

Dual Compensation Residual Networks for Class Imbalanced Learning

Ruibing Hou, *Student Member, IEEE*, Hong Chang, *Member, IEEE*, Bingpeng Ma, *Member, IEEE*, Shiguang Shan, *Fellow, IEEE*, and Xilin Chen, *Fellow, IEEE*

Abstract—Learning generalizable representation and classifier for class-imbalanced data is challenging for data-driven deep models. Most studies attempt to re-balance the data distribution, which is prone to overfitting on tail classes and underfitting on head classes. In this work, we propose Dual Compensation Residual Networks to better fit both tail and head classes. Firstly, we propose dual **Feature Compensation Module** (FCM) and **Logit Compensation Module** (LCM) to alleviate the overfitting issue. The design of these two modules is based on the observation: an important factor causing overfitting is that there is severe *feature drift* between training and test data on tail classes. In details, the test features of a tail category tend to drift towards feature cloud of multiple similar head categories. So FCM estimates a multi-mode feature drift direction for each tail category and compensate for it. Furthermore, LCM translates the deterministic feature drift vector estimated by FCM along intra-class variations, so as to cover a larger effective compensation space, thereby better fitting the test features. Secondly, we propose a **Residual Balanced Multi-Proxies Classifier** (RBMC) to alleviate the under-fitting issue. Motivated by the observation that re-balancing strategy hinders the classifier from learning sufficient head knowledge and eventually causes underfitting, RBMC utilizes uniform learning with a **residual path** to facilitate classifier learning. Comprehensive experiments on Long-tailed and Class-Incremental benchmarks validate the efficacy of our method.

Index Terms—Class Imbalance Learning, Class-Incremental Learning, Residual Path

1 INTRODUCTION

2 INTRODUCTION

Recently, many vision tasks have made significant progress with deep neural networks [1], driven by the development of deep neural networks as well as large-scale datasets [2]. However, these large-scale datasets are usually well-designed, and the number of samples in each class is balanced artificially. In real world, most data exhibits an *extremely class imbalanced* (long-tail) distribution [3], [4], where a few high-frequency (head) classes occupy most of the instances, while most low-frequency (tail) classes are under-represented by limited samples. The data imbalance poses great challenges to learning unbiased networks [5].

An intuitive solution to address class imbalance problem is to re-balance data distribution [6], [7], [8], [9], including *re-sampling* training data [9], [10], [11] and *re-weighting* loss functions [12], [13], [14]. However, the re-balanced distribution is easily fitted by the over-parameterized deep networks, increasing the risk of under-fitting the whole original data distribution while over-fitting the tail data [15]. Recent work [15] observes that uniform learning (*i.e.* training without re-balancing) actually learns more generalizable representation. Based on above observation, the work [15] proposed a decoupled scheme. In general, the decoupled scheme first learns feature representation under **Uniform**

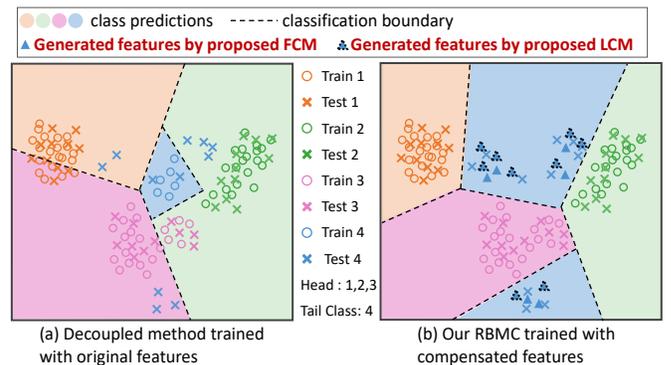


Fig. 1. The feature visualization with t-SNE [16]. (a) The classifier of decoupled method [15] *over-fits* the biased training features on tail classes and *under-fits* on head classes. (b) The proposed RBMC trained with compensated features by FCM and LCM could achieve a more accurate classification boundary, via better fitting both tail and head classes.

Sampling (US), and then re-adjusts the linear classifier on frozen features under **Class-Balanced Sampling (CBS)**. Nevertheless, such a two-stage decoupled scheme still tends to produce a biased classification boundary and fails to deal with imbalance issue well.

We argue that there are two main problems that make existing works [6], [7], [8], [9], [9], [10], [11], [12], [13], [14], [15], [17] produce *biased classification boundary*. The first is the **feature drift between training and test data on tail classes**. For a tail class with scanty training samples, the feature extractor easily overly memorizes the knowledge of the few samples, which cannot generalize to unseen test samples [18]. So the extracted features of training and test

- R. Hou, H. Chang, S. Shan and X. Chen are with Key Laboratory of Intelligent Information Processing, Institute of Computing Technology (ICT), Chinese Academy of Sciences (CAS), Beijing, 100190, China, and University of the Chinese Academy of Sciences, Beijing 100049, China.
- B. Ma is with the School of Computer Science and Technology, University of the Chinese Academy of Sciences, Beijing 100049, China.
E-mail: ruibing.hou@vipl.ict.ac.cn, bpma@ucas.ac.cn, {changhong, sgshan, xlchen}@ict.ac.cn

Manuscript received April 19, 2005; revised August 26, 2015.

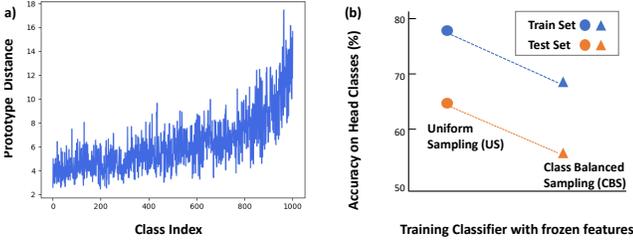


Fig. 2. Illustration of two problems of existing class-imbalanced methods [14], [15], [17]. (a) The distance of per-class prototype between training and test data of existing method [15]. As seen, the distance goes large as the number of training samples decreases, indicating a *large feature drift on tail classes*. (b) With frozen features, top-1 accuracy on head classes of classifier learning under US/CBS sampling [15]. As observed, re-balanced classifier performs poorly on both training and test compared to uniform classifier, proving that re-balanced classifier essentially suffers from *under-fitting the head classes*.

samples tend to deviate from each other. An example is illustrated in Fig. 1 (a) and Fig. 2 (a)¹. As a consequence, the classifier trained on the drifted features would fail to separate the test samples of tail classes, as shown in Fig. 1 (a). The second is the **re-balanced classifier under-fitting on head classes**. In the phase of classifier learning, existing methods typically utilize the re-balancing strategy to obtain a class-balanced classifier. However, the re-balancing strategy encourages the classifier to excessively focus on tail classes, which inevitably leads to under-optimize on head classes. An example is illustrated in Fig. 2 (b). As a result, the classifier trained only by re-balancing strategy is insufficient to recognize all head samples, as shown in Fig. 1 (a).

In this work, we propose **Dual Compensation Residual Networks** (DCRNets) to learn a more unbiased classifier on imbalanced data. DCRNets contains three key modules: dual compensation component consisting of **Feature Compensation Module** (FCM) and **Logit Compensation Module** (LCM) to alleviate the overfitting on tail classes, and **Residual Balanced Multi-Proxies Classifier** (RBMC) to alleviate the underfitting on head classes. The proposed modules can serve as plug-and-play modules, seamlessly integrating with existing methods with any network architecture, such as ResNet [1], ResNeXt [19] or Multi-expert framework [20], forming our DCRNets.

The first module, FCM, alleviates overfitting on tail data by compensating for feature drift on tail classes. FCM is designed based on the following observation as shown in Fig. 1 : (1) The test tail features tend to drift towards the feature cloud of **similar head categories**. (2) The feature drift direction presents a multi-mode, where the test features of one tail category could drift towards **multiple** similar head categories. Based on above observations, FCM estimates a **multi-mode feature drift** direction for each tail category based on **nearest head categories**. Along the estimated feature drift vector, original tail features can be moved closer to the test features, as shown in Fig. 1 (b).

Although FCM can reduce the feature drift of tail classes, it can be challenging to obtain entirely accurate feature drift vectors due to the complex distribution often presented by

test features of tail classes [21]. To this end, we design the second module, LCM, which translates the deterministic feature drift direction estimated by FCM along intra-class variations. In this way, the compensated features can cover a larger feature space and still live in true feature manifold, thus can better fit the test features. Furthermore, instead of performing explicit translation, LCM resorts to an efficient form of logits (classifier outputs) compensation that allows for end-to-end training. As shown in Fig. 1 (b), with FCM and LCM, the training features of a tail category can be compensated to corresponding test feature space, making the learned classification boundary more generalizable.

The last module, RBMC, alleviating underfitting on head data by introducing a residual path from uniform learning. Motivated by the fact that the classifier trained under uniform sampling (US) can better fit the head data [15], RBMC is designed to utilize the predictions of US classifier. In particular, RBMC learns a **residual** mapping from imbalanced predictions of US classifier to desired balanced predictions through class-balanced sampling (CBS) training. As a result, RBMC is optimized jointly under both US and CBS strategies, so that it can take full advantage of all training data to alleviate under-fitting, while also avoiding highly skewing towards head classes.

For evaluation, we conduct Class Imbalanced experiments on four long-tailed benchmarks. We demonstrate that various long-tailed methods [15], [20] can be directly incorporated into the architecture of DCRNets, yielding consistent performance improvements. Further, we conduct experiments on a closely connected field, Class Incremental Learning, to further show the generality of our method.

3 RELATED WORK

3.1 Long-tailed Classification

Long-tailed classification is a long-standing research problem in machine learning, where the key is to overcome the severe class imbalance issue [22]. With the great success deep neural networks have achieved in balanced classification tasks, increasing attention is being shifted to long-tailed classification. Recent studies can be roughly divided into the following six categories.

Re-sampling/Re-weighting. An intuitive solution to long-tailed task is to re-balance the data distribution. Re-sampling data is a common re-balancing strategy to artificially balance the imbalanced data. Typical re-sampling includes over-sampling [5], [12], [23], [24] by simply repeating data from tail classes, and under-sampling by discarding data from head classes [5], [9], [13]. However, duplicated tail data may cause over-fitting tail classes [14], [25]. In addition, it is revealed that over-sampling leads to a side-effect of making head classes under-represented [15], [17]. And discarding head data will certainly impair the generalization ability of deep networks [17].

Re-weighting loss is another prominent re-balancing strategy. The basic idea of Re-weighting is to upweight the tail samples and downweight the head samples in the loss function [6]. The class-sensitive loss [6], [26], [27] assigned the weight to each class inversely proportional to the number of samples. Cui *et al.* [14] proposed to adapt

1. The classes are re-indexed so that the number of training samples decreases as the class index increases

the effective number of samples instead of proportional frequency. A more fine-grained usage is at sample level, *e.g.* focal loss [28], using meta learning [29], [30] or Bayesian uncertainty [31]. The work [8] designed a label-distribution-aware loss to encourage a larger margin for the tail class. However, It is revealed that re-weighting is not capable to handle large-scale long-tailed data and easily causes optimization difficulty [15], [17].

Decoupled Methods. Kang *et al.* [15] firstly proposed the decoupled scheme. This work observed that an uniform sampling scheme actually generated more generalizable representations and achieved stronger performance after re-balancing the classifier. So [15] proposed a two-stage decoupled training pipeline that learned features using US at the first stage, then re-adjusted the classifier with frozen features using CBS at the second stage. However, such a two-stage design is not for end-to-end frameworks. The works [32], [33] improved the two-stage strategy by introducing a post-process to adjust the prediction score. Zhou *et al.* [17] simulated the two-stage training in a single stage by dynamically combining uniform sampling and class balanced sampling with cumulative learning strategy.

In spite of excellent results, above decoupled methods do not take into account the feature drift of tail categories, leading to a biased learning of classifier. On the contrary, our FCM and LCM effectively reduce the unfavorable feature drift. In addition, the decoupled methods typically only rely on class balanced sampling to optimize the classifier. So the classifier learning inevitably suffers from adverse effects of class balanced sampling, *i.e.*, a higher risk of under-fitting on head data. On the contrary, our RBMC utilizes the classifier trained with uniform sampling, which can take full advantage of all training data, thereby alleviating the under-fitting on head classes.

Self-supervised Contrastive Learning. Self-supervised contrastive learning [34], [35] trains the model in a pairwise way by aligning positive sample pairs and repulsing negative sample pairs. Recently, researchers have applied self-supervised learning to long-tailed recognition and demonstrated improved performance. SSP [36] applied self-supervised learning [34] for a good feature initialization. KCL [37] further designed a balanced contrastive loss. However, above works [36], [37] require a two-stage learning paradigm, which are not for end-to-end frameworks. Hybrid-SC [38] proposed a one-stage framework that used a supervised contrastive loss to learn better representation and a cross-entropy classification loss to learn a balanced classifier at the same time. Furthermore, PaCo [39] and GPaCo [40] applied re-balance strategy in supervised contrastive learning to tackle imbalance issue. Notably, self-supervised contrastive learning techniques are orthogonal to our method, which can be easily combined to learn better representation and further improve performance.

Transfer Learning. Another line is to transfer knowledge learned from head classes to tail classes. Some works [2], [41] exploited complex memory backs to transfer semantic features. Wang *et al.* [42] designed a meta network

to transfer the meta-knowledge of parameters evolution. Kim *et al.* [43] translated the head samples to a synthetic tail class via a complicated optimizing process. Chu *et al.* [44] generated tail features by mixing its class-specific features with the class-generic features transferred from head classes. Nevertheless, above methods required either large memory of historical features [2], [41] or complex optimization process [42], [43], [44], incurring high cost in time and memory. Recent methods [45], [46], [47] proposed to transfer the intra-class variance from head class to enlarge the diversity of tail classes. However, [45], [46], [47] made strong assumptions that each class has a shared variance, which could introduce harmful features that interfered with tail features [48]. And recent works [49], [50] proposed to transfer the mean and variance statistics for few-shot classification and long-tailed regression. But [49] focused on few-shot task, which could not distinguish the tail from head classes. And [50] aimed to create a continuity in feature space of continuous labels for regression task, which could not recognize discrete classes for recognition task.

Our FCM shares a similar idea that the head categories can be used to help the learning of tail categories. But to the best of our knowledge, our FCM is the first to point out the test features of tail category tend to drift towards the feature space of similar head categories. Therefore, FCM proposes to compensate the feature drift of tail categories with the information of similar head categories. In implementation, FCM is simple and waives the need for complex memory and optimization mechanism. So FCM is more easier to combine with other advanced long-tailed learning frameworks, such as multi-expert network [20].

Data Generation. In the context of long-tailed recognition, the tail class essentially lacks of sufficient diversity to learn a generalizable model. Therefore, one remedy is to generate more tail samples. In this part, we focus on the approaches that **performed generation without using head information**. For example, Zhang *et al.* [48] estimated a class-wise feature distribution with statistics calculated from observed samples, and then generated virtual features for tail classes based on the estimated distribution. Wang *et al.* [51] used training samples from tail classes and noise vectors to produce new hallucinated tail samples. Mullick *et al.* [52] used the convex combination of existing samples to generate new tail samples. Recent works [53], [54] used Implicit Semantic Data Augmentation (ISDA) algorithm [55] to perform implicit semantic augmentation.

Our LCM shares a similar idea of generating more diverse tail samples. However, different from most works [48], [51], [52] that explicitly generate samples, LCM performs an implicit generation which is highly efficient and easy to combine with other advanced long-tailed methods. Notably, ISDA [55] also implicitly generates samples. But ISDA generates new features around the training instances of the target class. In long-tail classification task, due to the feature drift of tail classes, the generated features would still be far away from the real test features. Thus ISDA has limited effectiveness in learning a generalized model for long-tailed data. Differently, LCM utilizes an extra feature drift compensation operation, which can generate more realistic features to improve the generalizability of

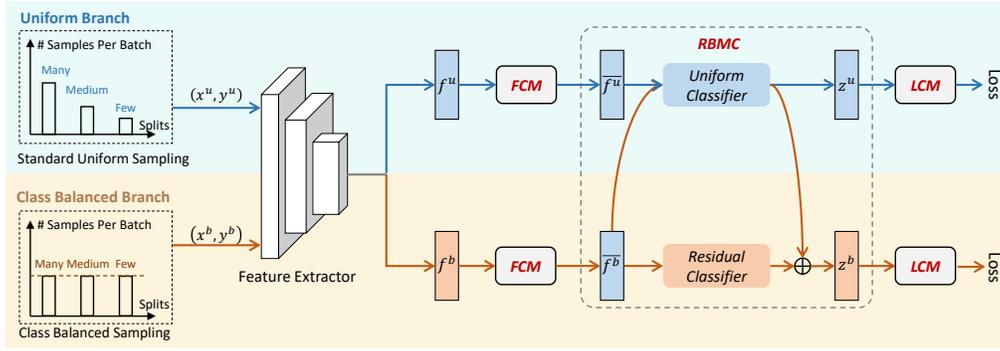


Fig. 3. Framework of Dual Compensation Residual Networks. It contains three key parts: *Feature Compensation Module* (FCM) to compensate the feature drift of tail categories, *Logit Compensation Module* (LCM) to implicitly enrich feature diversity, and *Residual Balanced Multi-proxies Classifier* (RBMC) to learn a class-balanced classification boundary.

long-tailed models.

Multi-Expert Networks. The recent trend in multi-expert networks shows their strong potential to effectively address long-tailed issue. LFME [56] and RIDE [20] are two typical multi-expert architectures. The two methods learned diverse models in parallel, combining with knowledge distillation and distribution-aware expert selection. Further, ACE [57] proposed to use complementary experts, where each expert is most knowledgeable in a specific subset and complementary to other experts. Similarity, ResLT [58] divided classes into several subsets, *i.e.*, all classes, middle and tail classes and tail classes. ResLT then trained multiple experts with various classes subsets under a residual learning mechanism. Notably, ResLT shares most parameters for different experts, which has less computational cost than other multi-expert networks [20], [20], [56]. NCL [59] collaboratively learned multiple experts concurrently and used knowledge distillation strategy [60] to enhance each single expert. TADE [61] extended the multi-expert network to test-agnostic long-tailed recognition task, where the training class distribution is long-tailed while the test class distribution is unknown and can be skewed arbitrarily. With the ensemble methods, the multi-expert networks set new state-of-the-art performance on standard long-tailed benchmarks. We validate in the experiments that our approach can be effectively inserted into the multi-expert networks to bring consistent performance improvements.

3.2 Class Incremental Learning

Incremental learning [62] aims to learn efficient models from the data that gradually comes in a sequence of training phases. Class Incremental Learning (CInL) [63] is a thriving subfield in incremental learning, where each phase has the data of a new set of classes coming from the identical dataset. Related methods mainly focus on how to solve the problems of forgetting old data, which can be roughly categorized into regularized-based and replay-based methods. Regularized-based methods [64], [65], [66] introduce regularization terms in the loss function to consolidate previous knowledge when learning from new data. Replay-based methods [63], [67] use the rehearsal strategy, which store a tiny subset of old data, and replay the model on them to reduce the forgetting.

For class incremental learning, data of old classes is generally not available when new classes appear. Even with the rehearsal strategy, the ratio of the number of new samples to that of old exemplars could be very high, leading to a very serious imbalance issue in class incremental learning [68]. In order to address the class imbalance issue in class incremental learning, Castro *et al.* [69] utilized a balanced fine tuning strategy. Belouadah *et al.* [70] stored the average confidence at each incremental phase to rectify the imbalanced logits. Wu *et al.* [71] added a bias correction layer to correct the output logits. Hou *et al.* [72] combined cosine normalization, less-forget constraint and inter-class separation to mitigate the impact of class imbalance. Zhao *et al.* [73] proposed weight aligning to correct the biased weight in classification layer after uniform training process, which is similar to the τ -norm classifier of decoupled method [15].

Above methods show that the techniques to address class imbalance can be applied to Class Incremental Learning. From the class imbalance view, most above methods [71], [72], [73] combine a uniform training stage on the imbalanced dataset and a balanced fine-tuning stage, similar to the two-phase decoupled strategy [15]. However, these methods neglect the issue of feature drift, which causes overfitting on the few stored old exemplars and poor performance on old classes. Our FCM and LCM can effectively alleviate the unfavorable feature drift, and improve classification performance under class imbalance, thereby benefiting class incremental learning.

4 DUAL COMPENSATION RESIDUAL NETWORKS

In this section, we first present the overall framework of DCRNets. Then, we elaborate on the dual compensation modules and the residual balanced classifier used as parts of DCRNets. We finally give the detailed training and test algorithm of DCRNets.

4.1 Overall Framework

Fig. 3 shows the overview of proposed DCRNets. The network consists of two branches: one Uniform Branch equipped with Uniform Sampling (US) and one Class Balanced Branch equipped with Class Balanced Sampling (CBS). A feature extractor (backbone) is shared between the two branches, and each branch has its own classifier. The

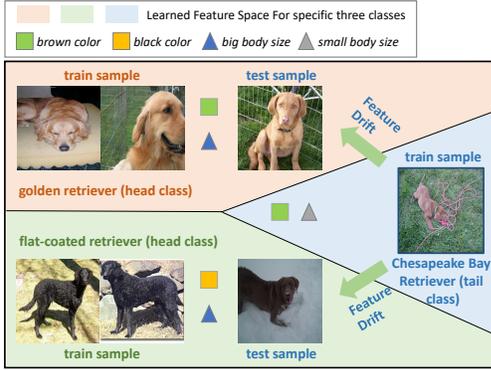


Fig. 4. Due to some common characteristics, the features of two test samples of tail category “Chesapeake Bay retriever” would drift towards the feature space of similar head categories “golden retriever” and “flat-coated retriever”, respectively.

Uniform Branch mainly aims to learn generalizable features. And the Class Balanced Branch is expected to learn a less biased classifier. Building on the two-branch architecture, DCRNets contains three novel parts: FCM, LCM and RBMC. FCM and LCM aims to alleviate overfitting on tail classes by compensating the feature drift between training and test data. RBMC aims to alleviate underfitting on head classes by adding a residual path from uniform branch to class balanced branch.

Formally, given an input image x^u/x^b sampled under US/CBS, the shared backbone is firstly used to extract the feature representation f^u/f^b for Uniform Branch/Class Balanced Branch. Motivated by the observation that there is severe feature drift between training and test data on tail classes, the original feature vectors (f^u and f^b) are sent into FCM to compensate the drift, obtaining compensated features \bar{f}^u and \bar{f}^b . After that, the compensated feature \bar{f}^u is sent into the uniform classifier of Uniform Branch for feature learning, and \bar{f}^b is sent into the RBMC of Class Balanced Branch for classifier learning. To prevent the classifier from fitting to FCM errors and enrich feature diversity especially for tail classes, the outputs of uniform classifier and RBMC are adjusted by LCM. At the last, DCRNets are optimized by applying cross-entropy classification loss on top of LCM.

4.2 Feature Compensation Module

In this section, we first analyze the feature drift phenomenon in detail, and then present the specific operations of FCM.

Analysis of Feature Drift. For imbalanced data, the tail categories typically cannot provide sufficient diversity to learn generalizable features. So there could be discrepancy in feature space between training and test data on tail categories. Fig. 2 (a) computes the Euclidean Distance of per-class feature prototype between training and test data. As seen, the distance increases as the number of training samples decreases, indicating a larger feature drift on tail categories. As a result, the classifier learned to differentiate the training samples can not generalize to differentiate test samples for tail classes.

Furthermore, we analyze the direction of feature drift. For one thing, we observe that the test features of tail

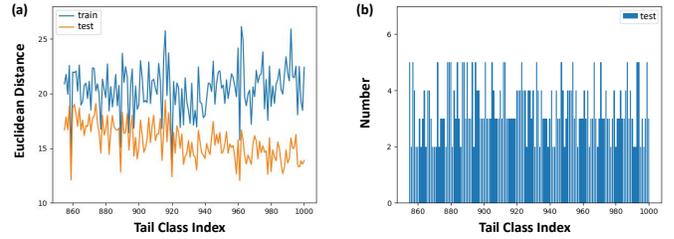


Fig. 5. (a) For each tail category on train/test set, the average distance of all of its features to the prototype of nearest head category. The tail features of **test set** become **closer** to the nearest head category, indicating a drift towards similar head features. (b) For each tail category on test set, the number of different nearest head categories across all samples. On test set, the samples of one tail category could be closest to different head categories. The result indicates that the tail features tend to drift towards **multiple** similar head categories.

category tend to drift towards the feature clouds of similar head categories. In detail, since tail category contains scanty training samples, the feature extractor typically learns **incomplete** features that over-fit to the scanty training samples. For example, as shown in Fig. 4, the *brown color* and *small size* features are enough to recognize the training samples of *Chesapeake Bay retriever* tail category. But the two features are insufficient to correctly recognize the unseen test samples of *Chesapeake Bay retriever*. Meanwhile, because test samples of tail category often share some characteristics with similar head categories, the test features of tail category tend to drift towards the feature space of similar head categories, as illustrated in Fig. 4. Fig. 5 (a) calculates the average distance between all train/test features of each tail category and the prototype of the nearest head category. As seen, as the number of samples in the tail category decreases, the test features of tail category become closer to the nearest head category. The result verifies that the test features of tail category could drift towards feature space of similar head categories.

For another, we observe that the direction of feature drift presents a multi-mode. In details, the test samples of one tail category usually have various characteristics, e.g., various shapes, sizes and poses. Thus, different test samples could share characteristics with different head categories, resulting in a drift towards multiple similar head categories. For example, in Fig. 4, the top test sample has the same *brown color* and *big body* as *golden retriever* category, while the bottom test sample has the same *black color* as *flat-coated retriever* category. As a result, the features of the two test samples drift towards *golden retriever* and *flat-coated retriever* respectively. Fig. 5 (b) counts the number of different nearest head categories across all test samples for each tail category. We observe that most tail categories have more than one nearest head category. The result indicates that the tail features tend to drift towards the feature space of multiple head categories.

Computing of Feature Drift. Aiming at reducing feature drift, we propose FCM to compensate for the training features. The key idea of FCM is to estimate a *multi-mode feature drift direction* for each tail category based on the class similarity.

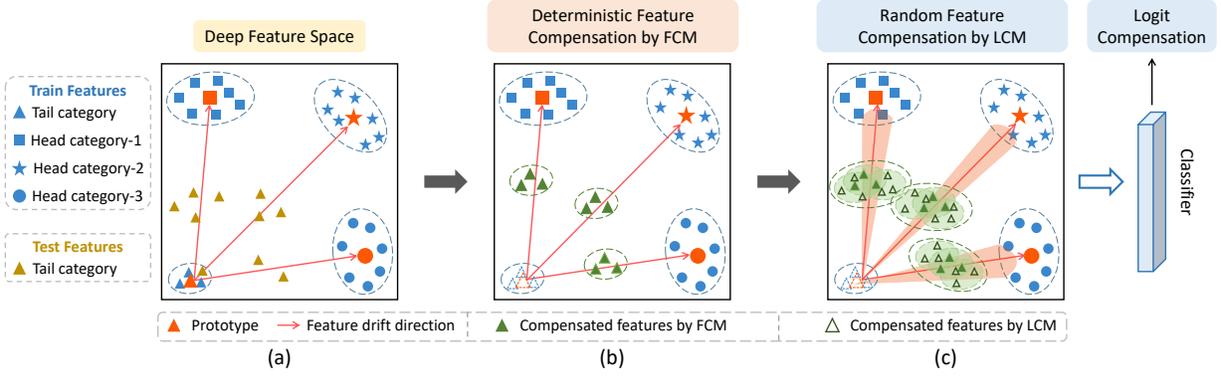


Fig. 6. Illustration of FCM and LCM. For each tail category (triangle), (a) FCM first computes multiple feature drift vectors (red lines) based on the distance vectors between the tail-category prototype and multiple similar head-category prototypes. (b) Then FCM shifts the original features (blue triangles) along the calculated drift vectors to obtain compensated features (green solid triangles). (c) LCM integrates uncertainty (red fan-shaped areas) into the deterministic drift vectors estimated by FCM, and generates richer compensated features (green hollow triangles).

Formally, FCM first computes a *prototype* $c_k \in \mathbb{R}^D$ for each category k . The prototype is defined as the mean vector of the training features belonging to its corresponding class [74]. Then for each tail category t , we select top m head categories with the closest distance to the prototype of this tail category:

$$g(c_t, c_j) = \frac{c_t^T c_j}{\|c_t\| \|c_j\|}, \quad t \in C_t, \quad (1)$$

$$\mathbb{D}_t = \{g(c_t, c_j) \mid j \in C_h\},$$

$$\mathbb{S}_t = \{j \mid g(c_t, c_j) \in \text{top}_m(\mathbb{D}_t)\},$$

where C_h and C_t denote the set of head and tail categories respectively, $\text{top}_m(\cdot)$ is an operator to select the top m elements from the input distance set \mathbb{D}_t . \mathbb{S}_t stores m nearest head categories with respect to the tail category t .

As above analysis, the test features of one tail category could drift towards multiple similar head categories. Therefore, we estimate a set of feature drift vectors for each tail category t based on the nearest head categories in \mathbb{S}_t . As shown in Fig. 6 (a), since the prototype refers to class representative points in feature space, the prototype distance of tail category to the similar head category in \mathbb{S}_t can be used to estimate corresponding feature drift vector:

$$\delta_{tj} = \alpha_t (c_j - c_t), \quad j \in \mathbb{S}_t, \quad (2)$$

where α_t is a positive coefficient to control the strength of drift compensation. As the generalization of learned features is often positively correlated with the number of training samples, the category with fewer training samples is more likely to be closer to similar head category, as shown in Fig. 5 (a). So α_t is set inversely proportional to the number of training samples of the tail category t . It is defined as:

$$\alpha_t = \alpha^0 (N_{\max} - N_t) / (N_{\max} - N_{\min}), \quad (3)$$

where α^0 is a hyperparameter, and N_{\max} and N_{\min} are the numbers of training samples for the most frequent and least frequent category in C_t respectively.

Compensation of Feature Drift. Based on the estimated feature drift vectors (Eq. 2), FCM compensates for original training features to reduce the feature drift. In details, for

each tail category t , we compute the probability s_{tj} of drifting towards similar head category j in \mathbb{S}_t , and the probability s_{tt} of retaining its feature:

$$s_{tj} = \frac{\exp(\tau g(c_t, c_j))}{\sum_{l \in \mathbb{S}_t} \exp(\tau g(c_t, c_l)) + \exp(\tau g(c_t, c_t))}, \quad j \in \mathbb{S}_t, \quad (4)$$

$$s_{tt} = \frac{\exp(\tau g(c_t, c_t))}{\sum_{l \in \mathbb{S}_t} \exp(\tau g(c_t, c_l)) + \exp(\tau g(c_t, c_t))}, \quad (5)$$

where τ is the temperature hyperparameter. Eq. 4 shows that the tail category has a higher probability of drifting towards more similar head category. Eq. 5 shows that the tail category is more likely to retain its features if it has a overall low similarity to all head categories in \mathbb{S}_t . This operation prevents the tail category from updating to a unrelated head category if this tail category has no similar characteristics to all head categories.

At last, for each training sample $(x, t)^2$ of tail category t , FCM shifts original feature $f \in \mathbb{R}^D$ after backbone along the estimated drift vector δ_{tj} with probability s_{tj} , to obtain a compensated feature $\bar{f} \in \mathbb{R}^D$. FCM is formulated as:

$$\bar{f} = f + r, \quad \text{where } p(r = \delta_{tj}) = s_{tj}, \quad j \in \mathbb{S}_t \cup \{t\}. \quad (6)$$

Eq. 6 is equivalent to:

$$\bar{f}_j = f + \delta_{tj}, \quad \text{where } p(\bar{f}_j) = s_{tj}, \quad j \in \mathbb{S}_t \cup \{t\}, \quad (7)$$

where $\delta_{tt} = \mathbf{0}$, indicating retaining original feature. Notably, we maintain training features of head categories, which are already fairly close to test features. As shown in Fig. 1 (b), when applying feature drift vectors to corresponding training features, the training features can be moved to be close to corresponding test space. In this way, the compensated features by FCM can better represent the test feature distribution, thereby achieving a more generalizable classifier.

Notably, although FCM uses prototypes, it does not employ contrastive learning. The prototypes on FCM are directly computed and unlearnable, which are only used to estimate the feature drift direction δ_{tj} (Eq. 2) for tail classes. With the updated features \bar{f}_j (Eq. 7) by FCM,

2. For simplicity, we omit superscript u and b of all symbols in Section 4.2 and 4.3.

we only employ *classification* learning with a cross-entropy classification loss.

4.3 Logit Compensation Module

While FCM can roughly estimate the feature drift direction, it is important to note that the test features of a tail class are unknown during training and typically exhibit complex distribution. As a result, the drift vectors estimated by FCM are prone to errors and difficult to cover all drift directions. To this end, we propose LCM which integrates uncertainty into drift estimation. Integrating uncertainty can prevent the subsequent classifier from over-fitting to FCM errors, thereby enhancing its robustness. Also, integrating uncertainty can expand estimated region to cover more drift directions, so that more diverse features for tail classes can be generated to improve the classification boundary.

As shown in Fig. 6 (c), LCM integrates uncertainty into the **deterministic** drift vectors estimated by FCM (δ_{tj} in Eq. 2), via translating δ_{tj} along intra-class variations. Concretely, we establish a zero-mean Gaussian Distribution $\mathcal{N}(0, \sigma_t^2)$ for each category t , where $\sigma_t \in \mathbb{R}^D$ is the class-conditional *sample standard deviation* estimated from all training features of category t .³ Then a **random** drift vector $\bar{\delta}_{tj}$ is obtained by translating δ_{tj} along a random direction sampled from $\mathcal{N}(0, \sigma_t^2)$. It is formulated as:

$$\bar{\delta}_{tj} = \delta_{tj} + \beta_t \mathbf{g}, \text{ where } \mathbf{g} \sim \mathcal{N}(0, \sigma_t^2), j \in \mathbb{S}_t \cup \{t\}, \quad (8)$$

where β_t is a positive coefficient to control the strength of translating $\bar{\delta}_{tj}$. Equivalently, we have $\bar{\delta}_{tj} \sim \mathcal{N}(\delta_{tj}, \beta_t \sigma_t^2)$. Based on $\bar{\delta}_{tj}$, the compensated feature $\bar{\mathbf{f}}$ can be obtained by applying it to original feature \mathbf{f} with probability s_{tj} ,

$$\bar{\mathbf{f}} = \mathbf{f} + \mathbf{v}, \text{ where } \mathbf{v} \sim \sum_{j \in \mathbb{S}_t \cup \{t\}} s_{tj} \mathcal{N}(\delta_{tj}, \beta_t \sigma_t^2), \quad (9)$$

where \mathbf{v} is the translation of the original feature \mathbf{f} , which is sampled from a *Gaussian Mixture Distribution*. With probability s_{tj} , the compensated feature will shift towards the j^{th} similar head class ($+\delta_{tj}$) and vary along intra-class variance ($\beta_t \sigma_t^2$). In this way, more diverse features that live in test manifold can be obtained, thereby achieving a more generalizable classifier. Notably, FCM (Eq. 6) is a special case of Eq. 9 where the standard deviation of Gaussian Distribution is equal to $\mathbf{0}$ ($\sigma_t^2 = 0$).

Class Adaptive Translation Coefficient. Previous works [45] observe that due to lack of intra-class diversity, the class with fewer training instances usually has a smaller *observed* variance. So the value of *observed standard deviation* (σ_t^2) for rarer class is relatively smaller, which harms the diversity of generated features by LCM (Eq. 9). To address this issue, existing works [45], [46] rely on a shared variance for all classes. But the strong assumption on a shared variance would introduce harmful features [48]. A recent work [53] uses a complex meta-learning strategy to learn variance of tail classes, which incurs excessive computation cost [18].

3. We assume the dimensions of the feature vector are independent of each other and consider the covariance matrix of the feature distribution to be diagonal. It can largely reduce the computation complexity and memory overload.

We use a simple and effective Class Adaptive Translation Coefficient mechanism. Eq. 9 shows β_t controls the variation range of generated feature $\bar{\mathbf{f}}$. Since the category with fewer training instances tends to have a smaller observed variance [45], β_t is set to inversely proportional to the number of training samples N_t . It is formulated as:

$$\beta_t = \beta^0 (N_{\max} - N_t) / (N_{\max} - N_{\min}), \quad (10)$$

where β^0 is a hyperparameter, and N_{\max} and N_{\min} are the numbers of training samples for the most frequent and least frequent category. In this way, class adaptive β_t can mitigate negative impact of smaller observed variance on rarer classes.

Logit Compensation Operation. In Eq. 9, a naive method is to explicitly translate original feature \mathbf{f} for M times, forming an augmented feature set $\{(\mathbf{f}^1, t), (\mathbf{f}^2, t), \dots, (\mathbf{f}^M, t)\}$. Then, the network can be trained with a classification loss on the augmented feature set. We consider the standard cross-entropy loss:

$$\mathcal{L}_M(\mathbf{f}, \mathbf{W}) = \frac{1}{M} \sum_{m=1}^M -\log \left(\frac{e^{\mathbf{w}_t^T \mathbf{f}^m}}{\sum_{k=1}^K e^{\mathbf{w}_k^T \mathbf{f}^m}} \right), \quad (11)$$

where $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_K] \in \mathbb{R}^{D \times K}$ is the weight matrix of the linear classification layer⁴, where K is the total class number. However, the naive implementation is computationally inefficient when M is large. Following ISDA [55], we consider the case when M grows to infinity, and derive an easy-to-compute logit compensation form for highly efficient implementation.

The augmented features can be divided into multiple groups $\{G_j\}_{j=1}^{|\mathbb{S}_t|+1}$ based on drift directions, where the features in G_j follow the Gaussian Distribution $\mathcal{N}(\mathbf{f} + \delta_{tj}, \beta_t \sigma_t^2)$. In the case of $M \rightarrow \infty$, as the probability of sampling from $\mathcal{N}(\delta_{tj}, \beta_t \sigma_t^2)$ is s_{tj} , the group size of G_j would be $s_{tj}M$. We can derive the expected classification loss over the augmented feature sets as:

$$\begin{aligned} & \lim_{M \rightarrow \infty} \mathcal{L}_M(\mathbf{f}, \mathbf{W}) \\ &= \lim_{M \rightarrow \infty} \frac{1}{M} \sum_{j \in \mathbb{S}_t \cup \{t\}} \sum_{\mathbf{f}^m \in G_j} -\log \frac{e^{\mathbf{w}_t^T \mathbf{f}^m}}{\sum_{k=1}^K e^{\mathbf{w}_k^T \mathbf{f}^m}} \\ &= \frac{s_{tj}M}{M} \sum_{j \in \mathbb{S}_t \cup \{t\}} \lim_{M \rightarrow \infty} \frac{1}{s_{tj}M} \sum_{\mathbf{f}^m \in G_j} -\log \frac{e^{\mathbf{w}_t^T \mathbf{f}^m}}{\sum_{k=1}^K e^{\mathbf{w}_k^T \mathbf{f}^m}} \quad (12) \\ &= s_{tj} \sum_{j \in \mathbb{S}_t \cup \{t\}} \mathbb{E}_{\bar{\mathbf{f}} \in \mathcal{N}(\mathbf{f} + \delta_{tj}, \beta_t \sigma_t^2)} -\log \frac{e^{\mathbf{w}_t^T \bar{\mathbf{f}}}}{\sum_{k=1}^K e^{\mathbf{w}_k^T \bar{\mathbf{f}}}} \\ &= s_{tj} \sum_{j \in \mathbb{S}_t \cup \{t\}} \mathbb{E}_{\bar{\mathbf{f}} \in \mathcal{N}(\mathbf{f} + \delta_{tj}, \beta_t \sigma_t^2)} \log \sum_{k=1}^K e^{(\mathbf{w}_k - \mathbf{w}_t)^T \bar{\mathbf{f}}}. \end{aligned}$$

It is infeasible to tackle Eq. 12 directly. Following [55], we use the Jensen's inequality $\mathbb{E}[\log X] \leq \log \mathbb{E}[X]$, and the upper bound of the expected loss can be derived as:

$$\lim_{M \rightarrow \infty} \mathcal{L}_M \leq s_{tj} \sum_{j \in \mathbb{S}_t \cup \{t\}} \log \sum_{k=1}^K \mathbb{E}_{\bar{\mathbf{f}} \in \mathcal{N}(\mathbf{f} + \delta_{tj}, \beta_t \sigma_t^2)} e^{(\mathbf{w}_k - \mathbf{w}_t)^T \bar{\mathbf{f}}}. \quad (13)$$

Then we leverage the fact that $(\mathbf{w}_k - \mathbf{w}_t)^T \bar{\mathbf{f}} \sim \mathcal{N}((\mathbf{w}_k - \mathbf{w}_t)^T (\mathbf{f} + \delta_{tj}), \beta_t ((\mathbf{w}_k - \mathbf{w}_t)^2)^T \sigma_t^2)$ and

4. We set the bias vector of the classifier to $\mathbf{0}$.

moment-generating function $\mathbb{E}[e^X] = e^{\mu + \frac{1}{2}\sigma}$ where $X \sim \mathcal{N}(\mu, \sigma)$ to obtain:

$$\begin{aligned} & \lim_{M \rightarrow \infty} \mathcal{L}_M \\ & \leq s_{tj} \sum_{j \in \mathbb{S}_t \cup \{t\}} \log \sum_{k=1}^K e^{(\mathbf{w}_k - \mathbf{w}_t)^T (\mathbf{f} + \delta_{tj}) + \frac{\beta_t}{2} ((\mathbf{w}_k - \mathbf{w}_t)^2)^T \sigma_t^2} \\ & = s_{tj} \sum_{j \in \mathbb{S}_t \cup \{t\}} -\log \frac{e^{\mathbf{w}_t^T (\mathbf{f} + \delta_{tj})}}{\sum_{k=1}^K e^{\mathbf{w}_k^T (\mathbf{f} + \delta_{tj}) + \frac{\beta_t}{2} ((\mathbf{w}_k - \mathbf{w}_t)^2)^T \sigma_t^2}} \end{aligned} \quad (14)$$

Let $\mathbf{z}_j(\mathbf{W}) = \mathbf{W}^T (\mathbf{f} + \delta_{tj})$ be the K -dimensional output logit of the linear classifier for compensated feature by FCM $\mathbf{f} + \delta_{tj}$ (Eq. 7). As shown in Eq. 14, with LCM, the k^{th} dimension of the logit is adjusted by $\frac{\beta_t}{2} ((\mathbf{w}_k - \mathbf{w}_t)^2)^T \sigma_t^2$. Therefore, **LCM is formulated as**:

$$\bar{\mathbf{z}}_j(\mathbf{W}) = \mathbf{z}_j(\mathbf{W}) + \frac{\beta_t}{2} ((\mathbf{W} - \mathbf{w}_t)^2)^T \sigma_t^2, \quad j \in \mathbb{S}_t \cup \{t\}. \quad (15)$$

Here, the compensated logits $\bar{\mathbf{z}}_j(\mathbf{W})$ can be used to compute the classification loss for end-to-end training. The r.h.s of Eq. 14 can be then denoted as:

$$\begin{aligned} & \lim_{M \rightarrow \infty} \mathcal{L}_M \\ & \leq s_{tj} \sum_{j \in \mathbb{S}_t \cup \{t\}} -\log \frac{e^{[\bar{\mathbf{z}}_j(\mathbf{W})]_t}}{\sum_{k=1}^K e^{[\bar{\mathbf{z}}_j(\mathbf{W})]_k}} \\ & = \mathcal{L}(\bar{\mathbf{z}}_j(\mathbf{W}) | j \in \mathbb{S}_t \cup \{t\}). \end{aligned} \quad (16)$$

4.4 Residual Balanced Multi-Proxies Classifier

Multi-Proxies Classifier. As above analysis, the test features of one tail class tend to drift towards **multiple** similar head classes. So the distribution of test features of one tail class could present a multi-mode [21]. In the naive linear classifier, the weight vector for each class acts as a single proxy [65], [75], which is difficult to optimize on a multi-mode setting. To this end, we consider a **multi-proxies** classifier to better capture the complex feature distribution.

In particular, the multi-proxies classifier allows for L proxies on each tail class during training. As pointed by [75], the proxies can be interpreted as the weight vectors in the classifier. Thus we design one weight vector for each head class, and L weight vectors $\{\mathbf{w}_{t,l}\}_{l=1}^L$ for each tail class t . Then, the class-wise logit $\mathbf{z} \in \mathbb{R}^K$ for each input feature $\bar{\mathbf{f}}$ can be computed as,

$$[\mathbf{z}]_k = \begin{cases} \mathbf{w}_k^T \bar{\mathbf{f}}, & \text{if } k \in C_h, \\ \sum_{l=1}^L \pi_{k,l} (\mathbf{w}_{k,l}^T \bar{\mathbf{f}}), & \text{if } k \in C_t, \end{cases} \quad \pi_{k,l} = \frac{e^{\mathbf{w}_{k,l}^T \bar{\mathbf{f}}}}{\sum_{i=1}^L e^{\mathbf{w}_{k,i}^T \bar{\mathbf{f}}}}, \quad (17)$$

where $\pi_{k,l}$ denotes the similarity of input feature to the l^{th} proxy of class k . Notably, the multi-proxies classifier can be regarded as a **sample-adaptive linear classifier**, where the equivalent weight matrix is $\widehat{\mathbf{W}} = [\widehat{\mathbf{w}}_1, \widehat{\mathbf{w}}_2, \dots, \widehat{\mathbf{w}}_K]$,

$$\widehat{\mathbf{w}}_k = \begin{cases} \mathbf{w}_k & \text{if } k \in C_h, \\ \sum_{l=1}^L \pi_{k,l} \mathbf{w}_{k,l} & \text{if } k \in C_t. \end{cases} \quad (18)$$

Notably, in the case of multi-proxies classifier, LCM simply needs to replace \mathbf{W} with $\widehat{\mathbf{W}}$ in Eq. 15.

Residual Balanced Multi-Proxies Classifier. In order to prevent the classifier predictions from highly biasing to head classes, existing methods [15], [17] use CBS to optimize the classifier. However, CBS leads to a side-effect of making head classes under-optimized, increasing the risk of classifier under-fitting the head data. As shown in Fig. 2 (b), CBS classifier performs poorly on both training and test set of head classes, which proves that CBS essentially suffers from under-fitting on head classes. To this end, we design a RBMC that uses a residual mechanism to learn a class-balanced classifier.

The structure of RBMC is illustrated in Fig. 3. RBMC consists of a Uniform Classifier and a Residual Classifier. Both classifiers adopt multi-proxies classifiers, which are trained under US and CBS respectively. In particular, given the feature vector $\bar{\mathbf{f}}^u \in \mathbb{R}^D$ after FCM on Uniform Branch with US, the **uniform classifier** can be formulated as:

$$\mathbf{z}^u = \widehat{\mathbf{W}}^u{}^T \bar{\mathbf{f}}^u, \quad (19)$$

where $\widehat{\mathbf{W}}^u$ denotes the weight matrix of the multi-proxies uniform classifier (Eq. 18). However, since US ignores the tail data [15], the prediction of uniform classifier is imbalanced and highly skewed to head classes.

Therefore, another **residual classifier** is used to produce a class balanced prediction. It learns the mapping from feature space to the residual between the imbalanced logits obtained by *uniform classifier* and desired balanced logits. Given the feature vector $\bar{\mathbf{f}}^b \in \mathbb{R}^D$ after FCM on Class Balanced Branch with CBS, **RBMC is formulated as**:

$$\mathbf{z}^b = \widehat{\mathbf{W}}^u{}^T \bar{\mathbf{f}}^b + \widehat{\mathbf{W}}^r{}^T \bar{\mathbf{f}}^b = (\widehat{\mathbf{W}}^u + \widehat{\mathbf{W}}^r)^T \bar{\mathbf{f}}^b, \quad (20)$$

where $\widehat{\mathbf{W}}^r$ is the weight matrix of the multi-proxies residual classifier. In such design, RBMC is jointly learned under both CBS (through *residual classifier*) and US (through *uniform classifier*) schemes. Therefore, the head information compromised by CBS can be recovered through uniform classifier. In this sense, RBMC can utilize the information from whole training dataset to alleviate under-fitting, while avoiding highly skewing to head classes.

4.5 Overall Algorithm

Training phase. The training algorithm of DCRNets is shown in Algorithm 1. We apply US and CBS to training set separately, and obtain sample $(\mathbf{x}^u, \mathbf{y}^u)$ for uniform branch and $(\mathbf{x}^b, \mathbf{y}^b)$ for class balanced branch. Then, the two samples are fed into the shared feature extractor to acquire the feature vectors $\bar{\mathbf{f}}^u \in \mathbb{R}^D$ and $\bar{\mathbf{f}}^b \in \mathbb{R}^D$. Since the training and test features of tail categories deviate from each other, the feature vectors are sent into FCM to compensate the drift. FCM generates a set of compensated features $\{\bar{\mathbf{f}}_j^u\}_{j \in \mathbb{S}_{y^u} \cup \{y^u\}}$ and $\{\bar{\mathbf{f}}_j^b\}_{j \in \mathbb{S}_{y^b} \cup \{y^b\}}$ (Eq. 7). After that, each $\bar{\mathbf{f}}_j^u$ is sent into the uniform classifier of weight matrix $\widehat{\mathbf{W}}^u$ to predict the class-wise logit $\mathbf{z}_j^u \in \mathbb{R}^K$ (Eq. 19). And each $\bar{\mathbf{f}}_j^b$ is sent into RBMC of weight matrix $\widehat{\mathbf{W}}^u + \widehat{\mathbf{W}}^r$ to predict the logit $\mathbf{z}_j^b \in \mathbb{R}^K$ (Eq. 20). In order to prevent the model from over-fitting to FCM errors, the output logits are adjusted by LCM and obtain compensated logits $\bar{\mathbf{z}}_j^u$ and $\bar{\mathbf{z}}_j^b$ (Eq. 15). Finally, the classification loss on Uniform Branch is

Algorithm 1 Learning algorithm of DCRNets.

Require: X is Training Dataset; $\mathcal{F}_{cnn}(\cdot; \Theta)$ is CNN feature extractor with parameter Θ ; $\widehat{\mathbf{W}}^u/\widehat{\mathbf{W}}^r$ is weight matrix of uniform/residual multi-proxies classifier; T is the total training epochs.

- 1: **for** $t = 1$ to T **do**
- 2: $(\mathbf{x}^u, y^u) \leftarrow \text{UniformSampler}(X)$
- 3: $(\mathbf{x}^b, y^b) \leftarrow \text{ClassBalancedSampler}(X)$
- 4: $\mathbf{f}^u \leftarrow \mathcal{F}_{cnn}(\mathbf{x}^u; \Theta)$
- 5: $\mathbf{f}^b \leftarrow \mathcal{F}_{cnn}(\mathbf{x}^b; \Theta)$
- 6: Send \mathbf{f}^u to FCM to obtain $\{\widehat{\mathbf{f}}_j^u\}_{j \in \mathbb{S}_{y^u} \cup \{y^u\}}$ with Eq. 7
- 7: Send \mathbf{f}^b to FCM to obtain $\{\widehat{\mathbf{f}}_j^b\}_{j \in \mathbb{S}_{y^b} \cup \{y^b\}}$ with Eq. 7
- 8: $\mathbf{z}_j^u = \widehat{\mathbf{W}}^u \widehat{\mathbf{f}}_j^u$ (Eq. 19)
- 9: $\mathbf{z}_j^b = (\widehat{\mathbf{W}}^u + \widehat{\mathbf{W}}^r)^T \widehat{\mathbf{f}}_j^b$ (Eq. 20)
- 10: send \mathbf{z}_j^u to LCM to obtain $\widehat{\mathbf{z}}_j^u(\widehat{\mathbf{W}}^u)$ with Eq. 15
- 11: send \mathbf{z}_j^b to LCM to obtain $\widehat{\mathbf{z}}_j^b(\widehat{\mathbf{W}}^u + \widehat{\mathbf{W}}^r)$ with Eq. 15
- 12: $\mathcal{L}_1 \leftarrow \mathcal{L}(\widehat{\mathbf{z}}_j^u(\widehat{\mathbf{W}}^u) | j \in \mathbb{S}_{y^u} \cup \{y^u\})$ with Eq. 16
- 13: $\mathcal{L}_2 \leftarrow \mathcal{L}(\widehat{\mathbf{z}}_j^b(\widehat{\mathbf{W}}^u + \widehat{\mathbf{W}}^r) | j \in \mathbb{S}_{y^b} \cup \{y^b\})$ with Eq. 16
- 14: $\mathcal{L} = \phi \mathcal{L}_1 + (1 - \phi) \mathcal{L}_2$
- 15: Update parameters $\{\Theta, \widehat{\mathbf{W}}^u, \widehat{\mathbf{W}}^r\}$ by minimizing \mathcal{L}
- 16: **end for**

denoted as $\mathcal{L}_1 = \mathcal{L}(\widehat{\mathbf{z}}_j^u(\widehat{\mathbf{W}}^u) | j \in \mathbb{S}_{y^u} \cup \{y^u\})$ (Eq. 16), and the classification loss on Class Balanced Branch is defined as $\mathcal{L}_2 = \mathcal{L}(\widehat{\mathbf{z}}_j^b(\widehat{\mathbf{W}}^u + \widehat{\mathbf{W}}^r) | j \in \mathbb{S}_{y^b} \cup \{y^b\})$ (Eq. 16). The overall loss for DCRNets is computed as $\phi \mathcal{L}_1 + (1 - \phi) \mathcal{L}_2$, where ϕ is a trade-off hyperparameter. DCRNets can be end-to-end optimized by minimizing the overall loss function.

Test Phase. In the test phase, only Class Balanced Branch is used for prediction and the compensation modules (FCM and LCM) are removed. Formally, given a test sample, we first use the feature extractor to obtain the feature vector. Then the feature is sent into RBMC to obtain a class balanced prediction (Eq. 20). So DCRNets only incur little computational burden, which is extremely efficient.

5 DISCUSSIONS

RBMC. The scheme of decoupling feature learning and classifier learning [15] has become a dominant direction for class-imbalanced task. Our RBMC uses uniform sampling to benefit **classifier** learning, which effectively alleviates the under-fitting of decoupled classifier [15] to head classes. Also, the feature learning and classifier learning can obey a unified way without decoupling. The jointly learning enhances the compatibility between features and classifier, and meets the end-to-end merit in deep learning. Considering its simplicity and effectiveness, RBMC could serve as a generic classifier in class-imbalanced learning literature.

FCM and LCM. Most studies on class imbalanced learning [15], [17], [76] attempt to deal with **label shift** between training and test phase. Existing works generally assume that the extracted features of training and test data follow a same distribution. However, this assumption may not hold. As shown in Fig. 2 (a), there is a severe **feature drift** between training and test especially on tail classes. To the best of our knowledge, we are the first to point out the test features of tail classes tend to drift towards similar head classes. And our FCM and LCM can effectively alleviate the feature

drift issue through compensation operations. Also, we combine FCM and LCM with various long-tailed methods and demonstrate that the two modules can consistently boost their performances. So we suggest that FCM and LCM can be used as general modules for class imbalanced learning to enhance generalization ability.

Combination of “FCM and LCM” and “RBMC”. In DCRNets, FCM and LCM use a set of hyper-parameters, $\{\alpha_t\}_{t=1}^K$ and $\{\beta_t\}_{t=1}^K$, to control a stronger data augmentation for tail class as compared to head class, while RBMC uses a residual path to alleviate underfitting on head class. The optimal combination of “FCM and LCM” and “RBMC” should achieve the best trade-off performance on both head and tail classes. As shown in Eq. 3 and Eq. 10, as the values of α^0 and β^0 of FCM and LCM decrease, the difference between augmentation strength applied to head and tail classes reduces. Thus, α^0 and β^0 control the ability of FCM and LCM to improve tail accuracy while potentially suppressing head accuracy. So adjusting α^0 and β^0 can modulate the ability of FCM and LCM, enabling them to work collaboratively with RBMC to achieve the best trade-off performance on both head and tail classes.

Comparison to BBN [17]. Our framework DCRNets share similarity to BBN [17] of using a two-branch architecture. However, DCRNets differ from BBN in the following ways:

(1) **Different Network Architecture.** From **classifier view**, the classifiers of the two branches of BBN are independent of each other. Differently, DCRNets add a *residual connection* from uniform classifier to class-balanced classifier. With the residual connection, the class-balanced classifier can effectively utilize the prediction of uniform classifier, thereby better fitting the head data; From **feature extractor view**, BBN uses *separate* parameters for part of feature extractors of the two branches. Instead, DCRNets *share* the backbone parameters, which largely reduces computation cost in the inference phase.

(2) **Different Learning Mechanism.** BBN adopts *cumulative learning* mechanism while our DCRNets adopt *residual learning* mechanism. The cumulative learning of BBN gradually assigns a bigger weight to class-balanced branch. However, as pointed by [15], over-emphasizing on class balanced learning even only at the later training stage damages the generalization of feature representation. Differently, DCRNets use a residual connection to bridge uniform learning and class-balanced learning, and always gives a fixed and dominant weight to uniform learning branch, which ensures high-quality representation.

(3) **DCRNets add FCM and LCM modules.** BBN suffers from feature drift between training and test data, which leads to overfitting on tail classes. On the contrary, FCM and LCM in DCRNets can effectively alleviate the feature drift, thereby achieving a more generalizable classifier. In addition, FCM and LCM can be applied to various long-tailed learning methods including one-branch [15] and two-branch methods [17]. Therefore, different from BBN, our proposed method is not limited to a two-branch architecture.

Comparison to ResLT [58]. Our RBMC and ResLT [58] share a similar idea of using a residual mechanism for long-tailed recognition. However, they have following sufficient differences:

(1) **The implementation and function of residual learning are fundamentally different.** ResLT [58] designs three branches equipped with many+medium+tail, medium+tail and tail data respectively. It adds the residual connections from medium+tail/tail branch to main branch, which aims to enhance performance on medium and tail data. On the contrary, our RBMC designs two branches equipped with uniform and class-balanced sampling respectively. RBMC adds a residual connection from uniform branch to class-balanced branch, which aims to enhance performance on head data.

(2) **The architecture of residual branches is different.** ResLT conducts re-balance in the aspect of parameter space of feature extractor. So ResLT allocates individual parameters for feature extractor and uses a shared classifier on residual branches. In contrast, RBMC conducts re-balance in the aspect of parameter space of classifier. So RBMC uses a shared feature extractor and allocates individual parameters for classifier on residual branch.

(3) **The learning strategy of residual branches is different.** ResLT optimizes residual branches using different sub-group data, i.e., medium+tail and tail data, while RBMC optimizes residual branch with a class-balanced sampling strategy. Therefore, RBMC does not need to artificially divide many, medium and tail data, which is more flexible and easier to implement

Comparison to ISDA [77]. Recently, [77] proposes a similar semantic data augmentation algorithm ISDA. We summarize the main differences between ISDA and LCM.

(1) Problem Setting: ISDA mainly focuses on a *balanced* classification task while our work considers the problem of learning from long-tailed *imbalanced* dataset.

(2) Insight: ISDA aims to *augment* training samples to regularize deep networks. Differently, our LCM aims to better *calibrate* the training features to the corresponding test feature space, which aims to learn a more generalizable classifier for tail classes.

(3) Technique: In order to better address long-tailed problem, LCM has the following different designs from ISDA. **Firstly, LCM performs augmentation around the feature drift vectors instead of original features.** As shown in Equation 4 of [77], i.e., $\tilde{\mathbf{a}}_i \sim \mathcal{N}(\mathbf{a}_i, \lambda \Sigma_{y_i})$, ISDA performs augmentation around the *original features*. But on long-tailed recognition, due to the scarce samples of tail classes, there is a large drift between training and test features. So the augmented features around original training features would probably not lie in true (test) feature manifold, especially for tail classes. Differently, as shown in Eq. 8, LCM performs augmentation around the *feature drift vector* $\delta_{t,j}$, in which $\delta_{t,j}$ is crucial to bridge the feature drift of tail classes to generate more realistic tail features.

Secondly, LCM derives derivation based on Gaussian Mixture Distribution instead of Gaussian Distribution. As shown in Equation 7 of [77], ISDA derives the derivation based on a *Gaussian Distribution*. It is built on the assumption that the feature representation of each class follows a *uni-modal* Gaussian distribution. However, as pointed by [21], the feature distribution of test data of a tail class usually presents *multi-modal* [21]. To this end, LCM takes a step further and extends the derivation to *Gaussian Mixture Distribution* (Eq. 12). Built on Gaussian Mixture Distribution,

LCM is more conducive to capture the complex multi-modal distributions of tail classes.

Thirdly, LCM adapts Class Adaptive Translation instead of same augmentation for each class. As shown in Equation 4 of [77], i.e., $\tilde{\mathbf{a}}_i \sim \mathcal{N}(\mathbf{a}_i, \lambda \Sigma_{y_i})$, ISDA uses the same augmentation strength (λ) for all classes. However, on long-tailed recognition, the observed variance of tail class is usually very small [45]. Therefore, the augmented tail features of ISDA would be in closer vicinity to the original ones, thereby limiting the tail data variety and even aggravating imbalance issue. Differently, as shown in Eq. 8, LCM designs a *class adaptive translation coefficient* (β_t) to adaptively increase the augmentation strength for tail classes, thereby ensuring the diversity of generated tail features.

6 EXPERIMENTS

We mainly evaluate the proposed DCRNets on long-tailed classification task. We use five long-tailed benchmarks, i.e., ImageNet-LT [2], iNaturalist18 [78], CIFAR100-LT [79], CIFAR10-LT [79] and PLACE-LT [2]. We incorporate DCRNets with multiple network architectures (ResNet and ResNeXt), and a multi-expert framework and achieve the best model performance for all settings. We also evaluate DCRNets on Class Incremental Learning task on CIFAR100 [79] to show its generality.

6.1 Datasets and Evaluation Protocol

ImageNet-LT dataset. ImageNet-LT is proposed by [2]. ImageNet-LT is a long-tailed subset of ImageNet-2012 [80] by sampling a subset following the Pareto distribution with power value of 6. It contains 105.8K training images and 50,000 test images from 1,000 categories. The number of samples per class ranges from 5 to 1280.

iNaturalist18 dataset. iNaturalist18 [78] is one species classification dataset with images collected from real world, which is on a large scale and suffers from extremely imbalanced class distribution. It is composed of 437.5K training images and 2.4K test images from 8,142 categories. The number of samples per class ranges from 2 to 1,000. In addition to the extreme imbalance, iNaturalist18 also confronts the fine-grained problem.

CIFAR100-LT and CIFAR10-LT dataset. CIFAR10 and CIFAR100 dataset have 60,000 images, 50,000 for training and 10,000 for validation with 100 and 10 categories. Following [8], [17], we use a long-tailed version of CIFAR10/CIFAR100 dataset, which follows an exponential decay in sample sizes across different classes with various imbalance factors. The imbalance factor is defined as the ratio between the numbers of training samples for the most frequent class and the least frequent class. We use the imbalance factor of 100 in our experiments.

Places-LT dataset. Places-LT is proposed by [2]. Places-LT is a long-tailed version of the large-scale scene classification dataset Places [81]. It consists of 184.5K images from 365 categories. The number of samples per class ranges from 5 to 4980.

Evaluation Protocol. After training on long-tailed dataset, we evaluate the models on a balanced test set, and report the top-1 accuracy over all classes, denoted as "All".

Following [2], we further report accuracy on three splits of classes: *Many-shot* (more than 100 samples), *Medium-shot* (with 20 to 100 samples) and *Few-shot* (less than 20 samples).

6.2 Implementation Details

Implementation Details on ImageNet-LT. For ImageNet-LT, we use standard ResNet-10, ResNet-50 [1] and ResNeXt-50 [19] as backbone. We follow the same training strategy in [15]. In details, we use SGD optimizer with momentum 0.9 and train for 90 and 200 epochs respectively. We adopt cosine learning rate schedule gradually decaying from 0.075 to 0, and batch size 192 and 64 for uniform branch and class balanced branch respectively. We define the classes with more than 100 training samples as head classes, and the remaining classes as tail classes. In FCM, the number of selected head classes m is 2 (Eq. 1). The parameter α^0 is set to 0.5/0.5/1.0 for ResNet-10/ResNet-50/ResNeXt-50. In LCM, the parameter β^0 is set to 1.5/6.0/9.0 for ResNet-10/ResNet-50/ResNeXt-50. In RBMC, the number of proxies on each tail class (L) is set to 2 (Eq. 17). The weighting factor of loss function ϕ is set to 0.8.

Implementation Details on iNaturalist18. Following [76], we utilize ResNet-50 as backbone and train it for 90 and 200 epochs respectively. We set α^0 to 0.5, β^0 to 3.0. Other settings are the same as those in ImageNet-LT.

Implementation Details on CIFAR100-LT. For CIFAR-100 dataset, the implementation details mainly follow [7]. We use ResNet32 as backbone. The networks are trained for 500 epochs. The initial learning rate is set to 0.05, and the batch size of uniform branch and class balanced branch is set to 384 and 128 respectively. For the hyperparameter, α^0 and β^0 are set to 0.5 and 1.0 respectively. Other setting are the same as those in ImageNet-LT.

Implementation Details on Places-LT. For Places-LT, following previous setting [17], we choose ResNet-152 pre-trained on ImageNet dataset as backbone network, and train the model for 30 epochs. We set α^0 to 0.5, β^0 to 1.0. Other settings are the same as those in ImageNet-LT.

6.3 Comparisons with State-of-the-art Methods

In this section, we compare our DCRNets with other state-of-the-arts on ImageNet-LT, iNaturalist18, CIFAR100-LT and Places-LT benchmarks. Numerical results can be found in Tab. 1-5.

The compared methods cover various categories of ideas for imbalanced classification:

- 1) **Re-weighting (RW).** Re-weighting the loss function, such as Focal Loss [28], BALMS [7] and LDAM [8].
- 2) **Decoupled scheme (DE).** Decoupling the feature learning and classification, such as Two-stage [15], [32], [76], [82], CBS+RRS [83] and BBN [17].
- 3) **Data Generation (DG).** Generating new samples for tail categories, such as MetaSAug [53] and RSG [84].
- 4) **Transfer Learning (TL).** Transferring knowledge learned from head classes to tail classes, such as OLTR [2] and M2m [43].
- 5) **Representation Learning (RL).** Using self-supervised for better representation: SSP [36], Hybrid-PSC [38] and PaCo [39].

TABLE 1
Long-tailed classification performance (overall top-1 accuracy %) compared to related methods on **ImageNet-LT**. † indicates our reproduced results with the released code. R-10, R-50 and X-50 means ResNet-10, ResNet-50 and ResNeXt-50.

Epoch	Method		R-10	R-50	RX-50
90	RW	Focal [28]	30.5	-	43.7
		BALMS [7]	41.8	48.8†	51.4
	DE	τ -norm [15]	40.6	46.7	49.4
		cRT [15]	41.8	47.3	49.5
		CBS+RRS [83]	41.9	47.3	-
		LWS [15]	41.4	47.7	49.9
		LADE [76]	41.6†	50.8†	51.9
		TDE [32]	-	51.1	51.8
		DisAlign [82]	-	51.3	52.6
	DG	MetaSAug [53]	-	47.4	-
		RSG [84]	-	-	51.8
	TL	M2m [43]	-	43.7	-
OLTR [2]		37.3	-	46.3	
		DCRNets	43.8	53.2	54.8
> 180	DE	cRT [15]	42.7†	50.8†	51.2†
		LADE [76]	43.1†	52.3†	53.0
		LWS [15]	43.5†	51.7†	51.9†
	RL	SSP [36]	43.2	51.3	-
	ME	ResLT [58]	43.8	-	52.9
		DCRNets	45.0	54.1	55.0
100	ME	RIDE [20](2 experts)	45.3	54.4	55.9
		DCRNets (2 experts)	46.6	55.7	56.8
		ACE [57](3 experts)	44.0	54.7	56.8
		RIDE [20](3 experts)	45.9	54.9	56.4
		DCRNets (3 experts)	47.6	56.9	57.2

- 6) **Multi-Expert Network (ME).** Using an ensemble network to extract features or extra distillation training procedures: RIDE [20], ACE [57], ResLT [58] and NCL [59].

Main Results on ImageNet-LT. Tab. 1 reports the main results on ImageNet-LT. Our method achieves the best performance under the same setting. It is noted that: **(1)** Decoupled methods [15], [32], [76], [82] (DE) are prone to under-fitting the head and over-fitting the tail classes. It is noteworthy that DCRNets have outperformed all the decoupled methods by about 2% top-1 accuracy, which validates its superiority. **(2)** Data-generation methods [53], [84] (DG) and Transfer-learning methods [2], [43] (TL) use complex generation/transfer modules and explicitly generate features. By contrast, DCRNets implicitly generate features by two shift operations. Our method puts much less overhead with a much better performance: about 3% top-1 accuracy improvement. **(3)** Representation-learning method [36] (RL) employs a self-supervised learning to obtain a good feature initialization. DCRNets still surpass SSP [36] by 1.8% top-1 accuracy with ResNet10 as backbone and 2.8% top-1 accuracy with ResNet50 as backbone. Notably, as a model initialization technique, SSP is orthogonal to our method and can be easily combined to further improve performance. **(4)** ResLT [58] and our method both use a residual mechanism for long-tailed recognition. Our method surpasses ResLT by 1.2% and 2.1% top-1 accuracy with ResNet10 and ResNeXt50 as backbone respectively, validating its superiority.

TABLE 2
Comparison to related methods on **iNaturalist18**.

Method		90E	>180E
RW	Focal [28]	61.1	-
	LDAM [8] 1	64.6	-
	BALMS [7]	-	69.8
DE	cRT [15]	65.2	67.6
	τ -norm [15]	65.6	69.3
	LWS [15]	65.9	69.5
	BBN [17]	66.3	69.3
	LADE [76]	-	70.0
	DisAlign [82]	67.8	70.6
	LDAM+DRW [8]	68.0	-
DG	MetaSAug [53]	66.3	-
	RSG [84]	67.9	70.3
RL	SSP [36]	-	68.1
	Hybrid-PSC [38]	68.1	70.4
ME	ResLT [58]	-	70.2
	DCRNets	70.3	72.3
ME	RIDE [20] (2 experts)	71.4	-
	DCRNets (2 experts)	72.0	-
	RIDE [20] (3 experts)	72.2	-
	ACE [57] (3 experts)	72.9	-
	DCRNets (3 experts)	73.0	-

TABLE 3
Comparison to related methods on **CIFAR100-LT** and **CIFAR10-LT** with an imbalance of 100.

Method		CIF100-LT	CIF10-LT
RW	Focal [28]	38.4	70.4
	LDAM [8]	42.1	77.0
	BALMS [7]	50.8	84.9
DW	BBN [17]	42.6	79.8
	Logit Adj. [33]	43.9	77.7
	cRT [15]	50.0	82.0
	LWS [15]	50.5	83.7
DG	RSG [84]	44.6	79.6
	MetaSAug [53]	48.0	80.7
TL	OLTR [2]	41.2	-
	M2m [43]	43.5	79.1
RL	Hybrid-PSC [38]	44.9	78.8
	Hybrid-PSC [38]	46.7	81.4
ME	ResLT [58]	45.3	80.4
	DCRNets	51.4	85.0
ME	RIDE [20] (2 experts)	47.0	-
	DCRNets (2 experts)	55.0	85.8
	RIDE [20] (3 experts)	48.0	-
	ACE [57] (3 experts)	49.4	81.2
	DCRNets (3 experts)	56.0	87.0

Recently, Multi-Expert networks [20], [57] (ME) achieve state-of-the-art performance on long-tailed recognition. For fair comparison, we further adopt the multi-expert network of RIDE [20] as backbone. Notably, the self-distillation and EA module in RIDE are not used for simplification. As shown in Tab. 1, under the same multi-expert network, our method outperforms RIDE [20] and ACE [57] by 1% – 2% top-1 accuracy, which shows the proposed modules can well incorporate with multi-expert networks.

Main Results on iNaturalist18. Tab. 2 reports the main results on a real world benchmark, iNaturalist18. Our method also achieves the best accuracy. Compared to

TABLE 4
Comparison to related methods on **Places-LT**.

Method		ResNet-152
RW	Focal [28]	34.6
	BALMS [7]	38.6
	LADE [8]	38.8
DE	cRT [15]	36.7
	LWS [15]	37.6
	τ -norm [15]	37.9
	LADE [76]	38.8
DG	RSG [84]	39.3
TL	OLTR [2]	35.9
ME	ResLT [58]	39.8
	DCRNets	40.0

non-Multi-Expert methods, DCRNets outperform the best Hybrid-PSC [38] by about 2% top-1 accuracy. Under the multi-expert network, DCRNets still outperform RIDE [20] and ACE [57]. The improvements prove that our proposed method consistently improves performance on both artificial and real-world large-scale datasets.

Main Results on CIFAR100-LT and CIFAR10-LT. Tab. 3 reports the main results on CIFAR100-LT and CIFAR10-LT. Our method significantly outperforms the existing methods. By introducing the powerful compensation modules without using multi-expert, DCRNets are even competitive compared with multi-expert methods. Combing with multi-expert networks, DCRNets surpass RIDE and ACE by up to 6% top-1 accuracy. The superior results on the two small-scale datasets further validate the effectiveness and robustness of our method.

Main Results on Places-LT. Tab. 4 reports the main results on Places-LT dataset. DCRNets yields 40.0% top-1 accuracy, with a notable performance gain at 0.2% over the prior methods. Places-LT has an extremely imbalance class distribution, where the imbalance factor is up to 990. The superior results on this highly long-tailed benchmark further validate that our method can effectively address the class imbalance issue.

6.3.1 Comparison with PaCo and NCL.

Recently, PaCo [39] and NCL [59] achieve state-of-the-art performance on long-tailed recognition. PaCo [39] improves supervised contrastive loss by adding a set of parametric learnable class centers to tackle imbalance issue. PaCo aims to utilize self-supervised contrastive learning to improve *feature representation*. Differently, our method employs classification learning, which uses residual connection, feature shift and logit shift operations to learn a more generalizable *classifier*. So PaCo is complementary to our method, which can be easily combined to learn better feature representation; NCL [59] is built on a multi-expert framework, which learns multiple experts concurrently and uses knowledge distillation to enhance each single expert. Our method can be integrated into NCL framework, where we use DCRNets as each expert of NCL framework. To sum up, the techniques of PaCo and NCL are complementary to our method, which can be easily incorporated to further improve performance.

TABLE 5

Comparison with PaCo [39] and NCL [59]. We use ResNet-50 as backbone network for ImageNet-LT, and use the same training strategies as [39], [59]: training all the models with RandAugment [85] for 400 epochs except models on Places-LT, which is 30 epochs. "RA" denotes RandAugment, "DIS" denotes Distillation and "SS" denotes self-supervised learning.

	Method	RA	DIS	SS	Img-LT	iNat18	CIF100-LT	CIF10-LT	Places-LT
Single Model	PaCo [39]	✓	×	✓	57.0	73.2	52.0	-	41.2
	DCRNets	✓	×	×	57.6	74.0	52.0	85.1	41.5
	DCRNets+PaCo	✓	×	✓	58.0	74.2	52.2	85.5	41.7
Multiple Models	NCL [59] (single)	✓	✓	✓	57.4	74.2	53.3	84.7	41.5
	DCRNets (single)	✓	✓	×	59.0	74.8	54.0	85.9	41.6
	NCL [59] (ensemble)	✓	✓	✓	59.5	74.9	54.2	85.5	41.8
	DCRNets (ensemble)	✓	✓	×	60.3	75.7	56.0	87.0	42.0

In implementation, PaCo and NCL both use the strong data augmentation RandAugment [85], which is well known to be a practical strategy to improve performance. And NCL uses multi-expert and knowledge distillation to further improve classification accuracy. So for a fair comparison, we use the same training strategies as [39], [59], *i.e.*, training all the models with RandAugment [85] for 400 epochs except models on Places-LT, which is 30 epochs⁵. To compare with NCL, we also apply multi-expert and knowledge distillation strategies to our method. The results are listed in Tab. 5.

Compared to PaCo [39]. As shown in Tab. 5, our method can be well incorporated with strong augmentation. Equipped with RandAugment [85], our method achieves 57.6% and 74.0% top-1 accuracy on ImageNet-LT and iNaturalist18, outperforming PaCo by 0.6% and 0.8% respectively. Furthermore, we add PaCo loss to our method. As shown in Tab. 5, integrating with PaCo can further lift the performance by 0.2% – 0.4%, which shows our method can well combine with self-supervised contrastive PaCo.

Compared to NCL [59]. NCL adopts a multi-expert framework. NCL maintains three experts and uses a **distillation** training phase, which increases the training computational cost by almost three times. For a fair comparison, we also train three separate DCRNet experts and add a simple Kullback-Leibler (KL) distillation loss [60] among any two experts. Notably, the Hard Category Mining and Self-supervised contrastive loss in NCL are not used in our method. Tab. 5 reports the performance of a single expert and an ensemble of multiple experts following [59]. As seen, our method can be well incorporated with distillation technique to further improve accuracy. **Under the same multi-expert distillation framework**, our method outperforms NCL by 1.6%, 0.6%, 0.7%, 1.2% (using a single expert) and 0.8%, 0.8%, 1.8%, 1.5% (using an ensemble of experts) on imageNet-LT, iNaturalist18, CIFAR100-LT and CIFAR10-LT benchmarks respectively, and achieves comparable performance on Places-LT benchmark, showing the superiority of our method.

6.4 Ablation Study

To investigate the effectiveness of each component in DCRNets, we conduct a series of ablation studies on ImageNet-LT and iNaturalist18. All models are trained for 90 epochs

⁵. Notably, NCL [59] did not use RandAugment [85] on Places-LT benchmark. For a fair comparison on Places-LT, we reimplement NCL with RandAugment on Places-LT, and report this result on Tab. 5.

TABLE 6

Ablation study on ImageNet-LT and iNaturalist18 for 90 epochs training.

Method	RBMC	FCM	LCM	ImageNet-LT			iNat
				R10	R50	RX50	
cRT [15]				41.8	47.3	49.5	65.2
	✓			42.8	51.4	52.9	67.7
	✓	✓		43.4	52.3	54.0	69.2
DCRNets	✓	✓	✓	43.8	53.2	54.8	70.3

in this section. Tab. 6 and Tab. 7 summary the comparison results for different settings.

Two-stage Learning *v.s.* Joint Learning. In order to verify the effectiveness of RBMC, we provide three variants of RBMC: **RBMC-wo-R**, **RBMC-wo-MC** and **RBMC-wo-R&MC**. RBMC-wo-R uses multi-proxies classifiers without residual connection; RBMC-wo-MC uses linear classifiers with a residual connection; RBMC-wo-R&MC uses linear classifiers without residual connection. We first investigate the necessity of joint learning of feature extractor and classifier. To verify this, we choose the two-stage decoupled cRT [15] as a baseline. As shown in Tab. 7, all variants of RBMC achieve obviously superior performance compared to cRT. We argue that the two-stage training of cRT harms the compatibility between features and classifier, resulting in poor adaptability of classifier. Therefore, it is necessary to jointly optimize the feature extractor and classifier for class imbalanced learning.

Effectiveness of Residual Connection in RBMC. We further assess the effectiveness of the residual connection in RBMC. As shown in Tab. 7, compared to RBMC-wo-R&MC, employing residual mechanism (RBMC-wo-MC) brings 1.1% overall gain with negligible computational overhead. Also, compared to RBMC-wo-R, RBMC achieves 1.0% overall gain, which validates the capability of the residual connection. Specifically, employing residual connection brings above 3% gain on many-shot split. The reason is that RBMC-wo-R directly applies CBS to train the classifier of class balanced branch, which leads to under-fitting on head classes. On the contrary, RBMC introduces a residual path from uniform classifier. It ensures the class balanced branch can exploit the rich information of head classes encoded by uniform classifier. Such that the under-fitting issue can be alleviated and a better accuracy on many-shot split could be achieved. However, RBMC-wo-R achieves better performance on few-shot split. It is reasonable since the classifier of RBMC-wo-R puts more focus on few-shot data.

TABLE 7

Ablation study and complexity comparison on ImageNet-LT. ResNet-50 is used as backbone for training 90 epochs. PN: the number of parameters. GFLOPs: the number of floating-point operation of a forward pass for an input image.

Method	PN	GFlops		ALL	Many	Med.	Few
		Train	Test				
cRT [15]	25.6M	4.100	4.100	47.3	58.8	44.0	26.1
RBMC-wo-R&MC	27.6M	4.100	4.100	49.9	59.7	47.7	30.1
RBMC-wo-R	28.1M	4.101	4.101	50.4	59.9	48.0	31.6
RBMC-wo-MC	27.6M	4.102	4.102	51.0	62.9	48.1	27.3
RBMC	28.1M	4.103	4.103	51.4	63.3	48.5	28.5
+FCM- $m=1$	-	-	-	51.7	61.1	49.6	32.8
+FCM	28.1M	4.113	4.103	52.3	61.2	50.5	33.6
+FCM-LCM-fix	-	-	-	52.6	62.0	50.5	33.1
+FCM-LCM-SV	-	-	-	52.7	61.8	50.8	33.4
+FCM-LCM	28.1M	4.115	4.103	53.2	62.3	51.2	33.8

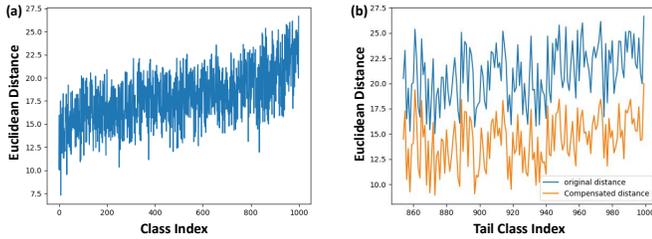


Fig. 7. (a) For each category, the average distance of all the test features to the closest original training feature of this category. (b) For each tail category, the average distance of all the test features to the closest original training feature (blue) / compensated feature by FCM (orange) of this category.

But the overall 1.0% accuracy gain manifests that RBMC derives a more balanced decision boundary which better trade-offs the head and tail classes.

Effectiveness of Multi-proxies Classifier in RBMC. We then investigate the influence of multi-proxies classifier (MC) in RBMC. As shown in Tab. 7, Compared to RBMC-wo-R&MC, employing MC (RBMC-wo-R) brings 0.5% overall gain. And compared to RBMC-wo-MC, RBMC achieves 0.4% overall gain with negligible computational overhead. The consistent improvements validate the effectiveness of MC. Specifically, MC achieves above 1.0% improvement on few-shot split, which indicates the robustness of MC on tail classes. Notably, RBMC shares a similar idea with Infinite Mixture Prototypes (IMP) [21] that represents each few-shot class by a set of clusters. However, IMP only considers the closest cluster of each class, while MC considers all clusters of a class. That allows MC to decrease the intra-class similarity when the training instances are insufficient, thereby alleviating the overfitting to the few training instances of tail classes.

Effectiveness of FCM. To show the separate influence of FCM and LCM, we provide a variant of DCRNets that adds FCM alone (+FCM). As shown in Tab. 7, FCM brings remarkable improvements on few-shot split by up to 5.1%. The significant improvements indicate that it is important to compensate the feature drift between training and test data on tail classes. Fig. 7 (a) visualizes the average distance of all test features of each category to the closest training feature

of this category. As shown, the distance goes large as the number of training instances decreases. The phenomenon reveals that the training and test features of a tail category do not occupy the same feature space. By contrast, as shown in Fig 7 (b), the compensated features of FCM are closer to the test features. As a consequence, FCM can make the classification boundary better adapt to the test data on tail classes.

We then investigate whether the *multi-mode* feature drift compensation is necessary. We explore a variant **FCM- $m=1$** that estimates a *uni-mode* feature drift direction for each tail category. It is equivalent to setting the number of selected head categories (m in Eq. 1) to 1. As shown in Tab. 7, FCM- $m=1$ leads to marginal improvements and FCM outperforms FCM- $m=1$ by 0.6%. The results verify the hypothesis that the features from one tail category could drift towards *multiple* head categories. Notably, [21] also observe the test data distribution of each few-shot category is clearly not uni-modal. So it is necessary for FCM to estimate a *multi-mode* feature drift vector to better fit the complex test distribution of tail category.

Effectiveness of LCM. Next, we assess the effectiveness of LCM. As shown in Tab. 7, LCM further achieves 0.9% overall top-1 gain over FCM. Also, we observe that LCM improves the accuracy of head and tail classes simultaneously. The improvements can be attributed to the effectiveness of integrating uncertainty into the estimation of feature drift directions. The uncertainty of LCM can prevent the subsequent classifier from over-fitting to FCM errors, thereby improving classifier’s generalization. Then, we investigate the influence of *Class Adaptive Translation Coefficient* (β_t in Eq. 8). We provide a variant **LCM-fix** that uses a fix translation coefficient for each class. LCM-fix is equivalent to setting β_t to 1. As shown in Tab. 7, LCM-fix brings negligible gain over FCM, which indicates the necessity of using a larger translation coefficient for the more minority category.

As discussed in Sec. 4.3, in order to deal with the smaller observed variance of a rarer class, previous work (SV) [45] uses a shared variance for all classes. Tab. 7 compares the proposed class adaptive strategy with SV. We observe that LCM outperforms **LCM-SV** by 0.5%. One possible explanation is that the unrelated classes, such as person and bag, generally have different intra-class variations. So for a specific tail class, the shared variance among all classes may lead to unreasonable translation directions, resulting in disturbed features. Differently, LCM uses the estimated variance from samples of the target class, which ensures the rationality of variation directions. In this way, the generated features still live in the true feature manifold, which can effectively benefit the learning of classifier.

Complexity Comparisons. To illustrate the cost of DCRNets, we report the number of parameters (PN) and the number of floating-point operations (GFLOPs). As shown in Tab. 7, RBMC only introduces 2.5M parameters compared to a naive linear classifier (cRT). In addition, RBMC requires 4.103 GFLOPs in a single forward pass for a 224×224 input image, corresponding to only 0.07% relative increase over cRT. FCM and LCM incurs 0.010 and 0.002 GFLOPs during training phase respectively. The increased

TABLE 8
 Combination FCM and LCM with various long-tailed methods on ImageNet-LT dataset. ResNet-50 is used as backbone trained for 90 epochs.

Method	FCM	LCM	ImageNet-LT
BBN [17]	✓		49.5
	✓	✓	51.9
BALMS [7]	✓		50.7
	✓	✓	53.0
LDAM [8]	✓		51.5
	✓	✓	53.3

cost of FCM mainly comes from the extra classification of generated compensated features. And the increased cost of LCM comes from the computation of logit compensation in Eq. 15. Both operations can be worked out by matrix multiplication thus occupy little training time in GPU libraries. During test phase, FCM and LCM are discarded and DCRNets only introduces negligible inference time.

Generalization of DCRNets across Different Backbones.

To investigate the effectiveness and generalization of DCRNets across different backbones, we also provide the ablation study based on ResNet-10 and ResNeXt-50 on ImageNet-LT dataset. As shown in Tab. 6, the proposed RBMC, FCM and LCM can consistently improve the results across different backbones. We also observe that: (1) the improvements of FCM and LCM are more obvious for deeper network. Specifically, the compensation modules, FCM and LCM, bring 1.0%/1.8%/1.9% gain with ResNet-10/ResNet-50/ResNeXt-50 as backbone. We argue that the deeper network has a higher risk of overfitting on tail classes, resulting in a larger feature drift between training and test data. Thus the compensation modules could be more effective on deeper networks. (2) The result of DCRNets based on ResNet-50 is much higher than the results of baseline cRT [15] based on ResNeXt-50 (+3.7%), where cRT based on ResNeXt-50 has more parameters and computational complexity. This comparison demonstrates that the improvement of DCRNets does not rely on the extra parameters and computational load.

Generalization of DCRNets across Different benchmarks.

We further provide the ablation study on iNaturalist18. As shown in Tab. 6, the proposed modules, RBMC, FCM and LCM, consistently improve the performance across different benchmarks. Notably, FCM and LCM bring 2.6% gain on iNaturalist18, which is higher than 1.8% gain on ImageNet-LT. The reason is that iNaturalist18 contains a large number of **fine-grained** classes. The tail classes would be more similar to head classes of the same superclass, resulting in a more serious feature drift. So FCM and LCM, which compensate the feature drift, could be more effective on the fine-grained iNaturalist18 benchmark.

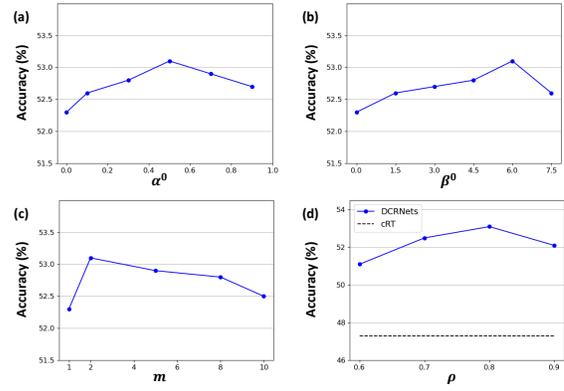


Fig. 8. Hyper-parameter analysis of DCRNets on ImageNet-LT. ResNet-50 is used as backbone for 90 epochs training. (a) FCM compensation strength parameter α^0 . (b) LCM translation strength parameter β^0 . (c) The number of selected head categories in FCM m . (d) Loss trade-off parameter ϕ .

6.5 Combination FCM/LCM with Various Methods

Most studies on class imbalanced learning attempt to deal with **label shift** between training and test phase. Existing works generally assume that the extracted features of training and test data follow the same distribution. In this work, we point out that this assumption does not hold on tail classes. And we propose FCM and LCM to alleviate the feature drift of tail classes between training and test phase. To demonstrate that addressing feature drift can consistently improve long-tailed performance, we combine FCM and LCM with several mainstream long-tailed methods, *i.e.*, BBN [17], BALMS [7] and LDAM [8], based on their open-source codes. As shown in Tab. 8, FCM and LCM consistently bring about 2% overall accuracy gains to various long-tailed methods, which validates the effectiveness and versatility of FCM and LCM. The consistent gains show that FCM and LCM can be used as general modules, which can combine with various long-tailed methods to enhance their feature generalization ability.

6.6 Parameter Analysis

The affects of hyper-parameters. Fig. 8 explores the hyper-parameter sensitivity of our method. Hyper-parameter α^0 and β^0 control the strength of feature drift compensation, m is the number of selected head categories for compensation and ϕ balances the two losses in DCRNets. We respectively evaluate each hyper-parameter, where we change its value and fix other hyper-parameters to the optimal values. Notably, we select the hyper-parameters based on the performance of the validation set. And the phenomenon *w.r.t* different hyper-parameters on validation set is consistent with that on test set. So we only present the performance on the test set in Fig. 8.

Firstly, we change α^0 from 0 to 1.0 and β^0 from 0 to 7.5. As shown in Fig. 8 (a) and (b), the performances of different $\alpha^0 > 0/\beta^0 > 0$ consistently outperform that of $\alpha^0 = 0/\beta^0 = 0$, which further verifies the effectiveness of FCM and LCM. In addition, the performances of different $\alpha^0 > 0$ and $\beta^0 > 0$ fluctuate in a small range ($\leq 0.5\%$),

TABLE 9

The performance of DCRNets with different batch sizes for uniform branch and class-balanced branch. ResNet-50 is used as backbone for training 90 epochs on ImageNet-LT. n_1/n_2 denote the batch size of uniform/class-balanced branch.

n_1/n_2	224/32	192/64	128/128	64/192
DCRNets	52.9	53.2	52.3	47.6

indicating the robustness of FCM and LCM to various α^0 and β^0 in some degree.

Secondly, we change m from 1 to 10. As shown in Fig. 8 (c), the performances of different $m > 1$ consistently outperform that of $m = 1$, which implies the effectiveness of the *multi-mode* drift vector in FCM. But there is an accuracy drop when m is too big. The main reason is that the tail features could compensate towards unrelated head categories if the number of selected head categories is too large.

Thirdly, we change ϕ from 0.6 to 0.9. As shown in Fig. 8 (d), the performances of different ϕ consistently outperform cRT baseline, and our method achieves the best performance when $\phi = 0.8$. Notice that there will be 2% accuracy drop when ϕ is too small ($\phi = 0.6$). The main reason is that the feature extractor would over-emphasis on class balanced learning when ϕ is too small, which damages the generalization of feature representations [15].

Overall, as shown in Fig. 8, the performances of DCRNets fluctuates within 2%. The phenomenons show that our method is not sensitive to the hyper-parameters in some degree, and can achieve competitive results under a wide range of hyper-parameters values.

The affects of batch sizes of two branches. We then explore how batch sizes of Uniform Branch and Class-Balanced Branch affect the performance. Denote the batch size of uniform branch and class-balanced branch as n_1 and n_2 respectively, Tab. 9 shows the performance of DCRNets with different n_1 and n_2 on ImageNet-LT. As shown, $n_1/n_2 = 192/64$ outperforms $n_1/n_2 = 128/128$ and $n_1/n_2 = 64/192$ by a large margin. We argue that using a too large n_2 would hurt feature representation. In particular, during each training iteration, a batch of n_1 examples and n_2 examples are sampled using uniform sampler and class-balanced sampler respectively. So using a too large n_2 would increase the risks of over-fitting the tail data (by over-sampling) and under-fitting the head data (by under-sampling). This could unexpectedly damage the representative ability of the learned deep features [15]. In addition, due to the small number of parameters in the classifier, a small n_2 is usually sufficient to learn a good balanced classifier. Therefore, it is reasonable to set $n_1 > n_2$. Tab. 9 also shows that $n_1/n_2 = 192/64$ achieves marginally better result than $n_1/n_2 = 224/32$. We argue that using a *too small* n_2 is inadequate to learn a class-balanced classifier. Therefore, we set $n_2 = n_1/3$ on all benchmarks, which is capable of simultaneously learning a good feature representation and class-balanced classifier.

6.7 Evaluation on Class Incremental Learning

For Class Incremental Learning (CInCL), an underlying problem is that the ratio of the number of new samples to

TABLE 10

Class incremental learning performance (top-1 accuracy %) on CIFAR100 with 2, 5 and 10 incremental steps. The average results over all incremental steps except the first step are reported. **CInCL** stands for Class Incremental Learning methods and **CImbL** indicates Class Imbalanced Learning inspired methods. We produce the results using their public code.

	Method	#incremental steps		
		2	5	10
CInCL	LwF [64]	52.6	47.1	39.7
	iCaRL [63]	62.0	63.3	61.6
	EEIL [69]	60.8	63.7	63.6
	IL2M [70]	65.0	65.5	63.7
	BiC [71]	65.3	66.0	63.7
	MDFCIL [73]	66.5	67.2	64.7
CImbL	τ -norm [15]	65.7	66.5	64.4
	cRT [15]	66.2	66.8	65.5
	LWS [15]	66.6	66.5	63.5
	LADE [76]	66.6	67.4	64.4
	DCRNets	68.2	68.5	67.4

that of old samples (preserved exemplars) could be very high, resulting in severe class imbalance. Thus DCRNets that address class imbalance can be leveraged for CInCL. In this section, we additionally evaluate DCRNets on CInCL.

Datasets and Implementations. We conduct CInCL experiments on CIFAR100 [79] following a closely related work [73]. CIFAR100 contains 60,000 samples of 32×32 images for 100 classes. Following the exact setting in [73], we deploy a 32-layer ResNet [1] as the baseline architecture for CIFAR100. We use SGD to train our model and set the batch size for uniform branch and class balanced branch to 128 and 32 respectively. The learning rate starts from 0.1 and reduces to 0.1 of the previous learning rate after 120, 180, 240 epochs (250 epochs in total). The hyper-parameters m , α^0 and β^0 are set to 2, 0.5 and 0.5 respectively. The weighting factor of loss function ϕ is set to 0.9.

Memory Budge. We follow the same data reply setting used in [73]. We store 2,000 samples for old classes, and select rehearsal exemplars based on herding selection [86]. More classes have been seen, fewer images can be retained per class. As a result, the problem of class imbalance becomes more serious.

Benchmark Protocol. We follow the common protocol used in [73]. On CIFAR100 benchmark, 100 classes are evenly split into 2/5/10 incremental steps. In each step, the model is evaluated on the test data for all seen classes. Since the first step is not related to class incremental learning actually, we report the average top-1 accuracy over all incremental phases except the first step.

Comparison to Other Methods. We compare DCRNets with several competitive or representative methods, including Class Incremental Learning methods (CInCL) and Class Imbalanced Learning inspired methods (CImbL). As for CInCL, we compare LwF [64], iCaRL [63], EEIL [69], IL2M [70], BiC [71] and MDFCIL [73]. Note that **EEIL**, **BiC**, **IL2M** and **MDFCIL** attempt to address the class imbalance problem in CInCL. As for CImbL-inspired methods, we compare cRT [15], τ -norm [15], LWS [15] and LADE [76]. We produce the results using their public codes. The compared results are shown in Tab. 10.

As for CInL methods, LwF [64] achieves much lower performance than others. It is reasonable since it incorporates no imbalance technique, making it vulnerable to class imbalance. EEIL [69], BiC [71], IL2M [70] and MDFCIL [73] propose different bias correction techniques to boost the performance. MDFCIL [73] is the most effective one. In MDFCIL, by normalizing the classifier weights to have similar norms, the bias induced by different numbers of samples can be removed. As for CImbL methods, cRT [15], τ -norm [15], LWS [15] and LADE [76] add a class-balanced step to correct the outputs, which achieve similar performance with MDFCIL. Our DCRNets consistently outperform the compared methods. Specifically, the overall performance on 10 incremental steps is improved by 2.7% accuracy compared to MDFCIL, and 1.9% compared to cRT. The reason is that all above methods assume that the extracted features of training and test data follow a same distribution, which is not hold on imbalanced data of CInL. Our method effectively reduces the unfavorable feature drift between training and test phases to better address imbalance issue, thereby benefiting incremental learning.

Overall, Tab. 10 indicates that the techniques in CImbL can be applied in CInL as well. And our approach is more effective to handle class imbalance in CInL. Our DCRNets can achieve better results compared to state-of-the-art approaches under different incremental settings.

7 CONCLUSION

We propose a framework, DCRNets, for class imbalanced learning. Firstly, we observe that there is a severe feature drift between training and test data, especially on tail class. So we design FCM and LCM to estimate the feature drift and compensate for it. Our study suggests that feature compensation is important to achieve a generalizable classifier on imbalanced data. Secondly, we observe that CBS essentially leads to underfitting on head classes. So RBMC is designed to add a residual path from Uniform Branch to facilitate classifier learning. Our study suggests that uniform learning is also useful for learning a class-balanced classifier. Our overall framework is *generic* and can be easily incorporated into existing Class Imbalanced methods to boost their performance. We also evaluate our approach on Class Incremental Learning to show its universality.

8 FURTHER WORK

In this part, we briefly discuss the feasibility of our approach on other tasks. We leave the specific application of our proposed method to other tasks as a further work.

Vehicle/Person ReID. This work mainly explores deep classification task. We argue our method can be also used to Vehicle/Person reID (*metric learning* task) when the identities present an imbalanced distribution.

Firstly, the overall two-branch framework is applicable for reID. The recently most successful approaches [87] show that the metric loss [87], [88] is superior to reID. However, on imbalanced setting, metric loss would concentrate more on head identities than tail ones since the comparison chances for tail identities reduce quadratically. So the metric training process is dominated by head identities, which could distort

overall feature space [45]. In our framework, the metric loss can be added on the class-balanced branch to perform balanced metric learning, to learn better feature representation for reID.

Secondly, our FCM and LCM which compensate tail identities with higher diversity are applicable for reID. As pointed by [45], on reID with imbalanced data, the head identity often occupies a large spatial span, while the tail identity often occupies a very small spatial span on feature space due to lack of intra-class diversity. This uneven distribution distorts overall feature space, consequentially compromising the discriminative ability of learned features. In our framework, FCM translates the tail features along meaningful nearest-neighbor directions, and LCM augments tail features along intra-class variations more strongly. Thus FCM and LCM can provide a larger space span with higher diversity for tail identities, thereby alleviating the distortion of feature space and improving representation learning on reID with imbalanced identities.

Unsupervised Learning. *Unsupervised Clustering Algorithm* [89], [90] is an important line of unsupervised learning. Unsupervised Clustering usually alternately learns the network parameters and the cluster assignments (pseudo-labels) for unlabeled data. We argue that our method could be applied to unsupervised clustering methods in the following two ways: **(1) Our two-branch framework could compensate for imbalanced clusters.** [89] observes that clustering algorithm, *e.g.*, k-means, is vulnerable to imbalanced clusters (pseudo-labels), where vast majority of examples are assigned to a few clusters. For unsupervised learning, our two-branch framework could use *pseudo-labels* to perform uniform and class-balanced sampling. In this way, the class-balanced branch could compensate for imbalanced clusters. **(2) Our FCM and LCM could reduce clustering errors.** [91], [92] observe that the clustering pseudo-labels are inherently noisy, *e.g.*, the samples sharing the same true class could be split into two or more clusters. The key of FCM and LCM is to synthesize new samples along the directions to its nearest neighbors. Therefore, FCM and LCM can be used to synthesize intermediate samples between two neighbor sub-clusters with the same true class. The synthesized samples can help correctly merge the two sub-clusters into one cluster. Therefore, we argue that FCM and LCM are able to reduce clustering errors and thus improve unsupervised learning.

Unsupervised Domain Adaptation (UDA). The key problem of UDA is the feature drift between source-domain and target-domain data. Our FCM and LCM are designed to deal with feature drift between training and test data, which can also address the feature drift issue on UDA. That is, FCM and LCM can translate the source features towards the target domain, by setting the feature drift vector of class c (δ_c) as the distance of class- c prototype between source data and target data. By adding the feature drift vector to original source features, the new generated source features will be closer to the target domain, which is conducive to train a more transferable classifier.

Continual UDA. The work [93] proposed Continual UDA where unlabeled target-domain samples are received in small batches and adaptation is performed continually. The target-domain data is collected *randomly* and *incrementally*.

tally in *small* batches, which easily causes the observed data to be class-imbalanced, especially in the early continual learning stage. Therefore, our framework can also be applied to Continual UDA to address its class imbalance issue.

ACKNOWLEDGMENTS

This work is partially supported by National Key R&D Program of China (No. 2018AAA0102402) and Natural Science Foundation of China (NSFC): 61976203, 62276246 and U19B2036, and in part by the National Postdoctoral Program for Innovative Talents under Grant BX20220310.

REFERENCES

- [1] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778, 2016.
- [2] Z. Liu, Z. Miao, X. Zhan, J. Wang, B. Gong, and S. X. Yu, "Large-scale long-tailed recognition in an open world," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2537–2546, 2019.
- [3] R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L.-J. Li, D. A. Shamma, et al., "Visual genome: Connecting language and vision using crowdsourced dense image annotations," *International Journal of Computer Vision*, vol. 123, no. 1, pp. 32–73, 2017.
- [4] P. D. Agrim Gupta and R. Girshick, "Lvis: A dataset for large vocabulary instance segmentation," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5356–5364, 2019.
- [5] M. Buda, A. Maki, and M. A. Mazurowski, "A systematic study of the class imbalance problem in convolutional neural networks," *Neural Networks*, vol. 106, pp. 249–259, 2018.
- [6] S. H. Khan, M. Hayat, M. Bennamoun, F. A. Sohel, and R. Togneri, "Cost-sensitive learning of deep feature representations from imbalanced data," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 29, no. 8, pp. 3573–3587, 2017.
- [7] R. Jiawei, C. Yu, X. Ma, H. Zhao, S. Yi, et al., "Balanced meta-softmax for long-tailed visual recognition," in *Advances in Neural Information Processing Systems*, pp. 4175–4186, 2020.
- [8] K. Cao, C. Wei, A. Gaidon, N. Arechiga, and T. Ma, "Learning imbalanced datasets with label-distribution-aware margin loss," in *Advances in Neural Information Processing Systems*, 2019.
- [9] C. Drummond, R. C. Holte, et al., "C4. 5, class imbalance, and cost sensitivity: why under-sampling beats over-sampling," in *Workshop on Learning from Imbalanced Datasets II*, vol. 11, pp. 1–8, Citeseer, 2003.
- [10] A. Menon, H. Narasimhan, S. Agarwal, and S. Chawla, "On the statistical consistency of algorithms for binary classification under class imbalance," in *International Conference on Machine Learning*, pp. 603–611, 2013.
- [11] H. Han, W.-Y. Wang, and B.-H. Mao, "Borderline-smote: a new over-sampling method in imbalanced data sets learning," in *International Conference on Intelligent Computing*, pp. 878–887, Springer, 2005.
- [12] J. Byrd and Z. Lipton, "What is the effect of importance weighting in deep learning," in *International Conference on Machine Learning*, pp. 872–881, 2019.
- [13] B. Krawczyk, "Learning from imbalanced data: open challenges and future directions," *Progress in Artificial Intelligence*, vol. 5, no. 4, pp. 221–232, 2016.
- [14] Y. Cui, M. Jia, T.-Y. Lin, Y. Song, and S. Belongie, "Class-balanced loss based on effective number of samples," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 9268–9277, 2019.
- [15] B. Kang, S. Xie, M. Rohrbach, Z. Yan, A. Gordo, J. Feng, and Y. Kalantidis, "Decoupling representation and classifier for long-tailed recognition," in *International Conference on Learning Representations*, 2020.
- [16] L. Van der Maaten and G. Hinton, "Visualizing data using t-sne.," *Journal of Machine Learning Research*, vol. 9, no. 11, 2008.
- [17] B. Zhou, Q. Cui, X.-S. Wei, and Z.-M. Chen, "BBN: Bilateral-branch network with cumulative learning for long-tailed visual recognition," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 9719–9728, 2020.
- [18] C. Finn, P. Abbeel, and S. Levine, "Model-agnostic meta-learning for fast adaptation of deep networks," in *International Conference on Machine Learning*, pp. 1126–1135, 2017.
- [19] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, "Aggregated residual transformations for deep neural networks," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1492–1500, 2017.
- [20] X. Wang, L. Lian, Z. Miao, Z. Liu, and S. X. Yu, "Long-tailed recognition by routing diverse distribution-aware experts," *arXiv preprint arXiv:2010.01809*, 2020.
- [21] K. Allen, E. Shelhamer, H. Shin, and J. Tenenbaum, "Infinite mixture prototypes for few-shot learning," in *International Conference on Machine Learning*, pp. 232–241, 2019.
- [22] G. Karakoulas and J. Shawe-Taylor, "Optimizing classifiers for imbalanced training sets," in *Advances in Neural Information Processing Systems*, 1999.
- [23] L. Shen, Z. Lin, and Q. Huang, "Relay backpropagation for effective learning of deep convolutional neural networks," in *European Conference on Computer Vision*, pp. 467–482, 2016.
- [24] Y. Zhong, W. Deng, M. Wang, J. Hu, J. Peng, X. Tao, and Y. Huang, "Unequal-training for deep face recognition with long-tailed noisy data," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7812–7821, 2019.
- [25] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer, "Smote: synthetic minority over-sampling technique," *Journal of Artificial Intelligence Research*, vol. 16, pp. 321–357, 2002.
- [26] C. Huang, Y. Li, C. C. Loy, and X. Tang, "Learning deep representation for imbalanced classification," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5375–5384, 2016.
- [27] C. Huang, Y. Li, C. C. Loy, and X. Tang, "Deep imbalanced learning for face recognition and attribute prediction," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 42, no. 11, pp. 2781–2794, 2019.
- [28] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2980–2988, 2017.
- [29] 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 *Advances in Neural Information Processing Systems*, 2019.
- [30] M. Ren, W. Zeng, B. Yang, and R. Urtasun, "Learning to reweight examples for robust deep learning," in *International Conference on Machine Learning*, pp. 4334–4343, 2018.
- [31] S. Khan, M. Hayat, S. W. Zamir, J. Shen, and L. Shao, "Striking the right balance with uncertainty," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 103–112, 2019.
- [32] K. Tang, J. Huang, and H. Zhang, "Long-tailed classification by keeping the good and removing the bad momentum causal effect," in *Advances in Neural Information Processing Systems*, pp. 1513–1524, 2020.
- [33] A. K. Menon, S. Jayasumana, A. S. Rawat, H. Jain, A. Veit, and S. Kumar, "Long-tail learning via logit adjustment," in *International Conference on Learning Representations*, 2021.
- [34] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," *arXiv preprint arXiv:1911.05722*, 2019.
- [35] X. Chen, H. Fan, R. Girshick, and K. He, "Improved baselines with momentum contrastive learning," *arXiv preprint arXiv:2003.04297*, 2020.
- [36] Y. Yang and Z. Xu, "Rethinking the value of labels for improving class-imbalanced learning," in *Advances in Neural Information Processing Systems*, pp. 19290–19301, 2020.
- [37] B. Kang, Y. Li, Z. Yuan, and J. Feng, "Exploring balanced feature spaces for representation learning," in *International Conference on Learning Representations*, 2021.
- [38] P. Wang, K. Han, X.-S. Wei, L. Zhang, and L. Wang, "Contrastive learning based hybrid networks for long-tailed image classification," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2021.
- [39] J. Cui, Z. Zhong, S. Liu, B. Yu, and J. Jia, "Parametric contrastive learning," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 715–724, 2021.
- [40] J. Cui, Z. Zhong, Z. Tian, S. Liu, B. Yu, and J. Jia, "Generalized parametric contrastive learning," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- [41] L. Zhu and Y. Yang, "Inflated episodic memory with region self-attention for long-tailed visual recognition," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4344–4353, 2020.

- [42] Y.-X. Wang, D. Ramanan, and M. Hebert, "Learning to model the tail," in *Advances in Neural Information Processing Systems*, pp. 7032–7042, 2017.
- [43] J. Kim, J. Jeong, and J. Shin, "M2m: Imbalanced classification via major-to-minor translation," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 13896–13905, 2020.
- [44] P. Chu, X. Bian, S. Liu, and H. Ling, "Feature space augmentation for long-tailed data," in *European Conference on Computer Vision*, pp. 694–710, Springer, 2020.
- [45] J. Liu, Y. Sun, C. Han, Z. Dou, and W. Li, "Deep representation learning on long-tailed data: A learnable embedding augmentation perspective," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2970–2979, 2020.
- [46] X. Yin, X. Yu, K. Sohn, X. Liu, and M. Chandraker, "Feature transfer learning for face recognition with under-represented data," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5704–5713, 2019.
- [47] B. Hariharan and R. Girshick, "Low-shot visual recognition by shrinking and hallucinating features," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3018–3027, 2017.
- [48] Y. Zang, C. Huang, and C. C. Loy, "Fasa: Feature augmentation and sampling adaptation for long-tailed instance segmentation," *arXiv preprint arXiv:2102.12867*, 2021.
- [49] S. Yang, L. Liu, and M. Xu, "Free lunch for few-shot learning: Distribution calibration," in *International Conference on Learning Representations*, 2021.
- [50] Y. Yang, K. Zha, Y.-C. Chen, H. Wang, and D. Katabi, "Delving into deep imbalanced regression," in *International Conference on Machine Learning*, 2021.
- [51] Y.-X. Wang, R. Girshick, M. Hebert, and B. Hariharan, "Low-shot learning from imaginary data," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7278–7286, 2018.
- [52] S. S. Mullick, S. Datta, and S. Das, "Generative adversarial minority oversampling," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1695–1704, 2019.
- [53] S. Li, K. Gong, C. H. Liu, Y. Wang, F. Qiao, and X. Cheng, "Metasaug: Meta semantic augmentation for long-tailed visual recognition," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 5212–5221, 2021.
- [54] S. Li, M. Xie, K. Gong, C. H. Liu, Y. Wang, and W. Li, "Transferable semantic augmentation for domain adaptation," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 11516–11525, 2021.
- [55] Y. Wang, X. Pan, S. Song, H. Zhang, G. Huang, and C. Wu, "Implicit semantic data augmentation for deep networks," in *Advances in Neural Information Processing Systems*, pp. 12635–12644, 2019.
- [56] L. Xiang, G. Ding, and J. Han, "Learning from multiple experts: Self-paced knowledge distillation for long-tailed classification," in *European Conference on Computer Vision*, pp. 247–263, Springer, 2020.
- [57] J. Cai, Y. Wang, and J.-N. Hwang, "Ace: Ally complementary experts for solving long-tailed recognition in one-shot," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 112–121, 2021.
- [58] J. Cui, S. Liu, Z. Tian, Z. Zhong, and J. Jia, "Reslt: Residual learning for long-tailed recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- [59] J. Li, Z. Tan, J. Wan, Z. Lei, and G. Guo, "Nested collaborative learning for long-tailed visual recognition," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 6949–6958, 2022.
- [60] Y. Zhang, T. Xiang, T. M. Hospedales, and H. Lu, "Deep mutual learning," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4320–4328, 2018.
- [61] Y. Zhang, B. Hooi, L. Hong, and J. Feng, "Test-agnostic long-tailed recognition by test-time aggregating diverse experts with self-supervision," *arXiv preprint arXiv:2107.09249*, 2021.
- [62] M. Delange, R. Aljundi, M. Masana, S. Parisot, X. Jia, A. Leonardis, G. Slabaugh, and T. Tuytelaars, "A continual learning survey: Defying forgetting in classification tasks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [63] S.-A. Rebuffi, A. Kolesnikov, G. Sperl, and C. H. Lampert, "icarl: Incremental classifier and representation learning," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2001–2010, 2017.
- [64] Z. Li and D. Hoiem, "Learning without forgetting," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 12, pp. 2935–2947, 2017.
- [65] A. Douillard, M. Cord, C. Ollion, T. Robert, and E. Valle, "Podnet: Pooled outputs distillation for small-tasks incremental learning," in *European Conference on Computer Vision*, pp. 86–102, 2020.
- [66] X. Tao, X. Chang, X. Hong, X. Wei, and Y. Gong, "Topology-preserving class-incremental learning," in *European Conference on Computer Vision*, pp. 254–270, Springer, 2020.
- [67] Y. Liu, Y. Su, A.-A. Liu, B. Schiele, and Q. Sun, "Mnemonics training: Multi-class incremental learning without forgetting," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 12245–12254, 2020.
- [68] C. He, R. Wang, and X. Chen, "A tale of two cils: The connections between class incremental learning and class imbalanced learning, and beyond," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3559–3569, 2021.
- [69] F. M. Castro, M. J. Marín-Jiménez, N. Guil, C. Schmid, and K. Alahari, "End-to-end incremental learning," in *European Conference on Computer Vision*, pp. 233–248, 2018.
- [70] E. Belouadah and A. Popescu, "Il2m: Class incremental learning with dual memory," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 583–592, 2019.
- [71] Y. Wu, Y. Chen, L. Wang, Y. Ye, Z. Liu, Y. Guo, and Y. Fu, "Large scale incremental learning," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 374–382, 2019.
- [72] S. Hou, X. Pan, C. C. Loy, Z. Wang, and D. Lin, "Learning a unified classifier incrementally via rebalancing," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 831–839, 2019.
- [73] B. Zhao, X. Xiao, G. Gan, B. Zhang, and S.-T. Xia, "Maintaining discrimination and fairness in class incremental learning," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 13208–13217, 2020.
- [74] J. Snell, K. Swersky, and R. Zemel, "Prototypical networks for few-shot learning," in *Advances in Neural Information Processing Systems*, pp. 4077–4087, 2017.
- [75] Y. Movshovitz-Attias, A. Toshev, T. K. Leung, S. Ioffe, and S. Singh, "No fuss distance metric learning using proxies," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 360–368, 2017.
- [76] Y. Hong, S. Han, K. Choi, S. Seo, B. Kim, and B. Chang, "Disentangling label distribution for long-tailed visual recognition," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2021.
- [77] Y. Wang, G. Huang, S. Song, X. Pan, Y. Xia, and C. Wu, "Regularizing deep networks with semantic data augmentation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [78] 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 *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 8769–8778, 2018.
- [79] A. Krizhevsky et al., "Learning multiple layers of features from tiny images," 2009.
- [80] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 248–255, 2009.
- [81] B. Zhou, A. Lapedriza, J. Xiao, A. Torralba, and A. Oliva, "Learning deep features for scene recognition using places database," in *Advances in Neural Information Processing Systems*, 2014.
- [82] S. Zhang, Z. Li, S. Yan, X. He, and J. Sun, "Distribution alignment: A unified framework for long-tail visual recognition," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2021.
- [83] J. Zhang, L. Liu, P. Wang, and C. Shen, "To balance or not to balance: A simple-yet-effective approach for learning with long-tailed distributions," in *European Conference on Computer Vision*, 2020.
- [84] J. Wang, T. Lukasiewicz, X. Hu, J. Cai, and Z. Xu, "Rsg: A simple but effective module for learning imbalanced datasets," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3784–3793, 2021.
- [85] E. D. Cubuk, B. Zoph, J. Shlens, and Q. V. Le, "Randaugment: Practical automated data augmentation with a reduced search space," in *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 702–703, 2020.
- [86] M. Welling, "Herding dynamical weights to learn," in *International Conference on Machine Learning*, pp. 1121–1128, 2009.
- [87] Y. Sun, C. Cheng, Y. Zhang, C. Zhang, L. Zheng, Z. Wang, and Y. Wei, "Circle loss: A unified perspective of pair similarity

optimization," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 6398–6407, 2020.

- [88] A. Hermans, L. Beyer, and B. Leibe, "In defense of the triplet loss for person re-identification," *arXiv preprint arXiv:1703.07737*, 2017.
- [89] M. Caron, P. Bojanowski, A. Joulin, and M. Douze, "Deep clustering for unsupervised learning of visual features," in *European Conference on Computer Vision*, pp. 132–149, 2018.
- [90] E. Amrani, L. Karlinsky, and A. Bronstein, "Self-supervised classification network," in *European Conference on Computer Vision*, pp. 116–132, Springer, 2022.
- [91] Y. Cho, W. J. Kim, S. Hong, and S.-E. Yoon, "Part-based pseudo label refinement for unsupervised person re-identification," in *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7308–7318, 2022.
- [92] Y. Ge, F. Zhu, D. Chen, R. Zhao, *et al.*, "Self-paced contrastive learning with hybrid memory for domain adaptive object re-id," in *Advances in Neural Information Processing Systems*, pp. 11309–11321, 2020.
- [93] A. M. N. Taufique, C. S. Jahan, and A. Savakis, "Conda: Continual unsupervised domain adaptation," *arXiv preprint arXiv:2103.11056*, 2021.



Ruibing Hou received the BS degree in Northwestern Polytechnical University, Xi'an, China, in 2016. She received PhD degree in computer science from the Institute of Computing Technology, Chinese Academy of Sciences, Beijing, China, in 2022. She is currently a post-doctorial researcher with the Institute of Computing Technology, Chinese Academy of Sciences. Her research interests are in machine learning and computer vision. She specially focuses on person re-identification, long-tailed learning and

few-shot learning.



Hong Chang received the Bachelor's degree from Hebei University of Technology, Tianjin, China, in 1998; the M.S. degree from Tianjin University, Tianjin, in 2001; and the Ph.D. degree from Hong Kong University of Science and Technology, Kowloon, Hong Kong, in 2006, all in computer science. She was a Research Scientist with Xerox Research Centre Europe. She is currently a Researcher with the Institute of Computing Technology, Chinese Academy of Sciences, Beijing, China. Her main research in-

terests include algorithms and models in machine learning, and their applications in pattern recognition and computer vision.



Bingpeng Ma received the BS degree in mechanics, in 1998 and the MS degree in mathematics, in 2003 from the Huazhong University of Science and Technology, respectively. He received the PhD degree in computer science from the Institute of Computing Technology, Chinese Academy of Sciences, P.R. China, in 2009. He was a post-doctorial researcher with the University of Caen, France, from 2011 to 2012. He joined the School of Computer Science and Technology, University of Chinese Academy of

Sciences, Beijing, in March 2013 and now he is a professor. His research interests cover computer vision, pattern recognition, and machine learning. He especially focuses on person re-identification, face recognition, and the related research topics.



Shiguang Shan (M'04-SM'15) received Ph.D. degree in computer science from the Institute of Computing Technology (ICT), Chinese Academy of Sciences (CAS), Beijing, China, in 2004. He has been a full Professor of this institute since 2010 and now the deputy director of CAS Key Lab of Intelligent Information Processing. His research interests cover computer vision, pattern recognition, and machine learning. He has published more than 300 papers, with totally more than 20,000 Google scholar citations. He served

as Area Chairs for many international conferences including CVPR, ICCV, AAAI, IJCAI, ACCV, ICPR, FG, etc. And he was/is Associate Editors of several journals including IEEE T-IP, Neurocomputing, CVIU, and PRL. He was a recipient of the China's State Natural Science Award in 2015, and the China's State S&T Progress Award in 2005 for his research work.



Xilin Chen is a professor with the Institute of Computing Technology, Chinese Academy of Sciences (CAS). He has authored one book and more than 400 papers in refereed journals and proceedings in the areas of computer vision, pattern recognition, image processing, and multimodal interfaces. He is currently an information sciences editorial board member of Fundamental Research, an editorial board member of Research, a senior editor of the Journal of Visual Communication and Image Representation, and

an associate editor-in-chief of the Chinese Journal of Computers, and Chinese Journal of Pattern Recognition and Artificial Intelligence. He served as an organizing committee member for multiple conferences, including general co-chair of FG 2013 / FG 2018, VCIP 2022, program co-chair of ICMI 2010 / FG 024, and an Area Chair / Senior PC of ICCV / CVPR / ECCV / ICMI for more than 10 times. He is a fellow of the ACM, IEEE, IAPR, and CCF.