

AFINet: Attentive Feature Integration Networks for Image Classification

Xinglin Pan^a, Jing Xu^b, Yu Pan^b, liangjian Wen^a, WenXiang Lin^c, Kun Bai^d, Zenglin Xu^{e,b,*}

^aUniversity of Electronic Science and Technology of China, Chengdu, China

^bSchool of Science and Technology, Harbin Institute of Technology, Shenzhen, China

^cBeijing Institute of Technology, Beijing, China

^dCloud and Smart Industries Group, Tencent, China

^eArtificial Intelligence Center, Peng Cheng Lab, Shenzhen, China

Abstract

Convolutional Neural Networks (CNNs) have achieved tremendous success in a number of learning tasks including image classification. Recent advanced models in CNNs, such as ResNets, mainly focus on the skip connection to avoid gradient vanishing. DenseNet designs suggest creating additional bypasses to transfer features as an alternative strategy in network design. In this paper, we design Attentive Feature Integration (AFI) modules, which are widely applicable to most recent network architectures, leading to new architectures named AFI-Nets. AFI-Nets explicitly model the correlations among different levels of features and selectively transfer features with a little overhead. AFI-ResNet-152 obtains a 1.24% relative improvement on the ImageNet dataset while decreases the FLOPs by about 10% and the number of parameters by about 9.2% compared to ResNet-152.

Keywords: CNN, Attention, Image Classification, Feature Integration

2010 MSC: 68T45

1. Introduction

Convolutional Neural Networks(CNNs) have achieved remarkable successes in various computer vision tasks, e.g., image classification, semantic segmentation, object

*Corresponding author

Email address: xuzenglin@hit.edu.cn (Zenglin Xu)

detection [1, 2, 3, 4], etc. A major drive to such successes is from the evolutionary design of network architectures, and such representative examples include AlexNet [5], VGG-Net [6], GoogleNet [7], ResNet [8], and DenseNet [9]. Notably, the skip connection introduced in ResNet has become a fundamental design strategy as an effective solution to the gradient vanishing problem, especially for very deep networks. And this strategy has been widely adopted in a variety of architectures, including ShuffleNet [10], ResNeXt [11], and Inception-ResNet [12], etc. Besides, as an alternative design strategy, DenseNet [9] suggests constructing extra bypasses to transfer previous layers of features for future reuse.

Constructing extra bypasses leads to many advantages in architecture designs, e.g., implicit deep supervision and diversified depths [9, 13]. Furthermore, preserved features, especially low-level features, are beneficial to overcome overfitting (which often leads to small training errors but large testing errors). As evidenced by a study of the class selectivity indices [14], the generalization gap (i.e., the difference between the training error and the testing error) of features is prone to increase with depths. Hence, high-level features are more useful to reduce the training error while low-level ones are in favor of closing the generalization gap. These observations suggest that aggregation of features at multiple levels is vital to the designs of architectures.

Despite the advantages of the current bypass design for classification, there are several issues to be addressed: (1) Some low-level features (e.g., edges of the background) may be irrelevant to classification and thus impair the learning performance; (2) The concatenation of all previous features involves a quadratic scale of memory usage which is inhibited in scenarios with limited computational and storage resources; (3) The correlation between features crossing layers is hard to model for the convolution operators.

To address these issues, we propose a lightweight and selective feature integration scheme for most residual-like networks, leading to the Attentive Feature Integration (AFI) module, as illustrated in Figure 1. Firstly, each of the input feature maps (i.e., raw features as shown in Figure 2) with C channels, is squeezed into a vector by a squeeze operation that captures information from a large spatial extent. Thus, the global context is embedded in the vectors. Secondly, the vectors are scored and nor-

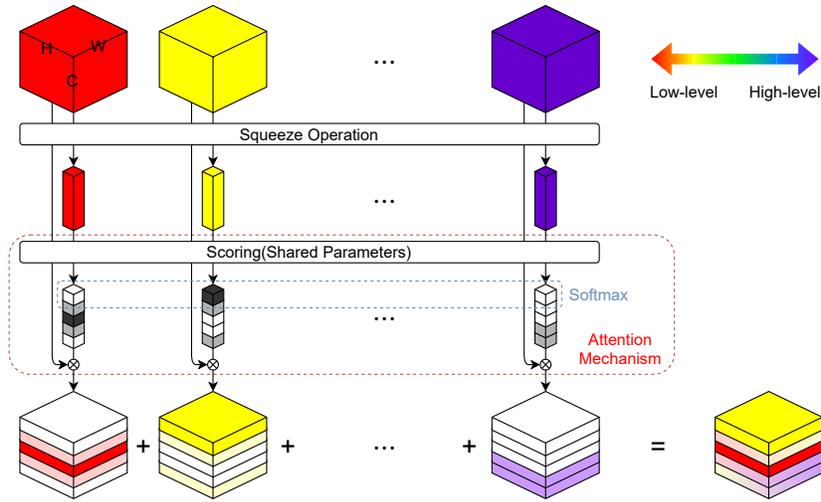


Figure 1: The AFI module can automatically extract important low-level features for the high-level features. Through two light-weight operations, i.e., the squeeze operation and the shared attention mechanism, every feature is re-calibrated along each channel dimension. Finally, the resulting feature is supplied to later layers.

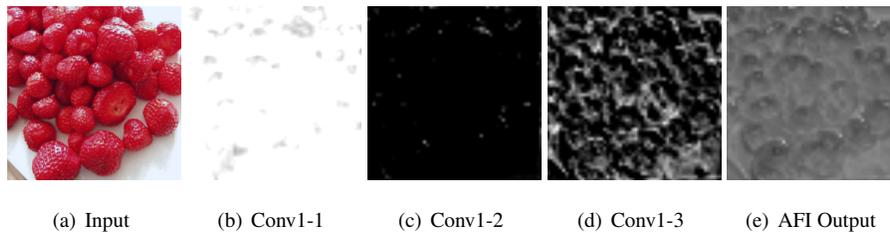


Figure 2: The subfigure (a) shows an image from the Imagenet dataset. The subfigure (b), (c) and (d) show features from conv1 stage in AFI-ResNet-50. All of these extract edge information. The subfigure (e) depicts the output of an AFI module whose inputs are the subfigure (b), (c), and (d). The subfigure (e) extracts the main edge information of the three subfigures.

35 malized sequentially via the attention mechanism (which consists of the shared scoring function and the channel-by-channel softmax function) in order to re-calibrate features. At last, we obtain the resulting feature via a summation of re-calibrated features, where each channel can be viewed as a convex combination over the raw features. In short, AFI modules rate the importance of features adaptively and have access to model correlations between distant layers.

40 AFI modules are not limited to a special backbone network. Instead, they can be

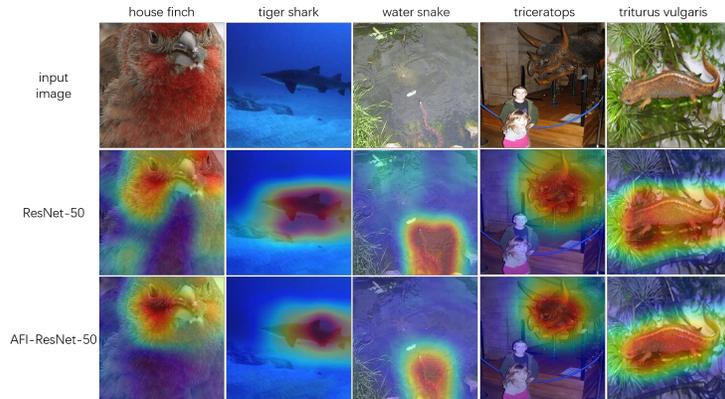


Figure 3: Illustration of the impact of the AFI module. All original images come from the ImageNet dataset. Heatmaps generated by Grad-CAM [19] illustrates which areas the network pays more attention to. Compared to the vanilla ResNet-50, the area that AFI-ResNet-50 paid most attention to is much smaller.

easily plugged to various backbones, leading to various architectures, namely AFI-Nets. To interpret the reason for no limit to a special backbone network, different backbone networks are uniformly viewed as instances of ordinary differential equations (ODEs) [15, 16, 17]. The application of AFI modules is very similar to the linear multistep method [18] apply to solve ODEs. (see Section 3.3 for more about this.)

In addition to the easy plug-and-go property, AFI modules also enjoy many advantages such as efficient utilizing of low-level features, and having lower overheads and higher accuracy. To illustrate the efficiency, we perform a series of indirect experiments. We firstly visual class activation mapping (CAM) as depicted in Figure 3. The activation area of AFI-ResNet-50 is much precise than ResNet-50. Then, the average of the class selectivity index [14] of features at different layers is lower than backbone networks, which argues that features AFI-Nets learnt are more common. At last, we obtain a 10.8% average relative improvement on various extremely easy to overfit datasets. The great performance of these experiments is usually attributed to the efficient utilization of low-level features.

AFI-Nets are trained on ImageNet with SGDM optimizers and regular data augmentations. As shown in Figure 4, AFI-ResNets are applicable with different depths. With the same number of layers, AFI-ResNet has fewer parameters, fewer FLOPs, and

60 more accuracy. Among AFI-ResNets, AFI-ResNet-152 increases the Top-1 accuracy rate in ImageNet by 1.24% compared to ResNet-152 while decreases the FLOPs by about 10% and the number of parameters by about 9.2%.

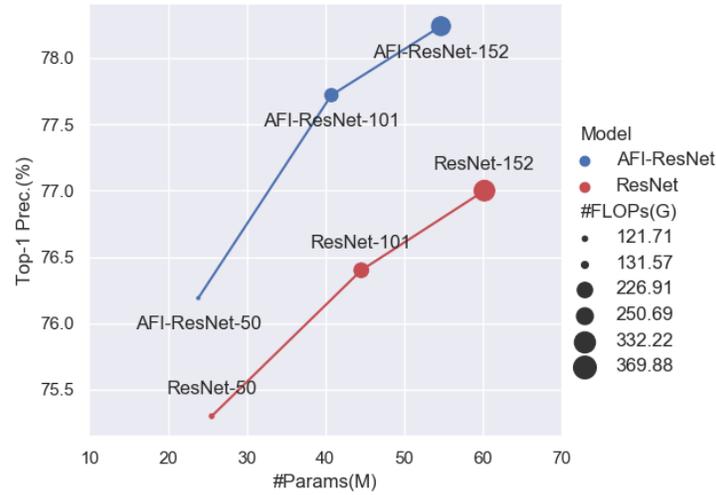


Figure 4: Illustration of the experiment results of AFI-ResNets. With the same number of layers, AFI-ResNets have fewer parameters, fewer FLOPs, and more accuracy. The FLOPs are calculated by assuming the batch size of 32.

Our major contributions can be summarized as follows:

- We propose lightweight AFI modules to selectively transfer features. From the perspective and inspiration of ODEs, AFI-Nets can be derived by applying AFI modules into residual-like networks.
- To avoid the difficulty of directly evaluating the efficiency in processing low-level features, we perform three indirect analytical experiments to verify AFI modules take more advantage of preserved low-level features.
- Experimental results show that our AFI-Nets significantly improve the representational power of the network.

2. Related Work

Modern CNN Architectures. Deep convolutional neural networks(CNNs) have dominated image classification since the AlexNet [5] and VGG-Net [6] are proposed. After that, substantial efforts have been made to improve the efficiency of CNNs. The modular design strategy in GoogleNet [7] simplifies the network architecture and the multi-path structure in each module shows a great success. ResNet introduces the short-cut connection alleviating the difficulty in deep network training [8]. DenseNet [9] densely connects all preceding layers to take full advantage of preceding feature maps. Based on those fundamental architectures, some advanced variants(e.g., ResNeSt [20]) have been proposed and have achieved impressive performance in many computer vision tasks.

Exploitation of Low-Level Features. The motivation to preserve low-level features is ample. As shown in Yosinski et al. [21], some shallow layer filters obtain similar effects as Gabor filters and color blobs. Such low-level features appear to be un-specific to particular datasets or tasks, and generalize well to other datasets or tasks. To maximize feature utilization, DenseNet [9] demonstrates that dense concatenating all the features in frontier layers can effectively alleviate the difficulty of training and improves the network performance with an increased calculation overhead. VoVNet [13], VoVNetV2 [22] overcome the inefficiency of dense connection by concatenating all features only once in the last feature map and achieve the state-of-the-art performance in instance segmentation. Dual Path Network [23] and Mixed link Network [24] try to transfer low-level feature based on ResNet [8]. However, few works focus on constructing bypass to transfer low-level features based on other residual-like networks.

Attention Mechanisms. The benefits of attention mechanism have been demonstrated across a range of tasks. Squeeze-and-Excitation block [25] highly appreciates attention mechanism and thus well improves the accuracy of varied CNNs. They use global average-pooled features to exploit the inter-channel relationship and to compute the channel-wise attention. Besides, there are several other researches to utilize the attention mechanism and improve the results of CNNs in various vision tasks. CBAM [26] further adds the spatial attention to the SE module and results in better

plug-and-play modules. Soft mask branches to refine the feature maps by adding attention knowledge proposed by [3]. Non-local Neural Networks [27] proposes non-local module to integrate the global attention information. Libra R-CNN [28] designs the balanced feature pyramid which refines the semantic feature from multi-level features. BASNet [29] pays more attention to the boundary of the mask by the boundary-aware loss function. However, few literature focuses on the mix of attention mechanism and transferring low-level features. The proposed AFI module thus aims to improve transferring based on attention mechanisms.

3. Our Model

In this section, we first introduce Attentive Feature Integration (AFI) modules and compare them with previous works. Next, we describe the ways of implementing AFI modules for residual-like networks with two examples. Specifically, we integrate AFI modules with backbone architectures, such as ResNet [8] and MobileNetV2 [30], to build new network architectures, i.e., AFI-ResNet and AFI-MobileNetV2, respectively. At last, inspired by the view of residual-like networks as instances of ordinary differential equations (ODEs), we point out that our modules is similar to Linear Multistep Method (LMM) in both motivation and formal.

3.1. Attentive Feature Integration Modules

AFI module aims to model the correlation among features at different levels and then produce a compact but comprehensive feature as an output. AFI module is composed of two operations: a squeeze operation and an attention mechanism. For convenient description, we denote N different-level features as $\mathbf{X}^{(i)} \in \mathbb{R}^{H \times W \times C}, i \in \{1, 2, \dots, N\}$.

Squeeze Operation. $F_{sq}(\cdot)$ represents the squeeze function(e.g., global average pooling), which gathers contextual long-range feature interactions, embedding global context into a vector descriptor. By shrinking $\mathbf{X}^{(i)}$ on its spatial dimensions $H \times W$, the channel-wise statistic $\mathbf{z}^{(i)} \in \mathbb{R}^C$ is generated, where C is the number of channels.

Attention Mechanism. The attention mechanism is a selective aggregation of information so that the resulting features are easily exploited for specific tasks. There are

two steps in the attention mechanism: score and normalization. The score step is based on a shared scoring function $F_{\text{sc}}(\cdot)$, which is applied into vector descriptors to produce an embedding of importance $\mathbf{s}^{(i)} \in \mathbb{R}^C$. $F_{\text{sc}}(\cdot)$ is formulated by two transformation matrices around the activation function

$$\mathbf{s}^{(i)} = F_{\text{sc}}\left(\mathbf{z}^{(i)}, \mathbf{W}\right) = \mathbf{W}_2 \text{ReLU}\left(\mathbf{W}_1 \mathbf{z}^{(i)}\right), \quad (1)$$

where the vector $\mathbf{s}^{(i)}$ is parameterized by forming a bottleneck layer with a weight $\mathbf{W}_1 \in \mathbb{R}^{\frac{C}{r} \times C}$, a dimensionality-increasing layer with a weight $\mathbf{W}_2 \in \mathbb{R}^{C \times \frac{C}{r}}$ and the reduction ratio r . Next, a matrix $\mathbf{S} = [\mathbf{s}^{(1)}, \mathbf{s}^{(2)}, \dots, \mathbf{s}^{(N)}] \in \mathbb{R}^{N \times C}$, is obtained by concatenating $\{\mathbf{s}^{(i)}\}_{i=1}^N$. Then, a matrix $\tilde{\mathbf{S}} \in \mathbb{R}^{N \times C}$ is formed by normalizing \mathbf{S} along the N dimension using a Softmax function.

In summary, AFI : $\{\mathbf{X}^{(i)}\}_{i=1}^N \rightarrow \mathbf{R}$ is formulated by

$$\begin{aligned} \mathbf{z}^{(i)} &= F_{\text{sq}}\left(\mathbf{X}^{(i)}\right), \quad \mathbf{s}^{(i)} = F_{\text{sc}}\left(\mathbf{z}^{(i)}, \mathbf{W}\right), \quad \mathbf{S} = [\mathbf{s}^{(1)}, \mathbf{s}^{(2)}, \dots, \mathbf{s}^{(N)}], \\ \tilde{\mathbf{S}}_{i,j} &= \frac{\exp(\mathbf{S}_{i,j})}{\sum_{k=1}^N \exp(\mathbf{S}_{k,j})}, \quad \mathbf{R}_{i,j,k} = \sum_{l=1}^N \tilde{\mathbf{S}}_{l,k} \mathbf{X}_{i,j,k}^{(l)}. \end{aligned} \quad (2)$$

where the resulting feature $\mathbf{R} \in \mathbb{R}^{H \times W \times C}$.

Remark on the connection to previous works. Recently, a series of literature attempts to incorporate the attention mechanism to improve the performance of CNNs. One of the most popular computational units is the Squeeze-and-Excitation(SE) module [25]. Compared to SE module, our AFI module focus on explicitly modeling the feature integration instead of channel-wise selection. Moreover, compared to DenseNet with self-attention modules like CAPR-DenseNet [31], our module can be applied to deeper and larger networks, while avoiding quadratic complexity of memory usage and running time by substituting the shared attention mechanism for the independent exciter. Compared to the SKNet [32] that selects the efficient kernel size, our AFI module can utilize different level(e.g., positional and semantic) information. Compare to the Mixed Link Network [24], our module can be applied to not only ResNet but also other residual-like networks.

3.2. AFI Modules for Residual-Like Networks

This subsection clarifies how to apply AFI module to residual-like networks and
 150 derive our AFI-Nets. In order to describe our method more clearly, we follow the definition of *stage* and *building blocks* in Hu et al. [33]. More specifically, a *stage* consists of several *building blocks* with features of the same shape stacking sequentially.

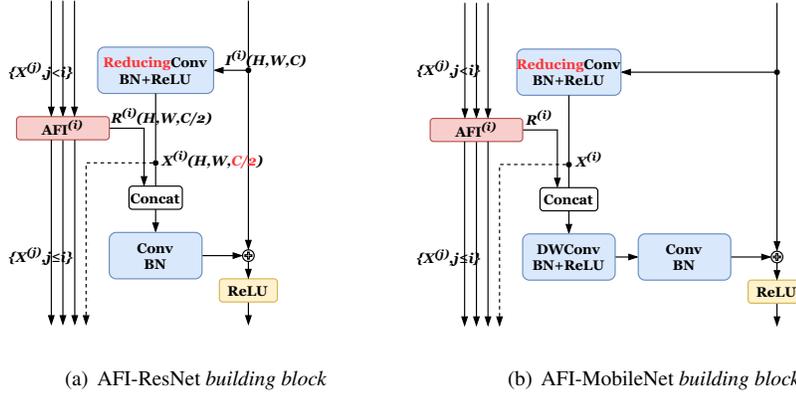


Figure 5: The architecture of AFI-MobileNet *building block* is shown in Figure 5(b). The architecture of AFI-ResNet *building block* is shown in Figure 5(a). The blue and yellow parts indicate the original backbone while the red parts indicate our revisions. The position denoted by $I^{(i)}$ is the input of the i -th *building block* of an original ResNet.

As shown in Figure 5, we half the output of the first convolution in i -th building blocks and denote it as $X^{(i)}$. For the i -th *building block* of each *stage*, $AFI^{(i)}$ denotes the AFI module corresponding to that *building block*. $\{X^{(1)}, X^{(2)}, \dots, X^{(i-1)}\}$ in
 155 previous *building blocks* are the inputs of $AFI^{(i)}$. We concatenate the output $R^{(i)}$ of $AFI^{(i)}$ with $X^{(i)}$ and then throw the concatenated result to the subsequent convolution layer. We name our networks as AFI-Nets.

Here, we would like to describe applications of our module in residual-like networks with an instance AFI-ResNet. AFI-ResNet *building block* is shown in Figure 5(a). We first review the architecture of ResNet *building block* [8]. For the i -th building block of ResNet

$$I^{(i+1)} = \text{ReLU} \left(I^{(i)} + \text{BN} \left(\text{ReLU} \left(\text{BN} \left(I^{(i)} * K_1^{(i)} \right) * K_2^{(i)} \right) \right) \right), \quad (3)$$

where $*$ denotes convolution operator, $\text{BN}(\cdot)$ denotes batch normalization, $\text{ReLU}(\cdot)$ is a rectified linear activation function. The input feature $\mathbf{I}^{(i)}$ is shown in Figure 5(a) and the kernels $\mathbf{K}_1^{(i)}, \mathbf{K}_2^{(i)} \in \mathbb{R}^{k \times k \times C \times C}$ are learnable filters where the kernel size $k = 3$ and the input channel and the output channel $C \in \{16, 32, 64\}$.

In AFI-ResNet, we replace $\mathbf{K}_1^{(i)}$ with a shrinking kernel $\widetilde{\mathbf{K}}_1^{(i)} \in \mathbb{R}^{k \times k \times C \times \frac{C}{2}}$ to reduce parameters. For the i -th *building block* of AFI-ResNet

$$\begin{aligned} \mathbf{X}^{(i)} &= \text{ReLU} \left(\text{BN} \left(\mathbf{I}^{(i)} * \widetilde{\mathbf{K}}_1^{(i)} \right) \right), \quad \mathbf{R}^{(i)} = \text{AFI}^{(i)} \left(\mathbf{X}^{(1)}, \mathbf{X}^{(2)}, \dots, \mathbf{X}^{(i-1)} \right) \\ \mathbf{I}^{(i+1)} &= \text{ReLU} \left(\mathbf{I}^{(i)} + \text{BN} \left([\mathbf{X}^{(i)}, \mathbf{R}^{(i)}] * \mathbf{K}_2^{(i)} \right) \right). \end{aligned} \quad (4)$$

Besides, we also provide an instance AFI-MobileNetV2 as shown in 5(b). We concatenate the output of the AFI module with the output feature of the first convolutional layer. As many building blocks (e.g., bottleneck *building blocks* [8]) use an 1×1 convolution at the front of building blocks to reduce channels, our AFI modules can neutralize the impairment of bottleneck convolution. It is elaborated in Section 4.1 that the parameters and FLOPs of AFI-Nets are usually smaller than the vanilla.

3.3. The Versatility of AFI modules

In this section, we will explain the versatility of AFI modules from the view of numerical analysis. Residual-like networks can be identified as instances of ordinary differential equations (ODEs), behaving like the forward Euler method with an initial value $y^{(1)}$ [15, 16, 17]. We firstly review the forward Euler method as background. Formally,

$$\begin{aligned} y^{(i)} - y^{(i-1)} &= \left(x^{(i)} - x^{(i-1)} \right) \frac{dy}{dx} + \left(x^{(i)} - x^{(i-1)} \right)^2 \frac{d^2y}{dx^2} + \dots \\ y^{(i)} - y^{(i-1)} &\approx \left(x^{(i)} - x^{(i-1)} \right) \cdot f \left(x^{(i)}, y^{(i)} \right), \end{aligned} \quad (5)$$

where $i \in \{2, 3, \dots, N\}$ and $f(x, y)$ denotes $\frac{dy}{dx}$. When satisfying $x^{(i)} - x^{(i-1)} = 1$, Equation (5) can be simplified as

$$y^{(i)} = y^{(i-1)} + f \left(x^{(i)}, y^{(i)} \right). \quad (6)$$

Euler method is a traditional first-order solution to the ODE. Therefore, this method can cause huge error in predicting the next value when ignoring higher order $\left\{ \frac{d^i y}{dx^i} \right\}_{i=2}^{\infty}$.

Linear Multi-step Method (LMM) [18] reuses the information in the previous steps by linear combination to fit them:

$$y^{(i)} = y^{(i-1)} + \sum_{j=1}^i \alpha_j^{(i)} f(x^{(j)}, y^{(j)}), \quad (7)$$

170 where $\{\alpha_j^{(i)}\}_{j=1}^i$ are coefficients of merging previous steps.

In residual-like networks, lower-level features play similar roles to the previous steps, thus it is also necessary to take these features into consideration. Thus, based on the consideration, we design AFI-Nets, which is the generalization of LMM. Specifically, we first replace the scalars in Equation (7) with the features in neural networks, namely

$$\mathbf{I}^{(i+1)} = \mathbf{I}^{(i)} + g^{(i)} \left(\sum_{j=1}^i \alpha_j^{(i)} f(\mathbf{I}^{(j)}) \right), \quad (8)$$

where $f(\mathbf{I}^{(j)})$ is a learnable function that acts similar to estimating $\frac{dy}{dx} \Big|_{x=x^{(j)}}$, and $g^{(i)}$ is a function to match the shape of $\mathbf{I}^{(i)}$ and the estimated residual. For instance, in AFI-ResNet, $f(\mathbf{I}^{(j)})$ is taken to be $\text{ReLU} \left(\text{BN} \left(\mathbf{I}^{(j)} * \widetilde{\mathbf{K}}_1^{(j)} \right) \right)$ with a kernel $\widetilde{\mathbf{K}}_1^{(j)}$, $g^{(i)}(\cdot)$ is taken to be $\text{BN}(\cdot * \hat{\mathbf{K}}_2^{(i)})$ with a kernel $\hat{\mathbf{K}}_2^{(i)}$, $\{\alpha_j^{(i)}, j < i\}$ is selected by AFI⁽ⁱ⁾, and $\alpha_i^{(i)} = 1$. Equation (8) can be simplified as:

$$\begin{aligned} & g^{(i)} \left(\sum_{j=1}^i \alpha_j^{(i)} f(\mathbf{I}^{(j)}) \right) \\ & g^{(i)} \left(\sum_{j=1}^{i-1} \alpha_j^{(i)} \mathbf{X}^{(j)} + \mathbf{X}^{(i)} \right) \\ & = \text{BN} \left(\left(\mathbf{R}^{(i)} + \mathbf{X}^{(i)} \right) * \hat{\mathbf{K}}_2^{(i)} \right) \\ & = \text{BN} \left(\mathbf{X}^{(i)} * \hat{\mathbf{K}}_2^{(i)} + \mathbf{R}^{(i)} * \hat{\mathbf{K}}_2^{(i)} \right). \end{aligned} \quad (9)$$

Compared with $\mathbf{K}_2^{(i)}$ in Equation (4), $\hat{\mathbf{K}}_2^{(i)}$ in Equation (9) is a special case of $\mathbf{K}_2^{(i)}$ when $\mathbf{X}^{(i)}$ and $\mathbf{R}^{(i)}$ is applied to convolution with the same kernel. Other AFI-Nets can be expressed by adopting different f and g .

4. Experiment

175 We firstly analyse overheads of our modules. Then, before presenting results on real-world datasets, we first illustrate the role of our AFI module by several experiments. If not otherwise specified, the reduction ratio r is 4. We show the ablation study to better understand the settings of the AFI modules in last.

4.1. Overheads of AFI modules

Table 1: The architecture details of AFI-ResNet-(6N+2) for CIFAR dataset. The operations and feature shapes are listed inside the brackets and the number of stacked blocks is shown outside. num_class depends on the categories.

Name	Output Size	AFI-ResNet-(6N+2)
Conv ₀	32×32	3×3, 16
Conv ₁	32×32	AFI + $\begin{matrix} 3 \times 3, 8 \\ 3 \times 3, 16 \end{matrix}$ × N
Conv ₂	16×16	AFI + $\begin{matrix} 3 \times 3, 16 \\ 3 \times 3, 32 \end{matrix}$ × N
Conv ₃	8×8	AFI + $\begin{matrix} 3 \times 3, 32 \\ 3 \times 3, 64 \end{matrix}$ × N
Output	num_class	AP, FC, Softmax

Additional AFI modules may raise a problem whether parameters and FLOPs are increased. To clarify it, we compare AFI-Nets with the vanilla about parameters and FLOPs. Table 1 shows AFI-ResNet architecture for instance. By setting $N = 5$ and $N = 18$ separately, AFI-ResNet-32 and AFI-ResNet-110 are acquired. In a $(N + 1)$ -block *stage*, the number of additional parameters is $2 \times (N - 1) \times C \times \frac{C}{r}$ and the times of extra multiplication in the scoring function is given by:

$$\sum_{i=3}^{N+1} 2 \times (i - 1) \times \frac{C}{r} \times C = \frac{C^2}{r} \times (N + 2)(N - 1). \quad (10)$$

180 The reason for enumerating i starting from 3 is that, there is no input to AFI⁽¹⁾ and there is only one input to AFI⁽²⁾. The input feature can skip AFI⁽²⁾ and is transferred immediately.

Table 2: Comparison between our AFI module and convolution to highlight the smartness of our block. No matter in which blocks, with AFI module, the number of parameters and FLOPs are always smaller compared with original ResNet. The bold shows the least parameters and FLOPs. The FLOPs are calculated by assuming the batch size of 32.

Stage	$H \times W$	C	Network	#Params	#FLOPs
Stage 1	32	16	ResNet-32	23.36K	765.46M
			AFI-ResNet-32	18.82K	612.38M
Stage 2	16	32	ResNet-32	88.77K	727.19M
			AFI-ResNet-32	70.72K	575.20M
Stage 3	8	64	ResNet-32	353.66K	724.30M
			AFI-ResNet-32	281.73K	573.02M

It is difficult to compare *building blocks* with AFI modules to the original because the FLOPs of the AFI module are related to the *building block* number N while the FLOPs of convolution operation is related to the size of images. For easier comparison, Table 2 shows the differences in parameters and FLOPs between ResNet and AFI-ResNet. All results is calculated by the tool¹. In each *stage*, the number of parameters and FLOPs in the AFI module are always smaller than the vanilla ones. Moreover, in most cases, our module assists the original block with decreasing parameters. For example, when the batch size is 32, the ResNet-110 costs 8.17G FLOPs while the AFI-ResNet-110 requires only 6.23G FLOPs, which corresponds to about 24% decrease.

4.2. Explore the Role of AFI Module

Class Activation Visualization. To illustrate whether the low-level features are collected by our AFI module, we generate the heatmap by Grad-CAM [19]. As we can see in the Figure 3, our network has extracted more details of detected objects, and meanwhile, the heatmap generated by ResNet-50 is too smooth to draw the specific borderline of detected objects. With the help of preserving low-level features, the tiny borderlines are more clearer shown in the figure. Besides, compared to ResNet-50, the

¹<https://github.com/Lyken17/pytorch-OpCounter>

area that AFI-ResNet-50 paid most attention to decreases obviously.

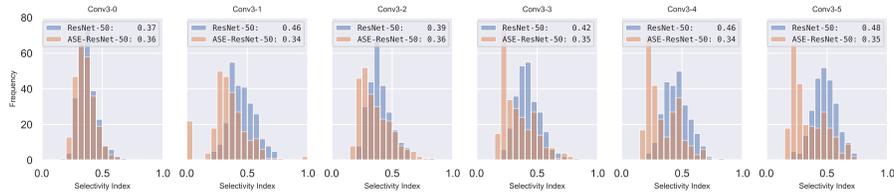


Figure 6: Each figure depicts the class selectivity index distribution for features in both the baseline ResNet-50 and AFI-ResNet-50 various building blocks in the *Conv3* stage of their architectures. The distributions come from the output of the AFI module and corresponding position of the vanilla in each building block. The mean value is reported after the label.

200 **Class Selectivity Index.** Besides, we adopt the class selectivity index metric [14] to analyze the semantic meaning of features. This metric computes, for each feature map, the difference between the highest class-conditional mean activity and the mean of all remaining class-conditional activities over the testing dataset. The measurement is normalized between zero and one where one indicates that a filter only fires for a single class and zero indicates that the filter produces the class-agnostic value. The less class selectivity index contend to the more generalization of channels in a degree. As shown in Figure 6, the AFI-ResNet-50 is able to learn tinier and more generalized features instead of features of a specific class than the vanilla in *Conv3 stage*. With the help of generalized features, the generalization gap between the training and the testing dataset will be closed.

205

210

Table 3: The table shows the accuracy rates (%) of networks on the Office-Home dataset. A:B denotes the task that a model fine-tunes on domain A and tests on domain B.

Model	Ar:Cl	Ar:Pr	Ar:Rw	Cl:Ar	Cl:Pr	Cl:Rw	Pr:Ar	Pr:Cl	Pr:Rw	Rw:Ar	Rw:Cl	Rw:Pr	Avg.
ResNet-50 [8]	34.9	50.0	58.0	37.4	41.9	46.2	38.5	31.2	60.4	53.9	41.2	59.9	46.1
AFI-ResNet-50	42.0	64.0	71.6	48.2	56.8	61.5	48.5	36.4	69.5	63.7	45.3	75.5	56.9

Performances on a Dataset with the Wide Generalization Gap. We consider the situation where the training dataset is quite different from the testing dataset to verify whether our module is able to bridge the generalization gap. We use another experiment to test whether our module is able to capture the task-agnostic feature which

215 can bridge the generalization gap. We consider the situation where the training dataset is quite different from the test dataset. The better performance indicates that the model has greater generalization ability.

To satisfy it, we choose the Office-Home dataset [34] which consists of four distinct domains (Artistic images (Ar), Clip Art (Cl), Product images (Pr), and Real-World
220 images (Rw)). We train the model on ImageNet on the pre-training process, fine-tune on the source domain (A) and then test on the target domain (B). This task is denoted as A:B as shown in Table 3. The pre-training details are similar to the training details in Section (4.3). The epoch of fine-tuning is 20. The learning rate of new full connection layer is $0.01 \times \left(1 + 10 \times \frac{\text{iteration}}{\text{total.iteration}}\right)^{-0.75}$. The learning rate of the other layers is
225 one-tenth of that of the new full connection layer.

Table 3 shows the comparison of various domain adaptation task accuracy between ResNet-50 and AFI-ResNet-50. When the generalization gap between the training dataset and the test dataset larger, the advantage of our network that adapts to learn more generalized features is remarked.

230 The experiment results show that our module can transfer low-level features more efficiently. In fact, our module provides prior knowledge that a convex combination of features can extract the main information. There are some previous methods that assuming additional prior knowledge regularizes the network. For instance, Mixup [35] extends the training distribution by incorporating the prior knowledge that linear inter-
235 polations of features should lead to linear interpolations of classified space.

4.3. Experiments on Real-World Datasets

CIFAR. The CIFAR [36] dataset consists of 60,000 RGB pictures, each with a size of 32×32 . 50,000 of them are used as the training set and 10,000 are used for testing. The CIFAR-10 task requires the network to correctly classify the pictures into
240 10 categories, such as airplanes and automobiles. CIFAR-100 requires the network to classify pictures into 100 categories. We train our network on the training dataset and evaluate it on the test dataset.

By integrating the AFI module with ResNet, ShuffleNetV2, MobileNetV2 and ResNext, we get their AFI counterparts. All the backbone networks have a residual

245 mapping. So we apply the AFI module to the first convolution layer in building blocks
to avoid affecting residual mapping and meanwhile neutralize the waste of the bottle-
neck convolution.

In this experiment, we set SGD with a momentum of 0.9 and a weight decay of
1e-4. We train the networks with the batch size to 64 for 300 epochs. The learning rate
250 is initialized to 0.1 and divided by 10 at 50%, 75% of training process, respectively.
Proportion is adopted in [9]. Data augmentation(mirroring/shifting) is used in training.
Because MobileNetV2 [30], ShuffleNetV2 [10], and other networks are not designed
for the CIFAR dataset, we adopt their variants from Github². All results are reproduced
and the same experiment settings are adopted for a fair comparison.

Table 4: Accuracy rates (%) on CIFAR-100 dataset. All results are reproduced by ourselves for a fair
comparison. Our network results are bold in the table. The FLOPs are calculated by setting the batch size to
32.

Model	C100	#Params	#FLOPs
ResNet-32 [8]	71.16	472.76K	2.23G
AFI-ResNet-32	71.09	378.23K	1.78G
ResNet-110 [8]	73.73	1.74M	8.17G
AFI-ResNet-110	74.03	1.35M	6.23G
ShuffleNetV2 [10]	70.71	0.94M	1.32G
AFI-ShuffleNetV2	72.06	0.95M	1.32G
MobileNetV2 [30]	75.20	2.41M	3.03G
AFI-MobileNetV2	75.94	2.25M	2.67G
DenseNet-40 [9]	75.58	1.06M	34.48G
ResNext-29(32×4d) [11]	78.44	4.87M	24.95G
AFI-ResNext-29(32×4d)	79.37	4.26M	21.74G
DenseNet-100 [9]	79.80	7.09M	230.119G

255 Table 4 shows the comparison of classification error between the original net-

²<https://github.com/kuangliu/pytorch-cifar>

works and their corresponding AFI counterparts. As we can see, most of networks work better in classification with assistance of our AFI module. For example, AFI-ResNext29(32×4d) increases the accuracy rate by 0.93% compared with ResNext29(32×4d) [11] on the CIFAR-100 dataset. Besides, as shown in the results, other AFI-Nets also increase the accuracy rates while decrease or remain at least the number of parameters and FLOPs. We additional compare DenseNets [9] with AFI-Nets. When the accuracy of both is comparable, DenseNets [9] require much more FLOPs than AFI-Nets.

ImageNet. The effect of the AFI module is also evaluated on the ImageNet 2012 dataset [6] which composes about 1.3 million training images and 50k validation images. Both top-1 and top-5 classification accuracy rates are reported on the validation dataset.

In this experiment, we use SGD with a momentum of 0.9 and a weight decay of $1e-4$. We train the networks with batch size 64 for 90 epoch. The initial learning rate is $0.1 * \text{batch_size} / 256$ and divided by 10 at 30, 60, and 80 epochs, respectively. 224×224 images serving as the inputs of the network are cropped from the resized raw images or their horizontal flips. Data augmentation in [37] is used in training. We evaluate our model by applying a center-crop with 224×224 .

The accuracy rates of baseline models and our network on the ImageNet validation set are shown in Table 5. With the same backbone architecture, our model always obtains a higher accuracy rate compared with ResNet. For instance, The AFI-ResNet-152 increases the accuracy rate by 1.24%, decreases the number of parameters by 9.2%, and meanwhile decreases FLOPs by 10% compared with ResNet-152.

Additionally, the memory usage of our model is smaller than the base model(ResNet-50) in both training and testing process as shown in Table 6. In contrast with DenseNet-121 [9] which requires an enormous running space for keeping features of the total 24-layer in DenseNet block (3), our AFI-ResNet-50 maintains features at most 6 layers, therefore the memory usage of our model is much smaller. Besides, reduction of the output of convolution also helps to optimize memory costs during running time. Furthermore, our AFI-Module could release more running space via reducing intermediate gradients in memory with the method proposed by [38].

Table 5: The table shows the accuracy rates (%) of networks on the ImageNet validation set. Our results are marked in bold. All results are reproduced for a fair comparison. The FLOPs are calculated by assuming the batch size of 32.

Model	Top-1 Prec.	Top-5 Prec.	#Params(M)	#FLOPs(G)
ResNet-50 [8]	75.3	92.2	25.56	131.57
AFI-ResNet-50	76.19 _(+0.89)	92.88 _(+0.68)	23.85 _(-1.71)	121.72 _(-9.85)
ResNet-101 [8]	76.4	92.9	44.55	250.69
AFI-ResNet-101	77.72 _(+1.32)	93.76 _(+0.86)	40.75 _(-3.80)	226.91 _(-23.78)
ResNet-152 [8]	77.0	93.3	60.19	369.88
AFI-ResNet-152	78.24 _(+1.24)	93.98 _(+0.68)	54.67 _(-5.52)	332.24 _(-37.64)
MobileNetV2 [30]	66.09	87.14	3.51	10.46
AFI-MobileNetV2	68.24 _(+2.15)	88.54 _(+1.40)	3.63 _(+0.12)	9.36 _(-1.10)
ResNeXt-50(32×4d) [11]	76.08	92.92	25.03	136.30
AFI-ResNeXt-50	77.16 _(+1.08)	93.40 _(+0.48)	21.84 _(-3.19)	116.61 _(-19.69)

Table 6: The comparison of the memory usage of AFI-ResNet-50 and ResNet-50 in the training and testing process. The batch size is 64.

Model	ResNet-50	AFI-ResNet-50
Training Memory (MiB)	7447	7001
Inference Memory (MiB)	3827	2671

4.4. Ablation Study

The Position of AFI Module. In this ablation study, we study whether low-level

Table 7: The accuracy rate (%) comparison of applying our AFI module to different residual stages of ResNet-32 and ResNet-110. Best results are marked in bold.

AFI Stage	ResNet-32	Stage 1	Stage 2	Stage 3
C100	71.16	71.38 _(+0.22)	70.63 _(-0.53)	70.75 _(-0.41)
AFI Stage	ResNet-110	Stage 1	Stage 2	Stage 3
C100	73.73	75.03 _(+1.30)	74.13 _(+0.40)	74.02 _(+0.29)

feature maps or high-level feature maps are more hospitable for exploitation by the

AFI module. Previous study [21] shows that the feature maps learned by the earlier
 290 convolutional layers are more general. According to Table 7, by only applying our
 AFI module to the *stage* 1 of ResNet-32 or ResNet-110, the accuracy rate increases
 obviously, which means low-level features have more practical impacts.

Table 8: The table shows the accuracy rates (%) of AFI-ResNet-50 with different r on the ImageNet validation set.

r	1	2	4	8
Top-1 Prec.	76.21	76.21	76.19	76.11

The Setting of Hyperparameter r . Tables 8 shows our AFI module with different
 reduction ratio r . The results show that our AFI module offers a trade-off between im-
 295 proved accuracy and increased model complexity for the real situations. For instance,
 for a mobile user, the company can adopt AFI-Nets with larger reduction ratio (i.e.
 $r = 4$) while a computer user can adopt AFI-Nets with smaller reduction ratio (i.e.
 $r = 2$) to get higher accuracy.

Table 9: The accuracy rate(%) of AFI-MobileNetV2 with various of the self-attention model in CIFAR-100.

MobileNetV2	Baseline	AFI	AFI-SE	AFI-GE	AFI-CBAM	AFI-ECA
C100	75.20	75.94	76.15	76.22	76.51	76.76

Compatibility with Other Self-attention Models. Intuitively, the AFI-module
 300 is a feature-wise attention mechanism; alternatively, the self-attention module aims to
 model interdependence between channels explicitly. We further conducted experiments
 on CIFAR-100 to demonstrate compatibility. In the experiments, we utilize the self-
 attention modules(SE [25], GE [33], CBAM [26], ECA [39]) to the resulting features.
 As shown in Table 9, all the AFI-SE, AFI-GE, AFI-CBAM, and AFI-ECA(+0.82%)
 305 models get better results, which indicates that our model is compatible with other self-
 attention modules.

5. Conclusion

In this paper, we proposed AFI modules that can adaptively select features for transferring and improve the representational power of networks. The output feature of AFI modules is aggregated by the channel-wise soft attention over a series of features. The structure of AFI modules is simple and can be used directly in existing state-of-the-art residual-like networks easily. Experimental results show the effectiveness of AFI-Networks, which achieves competitive performance on multiple datasets. Compared with baseline models, AFI counterparts can achieve better performance with lower computational complexity. We believe that the AFI modules are broadly applicable across various computer vision tasks, e.g., object detection, instance segmentation and semantic segmentation.

References

- [1] Y. LeCun, Y. Bengio, et al., Convolutional networks for images, speech, and time series, *The handbook of brain theory and neural networks* 3361 (10).
- [2] L. A. Gatys, A. S. Ecker, M. Bethge, Image style transfer using convolutional neural networks, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2414–2423.
- [3] F. Wang, M. Jiang, C. Qian, S. Yang, C. Li, H. Zhang, X. Wang, X. Tang, Residual attention network for image classification, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 3156–3164.
- [4] H. Jiang, J. Wang, Z. Yuan, Y. Wu, N. Zheng, S. Li, Salient object detection: A discriminative regional feature integration approach, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2013, pp. 2083–2090.
- [5] A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with deep convolutional neural networks, in: *Advances in neural information processing systems*, 2012, pp. 1097–1105.

- [6] K. Simonyan, A. Zisserman, Very deep convolutional networks for large-scale image recognition, arXiv preprint arXiv:1409.1556.
- 335 [7] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, A. Rabinovich, Going deeper with convolutions, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2015, pp. 1–9.
- [8] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: Proceedings of the IEEE conference on computer vision and pattern recognition,
340 2016, pp. 770–778.
- [9] G. Huang, Z. Liu, L. Van Der Maaten, K. Q. Weinberger, Densely connected convolutional networks, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 4700–4708.
- [10] N. Ma, X. Zhang, H.-T. Zheng, J. Sun, Shufflenet v2: Practical guidelines for
345 efficient cnn architecture design, in: Proceedings of the European Conference on Computer Vision (ECCV), 2018, pp. 116–131.
- [11] S. Xie, R. Girshick, P. Dollár, Z. Tu, K. He, Aggregated residual transformations for deep neural networks, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 1492–1500.
- 350 [12] C. Szegedy, S. Ioffe, V. Vanhoucke, A. A. Alemi, Inception-v4, inception-resnet and the impact of residual connections on learning, in: Thirty-first AAAI conference on artificial intelligence, 2017.
- [13] Y. Lee, J.-w. Hwang, S. Lee, Y. Bae, J. Park, An energy and gpu-computation efficient backbone network for real-time object detection, in: Proceedings of the
355 IEEE Conference on Computer Vision and Pattern Recognition Workshops, 2019, pp. 0–0.
- [14] A. S. Morcos, D. G. Barrett, N. C. Rabinowitz, M. Botvinick, On the importance of single directions for generalization, arXiv preprint arXiv:1803.06959.

- [15] B. Chang, L. Meng, E. Haber, F. Tung, D. Begert, Multi-level residual networks from dynamical systems view, arXiv preprint arXiv:1710.10348.
360
- [16] R. T. Chen, Y. Rubanova, J. Bettencourt, D. K. Duvenaud, Neural ordinary differential equations, in: Advances in neural information processing systems, 2018, pp. 6571–6583.
- [17] Q. Wang, B. Li, T. Xiao, J. Zhu, C. Li, D. F. Wong, L. S. Chao, Learning deep transformer models for machine translation, arXiv preprint arXiv:1906.01787.
365
- [18] U. M. Ascher, L. R. Petzold, Computer methods for ordinary differential equations and differential-algebraic equations, Vol. 61, Siam, 1998.
- [19] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, D. Batra, Grad-cam: Visual explanations from deep networks via gradient-based localization, in: Proceedings of the IEEE international conference on computer vision, 2017, pp. 618–626.
370
- [20] H. Zhang, C. Wu, Z. Zhang, Y. Zhu, Z. Zhang, H. Lin, Y. Sun, T. He, J. Mueller, R. Manmatha, et al., Resnest: Split-attention networks, arXiv preprint arXiv:2004.08955.
- [21] J. Yosinski, J. Clune, Y. Bengio, H. Lipson, How transferable are features in deep neural networks?, in: Advances in neural information processing systems, 2014, pp. 3320–3328.
375
- [22] Y. Lee, J. Park, Centermask: Real-time anchor-free instance segmentation, arXiv preprint arXiv:1911.06667.
- [23] Y. Chen, J. Li, H. Xiao, X. Jin, S. Yan, J. Feng, Dual path networks, in: Advances in neural information processing systems, 2017, pp. 4467–4475.
380
- [24] W. Wang, X. Li, J. Yang, T. Lu, Mixed link networks, arXiv preprint arXiv:1802.01808.

- [25] J. Hu, L. Shen, G. Sun, Squeeze-and-excitation networks, in: Proceedings of the
385 IEEE conference on computer vision and pattern recognition, 2018, pp. 7132–
7141.
- [26] S. Woo, J. Park, J.-Y. Lee, I. So Kweon, Cbam: Convolutional block atten-
tion module, in: Proceedings of the European Conference on Computer Vision
(ECCV), 2018, pp. 3–19.
- 390 [27] X. Wang, R. Girshick, A. Gupta, K. He, Non-local neural networks, in: Proceed-
ings of the IEEE conference on computer vision and pattern recognition, 2018,
pp. 7794–7803.
- [28] J. Pang, K. Chen, J. Shi, H. Feng, W. Ouyang, D. Lin, Libra r-cnn: Towards
balanced learning for object detection, in: Proceedings of the IEEE Conference
395 on Computer Vision and Pattern Recognition, 2019, pp. 821–830.
- [29] X. Qin, Z. Zhang, C. Huang, C. Gao, M. Dehghan, M. Jagersand, Basnet:
Boundary-aware salient object detection, in: Proceedings of the IEEE Confer-
ence on Computer Vision and Pattern Recognition, 2019, pp. 7479–7489.
- [30] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, L.-C. Chen, Mobilenetv2: In-
400 verted residuals and linear bottlenecks, in: Proceedings of the IEEE conference
on computer vision and pattern recognition, 2018, pp. 4510–4520.
- [31] K. Zhang, Y. Guo, X. Wang, J. Yuan, Z. Ma, Z. Zhao, Channel-wise and feature-
points reweights densenet for image classification, in: 2019 IEEE International
Conference on Image Processing (ICIP), IEEE, 2019, pp. 410–414.
- 405 [32] X. Li, W. Wang, X. Hu, J. Yang, Selective kernel networks, in: Proceedings of the
IEEE conference on computer vision and pattern recognition, 2019, pp. 510–519.
- [33] J. Hu, L. Shen, S. Albanie, G. Sun, A. Vedaldi, Gather-excite: Exploiting feature
context in convolutional neural networks, 2018.
- 410 [34] H. Venkateswara, J. Eusebio, S. Chakraborty, S. Panchanathan, Deep hashing
network for unsupervised domain adaptation, in: Proceedings of the IEEE con-
ference on computer vision and pattern recognition, 2017, pp. 5018–5027.

- [35] H. Zhang, M. Cisse, Y. N. Dauphin, D. Lopez-Paz, mixup: Beyond empirical risk minimization, arXiv preprint arXiv:1710.09412.
- [36] A. Krizhevsky, G. Hinton, et al., Learning multiple layers of features from tiny
415 images.
- [37] X. Li, J. Wu, Z. Lin, H. Liu, H. Zha, Recurrent squeeze-and-excitation context aggregation net for single image deraining, in: Proceedings of the European Conference on Computer Vision (ECCV), 2018, pp. 254–269.
- [38] G. Pleiss, D. Chen, G. Huang, T. Li, L. van der Maaten, K. Q.
420 Weinberger, Memory-efficient implementation of densenets, arXiv preprint arXiv:1707.06990.
- [39] Q. Wang, B. Wu, P. Zhu, P. Li, W. Zuo, Q. Hu, Eca-net: Efficient channel attention for deep convolutional neural networks, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020, pp. 11534–11542.