

Dynamic ReLU

Yinpeng Chen, Xiyang Dai, Mengchen Liu, Dongdong Chen,
Lu Yuan, Zicheng Liu

Microsoft

Abstract. Rectified linear units (ReLU) are commonly used in deep neural networks. So far ReLU and its generalizations (either non-parametric or parametric) are **static**, performing identically for all input samples. In this paper, we propose **Dynamic ReLU** (DY-ReLU), a dynamic rectifier whose parameters are input-dependent as *a hyper function over all input elements*. The key insight is that DY-ReLU *encodes the global context into the hyper function*, and adapts the piecewise linear activation function accordingly. Compared to its static counterpart, DY-ReLU has negligible extra computational cost, but significantly more representation capability, especially for light-weight neural networks. By simply using DY-ReLU for MobileNetV2, the top-1 accuracy on ImageNet classification is boosted from 72.0% to 76.2% with only 5% additional FLOPs.

Keywords: ReLU, Convolutional Neural Networks, Dynamic

1 Introduction

As the default recommendation, rectified linear unit (ReLU) [28,18] is one of the few milestones in the deep learning revolution. It is simple and powerful, greatly improving the performance of feed-forward networks. Thus, it has been widely used in many successful architectures (e.g. ResNet [12], MobileNet[14,31,13] and ShuffleNet [45,25]) for different vision tasks (e.g. recognition, detection, segmentation).

ReLU and its generalizations, either non-parametric (leaky ReLU [26]) or parametric(PReLU [11]) are **static**. They perform in the exactly same way for

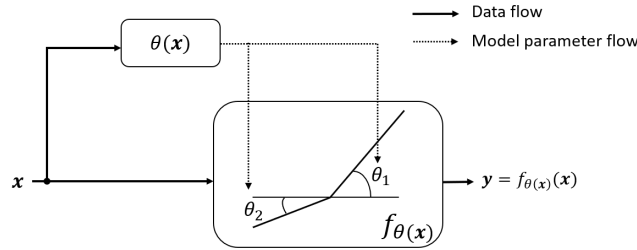


Fig. 1. Dynamic ReLU. The piecewise linear function is determined by the input x .

different inputs (e.g. images). This naturally raises an issue: *should rectifiers be fixed or adaptive to input (e.g. images)?* In this paper, we investigate *dynamic* rectifiers to answer this question.

We propose Dynamic ReLU (DY-ReLU), a piecewise linear function $f_{\theta(\mathbf{x})}(\mathbf{x})$ whose parameters are computed from a hyper function $\theta(\mathbf{x})$ over input \mathbf{x} . For example, Figure 1 illustrates that the slopes of two linear functions are determined by the hyper function. The key idea is that the *global context of all input elements* $\mathbf{x} = \{x_c\}$ is encoded in the hyper function $\theta(\mathbf{x})$ for adapting the piecewise linear activation function $f_{\theta(\mathbf{x})}(\mathbf{x})$. This enables significantly more representation capability especially for light-weight neural networks (e.g. MobileNet). Meanwhile it is computationally efficient as the hyper function $\theta(\mathbf{x})$ is simple with negligible extra computational cost.

Furthermore, we explore three variations of dynamic ReLU. They have different ways of sharing activation functions across spatial locations and channels: (a) spatial and channel-shared DY-ReLU-A, (b) spatial-shared and channel-wise DY-ReLU-B, and (c) spatial and channel-wise DY-ReLU-C. We have two findings. Firstly, channel-wise variations (DY-ReLU-B and DY-ReLU-C) are more suitable for image classification. Secondly, for keypoint detection, channel-wise variations (DY-ReLU-B and DY-ReLU-C) are more suitable for the backbone network while the spatial-wise variation (DY-ReLU-C) is more suitable for the head network.

We demonstrate the effectiveness of DY-ReLU on both image classification (ImageNet) and keypoint detection (COCO). Without bells and whistles, simply replacing static ReLU with dynamic ReLU in multiple networks (i.e. ResNet, MobileNet V2 and V3) achieves solid improvement with only a slight increase (5%) of computational cost. For instance, when using MobileNetV2, our method gains 4.2% top-1 accuracy on image classification and 3.5 AP on keypoint detection, respectively.

2 Related Work

Activation Functions: As a key factor in deep neural network, activation function introduces non-linearity. Among various activation functions, ReLU [10,28,18] is widely used. Three generalizations of ReLU are based on using a nonzero slopes α_i for negative input. Absolute value rectification [18] fixes $\alpha_i = -1$. LeakyReLU [26] fixes α_i to a small value, while PReLU [11] treats α_i as a learnable parameter. RReLU took a further step by making the trainable parameter a random number sampled from a uniform distribution [41]. Maxout [8] generalizes ReLU further, by dividing input into groups and outputs the maximum. One problem of ReLU is that it is not smooth. A number of smooth activation functions have been developed to address this, such as softplus [6], ELU [4], SELU [20], Mish [27]. PELU [34] introduced three trainable parameters into ELU. Recently, empowered by neural architecture search (NAS) techniques [46,30,47,23,40,2,33,36], Ramachandran et al. [29] found several novel activation functions, such as Swish function. Different to these static activation functions

that are input independent, our dynamic ReLU adapts the activation function to the input.

Dynamic Neural Networks: Our method is related to recent works of dynamic neural networks [21,24,35,38,43,16,15,42,3]. D²NN [24], SkipNet [35] and BlockDrop [38] learn an additional controller for skipping part of an existing model by using reinforcement learning. MSDNet [16] allows early-exit based on the prediction confidence. Slimmable Nets [43] learns a single neural network executable at different widths. Once-for-all [1] proposes a progressive shrinking algorithm to train one network that supports multiple sub-networks. Hypernetworks [9] generates network parameters using another hypernetwork. SENet [15] squeezes the global context and use it to reweight channels. Dynamic convolution [42,3] dynamically aggregates convolution kernels based on their attentions which are input dependent. Compared with these works, our method shifts the focus from kernel weights to activation functions, and shows dynamic ReLU is very powerful.

Efficient CNNs: Recently, designing efficient CNN architectures [17,14,31,13,45,25] has been an active research area. MobileNetV1 [14] decomposes 3×3 convolution to depthwise convolution and pointwise convolution. MobileNetV2 [31] introduces inverted residual and linear bottlenecks. MobileNetV3 [13] applies squeeze-and-excitation [15] and employs a platform-aware neural architecture search approach [33] to find the optimal network structure. ShuffleNet further reduces MAdds for 1×1 convolution by channel shuffle operations. ShiftNet [37] replaces expensive spatial convolution by the shift operation and pointwise convolution. Our method provides a new and effective component for efficient networks. It can be directly used in these networks by replacing static ReLU with our dynamic ReLU, with negligible extra computational cost.

3 Dynamic ReLU

We will describe dynamic ReLU (DY-ReLU) in this section. It is a **dynamic** piecewise linear function, whose parameters are input dependent. DY-ReLU does NOT increase either the depth or the width of the network, but increases the model capability efficiently with negligible extra computational cost.

This section is organized as follows. We firstly introduce the generic dynamic activation. Then, we present the mathematical definition of DY-ReLU and how to implement it. Finally, we compare it with prior work.

3.1 Dynamic Activation

For a given input vector (or tensor) \mathbf{x} , the dynamic activation is defined as a function $f_{\theta(\mathbf{x})}(\mathbf{x})$ with learnable parameters $\theta(\mathbf{x})$, which *adapt to the input \mathbf{x}* . As shown in Figure 1, it includes two functions:

1. *hyper function* $\theta(\mathbf{x})$: that computes parameters for the activation function.
2. *activation function* $f_{\theta(\mathbf{x})}(\mathbf{x})$: that computes the activation for the input. Its parameters are generated by the hyper function $\theta(\mathbf{x})$.

Note that the hyper function encodes the global context of all input elements ($x_c \in \mathbf{x}$) to determine the appropriate activation function. This enables significantly more representation power than its static counterpart (e.g. sigmoid, tanh, h-swish [13], ReLU [28,18], LeakyReLU [26], PRelu [11]), especially for light-weight models (e.g. MobileNet). Next, we will discuss dynamic ReLU.

3.2 Definition and Implementation of Dynamic ReLU

Definition: Let us denote the traditional or static ReLU as $\mathbf{y} = \max\{\mathbf{x}, 0\}$, where \mathbf{x} is the input vector. The activation of the c^{th} channel is computed as $y_c = \max\{x_c, 0\}$, where x_c is the input on the c^{th} channel. In contrast, DY-ReLU is defined as the maximum of multiple (K) linear functions as follows:

$$y_c = f_{\boldsymbol{\theta}(\mathbf{x})}(x_c) = \max_{1 \leq k \leq K} \{a_c^k x_c + b_c^k\}, \quad (1)$$

where the linear coefficients (a_c^k, b_c^k) are the output of the hyper function $\boldsymbol{\theta}(\mathbf{x})$ as:

$$\boldsymbol{\theta}(\mathbf{x}) = [a_1^1, \dots, a_C^1, \dots, a_1^K, \dots, a_C^K, b_1^1, \dots, b_C^1, \dots, b_1^K, \dots, b_C^K]^T \quad (2)$$

where C is the number of channels. Note that the activation parameters of each channel (a_c^k, b_c^k) are determined by considering all input channels, i.e. a_c^k and b_c^k are not only related to its corresponding input x_c , but also related to other input elements $x_{j, j \neq c}$.

Example (learning XOR): To make the idea of dynamic ReLU more concrete, we begin with a simple task, i.e. learning XOR function. In this example, we want our network to perform correctly on the four points $\mathbb{X} = \{[0, 0]^T, [0, 1]^T, [1, 0]^T, [1, 1]^T\}$. Compared with the solution [7] using *two* linear layers and one static ReLU, DY-ReLU only needs a *single* linear layer as follows:

$$y_1 = a_1^1 z_1 + b_1^1, \quad a_1^1(z_1) = z_1, \quad b_1^1(z_1) = 0, \\ z_1 = \mathbf{W}^T \mathbf{x}, \quad \mathbf{W} = \begin{bmatrix} 1 \\ -1 \end{bmatrix}, \quad (3)$$

where \mathbf{W} is the weight matrix of the first linear layer, which has a single output z_1 . Thus, the activation function only has one channel output y_1 . Here, we use subscript to make it consistent to Eq (1). Actually, DY-ReLU for this case only has one linear function ($K = 1$), with one non-zero parameter a_1^1 that equals to the input $a_1^1(z_1) = z_1$. Essentially, this is equivalent to a quadratic function $y_1 = (x_1 - x_2)^2$. This example demonstrates that dynamic ReLU has more representation power due to its hyper function.

Implementation: next, we show how to model the hyper function $\boldsymbol{\theta}(\mathbf{x})$ in CNNs, where the input is a 3D tensor. We use a light-weight network to model the hyper function that is similar to Squeeze-and-Excitation (SE) [15]. The global spatial information is firstly squeezed by global average pooling. It is then followed by two fully connected layers (with a ReLU between them) and a normalization layer. Different from SE, the output has $2KC$ elements, corresponding


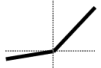
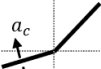
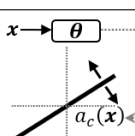
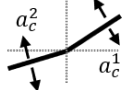
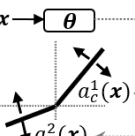
		Type	K	relation to DY-ReLU
ReLU [28,18]		static	2	special case $a_c^1(x) = 1, b_c^1(x) = 0$ $a_c^2(x) = 0, b_c^2(x) = 0$
LeakyReLU [26]		static	2	special case $a_c^1(x) = 1, b_c^1(x) = 0$ $a_c^2(x) = \alpha, b_c^2(x) = 0$
PReLU [11]		static	2	special case $a_c^1(x) = 1, b_c^1(x) = 0$ $a_c^2(x) = a_c, b_c^2(x) = 0$
SE [15]		dynamic	1	special case $a_c^1(x) = a_c(x), b_c^1(x) = 0$ $0 \leq a_c(x) \leq 1$
Maxout [8]		static	1,2,3,...	DY-ReLU is a dynamic and efficient Maxout.
DY-ReLU		dynamic	1,2,3,...	identical

Table 1. Relation to prior work. ReLU, LeakyReLU, PReLU and SE are special cases of DY-ReLU. DY-ReLU is a *dynamic* and *efficient* version of Maxout. α in LeakyReLU is a small number (e.g. 0.01). a_c in PReLU is a parameter to learn.

to the **residual** of $a_{1:C}^{1:K}$ and $b_{1:C}^{1:K}$, which are denoted as $\Delta a_{1:C}^{1:K}$ and $\Delta b_{1:C}^{1:K}$. We simply use $2\sigma(x) - 1$ to normalize the residual between -1 to 1, where $\sigma(x)$ denotes sigmoid function. The final output of the hyper function is computed as the sum of initialization and residual as follows:

$$a_c^k = \alpha^k + \lambda_a \Delta a_c^k, b_c^k = \beta^k + \lambda_b \Delta b_c^k, \quad (4)$$

where α^k and β^k are the initialization values of a_c^k and b_c^k , respectively. λ_a and λ_b are scalars which control the range of residual. $\alpha^k, \beta^k, \lambda_a$ and λ_b are hyper parameters. For the case of $K = 2$, the default initialization values are $\alpha^1 = 1, \alpha^2 = \beta^1 = \beta^2 = 0$, corresponding to the static ReLU. The default λ_a and λ_b are 1.0 and 0.5, respectively.

3.3 Relation to Prior Work

Table 1 shows the relationship between DY-ReLU and prior work. The three special cases of DY-ReLU are equivalent to ReLU [28,18], LeakyReLU [26] and PReLU [11], where the hyper function becomes static. Squeeze-and-Excitation

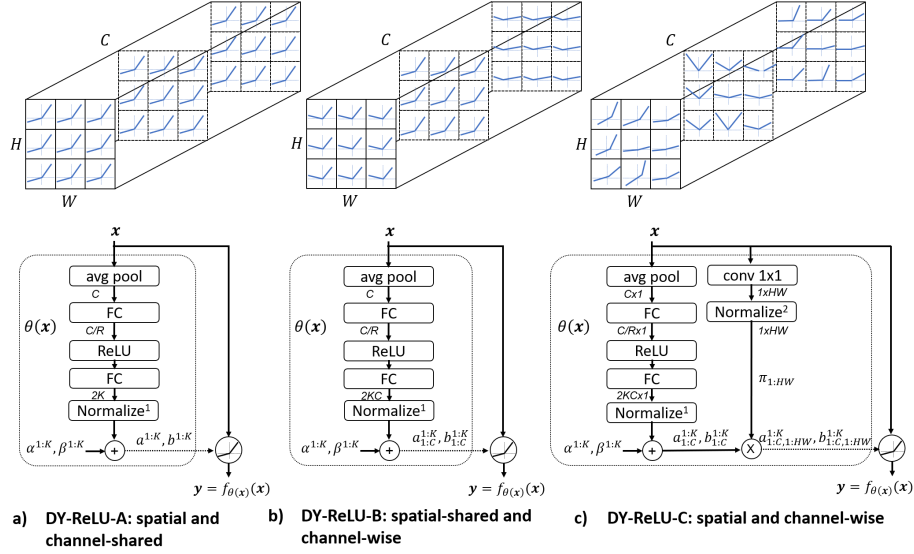


Fig. 2. Three DY-ReLU variations (from left to right). They have different ways of sharing activation functions. For each variation, the top row illustrates the piecewise linear function across spatial locations and channels, and the bottom row shows the network structure for the hyper function. Note that the first FC layer reduces the dimension by R , which is a hyper parameter.

[15] is another special case of DY-ReLU, with a *single* linear function $K = 1$ and zero intercept $b_c^1 = 0$.

DY-ReLU is a *dynamic* and *efficient* Maxout [8], with significantly less computations but even better performance. Literally, Maxout outputs the maximum of results for K convolutional kernels. In contrast, DY-ReLU applies K *dynamic* linear transforms on the results of a *single* convolutional kernel, and outputs the maximum of them. These dynamic linear transforms are powerful and computationally efficient.

4 Variations of Dynamic ReLU

In this section, we introduce another two variations of dynamic ReLU in addition to the option discussed in section 3.2. These three options have different ways of sharing activation functions as follows:

DY-ReLU-A: the activation function is *spatial and channel-shared*.

DY-ReLU-B: the activation function is *spatial-shared and channel-wise*.

DY-ReLU-C: the activation function is *spatial and channel-wise*.

DY-ReLU-B has been discussed in section 3.2.

Variation	Avg Pooling	FC-1	FC-2	Spatial Attention	$f_{\theta(\mathbf{x})}(\mathbf{x})$
DY-ReLU-A	$O(CHW)$	$O(C^2/R)$	$O(2KC/R)$	–	$O(CHW)$
DY-ReLU-B	$O(CHW)$	$O(C^2/R)$	$O(2KC^2/R)$	–	$O(CHW)$
DY-ReLU-C	$O(CHW)$	$O(C^2/R)$	$O(2KC^2/R)$	$O(CHW)$	$O(CHW)$

Table 2. Computational complexity (MAdds) of three DY-ReLU variations on an input tensor with dimension $C \times H \times W$. R is the reduction ratio of the first FC layer in the network for the hyper function (see Figure 2). Compared to a 1×1 convolution with complexity $O(C^2HW)$, DY-ReLU is cheaper by an order of magnitude.

4.1 Network Structure and Complexity

The network structures of three variations are shown in Figure 2. The detailed explanation is discussed as follows:

DY-ReLU-A (Spatial and Channel-shared): the same piecewise linear activation function is shared across all spatial positions and channels. Its hyper function has similar network structure (shown in Figure 2-(a)) to DY-ReLU-B, except the number of outputs is reduced to $2K$. Compared to DY-ReLU-B, DY-ReLU-A has less computational cost, but less representation capability.

DY-ReLU-B (Spatial-shared and Channel-wise): the implementation details are introduced in section 3.2 and its network structure is shown in Figure 2-(b). The activation function requires $2KC$ parameters ($2K$ per channel), which are computed by the hyper function.

DY-ReLU-C (Spatial and Channel-wise): as shown in Figure 2-(c), each input element $x_{c,h,w}$ has a unique activation function $\max_k \{a_{c,h,w}^k x_{c,h,w} + b_{c,h,w}^k\}$, where the subscript c,h,w indicates the c^{th} channel at the h^{th} row and w^{th} column of the feature map with dimension $C \times H \times W$. This introduces an issue that the output dimension is too large ($2KCHW$) to use a fully connected layer to generate. We address it by decoupling spatial locations from channels. Specifically, another branch for computing spatial attention $\pi_{h,w}$ is introduced. The final output is computed as the product of channel-wise parameters ($[a_{1:C}^{1:K}, b_{1:C}^{1:K}]^T$) and spatial attentions ($[\pi_{1:HW}]$). The spatial attention branch is simple, including an 1×1 convolution with a single output channel and a normalization function that is a softmax function with upper cutoff as follows:

$$\pi_{h,w} = \min\left\{\frac{\gamma \exp(z_{h,w}/\tau)}{\sum_{h,w} \exp(z_{h,w}/\tau)}, 1\right\}, \quad (5)$$

where $z_{h,w}$ is the output of 1×1 convolution, τ is the temperature, and γ is a scalar. The softmax is scaled up by γ is to prevent gradient vanishing. We empirically set $\gamma = \frac{HW}{3}$, making the average attention $\pi_{h,w}$ to $\frac{1}{3}$. A large temperature ($\tau = 10$) is used to prevent sparsity during the early training stage. The upper bound 1 constrains the attention between zero and one.

Computational Complexity: DY-ReLU is computationally efficient. For a feature map with dimension $C \times H \times W$, the computational complexities for the

	Top-1	Top-5
ReLU	60.3	82.9
DY-ReLU-A	63.3 _(3.0)	84.2 _(1.3)
DY-ReLU-B	66.4 _(6.1)	86.5 _(3.6)
DY-ReLU-C	66.3 _(6.0)	86.7 _(3.8)

Table 3. Comparing three DY-ReLU variations on Imagenet [5] classification. MobileNetV2 with width multiplier $\times 0.35$ is used. The numbers in brackets denote the performance improvement over the baseline. Channel-wise variations (DY-ReLU-B and DY-ReLU-C) are more effective than the channel-shared (DY-ReLU-A). Spatial-wise (DY-ReLU-C) does NOT introduce additional improvement.

Backbone	Head	AP	AP ^{0.5}	AP ^{0.75}	AP ^M	AP ^L	AR
ReLU	ReLU	59.2	84.3	66.4	56.2	65.0	65.6
DY-ReLU-A	ReLU	58.8 _(-0.4)	84.6	65.3	56.2	64.4	65.4
DY-ReLU-B	ReLU	61.5 _(+2.3)	85.8	68.9	58.5	67.5	67.9
DY-ReLU-C	ReLU	61.9 _(+2.7)	86.0	69.2	59.2	67.5	68.3
ReLU	DY-ReLU-A	57.0 _(-2.2)	83.7	63.5	54.1	62.6	63.6
ReLU	DY-ReLU-B	58.4 _(-0.8)	83.8	64.9	55.5	64.2	64.6
ReLU	DY-ReLU-C	61.0 _(+1.8)	85.2	68.6	58.0	67.1	67.3
DY-ReLU-C	DY-ReLU-C	63.3 _(+4.1)	86.3	71.4	60.3	69.2	69.4

Table 4. Comparing three DY-ReLU variations on COCO [22] keypoint detection. We use MobileNetV2 with width multiplier $\times 0.5$ as backbone and use up-sampling and inverted residual bottleneck blocks [3] in the head. The numbers in brackets denote the performance improvement over the baseline. Channel-wise variations (DY-ReLU-B and DY-ReLU-C) are more effective in the backbone and the spatial-wise variation (DY-ReLU-C) is more effective in the head.

three DY-ReLU variations are listed in Table 2. Compared to a 1×1 convolution, DY-ReLU has reduced complexity by an order of magnitude.

4.2 Ablations

Next, we investigate the three DY-ReLU variations by performing ablation studies on two tasks: image classification and keypoint detection. Here, we focus on the studies, which enable us to select the proper DY-ReLU variation for different tasks. The details of datasets, implementation and training setup will be shown later in the next section.

The comparison among three DY-ReLU variations on ImageNet [22] classification is shown in Table 3. Here we use MobileNetV2 with width multiplier $\times 0.35$. Although all three variations achieve improvement from the baseline, channel-wise DY-ReLUs (variation B and C) are clearly better than the channel-shared DY-ReLU (variation A). Variation B and C have similar accuracy, showing that spatial-wise is not critical for image classification.

Table 4 shows the ablation results on single person pose estimation (or keypoint detection). Similar to image classification, channel-wise DY-ReLUs (variation B and C) are better than the channel-shared one (variation A) in the backbone. In contrast, the spatial-wise variation DY-ReLU-C is critical in the head. This is because the keypoint detection task is spatially sensitive (distinguishing body joints in pixel level), especially in the head network with higher resolutions. The spatial attention allows different magnitudes of activation at different locations, which encourages better learning of DY-ReLU. Using spatial-wise DY-ReLU-C in both the backbone and the head achieves 4.1 AP improvement.

We also observe that the performance is even worse than the baseline if we use DY-ReLU-A in the backbone or use DY-ReLU-A and DY-ReLU-B in the head. We believe that DY-ReLU becomes sensitive if the hyper function is difficult to learn. Specifically, it is hard to learn spatially-shared hyper function for spatially sensitive task (i.e. distinguishing between keypoints and background in pixel level). Thus, learning DY-ReLU-A and DY-ReLU-B in the head with higher resolutions becomes sensitive. This sensitivity can be significantly reduced by introducing spatial attention, as learning spatial-wise hyper function to distinguish keypoints becomes easier. We also find that when spatial attention is involved in the head network, the training converges much faster at the beginning.

In summary, we have two findings from these ablations: (a) using channel-wise DY-ReLU (variation B or C) is important for image classification, and (b) for keypoint detection, channel-wise variations (DY-ReLU-B and DY-ReLU-C) are more suitable for the backbone, and spatial and channel-wise DY-ReLU-C is more suitable for the head. Thus, we use DY-ReLU-B for ImageNet classification and use DY-ReLU-C for keypoint detection in the next section.

5 Experimental Results

In this section, we present experimental results on image classification and single person pose estimation to demonstrate the effectiveness of DY-ReLU. We also report ablation studies to analyze different components of our approach.

5.1 ImageNet Classification

We use ImageNet [5] for all classification experiments. ImageNet has 1000 object classes, including 1,281,167 images for training and 50,000 images for validation. We evaluate DY-ReLU on three CNN architectures (MobileNetV2[31], MobileNetV3 [13] and ResNet [12]), by using DY-ReLU as the activation function after each convolution layer. Note that for MobileNetV3, we remove Squeeze-and-Excitation and replace ReLU and h-swish by DY-ReLU. The main results are obtained by using spatial-shared and channel-wise DY-ReLU-B with two piece-wise linear functions ($K = 2$). The batch size is 256. We use different training setups for the three architectures as follows:

Training setup for MobileNetV2: The initial learning rate is 0.05 and is scheduled to arrive at zero within a single cosine cycle. All models are trained

Network	Activation	#Param	MAdds	Top-1	Top-5
MobileNetV2 $\times 1.0$	ReLU	3.5M	300.0M	72.0	91.0
	DY-ReLU	7.5M	315.5M	76.2 _(4.2)	93.1 _(2.1)
MobileNetV2 $\times 0.75$	ReLU	2.6M	209.0M	69.8	89.6
	DY-ReLU	5.0M	221.7M	74.3 _(4.5)	91.7 _(2.1)
MobileNetV2 $\times 0.5$	ReLU	2.0M	97.0M	65.4	86.4
	DY-ReLU	3.1M	104.5M	70.3 _(4.9)	89.3 _(2.9)
MobileNetV2 $\times 0.35$	ReLU	1.7M	59.2M	60.3	82.9
	DY-ReLU	2.7M	65.0M	66.4 _(6.1)	86.5 _(3.6)
MobileNetV3-Large	ReLU/HS	5.4M	219.0M	75.2	92.2
	DY-ReLU	9.8M	230.5M	75.7 _(0.5)	92.5 _(0.3)
MobileNetV3-Small	ReLU/HS	2.9M	66.0M	67.4	86.4
	DY-ReLU	4.0M	68.7M	69.7 _(2.3)	88.3 _(1.9)
ResNet-50	ReLU	23.5M	3.8G	76.2	92.9
	DY-ReLU	27.6M	3.92G	77.2 _(1.0)	93.4 _(0.5)
ResNet-34	ReLU	21.3M	3.6G	73.3	91.4
	DY-ReLU	25.2M	3.71G	74.4 _(1.1)	92.0 _(0.6)
ResNet-18	ReLU	11.1M	1.81G	69.8	89.1
	DY-ReLU	12.8M	1.86G	71.8 _(2.0)	90.6 _(1.5)
ResNet-10	ReLU	5.2M	0.89G	63.0	84.7
	DY-ReLU	6.3M	0.91G	66.3 _(3.3)	86.7 _(2.0)

Table 5. Comparing DY-ReLU with its static counterpart (ReLU or h-swish, denoted as HS) on ImageNet [5] classification in three network architectures. We use spatial-shared and channel-wise DY-ReLU-B with $K = 2$ linear functions. Note that SE blocks are removed when using DY-ReLU in MobileNetV3. The numbers in brackets denote the performance improvement over the baseline. DY-ReLU outperforms its counterpart for all networks.

using SGD optimizer with 0.9 momentum for 300 epochs. The label smoothing (0.1) is used. The weight decay, dropout rate and data augmentation vary for different width multipliers. The details of hyper parameters are shown in appendix A.4.

Training setup for MobileNetV3: The initial learning rate is 0.1 and is scheduled to arrive at zero within a single cosine cycle. The weight decay is $3e-5$ and label smoothing is 0.1. We use SGD optimizer with 0.9 momentum for 300 epochs. We use dropout rate of 0.1 and 0.2 before the last layer for MobileNetV3-Small and MobileNetV3-Large respectively. We use more data augmentation (color jittering and Mixup [44]) for MobileNetV3-Large.

Training setup for ResNet: The initial learning rate is 0.1 and drops by 10 at epoch 30, 60. The weight decay is $1e-4$. All models are trained using SGD optimizer with 0.9 momentum for 90 epochs. We use dropout rate 0.1 before the last layer and label smoothing for ResNet-18, ResNet-34 and ResNet-50.

Main Results: We compare DY-ReLU with its static counterpart in three CNN architectures (MobileNetV2, MobileNetV3 and ResNet) in Table 5. Without bells and whistles, DY-ReLU outperforms its static counterpart by a clear margin

Activation	K	#Param	MAdds	Top-1	Top-5
ReLU	2	1.7M	59.2M	60.3	82.9
RReLU [41]	2	1.7M	59.2M	60.0 _(-0.3)	81.9 _(-1.0)
LeakyReLU [26]	2	1.7M	59.2M	60.9 _(+0.6)	82.3 _(-0.6)
PReLU (channel-wise) [11]	2	1.7M	59.2M	62.0 _(+1.7)	83.4 _(+0.5)
PReLU (channel-shared) [11]	2	1.7M	59.2M	63.1 _(+2.8)	84.0 _(+1.1)
SE[15]+ReLU	2	2.1M	60.9M	62.8 _(+2.5)	84.6 _(+1.7)
Maxout [8]	2	2.1M	118.3M	64.9 _(+4.6)	85.6 _(+2.7)
Maxout [8]	3	2.4M	177.4M	65.4 _(+5.1)	86.0 _(+3.1)
DY-ReLU-B	2	2.7M	65.0M	66.4 _(+6.1)	86.5 _(+3.6)
DY-ReLU-B	3	3.1M	67.8M	66.6 _(+6.3)	86.8 _(+3.9)

Table 6. Comparing DY-ReLU with related activation functions on ImageNet [5] classification. MobileNetV2 with width multiplier $\times 0.35$ is used. We use spatial-shared and channel-wise DY-ReLU-B with $K = 2, 3$ linear functions. The numbers in brackets denote the performance improvement over the baseline. DY-ReLU outperforms all prior work including Maxout, which has significantly more computations.

for all three architectures, with small extra computational cost ($\sim 5\%$). DY-ReLU gains more than 1.0% top-1 accuracy in ResNet and gains more than 4.2% top-1 accuracy in MobileNetV2. For the state-of-the-art MobileNetV3, our DY-ReLU outperforms the combination of Squeeze-and-Excitation and h-swish (key contributions of MobileNetV3). The top-1 accuracy is improved by 2.3% and 0.5% for MobileNetV3-Small and MobileNetV3-Large, respectively. Note that DY-ReLU achieves more improvement for smaller models (e.g. MobileNetV2 $\times 0.35$, MobileNetV3-Small, ResNet-10). This is because the smaller models are underfitted due to their model size, and dynamic ReLU significantly boosts their representation capability.

The comparison between DY-ReLU and prior work is shown in Table 6. Here we use MobileNetV2 $\times 0.35$, and replace ReLU with different activation functions in prior work. Our method outperforms all prior work with a clear margin. Compared to Maxout which has significantly more computational cost, DY-ReLU gains more than 1% top-1 accuracy. This demonstrates that DY-ReLU not only has more representation capability, but also is computationally efficient. The comparison between DY-ReLU and prior work using MobileNetV2 $\times 1.0$ is shown in appendix A.2. Similarly, our DY-ReLU outperforms all prior work. Note that channel-shared PReLU is better than channel-wise PReLU, which is different from the finding in [11]. This may be due to the different network usage (MobileNet vs VGG).

5.2 Ablation Studies on ImageNet

In this subsection, we run a number of ablations to analyze DY-ReLU. We focus on spatial-shared and channel-wise DY-ReLU-B, and use MobileNetV2 $\times 0.35$ for all ablations. By default, the number of linear functions in DY-ReLU is set as $K = 2$. The initialization values of slope and intercept are set as $\alpha^1 = 1$, $\alpha^2 = \beta^1 = \beta^2 = 0$. The range of slope and intercept are set as $\lambda_a = 1$ and

	K	intercept b_c^k	Activation Function	Top-1	Top-5
ReLU	2		$\max\{x_c, 0\}$	60.3	82.9
DY-ReLU	2		$\max\{a_c x_c, 0\}$	63.8	85.1
	2	✓	$\max\{a_c x_c + b_c, 0\}$	64.0	85.2
	2		$\max_{k=1}^2 \{a_c^k x_c\}$	65.7	86.2
	2	✓	$\max_{k=1}^2 (a_c^k x_c + b_c^k)$	66.4	86.5
	3		$\max_{k=1}^3 \{a_c^k x_c\}$	65.9	86.3
	3	✓	$\max_{k=1}^3 \{a_c^k x_c + b_c^k\}$	66.6	86.8

Table 7. Classification results on Imagenet [5] for using different piecewise linear activation functions in DY-ReLU.

	A1	A2	A3	Top-1	Top-5
ReLU	–	–	–	60.3	82.9
DY-ReLU	✓	–	–	64.2	84.9
	–	✓	–	65.3	85.9
	–	–	✓	62.7	83.8
	✓	✓	–	66.2	86.4
	✓	–	✓	64.5	85.3
	–	✓	✓	65.9	86.2
	✓	✓	✓	66.4	86.5

	R	#param	MAdds	Top-1	Top-5
ReLU	–	1.7M	59.2M	60.3	82.9
DY-ReLU	64	2.0M	64.3M	65.0	85.7
	32	2.1M	64.4M	65.5	86.0
	16	2.3M	64.6M	65.9	86.3
	8	2.7M	65.0M	66.4	86.5
	4	3.6M	65.9M	66.5	86.7

Table 8. Classification results on Imagenet [5] for two ablations of DY-ReLU. **Left:** using DY-ReLU at different layers in MobileNetV2 $\times 0.35$, where A1, A2 and A3 indicate the activation layers after three convolutional layers (1×1 conv, 3×3 depthwise conv, 1×1 conv) in an inverted residual block, respectively. **Right:** using different reduction ratio for the first FC layer in the hyper function (see Figure 2).

$\lambda_b = 0.5$, respectively. The reduction ratio of the first FC layer in the hyper function is set as $R = 8$.

Piecewise Linear Functions: Table 7 shows the classification accuracy using different piecewise linear functions. Compared to the static counterpart, all dynamic activation functions gain at least 3.5% on Top-1 accuracy. In addition, changing the second function from zero to a parametric linear function $a_c^2 x_c + b_c^2$ gains more improvement (1.9%+). The intercept b_c^k is helpful consistently. The gap between $K = 2$ and $K = 3$ is small.

Dynamic ReLU at Different Layers: Table 8-(Left) shows the classification accuracy for using DY-ReLU at three different layers (after 1×1 conv, 3×3 depthwise conv, 1×1 conv) of inverted residual block in MobileNetV2 $\times 0.35$. The accuracy is improved if DY-ReLU is used for more layers. Using DY-ReLU for all three layers yields the best accuracy. If only one layer is allowed to use DY-ReLU, using it after 3×3 depth-wise convolution yields the best performance.

Reduction Ratio R : The reduction ratio of the first FC layer in the hyper function $\theta(\mathbf{x})$ controls the representation capacity and computational cost of DY-ReLU. The comparison across different reduction ratios is shown in Table 8-(Right). Setting $R = 8$ achieves a good trade-off.

α^1	α^2	Top-1	Top-5
1.0	0.0	66.4	86.5
1.5	0.0	65.7	86.2
0.5	0.0	66.1	86.3
0.0	0.0	not converge	
1.0	-0.5	65.2	85.5
1.0	0.5	66.4	86.2
1.0	1.0	66.0	86.1

β^1	β^2	Top-1	Top-5
0.0	0.0	66.4	86.5
-0.1	0.0	66.4	86.5
0.1	0.0	66.2	86.4
0.0	-0.1	65.8	86.2
0.0	0.1	65.3	85.8

λ_a	Top-1	Top-5
0.5	65.3	86.0
1.0	66.4	86.5
2.0	66.3	86.5
3.0	65.5	86.1

Table 9. Classification results on Imagenet [5] for three ablation studies of hyper parameters in DY-ReLU. **Left:** initialization values of slopes (α^k in Eq (4)), **Middle:** initialization values of intercepts (β^k in Eq (4)), **Right:** ranges of slope (λ_a in Eq (4)).

Initialization of Slope (α^k in Eq (4)): As shown in Table 9-(Left), the classification accuracy is not sensitive to the initialization values of slopes if the first slope is not close to zero and the second slope is non-negative.

Initialization of Intercept (β^k in Eq (4)): the performance is stable (shown in Table 9-(Middle)) when both intercepts are close to zero. The second intercept is more sensitive than the first one, as it moves the interception of two lines further away from the origin diagonally.

Range of slope (λ_a in Eq (4)): Making slope range either too wide or too narrow is not optimal, as shown in Table 9-(Right). A good choice is to keep λ_a between 1 and 2.

5.3 COCO Single-Person Keypoint Detection

We use COCO 2017 dataset [22] to evaluate dynamic ReLU on single-person keypoint detection. All models are trained on `train2017`, including 57K images and 150K person instances labeled with 17 keypoints. These models are evaluated on `val2017` containing 5000 images by using the mean average precision (AP) over 10 object key point similarity (OKS) thresholds as the metric.

Implementation Details: We evaluate DY-ReLU on two backbone networks (MobileNetV2 and MobileNetV3) and one head network used in [3]. The head simply uses upsampling and four MobileNetV2’s inverted residual bottleneck blocks. We compare DY-ReLU with its static counterpart in both *backbone* and *head*. The spatial and channel-wise DY-ReLU-C is used here, as we show that the spatial attention is important for keypoint detection, especially in the head network (see section 4.2). Note that when using MobileNetV3 as backbone, we remove Squeeze-and-Excitation and replace either ReLU or h-swish by DY-ReLU. The number of linear functions in DY-ReLU is set as $K = 2$. The initialization values of slope and intercept are set as $\alpha^1 = 1$, $\alpha^2 = \beta^1 = \beta^2 = 0$. The range of slope and intercept are set as $\lambda_a = 1$ and $\lambda_b = 0.5$, respectively.

Training setup: We follow the training setup in [32]. All models are trained from scratch for 210 epochs, using Adam optimizer [19]. The initial learning rate is set as 1e-3 and is dropped to 1e-4 and 1e-5 at the 170th and 200th epoch,

Backbone	Activation	#Param	MAdds	AP	AP ^{0.5}	AP ^{0.75}	AP ^M	AP ^L	AR
MBNetV2 $\times 1.0$	ReLU	3.4M	993.7M	64.6	87.0	72.4	61.3	71.0	71.0
	DY-ReLU	9.0M	1026.9M	68.1 _(3.5)	88.5	76.2	64.8	74.3	73.9
MBNetV2 $\times 0.5$	ReLU	1.9M	794.8M	59.2	84.3	66.4	56.2	65.0	65.6
	DY-ReLU	4.6M	820.3M	63.3 _(4.1)	86.3	71.4	60.3	69.2	69.4
MBNetV3-Large	ReLU/HS	4.1M	896.4M	65.7	87.4	74.1	62.3	72.2	71.7
	DY-ReLU	10.1M	926.6M	67.2 _(1.5)	88.2	75.4	64.1	73.2	72.9
MBNetV3-Small	ReLU/HS	2.1M	726.9M	57.1	83.8	63.7	55.0	62.2	64.1
	DY-ReLU	4.8M	747.9M	60.7 _(3.6)	85.7	68.1	58.1	66.3	67.3

Table 10. Keypoint detection results on COCO validation set. All models are trained from scratch. The four groups use different backbones but share the same head. The first row of each group is the baseline using static activation function (ReLU or h-swish, denoted as HS). The second row of each group corresponds to our method using spatial and channel-wise DY-ReLU-C. Note that SE blocks are removed when using DY-ReLU in MobileNetV3. DY-ReLU outperforms its static counterpart by a clear margin. The numbers in brackets denote the performance improvement over the baseline.

respectively. All human detection boxes are cropped from the image and resized to 256×192 . The data augmentation includes random rotation ($[-45^\circ, 45^\circ]$), random scale ($[0.65, 1.35]$), flipping, and half body data augmentation.

Testing: We use the person detectors provided by [39] and follow the evaluation procedure in [39,32]. The keypoints are predicted on the average heatmap of the original and flipped images. The highest heat value location is then adjusted by a quarter offset from the highest response to the second highest response.

Main Results: Table 10 shows the comparison between DY-ReLU and its static counterpart in four different backbone networks (MobileNetV2 $\times 0.5$ and $\times 1.0$, MobileNetV3 Small and Large). The head network [3] is shared for these four experiments. DY-ReLU outperforms baselines by a clear margin. It gains 3.5 and 4.1 AP when using MobileNetV2 with width multiplier $\times 1.0$ and $\times 0.5$, respectively. It also gains 1.5 and 3.6 AP when using MobileNetV3-Large and MobileNetV3-Small, respectively. These results demonstrate that our method is also effective on keypoint detection.

6 Conclusion

In this paper, we introduce Dynamic ReLU (DY-ReLU), which adapts a piecewise linear activation function dynamically for each input. Compared to its static counterpart (ReLU and its generalizations), DY-ReLU significantly improves the representation capability with negligible extra computation cost, thus is more friendly to efficient CNNs. Our dynamic ReLU can be easily integrated into existing CNN architectures. By simply replacing ReLU (or h-swish) in ResNet and MobileNet (V2 and V3) with DY-ReLU, we achieve solid improvement for both image classification and human pose estimation. We hope DY-ReLU becomes a useful component for efficient network architecture.

A Appendix

In this appendix, we report additional analysis and experimental results for our dynamic ReLU (DY-ReLU) method.

A.1 Is DY-ReLU Dynamic?

In this section, we check if DY-ReLU is dynamic. Therefore, we inspect the input and output of DY-ReLU, and expect different activation values (y) across different images for a fixed input value (e.g. $x = 0.5$). In contrast, for a given input (e.g. $x = 0.5$), the output of ReLU is fixed ($y = 0.5$) regardless of channel or input image. Thus, the input-output pairs of ReLU fall into two lines ($y = x$ if $x > 0$, $y = 0$ otherwise).

Fig. 3 plots the input and output values of DY-ReLU at different blocks (from low level to high level) for 50,000 validation images in ImageNet [5]. We confirm

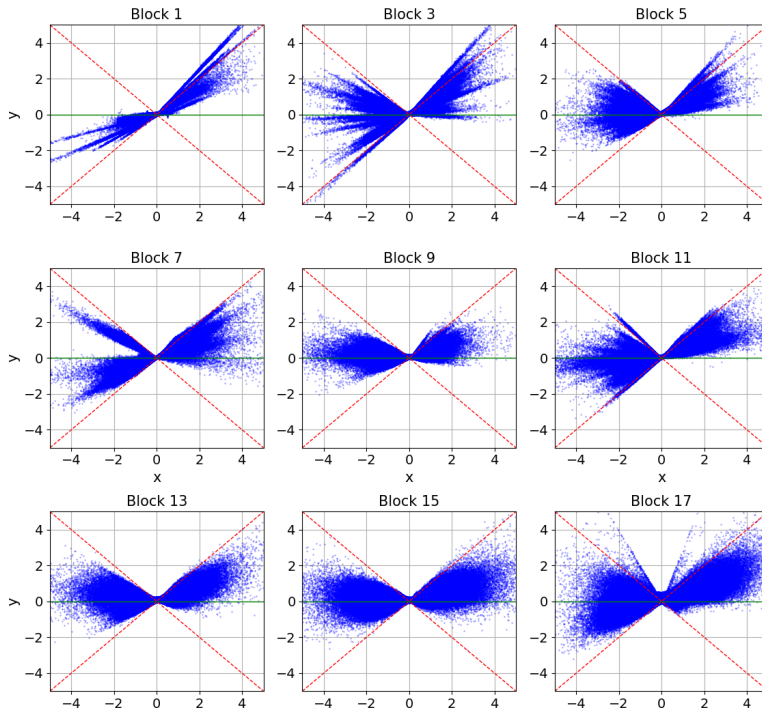


Fig. 3. Plots of input and output values of DY-ReLU in a well trained model (using MobileNetV2 $\times 0.35$) over 50,000 validation images in ImageNet [5]. We choose the dynamic ReLU after the depthwise convolution in every other mobile block. Block 1 (the top-left plot) corresponds to the lowest block, and Block 17 (the bottom-right plot) corresponds to the highest block. The two red lines correspond to $y = x$ and $y = -x$, respectively. Best viewed in color.

that the activation values (y) vary in a range (that blue dots in Fig. 3 cover) for a fixed input x . This demonstrates that the learnt DY-ReLU is dynamic to features. Furthermore, we observe that the distribution of input-output pairs varies across different blocks, indicating different dynamic functions learnt across levels.

A.2 Comparison between DY-ReLU with Prior Work

Table 11 shows the comparison between DY-ReLU and prior work, using MobileNetV2 $\times 1.0$. This is additional to results in Table 6, which are generated by using MobileNetV2 $\times 0.35$. The same conclusion holds for these two experiments: *our method outperforms all prior work*. Compared to Maxout [8] which has significantly more computational cost, DY-ReLU gains 1.1% and 0.4% top-1 accuracy for $K = 2$ and $K = 3$, respectively. This demonstrates that DY-ReLU not only has more representation capability, but also is computationally efficient.

A.3 Ablations of DY-ReLU Variations on Pose Estimation

In this section, we report additional results of comparing three DY-ReLU variations on COCO keypoint detection [22] (or pose estimation). The three variations are listed as follows:

DY-ReLU-A: the activation function is *spatial and channel-shared*.

DY-ReLU-B: the activation function is *spatial-shared and channel-wise*.

DY-ReLU-C: the activation function is *spatial and channel-wise*.

Activation	K	#Param	MAdds	Top-1	Top-5
ReLU	2	3.5M	300.0M	72.0	91.0
RReLU [41]	2	3.5M	300.0M	72.5 _(+0.5)	90.8 _(-0.2)
LeakyReLU [26]	2	3.5M	300.0M	72.7 _(+0.7)	90.8 _(-0.2)
PReLU (channel-wise) [11]	2	3.5M	300.0M	72.9 _(+0.9)	91.0 _(+0.0)
PReLU (channel-shared) [11]	2	3.5M	300.0M	73.3 _(+1.3)	91.2 _(+0.2)
SE[15]+ReLU	2	5.1M	304.8M	74.2 _(+2.2)	91.9 _(+0.9)
Maxout [8]	2	5.7M	579.1M	75.1 _(+3.1)	92.3 _(+1.3)
Maxout [8]	3	7.8M	866.4M	75.8 _(+3.8)	92.7 _(+1.7)
DY-ReLU-B	2	7.5M	315.5M	76.2 _(+4.2)	93.1 _(+2.1)
DY-ReLU-B	3	9.2M	322.8M	76.2 _(+4.2)	93.2 _(+2.2)

Table 11. Comparing DY-ReLU with related activation functions on ImageNet [5] classification. MobileNetV2 with width multiplier $\times 1.0$ is used. We use spatial-shared and channel-wise DY-ReLU-B with $K = 2, 3$ linear functions. The numbers in brackets denote the performance improvement over the baseline. DY-ReLU outperforms all prior work including Maxout, which has significantly more computations.

		Head			
		ReLU	DY-ReLU-A	DY-ReLU-B	DY-ReLU-C
Backbone	ReLU	59.2	57.0	58.4	61.0
	DY-ReLU-A	58.8	51.5	56.5	62.4
	DY-ReLU-B	61.5	54.3	58.6	63.2
	DY-ReLU-C	61.9	53.5	58.8	63.3

Table 12. Comparing three DY-ReLU variations on COCO keypoint detection [22]. The numbers in the table are average precision (AP) over 10 object key point similarity (OKS) thresholds. We use MobileNetV2 with width multiplier $\times 0.5$ as backbone and use up-sampling and inverted residual bottleneck blocks [3] as head. In the backbone, channel-wise variations (DY-ReLU-B and DY-ReLU-C) are more effective than channel-shared (DY-ReLU-A). In the head, spatial-wise variation (DY-ReLU-C) is more effective than spatial-shared (DY-ReLU-A and DY-ReLU-B).

Table 12 shows average precisions (AP) for all 16 combinations of using ReLU or DY-ReLU variations in both backbone and head. This is additional to the results in Table 4. The original conclusions hold: (a) spatial-wise (DY-ReLU-C) is critical in the head network as the last column in Table 12 has higher AP than the previous three columns, and (b) the optimal solution is to use channel-wise variation (variation B or C) in the backbone and use spatial-wise DY-ReLU-C in the head (see the last two rows in the last column of Table 12). Compared to the baseline that uses ReLU in both backbone and head, using DY-ReLU-C achieves 4.1 AP improvement.

The spatial-wise variation (DY-ReLU-C) is a better fit for keypoint detection, which is spatially sensitive (distinguishing body joints in pixel level). This is because the spatial attention allows different magnitudes of activation at different locations. Thus, it encourages better learning of DY-ReLU especially in the head network, which has higher resolutions.

A.4 Implementation Details of MobileNetV2

We now show the implementation details of MobileNetV2. Basically, we use larger weight decay, dropout rate and more data augmentation for higher width multipliers (e.g. $\times 1.0$) to prevent overfitting. We use weight decay $2e-5$ and dropout 0.1 for width $\times 0.35$ and increase weight decay ($3e-5$) and dropout (0.2) for width $\times 0.5$, $\times 0.75$, $\times 1.0$. Random cropping/flipping and color jitter are used for all width multipliers. Mixup [44] is used for width $\times 1.0$. Without using Mixup, the top-1 accuracy of DY-ReLU drops from 76.2% to 75.7%, which still outperforms the static counterpart (72.0%) by a clear margin.

References

1. Cai, H., Gan, C., Han, S.: Once for all: Train one network and specialize it for efficient deployment. ArXiv **abs/1908.09791** (2019)
2. Cai, H., Zhu, L., Han, S.: ProxylessNAS: Direct neural architecture search on target task and hardware. In: International Conference on Learning Representations (2019), <https://openreview.net/forum?id=Hy1VB3AqYm>
3. Chen, Y., Dai, X., Liu, M., Chen, D., Yuan, L., Liu, Z.: Dynamic convolution: Attention over convolution kernels. ArXiv **abs/1912.03458** (2019)
4. Clevert, D.A., Unterthiner, T., Hochreiter, S.: Fast and accurate deep network learning by exponential linear units (elus). arXiv preprint arXiv:1511.07289 (2015)
5. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: 2009 IEEE conference on computer vision and pattern recognition. pp. 248–255. Ieee (2009)
6. Dugas, C., Bengio, Y., Bélisle, F., Nadeau, C., Garcia, R.: Incorporating second-order functional knowledge for better option pricing. In: Advances in neural information processing systems. pp. 472–478 (2001)
7. Goodfellow, I., Bengio, Y., Courville, A.: Deep Learning. The MIT Press (2016)
8. Goodfellow, I.J., Warde-Farley, D., Mirza, M., Courville, A., Bengio, Y.: Maxout networks. arXiv preprint arXiv:1302.4389 (2013)
9. Ha, D., Dai, A.M., Le, Q.V.: Hypernetworks. ICLR (2017)
10. Hahnloser, R.H., Sarpeshkar, R., Mahowald, M.A., Douglas, R.J., Seung, H.S.: Digital selection and analogue amplification coexist in a cortex-inspired silicon circuit. Nature **405**(6789), 947–951 (2000)
11. He, K., Zhang, X., Ren, S., Sun, J.: Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In: ICCV (2015)
12. 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. pp. 770–778 (2016)
13. Howard, A., Sandler, M., Chu, G., Chen, L., Chen, B., Tan, M., Wang, W., Zhu, Y., Pang, R., Vasudevan, V., Le, Q.V., Adam, H.: Searching for mobilenetv3. CoRR **abs/1905.02244** (2019), <http://arxiv.org/abs/1905.02244>
14. 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 preprint arXiv:1704.04861 (2017)
15. Hu, J., Shen, L., Sun, G.: Squeeze-and-excitation networks. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2018)
16. Huang, G., Chen, D., Li, T., Wu, F., van der Maaten, L., Weinberger, K.: Multi-scale dense networks for resource efficient image classification. In: International Conference on Learning Representations (2018), <https://openreview.net/forum?id=Hk2aImxAAb>
17. Iandola, F.N., Moskewicz, M.W., Ashraf, K., Han, S., Dally, W.J., Keutzer, K.: Squeezenet: Alexnet-level accuracy with 50x fewer parameters and <1mb model size. CoRR **abs/1602.07360** (2016), <http://arxiv.org/abs/1602.07360>
18. Jarrett, K., Kavukcuoglu, K., Ranzato, M., LeCun, Y.: What is the best multi-stage architecture for object recognition? In: The IEEE International Conference on Computer Vision (ICCV) (2009)
19. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: International Conference on Learning Representations (ICLR) (2015)

20. Klambauer, G., Unterthiner, T., Mayr, A., Hochreiter, S.: Self-normalizing neural networks. In: Advances in neural information processing systems. pp. 971–980 (2017)
21. Lin, J., Rao, Y., Lu, J., Zhou, J.: Runtime neural pruning. In: Advances in Neural Information Processing Systems, pp. 2181–2191 (2017), <http://papers.nips.cc/paper/6813-runtime-neural-pruning.pdf>
22. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: European conference on computer vision. pp. 740–755. Springer (2014)
23. Liu, H., Simonyan, K., Yang, Y.: DARTS: Differentiable architecture search. In: International Conference on Learning Representations (2019), <https://openreview.net/forum?id=S1eYHoC5FX>
24. Liu, L., Deng, J.: Dynamic deep neural networks: Optimizing accuracy-efficiency trade-offs by selective execution. In: AAAI Conference on Artificial Intelligence (AAAI) (2018)
25. Ma, N., Zhang, X., Zheng, H.T., Sun, J.: Shufflenet v2: Practical guidelines for efficient cnn architecture design. In: The European Conference on Computer Vision (ECCV) (September 2018)
26. Maas, A.L., Hannun, A.Y., Ng, A.Y.: Rectifier nonlinearities improve neural network acoustic models. In: in ICML Workshop on Deep Learning for Audio, Speech and Language Processing (2013)
27. Misra, D.: Mish: A self regularized non-monotonic neural activation function. arXiv preprint arXiv:1908.08681 (2019)
28. Nair, V., Hinton, G.E.: Rectified linear units improve restricted boltzmann machines. In: ICML (2010)
29. Ramachandran, P., Zoph, B., Le, Q.V.: Searching for activation functions. arXiv preprint arXiv:1710.05941 (2017)
30. Real, E., Aggarwal, A., Huang, Y., Le, Q.V.: Regularized evolution for image classifier architecture search. In: AAAI Conference on Artificial Intelligence (AAAI) (2018)
31. 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. pp. 4510–4520 (2018)
32. Sun, K., Xiao, B., Liu, D., Wang, J.: Deep high-resolution representation learning for human pose estimation. In: CVPR (2019)
33. Tan, M., Chen, B., Pang, R., Vasudevan, V., Sandler, M., Howard, A., Le, Q.V.: Mnasnet: Platform-aware neural architecture search for mobile. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2019)
34. Trotter, L., Gigu, P., Chaib-draa, B., et al.: Parametric exponential linear unit for deep convolutional neural networks. In: 2017 16th IEEE International Conference on Machine Learning and Applications (ICMLA). pp. 207–214. IEEE (2017)
35. Wang, X., Yu, F., Dou, Z.Y., Darrell, T., Gonzalez, J.E.: Skipnet: Learning dynamic routing in convolutional networks. In: The European Conference on Computer Vision (ECCV) (September 2018)
36. 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: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2019)
37. Wu, B., Wan, A., Yue, X., Jin, P., Zhao, S., Golmant, N., Gholaminejad, A., Gonzalez, J., Keutzer, K.: Shift: A zero flop, zero parameter alternative to spatial convolutions (2017)

38. Wu, Z., Nagarajan, T., Kumar, A., Rennie, S., Davis, L.S., Grauman, K., Feris, R.: Blockdrop: Dynamic inference paths in residual networks. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2018)
39. Xiao, B., Wu, H., Wei, Y.: Simple baselines for human pose estimation and tracking. In: European conference on computer vision (04 2018)
40. Xie, S., Zheng, H., Liu, C., Lin, L.: SNAS: stochastic neural architecture search. In: International Conference on Learning Representations (2019), <https://openreview.net/forum?id=rylqooRqK7>
41. Xu, B., Wang, N., Chen, T., Li, M.: Empirical evaluation of rectified activations in convolutional network. CoRR (2015)
42. Yang, B., Bender, G., Le, Q.V., Ngiam, J.: Condconv: Conditionally parameterized convolutions for efficient inference. In: NeurIPS (2019)
43. Yu, J., Yang, L., Xu, N., Yang, J., Huang, T.: Slimmable neural networks. In: International Conference on Learning Representations (2019), <https://openreview.net/forum?id=H1gMCsAqY7>
44. Zhang, H., Cisse, M., Dauphin, Y.N., Lopez-Paz, D.: mixup: Beyond empirical risk minimization. In: International Conference on Learning Representations (2018), <https://openreview.net/forum?id=r1Ddp1-Rb>
45. Zhang, X., Zhou, X., Lin, M., Sun, J.: Shufflenet: An extremely efficient convolutional neural network for mobile devices. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2018)
46. Zoph, B., Le, Q.V.: Neural architecture search with reinforcement learning. CoRR **abs/1611.01578** (2017)
47. Zoph, B., Vasudevan, V., Shlens, J., Le, Q.V.: Learning transferable architectures for scalable image recognition. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2018)