

ResLT: Residual Learning for Long-tailed Recognition

Jiequan Cui, *Student Member, IEEE*, Shu Liu, Zhuotao Tian, *Student Member, IEEE*,
Zhisheng Zhong, *Student Member, IEEE*, Jiaya Jia, *Fellow, IEEE*

Abstract—Deep learning algorithms face great challenges with long-tailed data distribution which, however, is quite a common case in real-world scenarios. Previous methods tackle the problem from either the aspect of input space (re-sampling classes with different frequencies) or loss space (re-weighting classes with different weights), suffering from heavy over-fitting to tail classes or hard optimization during training. To alleviate these issues, we propose a more fundamental perspective for long-tailed recognition, *i.e.*, from the aspect of parameter space, and aims to preserve specific capacity for classes with low frequencies. From this perspective, the trivial solution utilizes different branches for the head, medium, tail classes respectively, and then sums their outputs as the final results is not feasible. Instead, we design the effective residual fusion mechanism – with one main branch optimized to recognize images from all classes, another two residual branches are gradually fused and optimized to enhance images from medium+tail classes and tail classes respectively. Then the branches are aggregated into final results by additive shortcuts. We test our method on several benchmarks, *i.e.*, long-tailed version of CIFAR-10, CIFAR-100, Places, ImageNet, and iNaturalist 2018. Experimental results manifest the effectiveness of our method. Our code is available at <https://github.com/jiequancui/ResLT>.

Index Terms—Residual Learning, Imbalanced Learning, Long-tailed Recognition.

1 INTRODUCTION

CONVOLUTIONAL neural networks (CNNs) have achieved impressive success on various tasks, including large-scale image classification [1], [2], [3], [4], [5], [6], [7], [8], [9], [10], object detection [11], [12], [13], [14] and semantic segmentation [15], [16], [17]. Especially, with the rise of neural network search [18], [19], [20], [21], [22], [23], performance of CNNs have further taken a big step. However, the incredible progress stems in part from high-quality and large-scale datasets, such as ImageNet [24], MS COCO [25] and Places [26]. These datasets are carefully designed with balanced distributions over different classes. In real-world applications, data could follow an unexpected long-tailed distribution where only a few head classes and a large number of tail classes exist. The long-tailed phenomenon may lead to severe degradation of performance for all models that do not take it into consideration. We in this paper address long-tailed recognition, *i.e.*, recognition on data with long-tailed distributions.

One of the greatest challenges in the long-tailed setting is data imbalance. It is also the principal reason for head classes dominating the training procedure, making head classes enjoy much higher accuracy than tail classes. In the literature, two kinds of methods tackle the long-tailed problem via re-sampling and re-weighting [27], [28], [29], [30], [31], [32]. It is found that over-sampling tail-class images [32], [33], [34], [35] may still suffer from heavy over-fitting to tail classes while under-sampling by dropping a

large number of head-class images [30], [34], [36] inevitably impairs the generalization ability of deep models. Also, for re-weighting methods, previous works [37], [38] have demonstrated that re-weighting strategies may cause optimization difficulties during training on large-scale and real-world datasets.

Regarding procedures, re-sampling rebalances head and tail classes in input space by constructing balanced mini-batches during training as shown in Fig. 1(a). Differently, re-weighting deals with the loss space by assigning different weights for classes according to the respective numbers of samples as shown in Fig. 1(b). Albeit the procedural difference, we note that these two lines, by nature, both *eventually make effects on model parameters* to adjust tail classes response. This finding leads to our key idea of *re-balancing in parameter space directly*.

Our direct operations on the head and tail classes regarding parameter space can avoid heavy over-fitting to tail classes and the difficult optimization problem mentioned above. To illustrate that designing these effective operations is nontrivial, we show in Fig. 1(c)(II) a naive solution. It preserves specialized parameters for the head, medium, and tail classes respectively via three branches of \mathcal{N}_h , \mathcal{N}_m and \mathcal{N}_t . The branches are optimized for respective recognition and the final result is obtained by summing or averaging outputs of them. Unfortunately, this naive solution does not work well in experiments.

Contrary to these naive operations, we model the long-tailed recognition problem as one residual learning process and propose a novel residual fusion mechanism, as shown in Fig. 1(c)(III). It has one main branch (\mathcal{N}_{h+m+t}) optimized to classify images of all classes while the other two residual branches (\mathcal{N}_{m+t} and \mathcal{N}_t) are optimized to recognize images in medium+tail and tail classes respectively. Outputs of

- J. Cui, Z. Tian, Z. Zhong and J. Jia are with the Department of Computer Science & Engineering, The Chinese University of Hong Kong, ShaTin, Hong Kong.
E-mail: {jqcui, zttian, zszhong21, leojia}@cse.cuhk.edu.hk; liushuhust@gmail.com
- J. Jia and S. Liu are with the SmartMore.

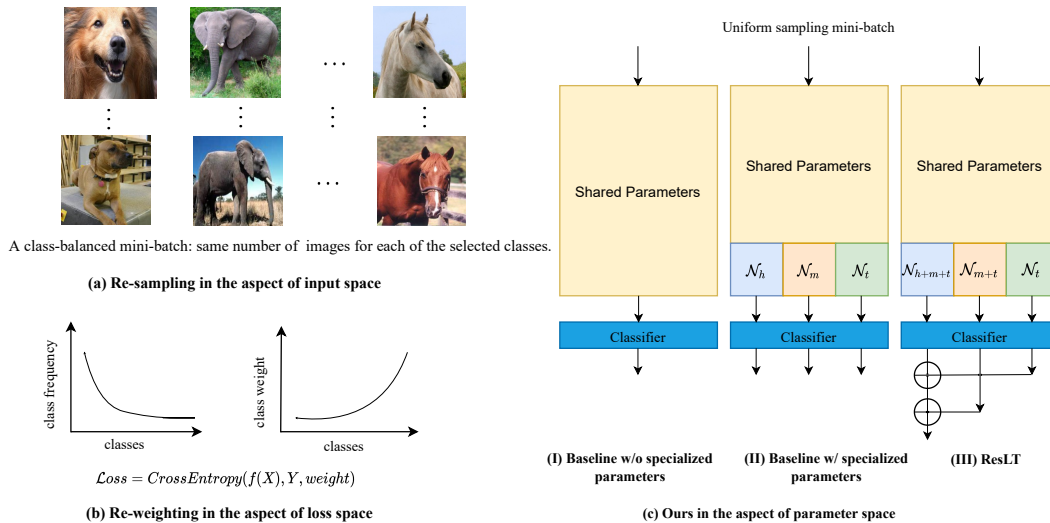


Fig. 1. Comparison between our method and previous ones. Previous re-sampling and re-weighting balance input space and loss space respectively, with effect on the model parameters. Our method is built regarding the most fundamental parameter space, individually preserves specialized parameters for the head, medium, and tail classes with three sub-branches. These sub-branches are combined finally to enhance classification results of the tail and medium classes by the proposed residual fusion mechanism. In (b), (X, Y) is a batch of images and their corresponding labels. $weight$ represents a vector of class-wise weights which usually is in reverse proportion to the number of samples of classes. The classifier is a fully connected (FC) layer that is shared among the three branches in (II) and (III) of (c).

these three branches are aggregated into the final result by additive shortcuts.

An intuitive explanation of our strategy is to gradually enhance classification results on tail classes with other learned residual branches. It is noteworthy that samples belonging to the head, medium, and tail classes still dominate N_{h+m+t} , N_{m+t} and N_t individually in the training procedure, and thus guarantee parameter specialization of the head, medium, and tail classes, which is coherent with the experimental phenomenon in 4.2.3.

We extensively validate our method on several long-tailed benchmark datasets using long-tailed versions of CIFAR-10, CIFAR-100, ImageNet, Places, and iNaturalist 2018 data. Experimental results manifest the effectiveness of our method for long-tailed recognition. Our key contributions are as follows.

- We study the problem from a new perspective of parameter space that leads to an effective re-balance between head and tail classes in the long-tailed setting.
- We propose the residual fusion mechanism, making re-balance in the aspect of parameter space feasible and alleviating limitations in previous works.
- We validate our method on representative benchmarks. Note that this is the first time that a one-stage method consistently surpasses two-stage methods on long-tailed CIFAR-10, CIFAR-100, ImageNet, iNaturalist 2018 all these datasets.

2 RELATED WORK

Re-sampling Strategy There are two groups of re-sampling strategies: over-sampling the tail class images [32], [33], [34], [35] and under-sampling the head class images [30], [34], [36]. Over-sampling is regularly useful on large datasets

and often suffers from heavy over-fitting to tail classes especially on small datasets. For under-sampling, it discards a large portion of data, which inevitably causes degradation of the generalization ability of deep models. Moreover, the under-sampling strategy is not suitable when data is largely imbalanced.

Re-weighting Strategy Re-weighting [37], [38], [39], [40], [41], [42] is another prominent strategy. It assigns different weights for classes and even samples. The vanilla re-weighting method gives class weights in reverse proportion to the number of samples of classes. However, with large-scale data, re-weighting makes the deep models difficult to optimize during training [37], [38]. Cui et al. [29] relieved the problem using “effective numbers” to calculate class weights. Tan et al. [43] conducted re-weighting from the perspective of gradients by ignoring the suppressed gradients of tail classes.

Another line of work is to adaptively re-weight each sample. Focal loss [44] assigned smaller weights for well-classified samples. Li et al. [45] down-weighted samples with very small or large gradients in response to the principle that samples with small gradients are usually well-classified and samples with large gradients are usually out of the distribution.

Two-stage Methods Cao et al. [46] first observed that re-weighting and re-sampling are inferior to the vanilla empirical risk minimization (ERM) algorithm before annealing the learning rate. A two-stage deferred re-balancing optimization schedule was proposed, which trains using vanilla ERM with the LDAM [46] loss before annealing learning rate, and then deploys a re-weighted LDAM loss with a much smaller learning rate.

Recently, Kang et al. [47] and Zhou et al. [48] concluded that although class re-balance strategies matter when

jointly training representation and classifier, uniform sampling gives more general representations. Based on this observation, Kang et al. [47] achieved state-of-the-art results on long-tailed ImageNet (ImageNet-LT) by decomposing representation and classifier learning, *i.e.*, first train the deep models with uniform sampling, then fine-tune the classifier with class-balanced sampling while keeping parameters of representation learning fixed. Similarly, Zhou et al. [48] integrated *mixup* into the proposed cumulative learning strategy with which they bridged the representation learning and classifier re-balancing and achieved state-of-the-art results on long-tailed CIFAR (CIFAR-LT) and iNaturalist 2018.

The two-stage design defies the end-to-end merit that we used to believe since the deep learning era. But why does the two-stage training outperform the end-to-end one largely in long-tailed classification? The work [49] by Tang et al. analyzed the reason from the perspective of causal graph and concluded that the bad momentum causal effects played a vital role. Differently, we address the problem in terms of residual learning.

Meta Learning and Transfer Learning MetaModelNet [39], by Wang et al., learned a meta regression network from head classes and used it to construct the classifier for tail classes. OLTR [50] used memory banks to store mid- and high-level features from head to tail classes. Then they transferred knowledge from head classes to tail classes with a dynamic-meta embedding mechanism. Abdullah Jamal et al. analyzed imbalanced learning from a domain adaptation perspective in [42] and proposed a meta learning framework for long-tailed recognition. Saurabh et al. [51] trained individual experts and transferred knowledge from these experts to the student model. However, the performance with meta learning or transfer learning methods still far from two-stage methods and the training procedure is usually much more complicated.

3 METHOD

Motivation As illustrated in Fig. 1(a), re-sampling strategies build class-balanced mini-batch data for training. However, over-sampling tail class images is easy to get into trouble of heavy over-fitting to tail classes while under-sampling head class images hurts the generalization ability of deep models due to discarding a large number of head-class images that are conducive in learning good representations. For re-weighting methods, previous works [37], [38] demonstrate that it makes deep models difficult to optimize. We, unlike these solutions, explore another way to re-balance regarding parameter space directly, avoiding above disadvantages and making further improvements.

3.1 Re-balancing in Parameter Space

We are the first to propose re-balancing regarding parameter space. To understand the difficulty, we show a naive way by preserving several specialized parameters for the head, medium, tail classes separately by three branches (*i.e.*, \mathcal{N}_h , \mathcal{N}_m , and \mathcal{N}_t). They are shown in Fig. 1(c)(II). After sorting regarding the number of images, all the classes are divided into 3 groups for the head, medium, tail classes, and we

TABLE 1
Top-1 accuracy (%) of baselines on CIFAR-10-LT with the imbalance factor 0.02. ResNet-32 is adopted.

Method	Many-shot	Medium-shot	Few-shot	All
Baseline (CE)	93.60	76.33	68.55	78.50
Baseline (1)	90.70	54.73	41.48	60.22
Baseline (2)	90.83	55.63	40.75	60.24
Baseline (3)	86.73	45.50	55.33	61.80
ResLT (Ours)	85.83	81.47	81.43	83.46

guarantee that the three groups have the same imbalance factor β . $\beta = \frac{N_{min}}{N_{max}}$ where N_{max} and N_{min} are the numbers of training samples for the most frequent class and the least frequent class in each group. Under this case, the core issue is the way to fuse the outputs of these branches into the final prediction because it is impossible to know whether one image belongs to head or medium or tail classes at inference time. Therefore, fusing the three branches is not trivial. Here we compare three straightforward fusion methods with the basic vanilla baseline using cross-entropy. They are set as described in Algorithm 1.

Algorithm 1 Pseudocode of baselines implementation in a PyTorch-like style.

- 1: **Input:** $x \in \mathcal{R}^{c \times h \times w}$ is an image feature map from the network backbone. Specifically, c is the number of channels, h and w are spatial dimensions. $\mathcal{N}_h(\cdot)$, $\mathcal{N}_m(\cdot)$, and $\mathcal{N}_t(\cdot)$ are three branches for the head, medium, and tail classes respectively. $cls(\cdot)$ is the shared linear classifier among the branches with the weight $w \in \mathcal{R}^{k \times c}$. k is the number of classes. $\sigma(\cdot)$ is the softmax activation function.
- 2: $out_h, out_m, out_t = cls(\mathcal{N}_h(x)), cls(\mathcal{N}_m(x)), cls(\mathcal{N}_t(x))$;
- 3: **Baseline (1):**
- 4: $out_a = \text{torch.stack}([out_h, out_m, out_t], \text{dim}=1)$;
- 5: $pred = \text{torch.argmax}(out_a.\text{max}(1), \text{dim}=1)$.
- 6: **Baseline (2):**
- 7: $out_a = out_h + out_m + out_t$;
- 8: $pred = \text{torch.argmax}(out_a, \text{dim}=1)$.
- 9: **Baseline (3):**
- 10: $out_a = \sigma(out_h) + \sigma(out_m) + \sigma(out_t)$;
- 11: $pred = \text{torch.argmax}(out_a, \text{dim}=1)$.

The experimental results are summarized in the Table 1. It shows that these straightforward fusion methods do not work satisfyingly and even yield inferior results compared to the vanilla baseline. We instead propose an effective residual fusion mechanism to handle this challenging issue.

3.2 Residual Fusion Mechanism

Under the long-tailed setting, models usually achieve high accuracy on head classes while the performance on tail classes is unsatisfactory. Based on this phenomenon, we model the long-tailed recognition problem as one residual learning process. It uses extra learned residual branches to enhance classification results on tail classes. Here we stress that the “residual” denotes the nested class assignments for different branches in our method.

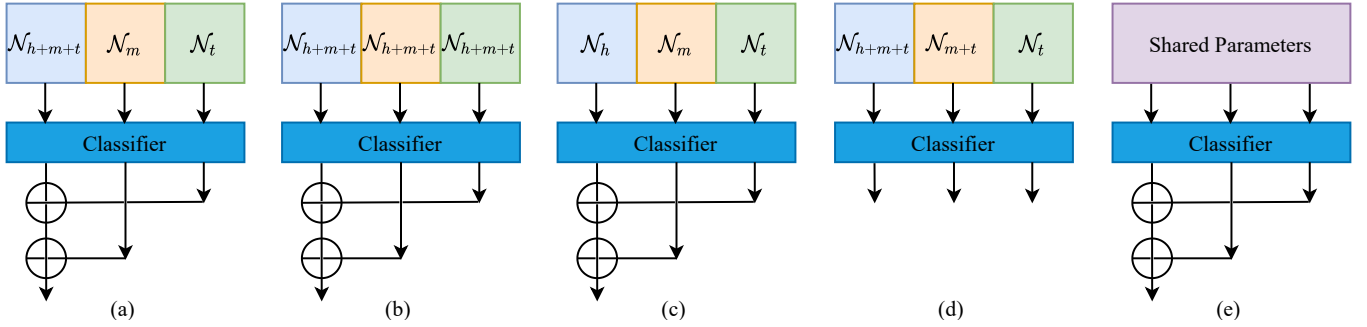


Fig. 2. Ablation study for the residual fusion and parameter specialization mechanisms. (a), (b), (c), (d), and (e) are variants of ResLT. (a) is a weak version of ResLT. For (b) and (c), the three branches have no residual relationship – nested class assignments, unlike ResLT. For (d), there is no additive shortcut, different from ResLT. (e) is without parameter specialization mechanism.

As shown in Fig. 1(c)(III), the model parameters are divided into two parts: shared and specialized parameters. Shared parameters are for common features across all the classes while specialized parameters are to preserve specific capacity for the head, medium, tail classes with three branches \mathcal{N}_{h+m+t} , \mathcal{N}_{m+t} and \mathcal{N}_t correspondingly. In fact, the three branch networks are implemented by only one extra 1x1 grouped convolution with group number 3. Though only a few parameters are introduced, this parameter specialization mechanism plays an important role as demonstrated in Sec. 3.3.1 with experimental verification.

To effectively aggregate the three branches, we propose the residual fusion module shown in Fig. 1(c)(III). We model the long-tailed recognition problem as such residual learning process where one main branch \mathcal{N}_{h+m+t} learns to recognize images of all the classes and the other two residual branches \mathcal{N}_{m+t} , \mathcal{N}_t are used to classify medium+tail classes and tail classes respectively. The outputs of the three branches are added together as the final result. It is worth noting that head, medium, and tail classes still dominate branches \mathcal{N}_{h+m+t} , \mathcal{N}_{m+t} , \mathcal{N}_t respectively, and thus the parameter specialization of these classes are preserved.

In our residual fusion mechanism, loss functions are deployed in training process as

$$\mathcal{L}_{fusion} = \mathcal{J}(\mathcal{N}_{h+m+t}(X) + \mathcal{N}_{m+t}(X) + \mathcal{N}_t(X), Y), \quad (1)$$

$$\mathcal{L}_{branch} = \sum_{i \in \{h+m+t, m+t, t\}} \mathcal{J}(\mathcal{N}_i(SX_i), SY_i), \quad (2)$$

$$\mathcal{L}_{all} = (1 - \alpha)\mathcal{L}_{fusion} + \alpha\mathcal{L}_{branch}, \quad (3)$$

where (X, Y) denotes the uniformly sampled images and labels of a mini-batch. (SX_{h+m+t}, SY_{h+m+t}) is the same as (X, Y) consisting of all class images. (SX_{m+t}, SY_{m+t}) is a subset of (X, Y) only containing images of medium and tail classes. (SX_t, SY_t) is a subset of (X, Y) only containing images belonging to tail classes. \mathcal{J} is cross-entropy loss and α is a hyper-parameter.

The fusion loss item \mathcal{L}_{fusion} in Eq. (1) optimizes for all the classes. The outputs of branches \mathcal{N}_{m+t} and \mathcal{N}_t are added to the outputs of main branch \mathcal{N}_{h+m+t} to obtain the fused outputs. Therefore, during inference, we just sum outputs of the three branches as the final result. The branch-independent loss item \mathcal{L}_{branch} in Eq. (2) is for \mathcal{N}_{h+m+t} , \mathcal{N}_{m+t} and \mathcal{N}_t respectively, further encouraging parameter specialization for head, medium, and tail classes respec-

tively. \mathcal{L}_{all} is our final loss which is a weighted sum of \mathcal{L}_{fusion} and \mathcal{L}_{branch} with a trade-off hyper-parameter α .

3.3 Analysis of Our Method

Our method depends on two key components: parameter specialization mechanism and residual learning mechanism. Here we analyze each of them with extensive experiments.

3.3.1 Importance of Parameter Specialization

Though the three sub-branches are implemented by only a 1x1 grouped convolution with group number 3, we show that it is important to reserve specialized parameters for head, medium and tail classes. We compare strategies with and without the parameter specialization mechanism on CIFAR-10-LT and CIFAR-100-LT with the imbalance factor 0.02. For the setting without parameter specialization, as shown in Fig. 2(e), we let outputs of the three branches be the same while keeping the same loss functions in Eqs. (1), (2) and (3) in implementation. As shown in Fig. 3(a) and Fig. 3(b), the model with the parameter specialization mechanism yields much higher accuracy than the model without it, which manifests the vast importance of our parameter specialization mechanism.

3.3.2 Analysis of Residual Fusion Module

To further study the proposed residual fusion mechanism, we conduct ablation study to understand the role of each component in our residual fusion module by comparing its variants on CIFAR-10-LT and CIFAR-100-LT with the imbalance factor 0.02:

- Possible variant (a), (b), (c): as shown in Fig. 2(a), (b), (c), we explore the necessary of the residual mechanism — the nested classes assignment for different branches in ResLT. \mathcal{N}_h , \mathcal{N}_m , \mathcal{N}_t and \mathcal{N}_{h+m+t} are optimized to classify images of head, medium, tail and all classes respectively.
- Possible variant (d): as shown in Fig. 2(d), we remove the additive shortcuts from ResLT. During inference, we take the output of main branch \mathcal{N}_{h+m+t} as the final results. The branches \mathcal{N}_{m+t} and \mathcal{N}_t only play the role of regularization during training.

As shown in Fig. 3(c) and Fig. 3(d), it is obvious that both the variants (b) and (c) yield lower performance than our

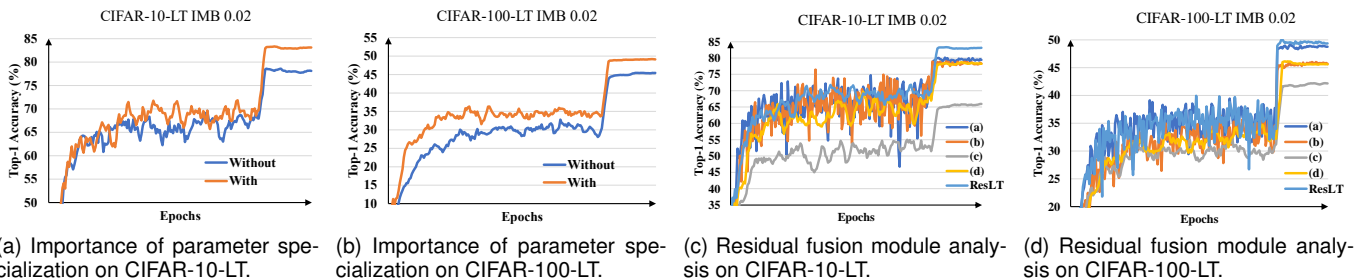


Fig. 3. Ablation studies for the importance of parameter specialization mechanism with (a), (b) and the effectiveness of residual fusion mechanism with (c), (d).

full residual fusion module. Variant (a) can enjoy better performance than (b) and (c) but is inferior to ResLT. This phenomenon demonstrates that the residual relationship (nested class assignments) among branch networks and the additive shortcuts are two core factors that make it feasible to solve long-tailed recognition in parameter space. Moreover, it shows the fundamental difference between our proposed residual learning mechanism and previous work [4].

Comparison with Re-weighting Our method may reminiscent attentive readers of re-weighting methods considering that the loss with respect to tail class images is calculated by three branches while the loss for head class images is calculated by only \mathcal{N}_{h+m+t} branch. In fact, the comparison between ResLT and the structure of Fig. 2(e) in Sec. 3.3.2 shows the importance of the parameter specialization mechanism. At the same time, it also implies that our high performance doesn't benefit from the naive re-weighting scheme. Specifically, for the structure of Fig. 2(e), the loss for tail class images is also calculated by three branches. However, there is a huge performance gap when compared with the structure of our residual fusion module, *i.e.*, ResLT, as indicated by Fig. 3(a) and Fig. 3(b).

Here we further clarify some fundamental differences between our method and re-weighting. We explicitly preserve specific capacity for the head, medium, tail classes. Besides, re-weighting independently pre-defines scalar factors as the class weights simply according to the number of images in the class while our method adaptively enhances classification results of tail classes with the learned residual branches of \mathcal{N}_{m+t} and \mathcal{N}_t . Further, our method captures the important relationship between different classes within each branch.

Discussion on ResLT and RIDE [52] Training with long-tailed data, deep models can easily overfit to tail classes and thus lead to high variance of model predictions. The motivation for RIDE proposed by Wang et al. is just to reduce the predictions variance by an improved model ensemble strategy. It constructs multiple experts and a distribution-aware diversity loss is deployed to promote high diversity among these experts. In contrast, different from re-weighting rebalancing from loss space and re-sampling rebalancing from input space, ResLT conducts rebalance in the aspect of parameter space by allocating individual parameters for tail classes and train the model with the designed residual mechanism.

TABLE 2

Top-1 accuracy (%) on iNaturalist 2018. Rows with † denote results directly copied from [47], [48]. * denotes model is trained with 360 epochs data. "§" denotes the model is trained with mixup and initialized with ImageNet pre-trained model by torchvision. Knowledge distillation is not applied to all methods for fair comparison.

Method	One-stage	ResNet-50
CB-Focal [53] †	✓	61.1
LDAM [54] †	✓	64.6
BBN [48] *	✓	69.6
τ -normalized [47]	✓	69.3
Causal Norm [49]	✓	64.0
LDAM+DRW [54] †	✗	68.0
cRT [47]	✗	67.6
LWS [47]	✗	69.5
LWS [47] §	✗	71.4
ResLT	✓	70.2
ResLT *	✓	70.5
ResLT §	✓	72.3

Reducing the variance of model predictions is a more general idea. A simple model ensemble strategy can even reduce the predictions variance of a model trained with balanced datasets, *e.g.*, full ImageNet, and thus further improve the performance. Differently, ResLT is just designed for long-tailed data and it should be compatible with techniques that aim to reduce the variance of model predictions. This analysis is coherent with our experimental results shown in Table 7. Further, the multiple experts structure in RIDE makes it hard to reuse public pre-trained models, *e.g.*, ImageNet pre-trained models that are usually adopted for Places-LT in the literature. Instead, ResLT is more compatible with previous work in this aspect.

TABLE 3

Top-1 accuracy (%) on ImageNet-LT for ResNeXt-50. † denotes results directly copied from [50]. "§" denotes the model is trained with strong data augmentation (mixup and autoaugment). Knowledge distillation is not applied to all methods for fair comparison.

Method	One-stage	Many	Medium	Few	All
Baseline(CE)	✓	67.6	42.8	13.4	47.6
Mixup [55]	✓	67.3	41.0	11.1	46.3
Lifted Loss [56] †	✓	41.1	35.4	24.0	35.2
Focal Loss [44] †	✓	41.1	34.8	22.4	34.6
Range Loss [57] †	✓	41.1	35.4	23.2	35.1
τ -normalized [47]	✓	60.9	49.2	36.8	51.5
Causal Norm [49]	✓	64.9	48.0	28.9	51.9
cRT [47]	✗	63.7	47.6	28.3	50.7
LWS [47]	✗	63.3	48.4	32.0	51.5
ResLT	✓	63.0	50.5	35.5	52.9
ResLT §	✓	63.6	55.7	38.9	56.1

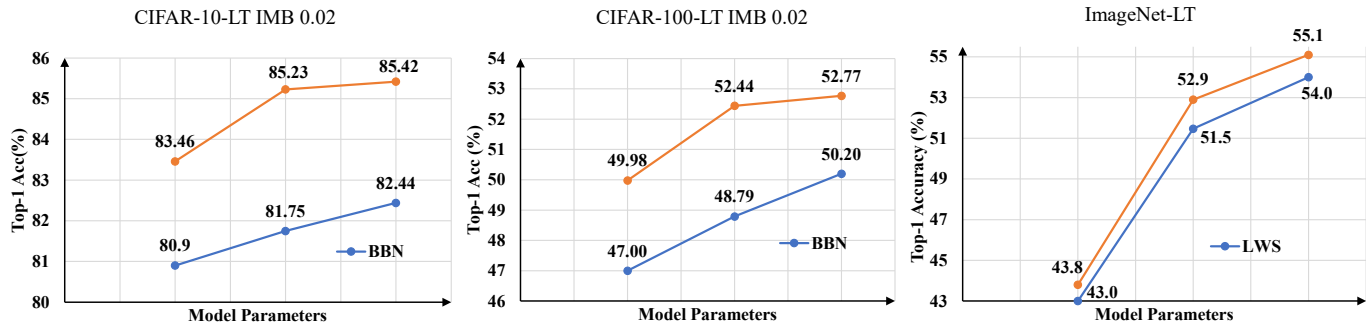


Fig. 4. Effects of model size on long-tailed recognition. ResNet-32(1x), ResNet-32(2x), ResNet-32(3x) are used for CIFAR-LT while ResNet-10, ResNeXt-50, ResNeXt-101 are applied for ImageNet-LT.

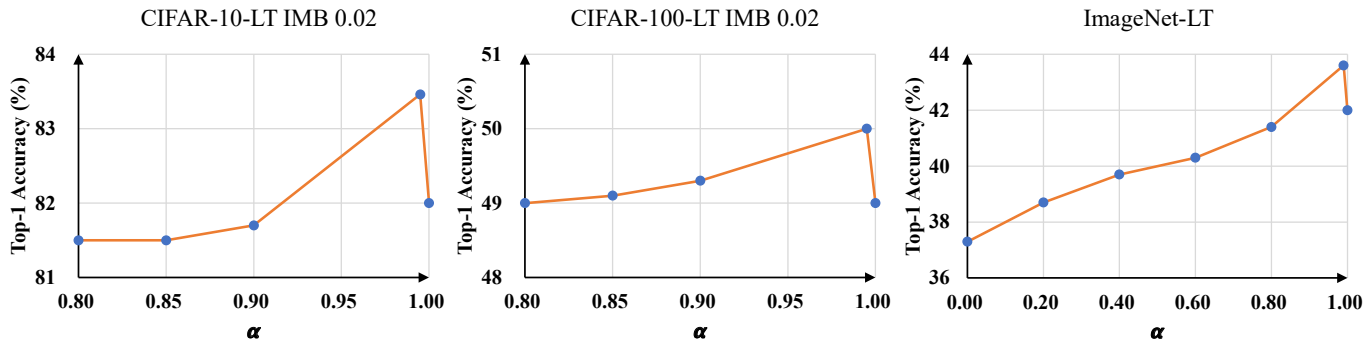


Fig. 5. We conduct ablation studies for choices of α on CIFAR-LT and ImageNet-LT. For CIFAR-LT, ResNet-32 is used, while we apply ResNet-10 on ImageNet-LT.

TABLE 4

Top-1 accuracy (%) on Places-LT for ResNet-152. † denotes results directly copied from [47]. ‡ denotes pre-trained model on ImageNet is trained with strong data augmentation (cutmix) provided by [58].

Method	One-stage	Many	Medium	Few	All
Baseline (CE)	✓	45.7	27.3	8.2	30.2
Lifted Loss [56] †	✓	41.1	35.4	24.0	35.2
Focal Loss [44] †	✓	41.1	34.8	22.4	34.6
Range Loss [57] †	✓	41.1	35.4	23.2	35.1
τ -normalized [47]	✓	37.8	40.7	31.8	37.9
OLTR [50]	✗	44.7	37.0	25.3	35.9
LFME [59]	✗	39.3	39.6	24.2	36.2
cRT [47]	✗	42.0	37.6	24.9	36.7
LWS [47]	✗	40.6	39.1	28.6	37.6
RIDE [52]	✗	-	-	-	-
ResLT	✓	39.8	43.6	31.4	39.8
ResLT ‡	✓	40.3	44.4	34.7	41.0

4 EXPERIMENTS

4.1 Datasets and Experimental Setting

CIFAR-10-LT and CIFAR-100-LT datasets CIFAR-10 and CIFAR-100 both have 60,000 images – 50,000 for training and 10,000 for validation with 10 categories and 100 categories respectively. For a fair comparison, we use the long-tailed versions of CIFAR datasets with the same setting as those used in [46], [48]. Following [48], we conduct experiments with imbalance factors 0.01, 0.02, and 0.1.

ImageNet-LT and Places-LT ImageNet-LT and Places-LT were proposed in [50] by Liu et al. ImageNet-LT is a long-tailed version of the large-scale object classification dataset ImageNet [24] by sampling a subset following the Pareto

distribution with power value $\alpha=6$. It contains 115.8K images from 1,000 categories, with class cardinality ranging from 5 to 1,280. Places-LT is a long-tailed version of the large-scale scene classification dataset Places [26]. It consists of 184.5K images from 365 categories with class cardinality ranging from 5 to 4,980.

iNaturalist 2018 The iNaturalist 2018 [61] is one species classification dataset, which is on a large scale and suffers from extremely imbalanced label distributions. It is composed of 437.5K images from 8,142 categories. In addition to the extreme imbalance, iNaturalist 2018 dataset also confronts the fine-grained problem [62].

Evaluation Protocol Following Liu et al. [50] and Kang et al. [47], on ImageNet-LT, Places-LT, and iNaturalist 2018, we report results on three splits of the set of classes: Many-shot (more than 100 images), Medium-shot (20 ~ 100 images) and Few-shot (less than 20 images). For CIFAR-10-LT and CIFAR-100-LT, we sort all the classes by the number of images in decreasing order. Specifically, we divide CIFAR-10-LT into 3 splits as Many-shot (class 0th ~ class 2th), Medium-shot (class 3th ~ class 5th) and Few-shot (class 6th ~ class 9th). For CIFAR-100-LT, we divided it into 3 splits as Many-shot (class 0th ~ class 34th), Medium-shot (class 35th ~ class 69th) and Few-shot (class 70th ~ class 99th).

Training Details For long-tailed CIFAR datasets, we follow [46], [48] to pre-process images. We randomly crop a 32x32 patch from the original image or its horizontal flip with 4 pixels padded on each side and normalize the pixel values into [0,1]. To be consistent with previous setting [26], [46],

TABLE 5

Top-1 accuracy (%) of ResNet-32 on long-tailed CIFAR-10 and CIFAR-100 (best results are marked in bold font). † denotes results directly copied from [48]. "§" denotes the model is trained with strong data augmentation (mixup and autoaugment). Knowledge distillation is not applied to all methods for fair comparison.

Dataset	One-stage	Long-tailed CIFAR-10			Long-tailed CIFAR-100		
Imbalance factor	-	0.01	0.02	0.1	0.01	0.02	0.1
Mixup [55] †	✓	73.06	77.82	87.10	39.54	44.99	58.02
Focal [44] †	✓	70.38	76.72	86.66	38.41	44.32	55.78
LDAM [54]	✓	73.35	-	86.96	39.60	-	56.91
BBN [48]	✓	79.82	82.18	88.32	42.56	47.02	59.12
CB-Focal [53] †	✓	74.57	79.27	87.10	39.60	45.17	57.99
Causal Norm [49]	✓	80.60	83.60	88.50	44.10	50.30	59.60
CE+DRW [54]	✗	76.34	79.97	87.56	41.51	45.29	58.12
CE+DRS [54]	✗	75.61	79.81	87.38	41.61	45.48	58.11
LDAM+DRW [54]	✗	77.03	81.03	88.16	42.04	46.62	58.71
ELF(LDAM)+DRW [60]	✗	78.10	82.40	88.00	43.10	47.50	58.90
RIDE (3 experts) [52]	✗	-	-	-	48.6	-	-
ResLT	✓	80.44	83.46	89.06	45.34	49.98	60.79
ResLT§	✓	82.40	85.17	89.70	48.21	52.71	62.01
ResLT (3 experts)	✓	-	-	-	49.73	54.51	63.73

we adopt ResNet-32 [4] as our backbone network for all experiments. SGD optimizer with momentum 0.9 is adopted. We train all models for 200 epochs for our ResLT method. The initial learning rate is set to 0.1 and the first five epochs are trained with the linear warm-up. The learning rate is decayed at the 160 and 180 epochs by 0.1 respectively. The batch size 128 is used through all the experiments.

For a fair comparison, we adopt the same experimental setting as [47] on ImageNet-LT, Places-LT and iNaturalist 2018. For Places-LT, following previous setting [47], [50], we choose ResNet-152 as the backbone network, pre-train it on the full ImageNet-2012 dataset (provided by torchvision), and finely tune it for 30 epochs on Places-LT. On ImageNet-LT, we report results with ResNet-10, ResNeXt-50 and ResNeXt-101. For consistency with previous setting, ResNet-50 and ResNet-152 are used for iNaturalist 2018 and we train 200 epochs following [47]. For all experiments, we use SGD optimizer with momentum 0.9 on 4 NVIDIA 2080Ti GPUs. For ResNet-10, ResNet-50, and ResNeXt-50, we adopt cosine learning rate schedule gradually decaying from 0.1 to 0. We use image resolution 224×224 and batch size 256. For ResNet-152 and ResNeXt-101, we adopt cosine learning rate schedule gradually decaying from 0.05 to 0, image resolution 224×224 and batch size 128 for the limited GPU memory.

4.2 Ablation Study

4.2.1 Grouped convolution or 3 separate convolution

We conduct ablation studies with an extra grouped convolution or 3 separate convolutions inserted in the top of the ResNet-32 on CIFAR-10-LT and CIFAR-100-LT. As shown in Fig. 7, our method significantly surpasses these two baselines, which indicates our high performance does not come from the small number of extra parameters.

Moreover, it is worthy to note that BBN [41] also has extra parameters. Our method significantly outperforms BBN as shown in Fig. 4.

4.2.2 How does the model size affect our method?

To go deeper in understanding the effects of model size on long-tailed recognition, we explore various models with dif-

ferent parameters on CIFAR100-LT, CIFAR10-LT with imbalance factor 0.02, and ImageNet-LT datasets. For CIFAR100-LT, CIFAR10-LT, we evaluate ResNet-32(1x), ResNet-32(2x) and ResNet-32(3x). Specifically 1x, 2x, 3x mean the multiplier of channels in each layer. For ImageNet-LT, ResNet-10, ResNeXt-50, ResNeXt-101 are used. As shown in Fig. 4, with the increase of model parameters, our method can consistently surpass BBN on CIFAR-LT and LWS on ImageNet-LT, which indicates the effectiveness of the proposed ResLT method.

4.2.3 Individual branch performance

Our parameter specialization mechanism reserves individual capacity for the head, medium, and tail classes. Along with the proposed residual fusion module, we make it possible to rebalance directly in the aspect of parameter space. To deeply understand the role of residual branches, we collect performance on Many-shot, Medium-shot, and Few-shot of each branch with experiments on CIFAR-100-LT, CIFAR-10-LT with an imbalance ratio of 0.01, and ImageNet-LT.

As shown in Figs. 6, head classes, medium classes and tail classes respectively dominate branch \mathcal{N}_{h+m+t} , \mathcal{N}_{m+t} , and \mathcal{N}_t , realizing the parameter specialization mechanism. Meanwhile, with the two residual branches \mathcal{N}_{m+t} and \mathcal{N}_t , the final performance on medium classes and tail classes is significantly enhanced compared with the single main branch \mathcal{N}_{h+m+t} , testifying the advantages of residual learning mechanism.

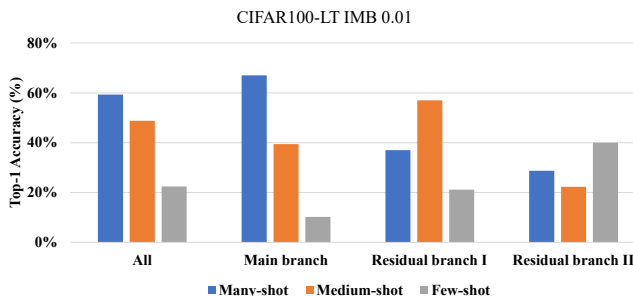
4.2.4 Accuracy on Many-shot, Medium-shot, and Few-shot

For ImageNet-LT, the detailed numbers of performance on Many-shot, Medium-shot and Few-shot are summarized in Table 6 while for iNaturalist 2018, the detailed numbers are put in Table 8. To highlight the advantages of ResLT, we plot the comparisons on ImageNet-LT with previous methods in Fig 8. We observe that ResLT usually can achieve higher performance on Medium-shot and Few-shot classes than previous methods, e.g, LWS, cRT. And the overall improvements just come from small performance degradation on Many-shot classes and such strong performance on Medium-shot and Few-shot classes, which demonstrates

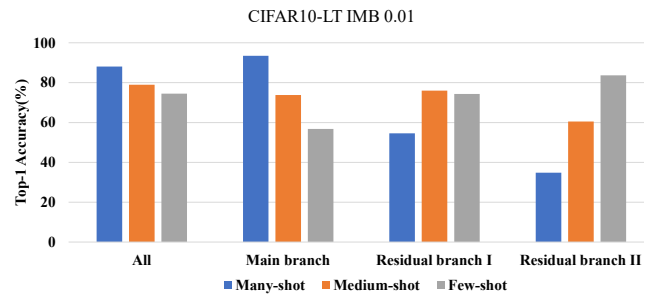
TABLE 6

Top-1 accuracy (%) of Many-shot, Medium-shot and Few-shot on ImageNet-LT with various ResNet backbones. † denotes results directly copied from their original paper.

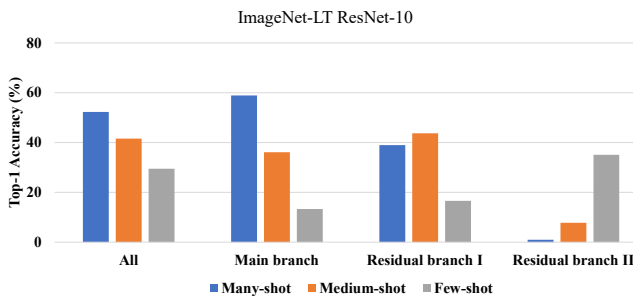
Method	Backbone Model	One-stage	Many-shot	Medium-shot	Few-shot	All
CE (baseline)	ResNet-10	✓	59.7	29.4	5.7	37.3
SEQL†	ResNet-10	✓	-	-	-	36.4
cRT	ResNet-10	✓	53.8	41.3	25.4	43.2
τ -normalize	ResNet-10	✓	50.4	42.1	26.7	42.7
LWS	ResNet-10	✗	51.8	41.6	27.6	43.0
ResLT	ResNet-10	✓	52.3	41.6	29.5	43.8
CE (baseline)	ResNeXt-50	✓	66.8	42.0	14.0	47.0
cRT	ResNeXt-50	✗	63.7	47.6	28.3	50.6
τ -normalize	ResNeXt-50	✓	62.3	48.9	33.7	51.5
LWS	ResNeXt-50	✗	63.3	48.4	32.0	51.5
ResLT	ResNeXt-50	✓	63.0	50.5	35.5	53.0
CE (baseline)	ResNeXt-101-32x4d	✓	69.6	44.6	15.6	49.6
cRT	ResNeXt-101-32x4d	✗	66.2	50.4	30.8	53.3
τ -normalize	ResNeXt-101-32x4d	✓	65.3	51.5	35.2	53.9
LWS	ResNeXt-101-32x4d	✗	65.7	51.4	34.7	54.0
ResLT	ResNeXt-101-32x4d	✓	63.3	53.3	40.3	55.1



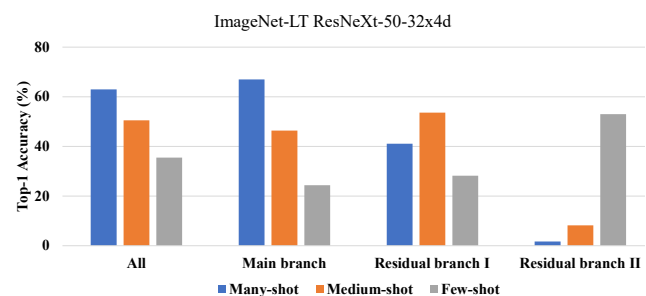
(a) Individual branch performance on CIFAR100-LT with imbalance ratio 0.01.



(b) Individual branch performance on CIFAR10-LT with imbalance ratio 0.01.

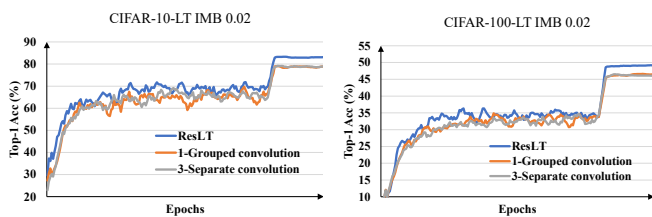


(c) Individual branch performance on ImageNet-LT with ResNet-10.



(d) Individual branch performance on ImageNet-LT with ResNeXt-50.

Fig. 6. Individual branch performance analysis. "All" denotes results are obtained with the aggregated final outputs. Main branch is the \mathcal{N}_{h+m+t} . "Residual branch I" represents the residual branch \mathcal{N}_{m+t} while "Residual branch II" denotes the residual branch \mathcal{N}_t .



(a) Ablation study for 1-Grouped convolution and 3-Separate convolutions on CIFAR-10-LT. (b) Ablation study for 1-Grouped convolution and 3-Separate convolutions on CIFAR-100-LT.

Fig. 7. Ablation study for 1-Grouped convolution and 3-Separate convolutions.

ResLT does a better trade-off among Many-shot, Medium-shot, and Few-shot classes.

4.2.5 How does the hyper-parameter α affect our method?

We conduct ablation study on choosing a proper α on CIFAR-10-LT, CIFAR-100-LT with imbalance factor 0.02, and ImageNet-LT datasets. For CIFAR-10-LT and CIFAR-100-LT, $\alpha = 0.995$ is adopted, while we use 0.99 and 0.90 for ImageNet-LT and iNaturalist 2018 respectively in our experiments. As shown in Fig. 5, even when $\alpha = 1$ that only with \mathcal{L}_{branch} loss in training, the performance is even much better than baselines in Table 1 mentioned in Sec. 3.1, demonstrating the importance of the residual mechanism – *nested class assignments for different branches*.

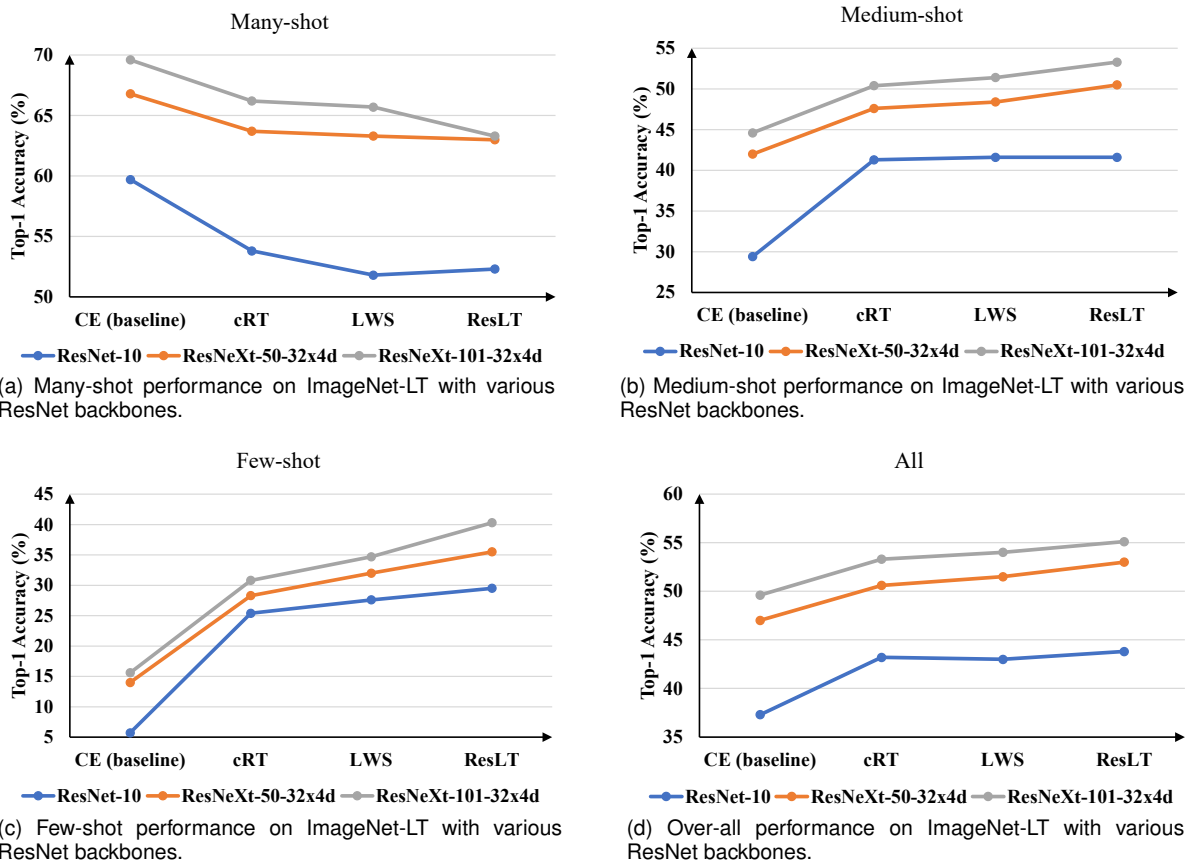


Fig. 8. Many-shot, Medium-shot, and Few-shot performance analysis on ImageNet-LT with various ResNet backbones.

TABLE 7

Top-1 accuracy with 180/200 training epochs and inference speed on ImageNet-LT and iNaturalist 2018. Inference time is calculated with a batch of 64 images on Nvidia GeForce 2080Ti GPU, CUDA11.6, Pytorch1.5, Python3.6.

Method	Epochs	Speed (ms)	Many-shot	Medium-shot	Few-shot	All
With RIDEResNeXt-50 on ImageNet-LT						
RIDE (2 experts) [52]	100	13.8	67.3	50.7	27.2	53.9
RIDE (3 experts) [52]	100	17.0	68.6	52.2	27.8	55.1
RIDE (2 experts) [52]	180	13.8	66.9	51.2	31.2	54.5
RIDE (3 experts) [52]	180	17.0	68.4	52.9	32.2	56.0
ResLT (2 experts)	180	13.9	63.3	55.4	44.3	56.4
ResLT (3 experts)	180	17.2	64.0	56.6	44.8	57.6
With RIDEResNet-50 on iNaturalist 2018						
RIDE (2 experts) [52]	100	12.4	56.5	69.9	70.9	69.0
RIDE (3 experts) [52]	100	15.9	62.0	71.6	72.0	70.8
RIDE (2 experts) [52]	200	12.4	62.5	71.0	70.9	70.1
RIDE (3 experts) [52]	200	15.9	68.3	72.6	71.8	71.7
ResLT (2 experts)	200	12.6	71.5	72.0	72.4	72.1
ResLT (3 experts)	200	16.2	73.0	72.6	73.1	72.9

4.2.6 How does the number of groups affect our method?

Here we have further explored the effects of the different number of groups on CIFAR-10-LT, CIFAR-100-LT with an imbalance factor 0.01. As shown in Table 9, when the number of groups is larger than 3, there are no obvious improvements.

4.2.7 Feature visualizations

To go deeper into understanding the effects of ResLT method, we visualize representation learned by vanilla cross-entropy and ResLT training method on CIFAR-10-LT

and CIFAR-100-LT with an imbalance ratio of 0.01. Specifically, for CIFAR-100-LT, we visualize the randomly selected 10 classes or 20 classes. As shown in Fig 9, The models trained with ResLT usually can produce more compact features, leading to more obvious separation and better performance.

4.2.8 More discussion on ResLT

Zhang et al. [63] identified the trade-off between model robustness and accuracy. For long-tailed classification, Cao et al. [54] derived a margin-based generalization bound for the optimal trade-off between a head class and a tail class

TABLE 8
Top-1 accuracy (%) of Many-shot, Medium-shot and Few-shot on iNaturalist 2018 with ResNet-50 backbone.

Method	Backbone Model	One-stage	Many-shot	Medium-shot	Few-shot	All
CE (baseline)	ResNet-50	✓	75.7	66.9	61.7	65.8
cRT	ResNet-50	✗	73.2	68.8	66.1	68.2
τ -normalize	ResNet-50	✓	71.1	68.9	69.3	69.3
LWS	ResNet-50	✗	71.0	69.8	68.8	69.5
ResLT	ResNet-50	✓	68.5	69.9	70.4	70.2

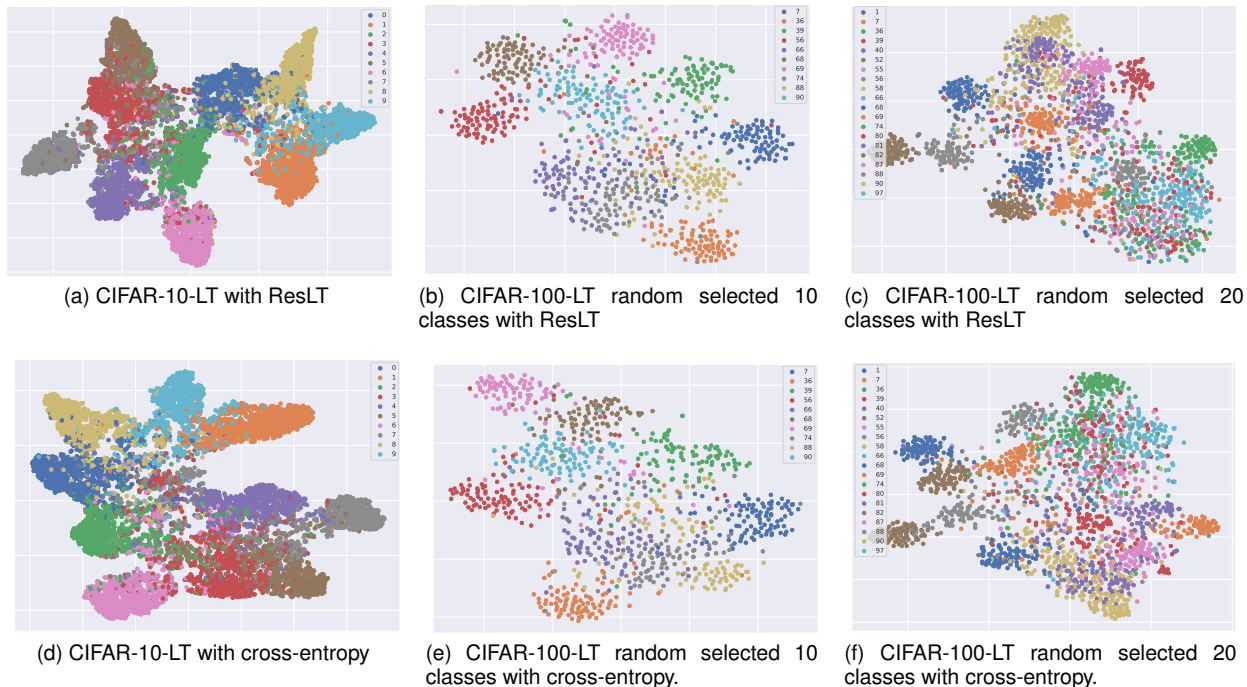


Fig. 9. Visualization comparison between ResLT trained model and cross-entropy trained model. Numbers in all the figures are class indexes.

TABLE 9
Ablation study for different number of groups. “#” represents the number of groups. Top-1 accuracy (%) are reported. We use ResNet-32 as backbone in experiments.

Dataset	#2	#3	#4	#5
CIFAR-10-LT	78.83	80.44	80.80	80.44
CIFAR-100-LT	45.11	45.34	45.34	45.14

margins in the binary classification problem, which also implied that there was a trade-off between head classes accuracy and tail classes accuracy. As discussed in Sec. 4.2.4, we observe ResLT can learn a better trade-off between head classes accuracy and tail classes accuracy with the proposed residual learning mechanism. We leave more theoretical analysis as our future work.

4.3 Comparisons with Previous Methods

Comparison methods In experiments, we compared with two kinds of methods: one-stage methods (with a consistent training strategy throughout the training) and state-of-the-art two-stage methods (input sampling strategy or loss function may differ in different stages).

- One-stage methods. For one-stage methods, we mainly compare with mixup [55], LDAM [46], BBN [48] and Causal Norm [49] methods.
- Two-stage methods. For two-stage methods, we mainly compare with the recently proposed decoupling representation and classifier learning [47] and RIDE [52]. For fair comparison with RIDE, we report the experimental results without knowledge distillation trick.

Comparison on ImageNet-LT Table 3 shows experimental results on ImageNet-LT. Here, we mainly compare with the recent SOTA method [47], [49], [52]. We observe that 90 epochs training is not enough to converge for deep models on ImageNet-LT. Under longer training with 180 epochs, Decouple method [47] can be significantly further improved while it seems to be helpless for Causal Norm [49]. We use this strong setting for comparisons. As shown in Table 3, ResNext-50 trained with our ResLT method in an end-to-end fashion enjoys 1.4% higher than [47] and 1.0% higher than [49]. Moreover, as shown in Fig. 4, we also test our method on ResNet-10 and ResNext-101. Our method consistently surpasses LWS [47] by 0.8% \sim 1.4% across small models to large models.

For fairly comparing with RIDE, we replace ResNeXt-50 with RIDEResNeXt-50 adopted by RIDE. As shown in Table 7, ResLT model achieves 57.6% top-1 accuracy with 3 experts in 180 epochs training, largely outperforming RIDE by 1.6%. With RIDEResNeXt-50 and 3 experts, our model only needs an extra 0.04G FLOPs computational cost at inference time. We also calculate the inference time with an image size 224×224 in Table 7. Our model costs 17.2ms when fed with a batch of 64 images while RIDE uses a comparable time of 17.0 ms.

Comparison on Places-LT Table 4 lists experimental results on Places-LT. Following previous setting [47], [50], we use ResNet-152 pre-trained on full ImageNet dataset as our backbone (provided by torchvision). As shown in Table 4, the model trained with our ResLT in an end-to-end manner achieves 39.8% Top-1 accuracy, which is 1.9% higher than previous Decoupling methods [47].

Comparison on iNaturalist 2018 Experimental results on iNaturalist 2018 are summarized in Table 7. We mainly compare our method with BBN and Kang et al. [47] with 200 training epochs. Note that, though Zhou et al. [48] claimed BBN(2x) trained with 180 epochs, they actually used 360 epochs data due to their dual sampler strategy. Under the same amount of data for training, our method outperforms [48] by 0.9%. Compared with [47], our ResLT model trained in an end-to-end manner achieves 70.2% top-1 accuracy, surpassing it by 0.7%. Further, with larger network ResNet-152, our method achieves 73.2% top-1 accuracy, surpassing [47] by 0.7% under same training setting. When compared with RIDE, ResLT model achieves 72.9% top-1 accuracy with 3 experts in 200 epochs training schedule, significantly surpassing RIDE by 1.2% as shown in Table 7.

Comparison on CIFAR-10-LT and CIFAR-100-LT We conduct extensive experiments on CIFAR-10 and CIFAR-100 with imbalance factors of 0.01, 0.02 and 0.1, same as previous setting [46], [48]. The experimental results are summarized in Table 5. Our method outperforms BBN by 2.78%, 2.96% and 1.76% on CIFAR-100-LT with imbalance factors 0.01, 0.02, and 0.1 respectively. Compared with Causal Norm by Tang et al., ResLT surpasses it by 1.24% and 1.19% separately with imbalance ratio 0.01 and 0.1 on CIFAR-100. With RIDEResNet, ResLT model boosts performance to 49.73% on CIFAR-100 under a imbalance ratio 0.01, surpassing RIDE.

5 CONCLUSION

In this paper, we provide a novel perspective to understand and address the long-tailed recognition problem. The proposed residual fusion mechanism effectively re-balances head and tail classes under this new point of view in parameter space. Extensive experimental results on various representative and challenging benchmarks manifest the effectiveness of our method.

To understand our method thoroughly, we conduct sound ablation studies to show that parameter specialization and residual learning mechanism are two key components to make it work. And we leave more theoretical analysis for ResLT as our future work.

REFERENCES

- [1] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *ICLR*, 2015.
- [2] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *NeurIPS*, 2012.
- [3] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *CVPR*, 2015.
- [4] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *CVPR*, 2016.
- [5] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "Mobilenetv2: Inverted residuals and linear bottlenecks," in *CVPR*, 2018.
- [6] N. Ma, X. Zhang, H.-T. Zheng, and J. Sun, "Shufflenet v2: Practical guidelines for efficient cnn architecture design," in *ECCV*, 2018.
- [7] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *CVPR*, 2018.
- [8] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "Mobilenets: Efficient convolutional neural networks for mobile vision applications," *arXiv preprint arXiv:1704.04861*, 2017.
- [9] X. Wang, R. Girshick, A. Gupta, and K. He, "Non-local neural networks," in *CVPR*, 2018.
- [10] M. Tan and Q. Le, "Efficientnet: Rethinking model scaling for convolutional neural networks," in *ICML*, 2019.
- [11] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *NeurIPS*, 2015.
- [12] S. Liu, L. Qi, H. Qin, J. Shi, and J. Jia, "Path aggregation network for instance segmentation," in *CVPR*, 2018.
- [13] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," in *CVPR*, 2017.
- [14] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *CVPR*, 2016.
- [15] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, "The cityscapes dataset for semantic urban scene understanding," in *CVPR*, 2016.
- [16] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, "Encoder-decoder with atrous separable convolution for semantic image segmentation," in *ECCV*, 2018.
- [17] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, "Pyramid scene parsing network," in *CVPR*, 2017.
- [18] B. Zoph, V. Vasudevan, J. Shlens, and Q. V. Le, "Learning transferable architectures for scalable image recognition," in *CVPR*, 2018.
- [19] H. Liu, K. Simonyan, and Y. Yang, "DARTS: differentiable architecture search," in *ICLR*, 2019.
- [20] M. Tan, B. Chen, R. Pang, V. Vasudevan, M. Sandler, A. Howard, and Q. V. Le, "Mnasnet: Platform-aware neural architecture search for mobile," in *CVPR*, 2019.
- [21] J. Cui, P. Chen, R. Li, S. Liu, X. Shen, and J. Jia, "Fast and practical neural architecture search," in *ICCV*, 2019.
- [22] H. Cai, C. Gan, T. Wang, Z. Zhang, and S. Han, "Once-for-all: Train one network and specialize it for efficient deployment," in *ICLR*, 2020.
- [23] A. Howard, M. Sandler, G. Chu, L.-C. Chen, B. Chen, M. Tan, W. Wang, Y. Zhu, R. Pang, V. Vasudevan *et al.*, "Searching for mobilenetv3," in *ICCV*, 2019.
- [24] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, and M. Bernstein, "ImageNet large scale visual recognition challenge," *IJCV*, 2015.
- [25] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft COCO: Common objects in context," in *ECCV*, 2014.
- [26] B. Zhou, A. Lapedriza, A. Khosla, A. Oliva, and A. Torralba, "Places: A 10 million image database for scene recognition," *IEEE TPAMI*, 2018.
- [27] M. Buda, A. Maki, and M. A. Mazurowski, "A systematic study of the class imbalance problem in convolutional neural networks," *Neural Networks*, 2018.
- [28] C. Huang, Y. Li, C. C. Loy, and X. Tang, "Learning deep representation for imbalanced classification," in *CVPR*, 2016.
- [29] Y. Cui, M. Jia, T. Lin, Y. Song, and S. J. Belongie, "Class-balanced loss based on effective number of samples," in *CVPR*, 2019.
- [30] H. He and E. A. Garcia, "Learning from imbalanced data," *IEEE TKDE*, 2009.

- [31] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, "SMOTE: Synthetic minority over-sampling technique," *JAIR*, 2002.
- [32] L. Shen, Z. Lin, and Q. Huang, "Relay backpropagation for effective learning of deep convolutional neural networks," in *ECCV*, 2016.
- [33] Q. Zhong, C. Li, Y. Zhang, H. Sun, S. Yang, D. Xie, and S. Pu, "Towards good practices for recognition & detection," in *CVPR workshops*, 2016.
- [34] M. Buda, A. Maki, and M. A. Mazurowski, "A systematic study of the class imbalance problem in convolutional neural networks," *Neural Networks*, 2018.
- [35] J. Byrd and Z. Lipton, "What is the effect of importance weighting in deep learning?" in *ICML*, 2019.
- [36] N. Japkowicz and S. Stephen, "The class imbalance problem: A systematic study," *Intelligent Data Analysis*, 2002.
- [37] C. Huang, Y. Li, C. Change Loy, and X. Tang, "Learning deep representation for imbalanced classification," in *CVPR*, 2016.
- [38] C. Huang, Y. Li, C. L. Chen, and X. Tang, "Deep imbalanced learning for face recognition and attribute prediction," *IEEE TPAMI*, 2019.
- [39] Y.-X. Wang, D. Ramanan, and M. Hebert, "Learning to model the tail," in *NeurIPS*, 2017.
- [40] M. Ren, W. Zeng, B. Yang, and R. Urtasun, "Learning to reweight examples for robust deep learning," in *ICML*, 2018.
- [41] J. Shu, Q. Xie, L. Yi, Q. Zhao, S. Zhou, Z. Xu, and D. Meng, "Meta-weight-net: Learning an explicit mapping for sample weighting," in *NeurIPS*, 2019.
- [42] M. A. Jamal, M. Brown, M.-H. Yang, L. Wang, and B. Gong, "Re-thinking class-balanced methods for long-tailed visual recognition from a domain adaptation perspective," in *CVPR*, 2020.
- [43] J. Tan, C. Wang, B. Li, Q. Li, W. Ouyang, C. Yin, and J. Yan, "Equalization loss for long-tailed object recognition," in *CVPR*, 2020.
- [44] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *ICCV*, 2017.
- [45] B. Li, Y. Liu, and X. Wang, "Gradient harmonized single-stage detector," in *AAAI*, 2019.
- [46] K. Cao, C. Wei, A. Gaidon, N. Arechiga, and T. Ma, "Learning imbalanced datasets with label-distribution-aware margin loss," in *NeurIPS*, 2019.
- [47] B. Kang, S. Xie, M. Rohrbach, Z. Yan, A. Gordo, J. Feng, and Y. Kalantidis, "Decoupling representation and classifier for long-tailed recognition," *arXiv preprint arXiv:1910.09217*, 2019.
- [48] B. Zhou, Q. Cui, X.-S. Wei, and Z.-M. Chen, "Bbn: Bilateral-branch network with cumulative learning for long-tailed visual recognition," *arXiv preprint arXiv:1912.02413*, 2019.
- [49] K. Tang, J. Huang, and H. Zhang, "Long-tailed classification by keeping the good and removing the bad momentum causal effect," in *NeurIPS*, 2020.
- [50] Z. Liu, Z. Miao, X. Zhan, J. Wang, B. Gong, and S. X. Yu, "Large-scale long-tailed recognition in an open world," in *CVPR*, 2019.
- [51] S. Sharma, N. Yu, M. Fritz, and B. Schiele, "Long-tailed recognition using class-balanced experts," in *DAGM German Conference on Pattern Recognition*, 2020.
- [52] X. Wang, L. Lian, Z. Miao, Z. Liu, and S. Yu, "Long-tailed recognition by routing diverse distribution-aware experts," in *ICLR*, 2021.
- [53] Y. Cui, M. Jia, T.-Y. Lin, Y. Song, and S. Belongie, "Class-balanced loss based on effective number of samples," in *CVPR*, 2019.
- [54] K. Cao, C. Wei, A. Gaidon, N. Arechiga, and T. Ma, "Learning imbalanced datasets with label-distribution-aware margin loss," in *NeurIPS*, 2019.
- [55] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz, "mixup: Beyond empirical risk minimization," in *ICLR*, 2018.
- [56] H. O. Song, Y. Xiang, S. Jegelka, and S. Savarese, "Deep metric learning via lifted structured feature embedding," in *CVPR*, 2016.
- [57] X. Zhang, Z. Fang, Y. Wen, Z. Li, and Y. Qiao, "Range loss for deep face recognition with long-tailed training data," in *ICCV*, 2017.
- [58] S. Yun, D. Han, S. J. Oh, S. Chun, J. Choe, and Y. Yoo, "Cutmix: Regularization strategy to train strong classifiers with localizable features," in *ICCV*, 2019.
- [59] L. Xiang, G. Ding, and J. Han, "Learning from multiple experts: Self-paced knowledge distillation for long-tailed classification," in *ECCV*, 2020.
- [60] R. Duggal, S. Freitas, S. Dhamnani, D. H. Chau, and J. Sun, "ELF: an early-exiting framework for long-tailed classification," *arXiv preprint arXiv:2006.11979*, 2020.
- [61] G. Van Horn, O. Mac Aodha, Y. Song, Y. Cui, C. Sun, A. Shepard, H. Adam, P. Perona, and S. Belongie, "The iNaturalist species classification and detection dataset," in *CVPR*, 2018.
- [62] X.-S. Wei, P. Wang, L. Liu, C. Shen, and J. Wu, "Piecewise classifier mappings: Learning fine-grained learners for novel categories with few examples," *IEEE TIP*, 2019.
- [63] H. Zhang, Y. Yu, J. Jiao, E. Xing, L. El Ghaoui, and M. Jordan, "Theoretically principled trade-off between robustness and accuracy," in *ICML*, 2019.



Jiequan Cui received the B.Eng. degree (Honors) in Computer Science from the School of Computer Science and Technology, Shandong University (SDU) in 2018. He is currently a Ph.D. Candidate at the Chinese University of Hong Kong (CUHK), under the supervision of Prof. Ji-aya Jia. He serves as a reviewer for IJCV, CVPR, ICCV, ECCV, ICLR, NeurIPS. His research interests include model generalization, neural architecture search, adversarial robustness, imbalanced learning, and image segmentation.



Shu Liu now serves as Co-Founder and Technical Head in SmartMore. He received the BS degree from Huazhong University of Science and Technology and the PhD degree from the Chinese University of Hong Kong. He was the winner of 2017 COCO Instance Segmentation Competition and received the Outstanding Reviewer of ICCV in 2019. He continuously served as a reviewer for TPAMI, CVPR, ICCV, NeurIPS, ICLR and etc. His research interests lie in deep learning and computer vision.



Zhuotao Tian received the B.Eng. degree (Honors) in Computer Science from the School of Computer Science and Technology, Harbin Institute of Technology (HIT) in 2018. He is currently a 3rd year Ph.D. student at the Chinese University of Hong Kong (CUHK), under the supervision of Prof. Jiaya Jia. He serves as a reviewer for IJCV, CVPR, ICCV, ECCV, AAAI. His research interests include few-shot learning, semi-supervised learning, semantic segmentation and scene text detection.



Zhisheng Zhong received the B.Eng. Degree in Communication Engineering from Beijing University of Posts and Telecommunications (BUPT) in 2016. He received the Master Degree in Computer Science from Peking University (PKU) in 2019. Now he is a Ph.D. student at the Department of Computer Science Engineering (CSE), the Chinese University of Hong Kong (CUHK). He serves as a reviewer for NeurIPS, CVPR, ICCV, ICLR and etc. His research interests lie in deep learning and computer vision.



Jiaya Jia received the Ph.D. degree in Computer Science from Hong Kong University of Science and Technology in 2004 and is currently a full professor in Department of Computer Science and Engineering at the Chinese University of Hong Kong (CUHK). He assumes the position of Associate Editor-in-Chief of IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI) and is in the editorial board of International Journal of Computer Vision (IJCV). He continuously served as area chairs for ICCV, CVPR, AAAI, ECCV, and several other conferences for the organization. He was on program committees of major conferences in graphics and computational imaging, including ICCP, SIGGRAPH, and SIGGRAPH Asia. He is a Fellow of the IEEE.