Regularization

Yi Wang, *Student Member, IEEE*, Zhen-Peng Bian, Junhui Hou, *Senior Member, IEEE*,  
and Lap-Pui Chau, *Fellow, IEEE*

**Abstract**—Regularization is commonly used for alleviating overfitting in machine learning. For convolutional neural networks (CNNs), regularization methods, such as DropBlock and Shake-Shake, have illustrated the improvement in the generalization performance. However, these methods lack a self-adaptive ability throughout training. That is, the regularization strength is fixed to a predefined schedule, and manual adjustments are required to adapt to various network architectures. In this paper, we propose a dynamic regularization method for CNNs. Specifically, we model the regularization strength as a function of the training loss. According to the change of the training loss, our method can dynamically adjust the regularization strength in the training procedure, thereby balancing the underfitting and overfitting of CNNs. With dynamic regularization, a large-scale model is automatically regularized by the strong perturbation, and vice versa. Experimental results show that the proposed method can improve the generalization capability on off-the-shelf network architectures and outperform state-of-the-art regularization methods.

**Index Terms**—CNN, image classification, regularization, overfitting, generalization.

## I. INTRODUCTION

CONVOLUTIONAL neural networks (CNNs), which use a stack of convolution operations followed by non-linear activation (e.g., Rectified Linear Unit, ReLU) to extract high-level discriminative features, have achieved considerable improvements for visual tasks [1], [2], [3]. Recent advances of the CNN architectures, such as ResNet [2], DenseNet [4], ResNeXt [5], and PyramidNet [6], ease the vanishing gradient problem and boost the performance. However, CNNs still suffer from the overfitting problem, which reduces their generalization capability.

A wide variety of regularization strategies were exploited to alleviate overfitting and decrease the generalization error. Data augmentation [1] is a simple yet effective manner to improve the diversity of training data. Batch normalization [7] standardizes the mean and variance of features of each mini-batch, which makes the optimization landscape smoother [8]. Dropout [9] aims to train an ensemble of sub-networks,

weakening the effect of “co-adaptions” on training data. DropBlock [10] introduces a structured dropout approach, which drops the contiguous regions of a feature map. Shake-Shake regularization [11] was proposed to randomly interpolate two complementary features in the two residual branches of ResNeXt, achieving state-of-the-art classification performance. ShakeDrop [12] incorporates the idea of stochastic depth [13] with Shake-Shake regularization to stabilize the training process for ResNet-like architectures. Despite the impressive improvement of the regularization methods, there are two main drawbacks with these methods.

- 1) The regularization strength (or amplitude) is not flexible for different network architectures. For example, the ShakeDrop was designed for deep networks but not suitable for shallow networks. Instead of improving classification performance, it even worsens the performance of shallow networks (see the TABLE I).
- 2) The regularization strength is unchangeable over the whole training process. The fixed strong regularization is beneficial to reduce overfitting, but it causes difficulties to fit data at the beginning of training. From the perspective of curriculum learning [14], the learner should begin with easy examples.

In view of these issues, we propose a dynamic regularization method for CNNs, in which the regularization strength is adaptable to the change of the training loss. During training, the regularization strength is gradually increased with respect to the training status. Analogous to human education, the regularizer is regarded as an instructor who gradually increases the difficulty of training examples in a form of feature perturbation. The dynamic regularization can adapt to different model sizes. It provides a strong regularization for large-scale models, and vice versa. (See Fig. 5 (b)). That is, the regularization strength grows faster and achieves a higher value for a large-scale model than that of a light model.

Fig. 1 shows the proposed dynamic regularization in the ResNet structure. The training loss is not only used to perform backpropagation but also exploited to update the amplitude of the regularization. The features are multiplied by the regularizer in the residual branch. The regularizer works as a perturbation which introduces an augmentation in feature space, so CNNs are trained by the diversity of augmented features. Additionally, the regularization amplitude is changeable with respect to the change of the training loss. We conduct experiments on the image classification task to evaluate our regularization strategy. Experimental results show

Yi Wang and Lap-Pui Chau are with School of Electrical and Electronics Engineering, Nanyang Technological University, Singapore, 639798 (e-mail: wang1241@e.ntu.edu.sg, elpchau@ntu.edu.sg).

Zhen-Peng Bian is with Singapore Telecommunications Limited, Singapore, 239732 (e-mail: zbian1@ntu.edu.sg).

Junhui Hou is with Department of Computer Science, City University of Hong Kong (e-mail: jh.hou@cityu.edu.hk).

Yi Wang and Zhen-Peng Bian contributed to this work equally.

Corresponding author: Lap-Pui Chau.

This work was supported in part by the Hong Kong Research Grants Council under grants 9048123 (CityU 21211518) and 9042820 (CityU 11219019).

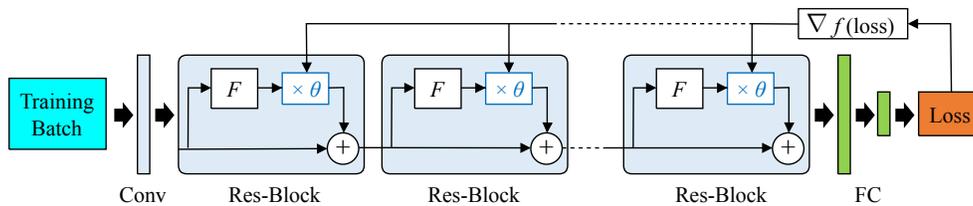


Fig. 1. The proposed dynamic regularization in the ResNet structure. Conv denotes the convolutional layer. FC denotes the fully connected layer.  $F$  denotes the residual function.  $\nabla f(\text{loss})$  denotes a backward difference of the training loss. The dynamic regularization aims to make a self-adaptive schedule throughout training for various network sizes by adjusting the strength of the random perturbation  $\theta$ . As a manner of feature augmentation, the  $\theta$  introduces noises for the residual branch in the forward and backward process.

that the proposed dynamic regularization outperforms state-of-the-art regularization methods, i.e., PyramidNet, ResNeXt, and DenseNet equipped with our dynamic regularization improve the classification accuracy in various model settings, when compared with the same networks with ShakeDrop [12], Shake-Shake [11], and DropBlock [10], respectively.

The rest of this paper is organized as follows. We first briefly introduce the related work on deep CNNs and regularization methods in Section II. Then, the proposed dynamic regularization is presented in Section III. Experimental results and discussion are given in Section IV. Finally, Section V concludes this paper.

## II. RELATED WORK

### A. Deep CNNs

CNNs have become deeper and wider with a more powerful capacity [2], [4], [6], [15], [16]. As our proposed regularization is based on ResNet and its variants, we briefly review the basic structure of ResNet, i.e., the residual block.

**Residual block.** The residual block (Res-Block, shown in Fig. 1) is formulated as

$$x_{l+1} = x_l + F(x_l, \mathcal{W}_l), \quad (1)$$

where an identity branch  $x_l$  is the input features of the  $l^{\text{th}}$  Res-Block, which is added by a residual branch  $F$  that is a non-linear transformation between  $x_l$  and a set of parameters  $\mathcal{W}_l$  ( $\mathcal{W}_l$  will be omitted for simplicity in the following).  $F$  consists of two Conv-BN-ReLU or Bottleneck Architectures in the original ResNet structure [2]. In recent works,  $F$  was also designed to other forms, e.g. Wide-ResNet [17], Inception module [18], PyramidNet [6], and ResNeXt [5]. PyramidNet gradually increases the number of channels in the Res-Blocks as the layers go deep. ResNeXt has multiple aggregated residual branches expressed as

$$x_{l+1} = x_l + F_1(x_l) + F_2(x_l), \quad (2)$$

where  $F_1$  and  $F_2$  are two residual branches. The number of branches (namely cardinality) is not limited.

### B. Regularization

In addition to the advances of network architectures, many regularization techniques, e.g., data augmentation [1], [19], stochastic dropping [9], [13], [20], [21], [10], and Shake-based regularization methods [11], [12], have been successfully applied to avoid overfitting of CNNs.

Data augmentation (e.g., random cropping, flipping, and color adjusting [1]) is a simple yet effective strategy to increase the diversity of data. DeVries and Taylor [19] introduced an image augmentation technique, where augmented images are generated by randomly cutting out square regions from input images (called Cutout). Dropout [9] is a widely-used technique which stochastically drops out the hidden nodes from the networks during the training process. Following this idea, Maxout [22], Continuous Dropout [23], DropPath [20], and stochastic depth [13] were proposed. Stochastic depth randomly drops a certain number of residual branches of ResNet so that the network is shrunk in training. By incorporating Dropout with Cutout, DropBlock [10] drops the contiguous regions in a feature map. Adding a parameter norm penalty to the loss function, the weight decay (or Tikhonov regularization) is commonly used for neural networks and linear inverse problems [24]. DisturbLabel [25] imposes noisy labels in the loss function. Shake-based regularization approaches [11], [12] were recently proposed to augment features inside CNNs, which achieve appealing classification performance.

**Shake-based regularization approaches.** Gastaldi [11] proposed a Shake-Shake regularization method, as shown in Fig. 2 (a). A random variable  $\alpha$  is used to control the interpolation of the two residual branches (i.e.,  $F_1(x)$  and  $F_2(x)$  in 3-branch ResNeXt). It is given by:

$$x_{l+1} = x_l + \alpha F_1(x_l) + (1 - \alpha) F_2(x_l), \quad (3)$$

where  $\alpha \in [0, 1]$  follows the uniform distribution in the forward pass. For the backward pass,  $\alpha$  is replaced by another uniform random variable  $\beta \in [0, 1]$  to disturb the learning process. The regularization amplitude of each branch is fixed to 1.

To extend the use of Shake-Shake regularization, Yamada *et al.* [12] introduced a single Shake in 2-branch architectures (e.g., ResNet or PyramidNet) as shown in Fig. 2 (b). Stochastic depth [13] was adopted to stabilize the learning:

$$x_{l+1} = x_l + (b_l + \alpha - b_l \alpha) F(x_l), \quad (4)$$

where  $\alpha \in [-1, 1]$  is a uniform random variable and  $b_l \in \{0, 1\}$  is the Bernoulli random variable determining when to perform the original network (i.e.,  $x_{l+1} = x_l + F(x_l)$ , if  $b_l = 1$ ) or the perturbed one (i.e.,  $x_{l+1} = x_l + \alpha F(x_l)$ , if  $b_l = 0$ ). In the backward pass,  $\alpha$  is replaced by  $\beta \in [0, 1]$ . The probability of  $b_l$  denotes  $p_l = P(b_l = 1)$ , which follows a linear decay rule, i.e.,  $p_l = 1 - \frac{l}{L}(1 - p_L)$ , where  $L$  is the

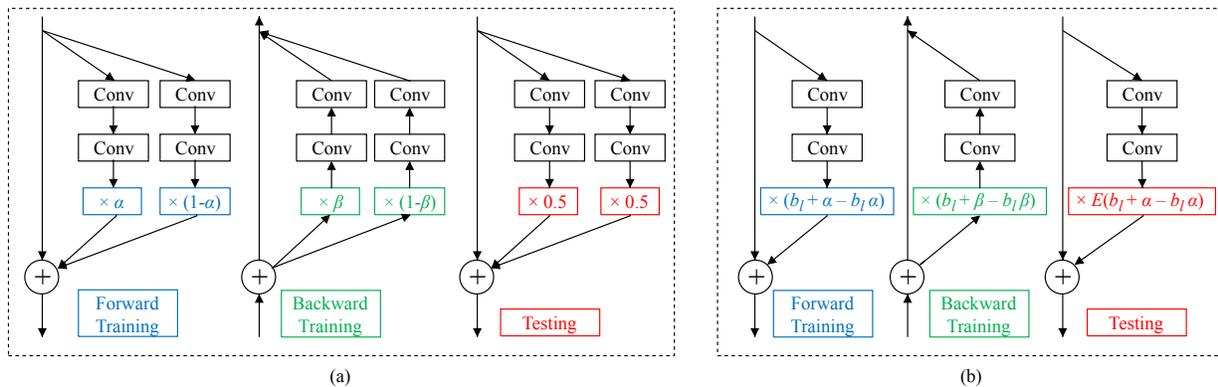


Fig. 2. Shake-based regularization methods in the Res-Block. Some layers (e.g., batch normalization and ReLU) in the residual branch is omitted for simplicity. (a) 3-branch architecture with Shake-Shake regularization [11]. (b) 2-branch architecture with ShakeDrop [12].

total number of Res-Blocks and  $p_L = 0.5$ . The regularization amplitude of the branch is also fixed to 1. We argue that this heavy regularization overemphasizes the overfitting and the fixed regularization amplitude cannot fit the dynamics of the training process and different model sizes well.

### III. THE PROPOSED METHOD

As aforementioned, the fixed regularization strength in the existing regularization methods, such as DropPath [20], Stochastic depth [13], Shake-Shake [11], and Shakedrop [12], departs from the human learning paradigm (e.g., the curriculum learning [14], [21] or self-paced learning [26]). A naive way is to predefine the schedule for updating the regularization strength, such as the linear increment scheme in [10], [27], which linearly increases the regularization strength from low to high. We argue that the predefined schedule is not flexible enough to reveal the learning process. Based on the fact that the loss of the learning system can fully provide the learning status, we propose a dynamic regularization, which is capable of adjusting the regularization strength adaptively.

Our dynamic regularization for CNNs leverages the dynamics of the training loss. That is, at the beginning of training, both the training and testing losses keep decreasing. Through a certain number of iterations, the network overfits the training data, resulting in that the training loss decreases more rapidly than the testing loss. We design a regularization strategy to follow this dynamics. If the training loss drops in an iteration, the regularization strength should increase against the overfitting in the next iteration; otherwise, the regularization strength should decrease against the underfitting. In what follows, we first introduce the dynamic regularization in the residual architectures and then deliberate the update of the regularization strength in each iteration of the training process. We finally extend our dynamic regularization in the densely-connected networks.

#### A. Residual Architectures with Dynamic Regularization

We apply the dynamic regularization method in two residual network architectures: the 2-branch architecture (e.g., PyramidNet [6]) and the 3-branch architecture (e.g., ResNeXt [5]).

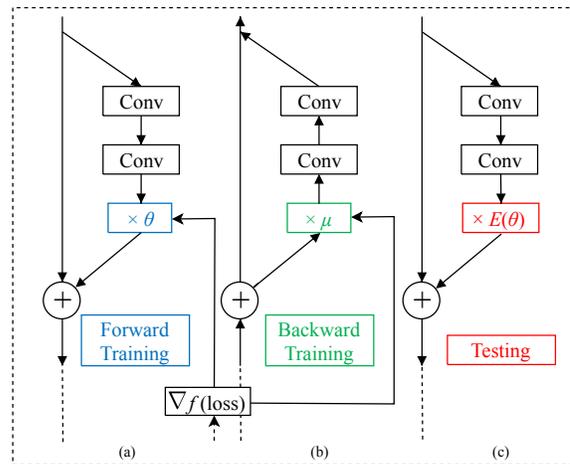


Fig. 3. The 2-branch Res-Block with dynamic regularization.

1) *The 2-branch architecture with dynamic regularization: Training phase.* The dynamic regularization adopted in a Res-Block is shown in Figs. 3 (a) and (b). Specifically, a dynamic regularization unit (called random perturbation) is embedded into the residual branch of a Res-Block. The random perturbation  $\theta$  is achieved by

$$\theta = A + s_i \cdot r, \quad (5)$$

where  $A$  is the basic constant amplitude,  $s_i$  is the dynamic factor at the  $i^{\text{th}}$  iteration, and  $r \in [-R, R]$  is the uniform random noise with the expected value  $E(r) = 0$ . The value of  $s_i$  is updated via the backward difference of the training loss (See Section III.B). The regularization amplitude is proportional to  $A + s_i \cdot R$ . In the forward pass, the output of the  $(l+1)^{\text{th}}$  Res-Block can be expressed as:

$$x_{l+1} = x_l + (A + s_i \cdot r)F(x_l). \quad (6)$$

In the backward pass,  $\theta$  has a different value (represented by  $\mu$  in Fig. 3 (b)) due to the random noise  $r$ .

**Random noise.** The range of  $r$ , i.e.,  $R$ , is a hyper-parameter in the training phase. A straightforward way is to set  $R$  to be uniform inside all Res-Blocks. According to [13], the features of the bottom Res-Blocks should remain more than those of

the top Res-Blocks. Hence, we propose a linear enhancement rule to configure this range inside Res-Blocks. For the  $l^{\text{th}}$  Res-Block, the range denoted as  $R_l$  is given by

$$R_l = l/L, \quad (7)$$

where  $L$  is the total number of Res-Blocks. With the linearly increased  $R$ , the regularization strength is gradually raised from the bottom layers to the top layers. We conducted comparative experiments on different settings of  $R$  in Section IV.C.3.

**Inference phase.** As shown in Fig. 3 (c), we calculate the expected value of  $\theta$  as

$$E(\theta) = E(A + s_i \cdot r) = A, \quad (8)$$

and obtain a forward pass for inference:

$$x_{l+1} = x_l + A \cdot F(x_l). \quad (9)$$

Since  $A$  is a constant, Eq. (9) is equivalent to the standard Res-Block.

2) *The 3-branch architecture with dynamic regularization:* As shown in Fig. 2 (a), we apply the dynamic regularization in the 3-branch architecture. Formally, we use the proposed random perturbation  $\theta$  of Eq. (5) to replace  $\alpha$  of Eq. (3) in Shake-Shake regularization. Hence, the Res-Block with dynamic regularization can be defined as

$$x_{l+1} = x_l + (A + s_i \cdot r)F_1(x_l) + (1 - A - s_i \cdot r)F_2(x_l). \quad (10)$$

If we set  $A = 0.5$ ,  $r \in [-0.5, 0.5]$ , and  $s_i = 1$ , then  $\theta$  ranges from 0 and 1, which is equivalent to  $\alpha$  of Eq. (3). The Shake-Shake regularization can be thought of as a special case of our dynamic regularization with a fixed strength.

### B. Update of the Regularization Strength

The proposed updating solution for the dynamic regularization strength is achieved by the dynamics of the training loss. In particular, the dynamic characteristic of the training loss can be model as the backward difference between the training losses at successive iterations:

$$\nabla \text{loss}_i = \text{loss}_i - \text{loss}_{i-1}, \quad (11)$$

where  $\text{loss}_i$  denotes the training loss at the  $i^{\text{th}}$  iteration. Although the training loss shows a downtrend in training, large fluctuations appear when sequential mini-batches are fed. To eliminate the fluctuations and obtain the overall trend of the loss, we apply a Gaussian filter to smooth it. The filtered backward difference can be rewritten as

$$\nabla f(\text{loss}_i) = f(\text{loss}_i) - f(\text{loss}_{i-1}), \quad (12)$$

where  $f(\cdot)$  is the filtering operation defined as

$$f(\text{loss}_i) = \sum_{n=0}^N w[n] \cdot \text{loss}_{i-n}. \quad (13)$$

The filter length is  $N+1$ . Here we use the normalized Gaussian window and formulate  $w[n]$  as

$$w[n] = \frac{1}{\sqrt{2\pi}(\sigma N/2)} e^{-\frac{1}{2} \left( \frac{n-N/2}{\sigma N/2} \right)^2}, \quad (14)$$

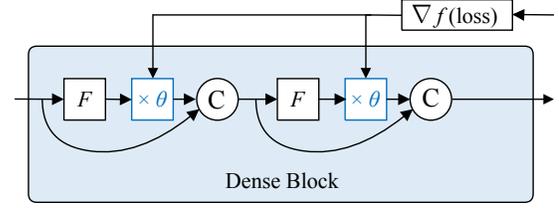


Fig. 4. Dense block with dynamic regularization.  $F$  denotes a convolution operation.  $\theta$  is a random perturbation. 'C' means a concatenation operation.  $\nabla f(\text{loss})$  denotes a backward difference of the training loss.

where  $\sigma = 0.4$ , and  $0 \leq n \leq N$ . The standard deviation is determined by  $\sigma \cdot N/2$ . We will discuss the effectiveness of the Gaussian filter in Section IV.C.4. The dynamic factor in Eqs. (6) and (10) is updated with respect to  $\nabla f(\text{loss}_i)$ , i.e.,

$$s_{i+1} = \begin{cases} s_i + \Delta s, & \nabla f(\text{loss}_i) \leq 0 \\ s_i - \Delta s, & \nabla f(\text{loss}_i) > 0 \end{cases} \quad (15)$$

where  $\Delta s$  is a small constant step for changing the regularization amplitude. From Eq. (15), it can be observed that if the training loss decreases ( $\nabla f(\text{loss}_i) \leq 0$ ), the regularization amplitude increases to avoid overfitting; otherwise, it decreases to prevent underfitting. The dynamic factor keeps updating to reflect the dynamics of the training loss.

**Remark.** Some methods have been proposed to change the regularization strength. For instance, Zoph *et al.* [27] introduced a linear increment scheme, ScheduledDropPath, to regularize NASNets. The probability of dropping out a path is increased linearly throughout training. Following this, DropBlock [10] employs a linearly-increased dropping rate. However, the constant or linear scheme is still a predefined rule, which cannot adapt to the training procedure and different model size. Different from them, our dynamic scheduling exploits the dynamics of the training loss, which is applicable to different network architectures. In Section IV.C.2, we conducted comparisons between them.

### C. Extension to Densely-Connected Networks

To illustrate the flexibility of our method, we further incorporate it into DenseNet [4] by assigning the random perturbations inside the dense block. Fig. 4 shows a two-layer dense block with dynamic regularization, where the perturbations are inserted behind the output features of the convolutional layers. This manner accumulates the noise from all preceding layers to the current layer. A small perturbation could lead to serious noise for the subsequent layers. In experiments, we found that ShakeDrop and DropBlock with default hyper-parameters yield worse results. This is caused by the strong regularization. To decrease the regularization strength, we increase the probability of the Bernoulli random variable (i.e.,  $p_L$ , the rate of keeping original features rather than shaking features) in Eq. (4) for ShakeDrop and increase the *keep\_prob* for DropBlock. Note that for our dynamic regularization, it is not needed to adjust the hyper-parameters used in the residual structure. Our method obtains consistent better results. More details can be seen in Section IV.

TABLE I

COMPARISON OF REGULARIZATION METHODS ON CIFAR100 IN THE 2-BRANCH ARCHITECTURE (I.E., PYRAMIDNET) AND DENSENET ARCHITECTURE. TOP-1 ERROR RATES (%) ARE SHOWN. DYNAMIC DENOTES THE PROPOSED REGULARIZATION METHOD. HP: HYPER-PARAMETERS. THE BEST RESULT UNDER EACH CASE IS BOLD.

Network Architecture	Params	Regularization	Top-1 Error
PyramidNet-110-a48	1.8M	Baseline [6]	23.40
		ShakeDrop [12]	21.60
		DropBlock [10]	21.50
		Dynamic (ours)	<b>21.32</b>
PyramidNet-26-a84	0.9M	Baseline [6]	26.30
		ShakeDrop [12]	31.83
		DropBlock [10]	23.88
		Dynamic (ours)	<b>23.83</b>
PyramidNet-26-a200	3.8M	Baseline [6]	22.53
		ShakeDrop [12]	26.11
		DropBlock [10]	21.22
		Dynamic (ours)	<b>20.34</b>
DenseNet-BC-100-k12 (default HP)	0.8M	Baseline [6]	22.26
		ShakeDrop [6] (pL=0.5)	25.29
		DropBlock [10] (kp=0.9)	23.57
		Dynamic (ours)	<b>20.59</b>
DenseNet-BC-100-k12 (optimized HP)	0.8M	ShakeDrop [6] (pL=0.9)	21.41
		DropBlock [10] (kp=0.95)	21.20

#### IV. EXPERIMENTAL RESULTS

In this section, we evaluate the proposed dynamic regularization on the classification benchmark: CIFAR100 [28] and ImageNet [29], in comparison with three state-of-the-art approaches: Shake-Shake [11], ShakeDrop [12], and DropBlock [10]. Then we conduct ablation studies to compare the fixed or linear-increment scheme of the regularization strength, and discuss the effectiveness of the Gaussian filter and the random noise.

##### A. Implementation Details

1) *CIFAR100*: The following settings are used throughout the experiments. We set the training epoch to 300 and the batch size to 128. The learning rate was initialized to 0.1 for the 2-branch architecture [12], and 0.2 for the 3-branch architecture [11] and DenseNet [4]. We used the cosine learning schedule to gradually reduce the learning rate to 0. The weight decay and momentum was set to 0.0001 and 0.9, respectively. PyramidNet [6], ResNeXt [5], and DenseNet [4] were used as baselines. We employed the standard translation, flipping [1] and Cutout [19] as the data augmentation scheme. Therefore, the regularizer is the only factor to affect experiments. All experimental results are presented by the average of 3 runs at the 300-th epoch.

2) *ImageNet*: ResNet-18 was trained on ImageNet-1k [1] with 120 epochs. The learning rate was initialized to 0.1. Other settings were the same as those in CIFAR100. We reported the single-crop testing results.

3) *Regularizer*: For the dynamic regularization, we set the initial dynamic factor  $s_0 = 0$  and  $A = 0.5$ . We used  $\Delta s = 0.0003$  for the 2-branch architecture and DenseNet, and  $\Delta s = 0.00025$  for the 3-branch architecture. The length of the Gaussian filter was 501. In ShakeDrop [12] and Shake-Shake [11], the default hyper-parameters were employed. In

TABLE II

COMPARISON OF REGULARIZATION METHODS ON CIFAR100 IN THE 3-BRANCH ARCHITECTURE (I.E., RESNEXT). TOP-1 ERROR RATES (%) ARE SHOWN. THE BEST RESULT UNDER EACH CASE IS BOLD.

Network Architecture	Params	Regularization	Top-1 Error (%)
ResNeXt-26-2x32d	2.9M	Baseline [5]	22.95
		Shake-Shake [11]	21.45
		DropBlock [10]	21.20
		Dynamic (ours)	<b>20.91</b>
ResNeXt-26-2x64d	11.7M	Baseline [5]	20.59
		Shake-Shake [11]	19.19
		DropBlock [10]	19.26
		Dynamic (ours)	<b>18.76</b>

DropBlock [10], we used the default hyper-parameters in their paper, i.e., the *keep\_prob* and *block\_size* were set to 0.9 and 7, respectively. To prevent high regularization strength in DenseNet, we increased the value of  $p_L$  of ShakeDrop from 0.5 to 0.9 and increased the value of *keep\_prob* of DropBlock to 0.95. To prevent underfitting of ImageNet, we applied a small  $\Delta s = 5 \times 10^{-7}$  in the dynamic regularization and 0.99 *keep\_prob* in DropBlock.

##### B. Comparison with State-of-the-Art Regularization Methods

1) *2-branch architecture*: We start with comparing the proposed dynamic regularization with ShakeDrop [12] and DropBlock [10] in the 2-branch architecture on CIFAR100. Following the ShakeDrop, we used PyramidNet [6] as our baseline (denoted as Baseline in Table I) and chose different architectures including: 1) PyramidNet-110-a48 (i.e., the network has a depth of 110 layers and a widening factor of 48) which is a deep and narrow network, 2) PyramidNet-26-a84 which is a shallow network, and 3) PyramidNet-26-a200 which is a shallow and wide network.

The first three entries of Table I are the results of PyramidNet. From Table I, it can be observed that our dynamic regularization outperforms the counterparts of ShakeDrop and DropBlock in various architectures. The error rates of ShakeDrop are even worse than those of Baseline in the shallow architectures, i.e., PyramidNet-26-a84 and PyramidNet-26-a200, which means ShakeDrop with a fixed regularization strength fails in this case. This issue comes from stochastic depth [13], where stochastic depth is designed for deep networks. With a linearly-increased dropping rate, DropBlock gains lower error rates than Baseline. However, the predefined schedule of dropping rate is inferior to our dynamic schedule. Regardless of the depth of networks, the dynamic regularization method obtains a consistent improvement.

2) *3-branch architecture*: For the 3-branch architecture, we compare the dynamic regularization with Shake-Shake [11] and DropBlock [10] in ResNeXt-26-2x32d (i.e., the network has a depth of 26 layers and 2 residual branches, and the first residual block has a width of 32 channels) and ResNeXt-26-2x64d. The results are shown in Table II. We can see that the error rates of dynamic regularization are lower than those of Shake-Shake and DropBlock. The results from Tables I and II show that our dynamic regularization can adapt to various network architectures. Our method can decrease the errors by more than 2% on average in comparison with Baseline.

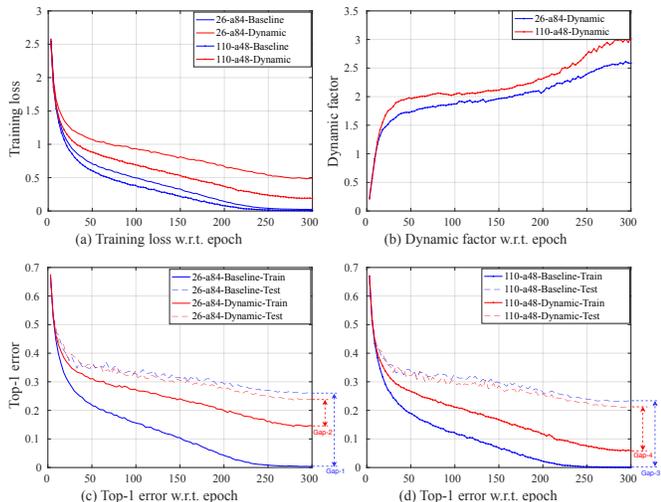


Fig. 5. Illustration of the training loss, dynamic factor, and Top-1 error with respect to epoch for PyramidNet. Gap stands for the difference between training and testing errors. Zoom in the figure for better viewing.

TABLE III

COMPARISON OF REGULARIZATION METHODS ON IMAGENET, WITH SINGLE-CROP TESTING. TOP-1 ERROR RATES (%) ARE SHOWN. THE BEST RESULT IS BOLD.

Network Architecture	Params	Regularization	Top-1 Error (%)
ResNet-18	11.7M	Baseline [2]	29.05
		ShakeDrop [12]	-
		DropBlock [10]	29.06
		Dynamic (ours)	<b>28.82</b>

3) *Densely-connected architecture*: Moreover, we evaluate the regularization methods in DenseNet-BC-100-k12 (i.e., the network uses bottleneck layers and compression with a depth of 100 layers and a growth rate of 12 [4]). The results are shown in the bottom of Table I. Default HP means that all regularizers employed the same hyper-parameters as the ones in PyramidNet. Optimized HP means we adjusted the hyper-parameters in terms of the regularization strength for DenseNet. With the default HP, ShakeDrop and DropBlock damaged the performance of Baseline. We found that the training errors of the two methods were much higher than the testing errors, which means the model underfitted data due to the high regularization strength. With the Optimized HP, we decreased the regularization strengths. We set a larger *keep\_rate* for DropBlock and a larger  $p_L$  for ShakeDrop, so the performance was optimized accordingly. On the contrary, without adjusting the hyper-parameters, our dynamic regularization was stable and reduced the Top-1 error by 1.67% (from 22.26% to 20.59%).

4) *Results on ImageNet*: We evaluate ResNet-18 with the dynamic regularization on ImageNet. The classification results are shown in Table III. Due to a large amount of data trained by a light model, ResNet-18 Baseline underfitted the training data, leading to worse performance when strong regularizers were used. ShakeDrop cannot converge on such a shallow network, so we did not report it. DropBlock also cannot work well, even though we reduced the regularization strength (i.e., the value

TABLE IV  
COMPARISON OF REGULARIZATION SCHEDULES.

PyramidNet-26-a84	Top-1 Error(%)	PyramidNet-26-a84	Top-1 Error(%)
Fix-1	25.45	Linear-1	25.76
Fix-2	24.75	Linear-2	25.09
Fix-3	25.52	Linear-3	24.28
Fix-4	30.52	Linear-4	25.80
Baseline	26.30	Dynamic	<b>23.83</b>

of *keep\_rate* was set to 0.99). Our dynamic regularization performed well in this situation and produced the best result (28.82%).

### C. Ablation Study and Discussion

1) *Effectiveness of dynamic regularization*: Fig. 5 shows the training loss, dynamic factor, and Top-1 error with respect to the epoch in the two networks, i.e., PyramidNet-26-a84 and PyramidNet-110-a48. As shown in Fig. 5 (a), one property of the dynamic regularization is that it prevents the training loss from a rapid descent. In other words, the networks are not easy to fit the training data by rote. Fig. 5 (b) illustrates that the dynamic factors of two networks gradually increase throughout training. Instead of using a predefined scheduling function in [27], our dynamic scheduling is self-adaptive according to the backward difference of the training loss. Another important property of the dynamic regularization is that a low regularization strength is generated for a light model (e.g., PyramidNet-26-a84), and a high strength is for a heavy model (e.g., PyramidNet-110-a48). Figs. 5 (c) and (d) show the networks with dynamic regularization could narrow the gap between the training and testing errors (from Gap-1 to Gap-2 for PyramidNet-26-a84 and from Gap-3 to Gap-4 for PyramidNet-110-a48, respectively) and achieve lower testing errors than the Baselines.

2) *Schedules of the regularization strength*: In [27], [10], the regularization strength is adjusted by a linear-increment schedule, where ScheduledDropPath is used to linearly increase the probability of dropped path (that can also be considered as the regularization strength) in training. Besides, the fixed regularization schedule is commonly used in previous methods [20], [13], [11], [12]. We used PyramidNet-26-a84 as a backbone to compare different regularization schedules.

Table IV illustrates six different configurations of the regularization strength. ‘Fix- $x$ ’ means the dynamic factor is fixed to  $x$  and ‘Linear- $x$ ’ means the dynamic factor is linearly scheduled from 0 to  $x$  over the course of training steps. ‘Fix-2’ and ‘Linear-3’ achieve the best results in fixed and linear schedules, respectively. Compared with them, the dynamic setting with 23.83% error rate achieved the best performance, which shows the effectiveness of our dynamic regularization schedule.

3) *Random noise*: As mentioned in Section III, the range of the random noise involved in our dynamic regularization, i.e.,  $R$ , is designed to grow from bottom Res-Blocks to top Res-Blocks linearly. To evaluate this setting, we performed the dynamic regularization with uniform  $R$  and linearly growing

TABLE V  
EFFECTIVENESS OF LINEARLY GROWING  $R$  AND GAUSSIAN FILTERING.

PyramidNet-26-a84	Top-1 Error (%)
Baseline	26.30
Dynamic-Uniform $R$	25.28
Dynamic-Linear growth $R$	<b>23.83</b>
Dynamic-No filter	25.21
Dynamic-Gaussian filter	<b>23.83</b>

$R$  in PyramidNet-26-a84. From the 2nd and 3rd entries of Table V, we can see the model with uniform  $R$  is inferior to the model with a linearly growing  $R$  (25.28% v.s. 23.83%).

4) *Gaussian Filtering*: In the process of updating the dynamic factor, we employed a Gaussian filter to remove the instant change of the training loss in a mini-batch mode. That is, we refer to the Eq. (12) rather than the Eq. (11) to update the dynamic factor. To study the effectiveness of the Gaussian filter, we conducted comparative experiments between the Dynamic with and without the Gaussian filter. The last two entries of Table V show that if we remove the Gaussian filter, the error rate increases by 1.38%, which validates that the Gaussian filter also plays an important role in dynamic regularization.

## V. CONCLUSION

In this paper, we have presented a dynamic schedule to adjust the regularization strength to fit various network architectures and the training process. Our dynamic regularization is self-adaptive in accordance with the change of the training loss. It produces a low regularization strength for light network architectures and high regularization strength for heavy ones. Furthermore, the strength is self-paced grown to avoid overfitting. Experimental results demonstrate that the proposed dynamic regularization outperforms state-of-the-art ShakeDrop, Shake-Shake, and DropBlock regularization methods. In future, we will investigate the potential of the dynamic regularization in data augmentation and Dropout-based methods.

## REFERENCES

- [1] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in Neural Information Processing Systems*, 2012, pp. 1097–1105.
- [2] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [3] Z.-Q. Zhao, P. Zheng, S.-t. Xu, and X. Wu, "Object detection with deep learning: a review," *IEEE Transactions on Neural Networks and Learning Systems*, 2019.
- [4] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 4700–4708.
- [5] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, "Aggregated residual transformations for deep neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1492–1500.
- [6] D. Han, J. Kim, and J. Kim, "Deep pyramidal residual networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 5927–5935.
- [7] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *International Conference on Machine Learning*, 2015, pp. 448–456.

- [8] S. Santurkar, D. Tsipras, A. Ilyas, and A. Madry, "How does batch normalization help optimization?" in *Advances in Neural Information Processing Systems*, 2018, pp. 2483–2493.
- [9] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, "Dropout: a simple way to prevent neural networks from overfitting," *The Journal of Machine Learning Research*, vol. 15, no. 1, pp. 1929–1958, 2014.
- [10] G. Ghiasi, T.-Y. Lin, and Q. V. Le, "Dropblock: A regularization method for convolutional networks," in *Advances in Neural Information Processing Systems*, 2018, pp. 10 727–10 737.
- [11] X. Gastaldi, "Shake-shake regularization," *CoRR*, vol. abs/1705.07485, 2017.
- [12] Y. Yamada, M. Iwamura, T. Akiba, and K. Kise, "Shakedrop regularization for deep residual learning," *IEEE Access*, vol. 7, pp. 186 126–186 136, 2019.
- [13] G. Huang, Y. Sun, Z. Liu, D. Sedra, and K. Q. Weinberger, "Deep networks with stochastic depth," in *European Conference on Computer Vision*. Springer, 2016, pp. 646–661.
- [14] Y. Bengio, J. Louradour, R. Collobert, and J. Weston, "Curriculum learning," in *Proceedings of the Annual International Conference on Machine Learning*. ACM, 2009, pp. 41–48.
- [15] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *International Conference on Learning Representations*, 2014.
- [16] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 1–9.
- [17] S. Zagoruyko and N. Komodakis, "Wide residual networks," in *British Machine Vision Conference*, 2016.
- [18] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi, "Inception-v4, inception-resnet and the impact of residual connections on learning," in *AAAI Conference on Artificial Intelligence*, 2017.
- [19] T. DeVries and G. W. Taylor, "Improved regularization of convolutional neural networks with cutout," *CoRR*, vol. abs/1708.04552, 2017.
- [20] G. Larsson, M. Maire, and G. Shakhnarovich, "Fractalnet: Ultra-deep neural networks without residuals," in *International Conference on Learning Representations*, 2017.
- [21] P. Morerio, J. Cavazza, R. Volpi, R. Vidal, and V. Murino, "Curriculum dropout," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 3544–3552.
- [22] I. J. Goodfellow, D. Warde-Farley, M. Mirza, A. Courville, and Y. Bengio, "Maxout networks," in *International Conference on Machine Learning*, 2013.
- [23] X. Shen, X. Tian, T. Liu, F. Xu, and D. Tao, "Continuous dropout," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 29, no. 9, pp. 3926–3937, 2017.
- [24] G. De Nicolao and G. Ferrari-Trecate, "Regularization networks for inverse problems: A state-space approach," *Automatica*, vol. 39, no. 4, pp. 669–676, 2003.
- [25] L. Xie, J. Wang, Z. Wei, M. Wang, and Q. Tian, "Disturblabel: Regularizing cnn on the loss layer," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 4753–4762.
- [26] M. P. Kumar, B. Packer, and D. Koller, "Self-paced learning for latent variable models," in *Advances in Neural Information Processing Systems*, 2010, pp. 1189–1197.
- [27] B. Zoph, V. Vasudevan, J. Shlens, and Q. V. Le, "Learning transferable architectures for scalable image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8697–8710.
- [28] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," Citeseer, Tech. Rep., 2009.
- [29] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2009, pp. 248–255.