# Dynamic ReLU

Yinpeng Chen[0000−0003−1411−225X], Xiyang Dai[0000−0003−1761−8715], Mengchen Liu[0000−0002−9040−1013], Dongdong Chen[0000−0002−4642−4373], Lu Yuan[0000−0001−7879−0389], and Zicheng Liu[0000−0001−5894−7828]

Microsoft Corporation, Redmond WA 98052, USA
{yiche,xidai,mengcliu,dochen,luyuan,zliu}@microsoft.com

**Abstract.** Rectified linear units (ReLU) are commonly used in deep neural networks. So far ReLU and its generalizations (non-parametric or parametric) are **static**, performing identically for all input samples. In this paper, we propose **Dynamic ReLU** (DY-ReLU), a dynamic rectifier of which parameters are generated by 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

Rectified linear unit (ReLU) [27,17] 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 [11], MobileNet[13,30,12] and ShuffleNet [44,24]) for different vision tasks (e.g. recognition, detection, segmentation). ReLU and its generalizations, either non-parametric (leaky ReLU [25]) or parametric(PReLU [10]) are **static**. They perform in the exactly same way for 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.
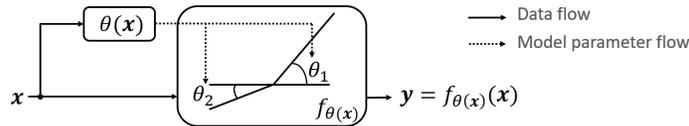


**Fig. 1.** Dynamic ReLU. The piecewise linear function is determined by the input $x$.

We propose dynamic ReLU (DY-ReLU), a piecewise function $f_{\boldsymbol{\theta}(\boldsymbol{x})}(\boldsymbol{x})$ whose parameters are computed from a hyper function $\boldsymbol{\theta}(\boldsymbol{x})$ over input $\boldsymbol{x}$. Figure 1 shows an example 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* $\boldsymbol{x} = \{x_c\}$ is encoded in the hyper function $\boldsymbol{\theta}(\boldsymbol{x})$ for adapting the activation function $f_{\boldsymbol{\theta}(\boldsymbol{x})}(\boldsymbol{x})$. This enables significantly more representation capability, especially for light-weight neural networks (e.g. MobileNet). Meanwhile, it is computationally efficient as the hyper function $\boldsymbol{\theta}(\boldsymbol{x})$ is simple with negligible computational cost.

Furthermore, we explore three variations of dynamic ReLU, which share activation functions across spatial locations and channels differently: (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. They perform differently at different tasks. Channel-wise variations (DY-ReLU-B and DY-ReLU-C) are more suitable for image classification. When dealing with keypoint detection, DY-ReLU-B and DY-ReLU-C are more suitable for the backbone network while the spatial-wise DY-ReLU-C is more suitable for the head network.

We demonstrate the effectiveness of DY-ReLU on both ImageNet classification and COCO keypoint detection. Without bells and whistles, simply replacing static ReLU with dynamic ReLU in multiple networks (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:** activation function introduces non-linearity in deep neural networks. Among various activation functions, ReLU [9,27,17] is widely used. Three generalizations of ReLU are based on using a nonzero slopes $\alpha$ for negative input. Absolute value rectification [17] fixes $\alpha = -1$. LeakyReLU [25] fixes $\alpha$ to a small value, while PReLU [10] treats $\alpha$ as a learnable parameter. RReLU took a further step by making the trainable parameter a random number sampled from a uniform distribution [40]. Maxout [7] 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 [19], Mish [26]. PELU [33] introduced three trainable parameters into ELU. Recently, empowered by neural architecture search (NAS) techniques [45,29,46,22,39,2,32,35], Ramachandran et al. [28] found several novel activation functions, such as Swish function. Different with 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 work of dynamic neural networks [20,23,34,37,42,15,14,41,3]. $D^2NN$ [23], SkipNet [34] and BlockDrop [37] learn a controller for skipping part of an existing model by using reinforcement learning. MSDNet [15] allows early-exit based on the prediction confidence. Slimmable Net [42] learns a single network executable at different

widths. Once-for-all [1] proposes a progressive shrinking algorithm to train one network that supports multiple sub-networks. Hypernetworks [8] generates network parameters using anther network. SENet [14] squeezes global context to reweight channels. Dynamic convolution [41,3] adapts convolution kernels based on their attentions that are input dependent. Compared with these works, our method shifts the focus from kernel weights to activation functions.

**Efficient CNNs:** Recently, designing efficient CNN architectures [16,13,30,12,44,24] has been an active research area. MobileNetV1 [13] decomposes a $3 \times 3$ convolution to a depthwise convolution and a pointwise convolution. MobileNetV2 [30] introduces inverted residual and linear bottlenecks. MobileNetV3 [12] applies squeeze-and-excitation [14], and employs a platform-aware neural architecture search approach [32] to find the optimal network structure. ShuffleNet further reduces MAdds for $1 \times 1$ convolution by group convolution. ShiftNet [36] replaces expensive spatial convolution by the shift operation and pointwise convolution. Our method provides an effective activation function, which can be easily used in these networks (by replacing ReLU) to improve representation capability with low computational cost.

## 3 Dynamic ReLU

Dynamic ReLU (DY-ReLU) is a **dynamic** piecewise function, of which parameters are input dependent. It 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) $\boldsymbol{x}$, the dynamic activation is defined as a function $f_{\boldsymbol{\theta}(\boldsymbol{x})}(\boldsymbol{x})$ with learnable parameters $\boldsymbol{\theta}(\boldsymbol{x})$, which *adapt to the input* $\boldsymbol{x}$. As shown in Figure 1, it includes two functions:

1. *hyper function* $\boldsymbol{\theta}(\boldsymbol{x})$: that computes parameters for the activation function.
2. *activation function* $f_{\boldsymbol{\theta}(\boldsymbol{x})}(\boldsymbol{x})$: that uses the parameters $\boldsymbol{\theta}(\boldsymbol{x})$ to generate activation for all channels.

Note that the hyper function encodes the global context of all input elements $(x_c \in \boldsymbol{x})$ to determine the appropriate activation function. This enables significantly more representation power than its static counterpart, 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 $\boldsymbol{y} = \max\{\boldsymbol{x}, 0\}$, where $\boldsymbol{x}$ is the input vector. For the input $x_c$ at the $c^{th}$ channel, the activation

is computed as $y_c = \max\{x_c, 0\}$. ReLU can be generalized to a parametric piecewise linear function $y_c = \max_k\{a_c^k x_c + b_c^k\}$. We propose dynamic ReLU to further extend this piecewise linear function from static to dynamic by adapting $a_c^k$, $b_c^k$ based upon all input elements $\boldsymbol{x} = \{x_c\}$ as follows:

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

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

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

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

**Implementation of hyper function $\boldsymbol{\theta}(\boldsymbol{x})$:** We use a light-weight network to model the hyper function that is similar to Squeeze-and-Excitation (SE) [14]. For an input tensor $\boldsymbol{x}$ with dimension $C \times H \times W$, the 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. The output has $2KC$ elements, corresponding 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 is computed as the sum of initialization and residual as follows:

$$a_c^k(\boldsymbol{x}) = \alpha^k + \lambda_a \Delta a_c^k(\boldsymbol{x}), \; b_c^k(\boldsymbol{x}) = \beta^k + \lambda_b \Delta b_c^k(\boldsymbol{x}), \tag{3}$$

where $\alpha^k$ and $\beta^k$ are initialization values of $a_c^k$ and $b_c^k$, respectively. $\lambda_a$ and $\lambda_b$ are scalars that control the range of residual. $\alpha^k$, $\beta^k$, $\lambda_a$ and $\lambda_b$ are hyper parameters. For the case of $K = 2$, the default values are $\alpha^1 = 1$, $\alpha^2 = \beta^1 = \beta^2 = 0$, corresponding to static ReLU. The default $\lambda_a$ and $\lambda_b$ 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 [27,17], LeakyReLU [25] and PReLU [10], where the hyper function becomes static. SE [14] 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 [7], with significantly less computations but even better performance. Different with Maxout that requires multiple $(K)$ convolutional kernels, DY-ReLU applies $K$ *dynamic* linear transforms on the results of a *single* convolutional kernel, and outputs the maximum of them. This results in much less computations and even better performance.

## 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:

| | | Type | $K$ | relation to DY-ReLU |
|---|---|---|---|---|
| ReLU [27,17] | | 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 [25] | | 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 [10] | $a_c$ | 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 [14] | $x \rightarrow \boxed{\theta}$ $a_c(x)$ | dynamic | 1 | special case $a_c^1(x) = a_c(x),\ b_c^1(x) = 0$ $0 \leq a_c(x) \leq 1$ |
| Maxout [7] | $a_c^2$ $a_c^1$ | static | 1,2,3,... | DY-ReLU is a dynamic and efficient Maxout. |
| DY-ReLU | $x \rightarrow \boxed{\theta}$ $a_c^1(x)$ $a_c^2(x)$ | 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. $\alpha$ in LeakyReLU is a small number (e.g. 0.01). $a_c$ in PReLU is a parameter to learn.

**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.

### 4.1   Network Structure and Complexity

The network structures of three variations are shown in Fig. 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 Fig. 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.
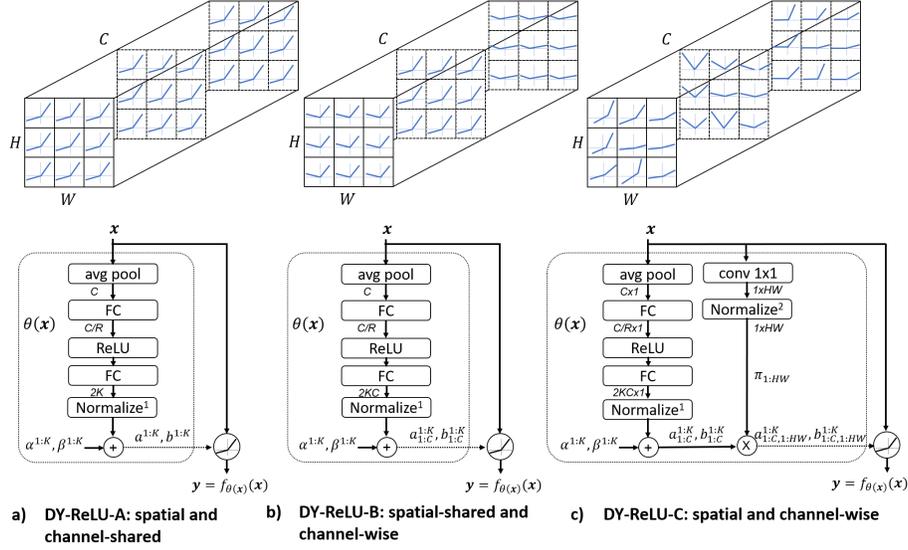
**Fig. 2.** Three DY-ReLU variations. They have different ways of sharing activation functions. The top row illustrates the piecewise linear function across spatial locations and channels, and the bottom row shows the network structure for the hyper function. Note that the first FC layer reduces the dimension by $R$, which is a hyper parameter.

**DY-ReLU-B (Spatial-shared and Channel-wise)**: the implementation details are introduced in section 3.2 and the network structure is shown in Fig. 2-(b). The hyper function outputs $2KC$ parameters ($2K$ per channel).

**DY-ReLU-C (Spatial and Channel-wise)**: as shown in Fig. 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 that has dimension $C \times H \times W$. This introduces an issue that the output dimension is too large ($2KCHW$), resulting in significantly more parameters in the fully connected layer. 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 a $1 \times 1$ convolution with a single output channel and a normalization that is a softmax function with upper cutoff as follows:

$$\pi_{h,w} = \min\{\frac{\gamma \exp(z_{h,w}/\tau)}{\sum_{h,w} \exp(z_{h,w}/\tau)}, 1\}, \tag{4}$$

where $z_{h,w}$ is the output of $1 \times 1$ convolution, $\tau$ is the temperature, and $\gamma$ is a scalar. The softmax is scaled up by $\gamma$ 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

|  | Top-1 |
|---|---|
| ReLU | $60.32 \pm 0.13$ |
| DY-ReLU-A | $63.28 \pm 0.12_{(2.96)}$ |
| DY-ReLU-B | $\mathbf{66.36 \pm 0.12_{(6.04)}}$ |
| DY-ReLU-C | $66.31 \pm 0.14_{(5.99)}$ |

**Table 2.** Comparing three DY-ReLU variations on Imagenet [5] classification. MobileNetV2 with width multiplier $\times 0.35$ is used. The mean and standard deviations of three runs are shown. 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 |
|---|---|---|
| ReLU | ReLU | $59.26 \pm 0.21$ |
| DY-ReLU-A | ReLU | $58.97 \pm 0.15_{(-0.29)}$ |
| DY-ReLU-B | ReLU | $61.76 \pm 0.27_{(+2.50)}$ |
| DY-ReLU-C | ReLU | $62.23 \pm 0.32_{(+2.97)}$ |
| ReLU | DY-ReLU-A | $57.12 \pm 0.25_{(-2.14)}$ |
| ReLU | DY-ReLU-B | $58.72 \pm 0.35_{(-0.54)}$ |
| ReLU | DY-ReLU-C | $61.03 \pm 0.11_{(+1.77)}$ |
| DY-ReLU-C | DY-ReLU-C | $\mathbf{63.27 \pm 0.15_{(+4.01)}}$ |

**Table 3.** Comparing three DY-ReLU variations on COCO [21] keypoint detection. We use MobileNetV2 $\times 0.5$ as backbone and use upsampling and inverted residual bottleneck blocks [3] in the head. The mean and standard deviations of three runs are shown. 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.

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. It includes four components: (a) average pooling, (b) the first FC layer (with ReLU), (c) the second FC layer (with normalization), and (d) piecewise function $f_{\boldsymbol{\theta}(\boldsymbol{x})}(\boldsymbol{x})$. For a feature map with dimension $C \times H \times W$, all three DY-ReLU variations share complexity for average pooling $O(CHW)$, the first FC layer $O(C^2/R)$ and piecewise function $O(CHW)$. The second FC layer has complexity $O(2KC/R)$ for DY-ReLU-A and $O(2KC^2/R)$ for DY-ReLU-B and DY-ReLU-C. Note that DY-ReLU-C spends additional $O(CHW)$ on computing spatial attentions. In most of the layers of MobileNet and ResNet, DY-ReLU has much less computation than a $1 \times 1$ convolution, which has complexity $O(C^2HW)$.

### 4.2    Ablations

Next, we study the three DY-ReLU variations on image classification and keypoint detection. Our goal is to understand their differences when performing 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 [21] classification is shown in Table 2. MobileNetV2 $\times 0.35$ is used. 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 3 shows the comparison on COCO keypoint detection. Similar to image classification, ***channel-wise variations (B and C) are better than channel-shared variation A in the backbone***. In contrast, ***the spatial-wise variation C is critical in the head***. Using DY-ReLU-C in both backbone and head achieves 4 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 the spatially-shared hyper function in DY-ReLU-A or DY-ReLU-B is difficult to learn when dealing with spatially sensitive task (e.g. distinguishes body joints in pixel level), especially in the head that has higher resolutions. This difficulty can be effectively alleviated by making hyper function spatial-wise, which encourages learning different activation functions at different positions. We observe that the training converges much faster when using spatial attention in the head network.

Base upon these ablations, we use DY-ReLU-B for ImageNet classification and use DY-ReLU-C for COCO 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 classes, including 1,281,167 images for training and 50,000 images for validation. We evaluate DY-ReLU on three CNN architectures (MobileNetV2 [30], MobileNetV3 [12] and ResNet [11]). We replace their default activation functions (ReLU in ResNet and MobileNetV2, ReLU/hswish/SE in MobileNetV3) with DY-ReLU. The main results are obtained by using spatial-shared and channel-wise DY-ReLU-B with two piecewise linear functions ($K = 2$). Note that MobileNetV2 and V3 have no activation after the last convolution layer in each block, where we add DY-ReLU with $K = 1$. 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 using SGD optimizer with 0.9 momentum for 300 epochs. The label smoothing 0.1 is used. We use weight decay 2e-5 and dropout 0.1 for width ×0.35, and increase weight decay 3e-5 and dropout 0.2 for width ×0.5, ×0.75, ×1.0. Random cropping/flipping and color jittering are used for all width multipliers. Mixup [43] is used for width ×1.0 to prevent overfitting.

**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

| Network | Activation | #Param | MAdds | Top-1 | Top-5 |
|---|---|---|---|---|---|
| MobileNetV2 ×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 ×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 ×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 ×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/SE/HS | 5.4M | 219.0M | 75.2 | 92.2 |
|  | DY-ReLU | 9.8M | 230.5M | $75.9_{(0.7)}$ | $92.7_{(0.5)}$ |
| MobileNetV3-Small | ReLU/SE/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 | 3.86G | 76.2 | 92.9 |
|  | DY-ReLU | 27.6M | 3.88G | $77.2_{(1.0)}$ | $93.4_{(0.5)}$ |
| ResNet-34 | ReLU | 21.3M | 3.64G | 73.3 | 91.4 |
|  | DY-ReLU | 24.5M | 3.65G | $74.4_{(1.1)}$ | $92.0_{(0.6)}$ |
| ResNet-18 | ReLU | 11.1M | 1.81G | 69.8 | 89.1 |
|  | DY-ReLU | 12.8M | 1.82G | $71.8_{(2.0)}$ | $90.6_{(1.5)}$ |
| ResNet-10 | ReLU | 5.2M | 0.89G | 63.0 | 84.7 |
|  | DY-ReLU | 6.3M | 0.90G | $66.3_{(3.3)}$ | $86.7_{(2.0)}$ |

**Table 4.** Comparing DY-ReLU with baseline activation functions (ReLU, SE or h-swish, denoted as HS) on ImageNet [5] classification in three network architectures. DY-ReLU-B with $K = 2$ linear functions is used. 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.

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 [43]) 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 4. Without bells and whistles, DY-ReLU outperforms its static counterpart by a clear margin 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 SE and h-swish (key contributions of MobileNetV3). The top-1 accuracy is improved by 2.3% and 0.7% for MobileNetV3-Small and MobileNetV3-Large, respectively. Note that DY-ReLU achieves more

| Activation | K | MobileNetV2 ×0.35 | | | MobileNetV2 ×1.0 | | |
|---|---|---|---|---|---|---|---|
| | | #Param | MAdds | Top-1 | #Param | MAdds | Top-1 |
| ReLU | 2 | 1.7M | 59.2M | 60.3 | 3.5M | 300.0M | 72.0 |
| RReLU [40] | 2 | 1.7M | 59.2M | $60.0_{(-0.3)}$ | 3.5M | 300.0M | $72.5_{(+0.5)}$ |
| LeakyReLU [25] | 2 | 1.7M | 59.2M | $60.9_{(+0.6)}$ | 3.5M | 300.0M | $72.7_{(+0.7)}$ |
| PReLU [10] | 2 | 1.7M | 59.2M | $63.1_{(+2.8)}$ | 3.5M | 300.0M | $73.3_{(+1.3)}$ |
| SE[14]+ReLU | 2 | 2.1M | 62.0M | $62.8_{(+2.5)}$ | 5.1M | 307.5M | $74.2_{(+2.2)}$ |
| Maxout [7] | 2 | 2.1M | 106.6M | $64.9_{(+4.6)}$ | 5.7M | 575.8M | $75.1_{(+3.1)}$ |
| Maxout [7] | 3 | 2.4M | 157.6M | $65.4_{(+5.1)}$ | 7.8M | 860.2M | $75.8_{(+3.8)}$ |
| DY-ReLU-B | 2 | 2.7M | 65.0M | $\mathbf{66.4}_{(+6.1)}$ | 7.5M | 315.5M | $\mathbf{76.2}_{(+4.2)}$ |
| DY-ReLU-B | 3 | 3.1M | 67.8M | $\mathbf{66.6}_{(+6.3)}$ | 9.2M | 322.8M | $\mathbf{76.2}_{(+4.2)}$ |

**Table 5.** Comparing DY-ReLU with related activation functions on ImageNet [5] classification. MobileNetV2 with width multiplier ×0.35 and ×1.0 are 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.

improvement for smaller models (e.g. MobileNetV2 ×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 5. Here we use MobileNetV2 (×0.35 and ×1.0), and replace ReLU with different activation functions in prior work. Our method outperforms all prior work with a clear margin, including Maxout that has significantly more computational cost. This demonstrates that DY-ReLU not only has more representation capability, but also is computationally efficient.

### 5.2   Inspecting DY-ReLU: Is It Dynamic?

We check if DY-ReLU is dynamic by examining its input and output over multiple images. Different activation values ($y$) across different images for a given input value (e.g. $x = 0.5$) is expected to differentiate from static ReLU, which has a fixed output (e.g. $y = 0.5$ when $x = 0.5$).

Fig. 3-(Top) 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]. Clearly, the learnt DY-ReLU is dynamic over features as activation values ($y$) vary in a range (that blue dots cover) for a given input $x$. The dynamic range varies across different blocks, indicating different dynamic functions learnt across levels. We also observe many positive activations for negative inputs. Statistically, 51% of DY-ReLU have segments with either negative slope or slope above 1, and 37% of DY-ReLU have at least one segment with intercept more than 0.05. These cases cannot be handled by ReLU, SE or MaxOut of two SEs.

We also analyzed the angle between two segments in DY-ReLU (i.e. slope difference $|a_c^1 - a_c^2|$). The slope difference decreases from lower to higher levels
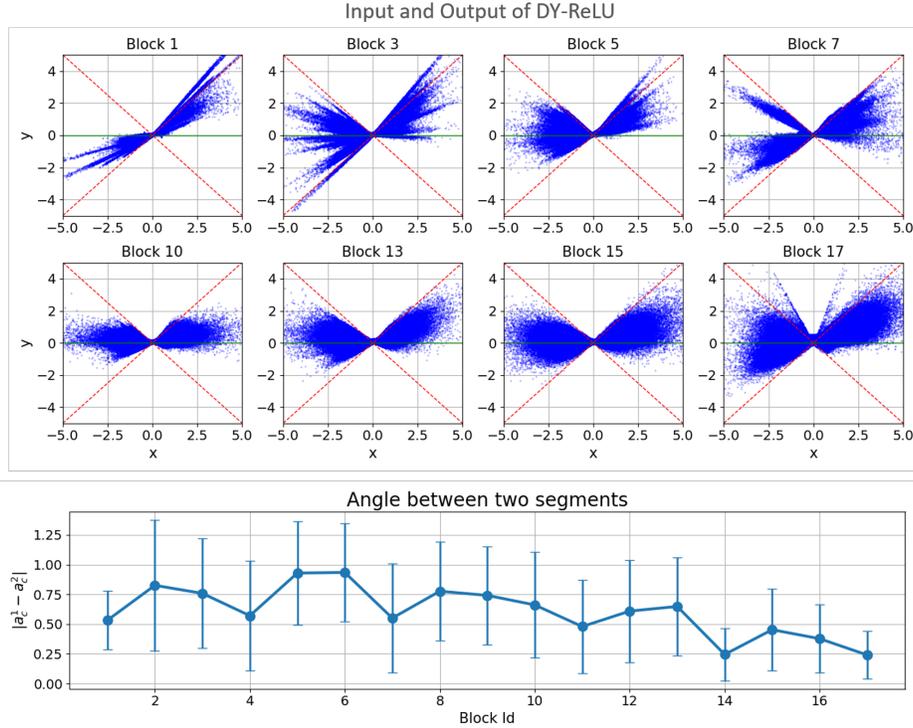
**Fig. 3. Top:** plots of input and output values of DY-ReLU in a well trained model (using MobileNetV2 ×0.35) over 50,000 validation images in ImageNet [5]. We choose the dynamic ReLU after the depthwise convolution in every other mobile block. Block 1 corresponds to the lowest block, and Block 17 corresponds to the highest block. The two red lines correspond to $y = x$ and $y = -x$, respectively. **Bottom:** Angle (or slope difference $|a_c^1 - a_c^2|$) between two segments in DY-ReLU across blocks. The bending of the activation functions decreases from low levels to high levels. Best viewed in color.

(shown in Fig. 3-(Bottom)). This indicates that the activation functions tend to have lower bending in higher levels.

### 5.3    Ablation Studies on ImageNet

We run a number of ablations to analyze DY-ReLU. We focus on spatial-shared and channel-wise DY-ReLU-B, and use MobileNetV2 ×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 $\lambda_b = 0.5$, respectively. The reduction ratio of the first FC layer in the hyper function is set as $R = 8$.

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

**Table 6. Different dynamic piecewise functions** evaluated on ImageNet classification. MobileNetV2 ×0.35 is used.

| | $A_1$ | $A_2$ | $A_3$ | 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 |

**Table 7. DY-ReLU at different layers** evaluated on ImageNet. MobileNetV2 ×0.35 is used. $A_1, A_2, A_3$ indicate activations after three convolution layers in an inverted residual block.

| | $R$ | #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. Different reduction ratios $R$** for the first fully connected layer in the hyper function (see Fig. 2). Evaluation is on ImageNet classification. MobileNetV2 ×0.35 is used. Setting $R = 8$ achieves a good trade-off.

**Dynamic Piecewise Functions:** Table 6 shows the classification accuracy using different piecewise functions. The major gain is due to making ReLU dynamic. Specifically, making the first segment dynamic boosts top-1 accuracy from 60.3% to 63.8%. Making the second segment dynamic gains additional 1.9%. The intercept $b_c^k$ is helpful consistently. The gap between $K = 2$ and 3 is small. In most of DY-ReLU with $K = 3$ segments, 2 of the 3 segments have similar slopes.

**Dynamic ReLU at Different Layers:** Table 7 shows the classification accuracy for using DY-ReLU at three different layers (after $1 \times 1$ conv, $3 \times 3$ depthwise conv, $1 \times 1$ conv) in an inverted residual block in MobileNetV2 ×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 $3 \times 3$ depth-wise convolution yields the best performance.

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

| $\alpha^1$ | $\alpha^2$ | 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 |

(a) **Initialization of $\alpha^k$.**

| $\beta^1$ | $\beta^2$ | 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 |

(b) **Initialization of $\beta^k$.**

| $\lambda_a$ | 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 |

(c) **Range of slope $\lambda_a$.**

**Table 9.** Ablations of three hyper parameters in DY-ReLU on Imagenet classification.

**Initialization of Slope** ($\alpha^k$ in Eq (3)): As shown in Table 9-(a), 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** ($\beta^k$ in Eq (3)): the performance is stable (shown in Table 9-(b)) 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** ($\lambda_a$ in Eq (3)): Making slope range either too wide or too narrow is not optimal, as shown in Table 9-(c). A good choice is to keep $\lambda_a$ between 1 and 2.

### 5.4   COCO Single-Person Keypoint Detection

We use COCO 2017 dataset [21] 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 [31]. All models are trained from scratch for 210 epochs, using Adam optimizer [18]. The initial learning rate is set as 1e-3 and is dropped to 1e-4 and 1e-5 at the $170^{th}$ and $200^{th}$ epoch,

| Backbone | Activation | Param | MAdds | AP | $AP^{0.5}$ | $AP^{0.75}$ | $AP^M$ | $AP^L$ |
|---|---|---|---|---|---|---|---|---|
| MBNetV2 ×1.0 | ReLU | 3.4M | 993.7M | 64.6 | 87.0 | 72.4 | 61.3 | 71.0 |
| | DY-ReLU | 9.0M | 1026.9M | $68.1_{(3.5)}$ | 88.5 | 76.2 | 64.8 | 74.3 |
| MBNetV2 ×0.5 | ReLU | 1.9M | 794.8M | 59.3 | 84.3 | 66.4 | 56.2 | 65.0 |
| | DY-ReLU | 4.6M | 820.3M | $63.3_{(4.0)}$ | 86.3 | 71.4 | 60.3 | 69.2 |
| MBNetV3 Large | ReLU/SE/HS | 4.1M | 896.4M | 65.7 | 87.4 | 74.1 | 62.3 | 72.2 |
| | DY-ReLU | 10.1M | 926.6M | $67.2_{(1.5)}$ | 88.2 | 75.4 | 64.1 | 73.2 |
| MBNetV3 Small | ReLU/SE/HS | 2.1M | 726.9M | 57.1 | 83.8 | 63.7 | 55.0 | 62.2 |
| | DY-ReLU | 4.8M | 747.9M | $60.7_{(3.6)}$ | 85.7 | 68.1 | 58.1 | 66.3 |

**Table 10.** Comparing DY-ReLU with baseline activation functions (ReLU, SE or h-swish, denoted as HS) on COCO Keypoint detection. The evaluation is on validation set. The head structure in [3] is used. DY-ReLU-C with $K = 2$ is used in both backbone and head. 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 static counterpart by a clear margin.

respectively. All human detection boxes are cropped from the image and resized to $256 \times 192$. The data augmentation includes random rotation ($[-45°, 45°]$), random scale ($[0.65, 1.35]$), flipping, and half body data augmentation.
**Testing:** We use the person detectors provided by [38] and follow the evaluation procedure in [38,31]. 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 ×0.5 and ×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 4.0 AP when using MobileNetV2 with width multiplier ×1.0 and ×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.

# 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=HylVB3AqYm`
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.J., Warde-Farley, D., Mirza, M., Courville, A., Bengio, Y.: Maxout networks. arXiv preprint arXiv:1302.4389 (2013)
8. Ha, D., Dai, A.M., Le, Q.V.: Hypernetworks. ICLR (2017)
9. 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)
10. He, K., Zhang, X., Ren, S., Sun, J.: Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In: ICCV (2015)
11. 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)
12. 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`
13. 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)
14. Hu, J., Shen, L., Sun, G.: Squeeze-and-excitation networks. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2018)
15. 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=Hk2aImxAb`
16. 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`
17. Jarrett, K., Kavukcuoglu, K., Ranzato, M., LeCun, Y.: What is the best multistage architecture for object recognition? In: The IEEE International Conference on Computer Vision (ICCV) (2009)
18. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. In: International Conference on Learning Representations (ICLR) (2015)
19. Klambauer, G., Unterthiner, T., Mayr, A., Hochreiter, S.: Self-normalizing neural networks. In: Advances in neural information processing systems. pp. 971–980 (2017)

20. 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`
21. 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)
22. Liu, H., Simonyan, K., Yang, Y.: DARTS: Differentiable architecture search. In: International Conference on Learning Representations (2019), `https://openreview.net/forum?id=S1eYHoC5FX`
23. Liu, L., Deng, J.: Dynamic deep neural networks: Optimizing accuracy-efficiency trade-offs by selective execution. In: AAAI Conference on Artificial Intelligence (AAAI) (2018)
24. 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)
25. 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)
26. Misra, D.: Mish: A self regularized non-monotonic neural activation function. arXiv preprint arXiv:1908.08681 (2019)
27. Nair, V., Hinton, G.E.: Rectified linear units improve restricted boltzmann machines. In: ICML (2010)
28. Ramachandran, P., Zoph, B., Le, Q.V.: Searching for activation functions. arXiv preprint arXiv:1710.05941 (2017)
29. Real, E., Aggarwal, A., Huang, Y., Le, Q.V.: Regularized evolution for image classifier architecture search. In: AAAI Conference on Artificial Intelligence (AAAI) (2018)
30. 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)
31. Sun, K., Xiao, B., Liu, D., Wang, J.: Deep high-resolution representation learning for human pose estimation. In: CVPR (2019)
32. 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)
33. Trottier, 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)
34. 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)
35. 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)
36. 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)
37. 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)

38. Xiao, B., Wu, H., Wei, Y.: Simple baselines for human pose estimation and tracking. In: European conference on computer vision (04 2018)
39. 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`
40. Xu, B., Wang, N., Chen, T., Li, M.: Empirical evaluation of rectified activations in convolutional network. CoRR (2015)
41. Yang, B., Bender, G., Le, Q.V., Ngiam, J.: Condconv: Conditionally parameterized convolutions for efficient inference. In: NeurIPS (2019)
42. 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`
43. 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`
44. 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)
45. Zoph, B., Le, Q.V.: Neural architecture search with reinforcement learning. CoRR **abs/1611.01578** (2017)
46. 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)