# Deep Long-Tailed Learning: A Survey

Yifan Zhang, Bingyi Kang, Bryan Hooi, Shuicheng Yan, *Fellow, IEEE*, and Jiashi Feng

**Abstract**—Deep long-tailed learning, one of the most challenging problems in visual recognition, aims to train well-performing deep models from a large number of images that follow a long-tailed class distribution. In the last decade, deep learning has emerged as a powerful recognition model for learning high-quality image representations and has led to remarkable breakthroughs in generic visual recognition. However, long-tailed class imbalance, a common problem in practical visual recognition tasks, often limits the practicality of deep network based recognition models in real-world applications, since they can be easily biased towards dominant classes and perform poorly on tail classes. To address this problem, a large number of studies have been conducted in recent years, making promising progress in the field of deep long-tailed learning. Considering the rapid evolution of this field, this paper aims to provide a comprehensive survey on recent advances in deep long-tailed learning. To be specific, we group existing deep long-tailed learning studies into three main categories (*i.e.,* class re-balancing, information augmentation and module improvement), and review these methods following this taxonomy in detail. Afterward, we empirically analyze several state-of-the-art methods by evaluating to what extent they address the issue of class imbalance via a newly proposed evaluation metric, *i.e.,* relative accuracy. We conclude the survey by highlighting important applications of deep long-tailed learning and identifying several promising directions for future research.

**Index Terms**—Long-tailed Learning, Deep Learning, Imbalanced Learning, Convolutional Neural Networks

✦

## 1 INTRODUCTION

DEEP learning allows computational models, composed of multiple processing layers, to learn data representations with multiple levels of abstraction [1], [2] and has made incredible progress in computer vision [3], [4], [5], [6], [7], [8]. The key enablers of deep learning are the availability of large-scale datasets, the emergence of GPUs, and the advancement of deep network architectures [9]. Thanks to the strong ability of learning high-quality data representations, deep neural networks have been applied with great success to many visual discriminative tasks, including image classification [6], [10], object detection [7], [11] and semantic segmentation [8], [12].

In real-world applications, training samples typically exhibit a long-tailed class distribution, where a small portion of classes have massive sample points but the others are associated with only a few samples [13], [14], [15], [16]. Such class imbalance of training sample numbers, however, makes the training of deep network based recognition models very challenging. As shown in Fig. 1, the trained model can be easily biased towards head classes with massive training data, leading to poor model performance on tail classes that have limited data [17], [18], [19]. Therefore, the deep models trained by the common practice of empirical risk minimization [20] cannot handle real-world applications with long-tailed class imbalance, *e.g.,* face recognition [21], [22], species classification [23], [24], medical image diagnosis [25], urban scene understanding [26] and unmanned aerial vehicle detection [27].

To address long-tailed class imbalance, massive deep long-tailed learning studies have been conducted in recent years [15], [16], [28], [29], [30]. Despite the rapid evolution in this field, there is still no systematic study to review and discuss existing progress. To fill this gap, we aim to provide a comprehensive survey for recent long-tailed learning studies conducted before mid-2021.
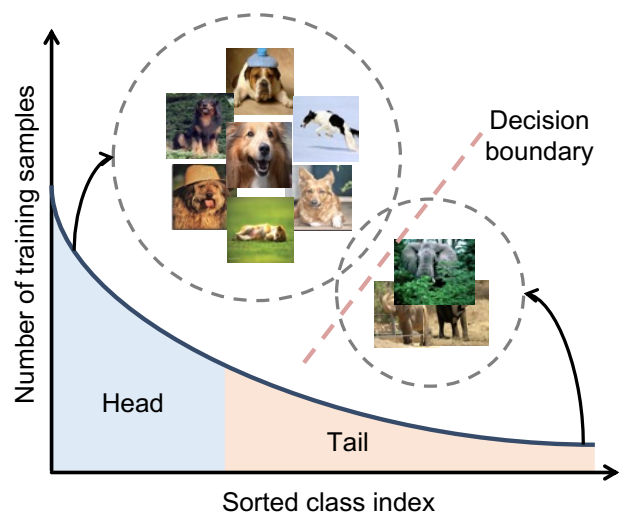
- *Y. Zhang and B. Hooi are with School of Computing, National University of Singapore, Singapore. E-mail: yifan.zhang@u.nus.edu, dcsbhk@nus.edu.sg.*
- *B. Kang, S. Yan, J. Feng are with SEA AI Lab, Singapore. E-mail: bingykang@gmail.com, yansc@sea.com, fengjs@sea.com.*

Fig. 1. The label distribution of a long-tailed dataset (*e.g.,* the iNaturalist species dataset [23] with more than 8,000 classes). The head-class feature space learned on these sampled is often larger than tail classes, while the decision boundary is usually biased towards dominant classes.

As shown in Fig. 2, we group existing methods into three main categories based on their main technical contributions, *i.e.,* class re-balancing, information augmentation and module improvement; these categories can be further classified into nine sub-categories: re-sampling, cost-sensitive learning, logit adjustment, transfer learning, data augmentation, representation learning, classifier design, decoupled training and ensemble learning. According to this taxonomy, we provide a comprehensive review of existing methods, and also empirically analyze several state-of-the-art methods by evaluating their abilities of handling class imbalance using a new evaluation metric, namely *relative accuracy*. We conclude the survey by introducing several real-world application scenarios of deep long-tailed learning and identifying several promising research directions that can be explored by the community in the future.

Fig. 2 taxonomy (left column, a tree diagram rendered as text):

**Long-tailed Learning**

- **Class Re-balancing** (Sec. 3.1)
  - **Re-sampling** (Sec. 3.1.1)
  - **Cost-sensitive Learning** (Sec. 3.1.2)
  - **Logit Adjustment** (Sec. 3.1.3)
- **Information Augmentation** (Sec. 3.2)
  - **Transfer Learning** (Sec. 3.2.1)
  - **Data Augmentation** (Sec. 3.2.2)
- **Module Improvement** (Sec. 3.3)
  - **Representation Learning** (Sec. 3.3.1)
  - **Classifier Design** (Sec. 3.3.2)
  - **Decoupled Training** (Sec. 3.3.3)
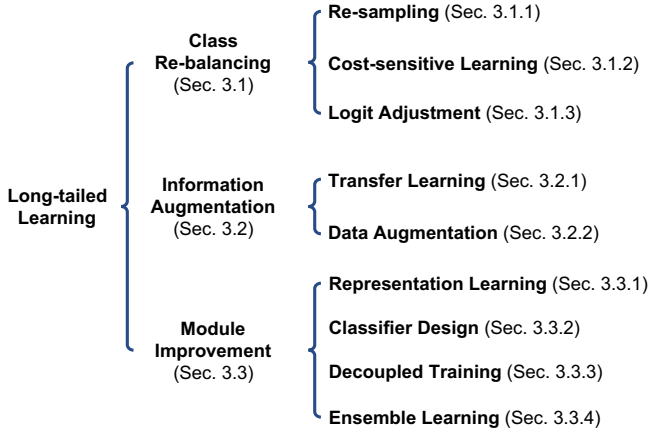  - **Ensemble Learning** (Sec. 3.3.4)

Fig. 2. Taxonomy of existing deep long-tailed learning methods.

We summarize the key contributions of this survey as follows.

- To the best of our knowledge, this is the first comprehensive survey of deep long-tailed learning, which will provide a better understanding of long-tailed visual learning with deep neural networks for researchers and the community.
- We provide an in-depth review of advanced long-tailed learning studies, and empirically study state-of-the-art methods by evaluating to what extent they handle long-tailed class imbalance via a new relative accuracy metric.
- We identify four potential directions for method innovation as well as eight new deep long-tailed learning task settings for future research.

The rest of this survey will be organized as follows: Section 2 presents the problem definition and introduces widely-used datasets, metrics and network backbones. Section 3 provides a comprehensive review of advanced long-tailed learning methods and Section 4 empirically analyzes several state-of-the-art methods based on a new evaluation metric. Section 5 presents the application scenarios of deep long-tailed learning, while Section 6 identifies future research directions. We conclude the survey in Section 7.

# 2 PROBLEM DEFINITION AND BASIC CONCEPTS

## 2.1 Problem Definition

Deep long-tailed learning seeks to learn a deep neural network model from a training dataset with a long-tailed class distribution, where a small fraction of classes have massive samples and the rest classes are associated with only a few samples (c.f. Fig. 1). Let $\{x_i, y_i\}_{i=1}^n$ be the long-tailed training set, where each sample $x_i$ has a corresponding class label $y_i$. The total number of training set over $K$ classes is $n = \sum_{k=1}^{K} n_k$, where $n_k$ denotes the data number of class $k$; let $\pi$ denote the vector of label frequencies, where $\pi_k = n_k/n$ indicates the label frequency of class $k$. Without loss of generality, a common assumption in long-tailed learning [31], [32] is that the classes are sorted by cardinality in decreasing order (*i.e.*, if $i_1 < i_2$, then $n_{i_1} \geq n_{i_2}$, and $n_1 \gg n_K$), and then the imbalance ratio is defined as $n_1/n_K$.

This task is challenging due to two difficulties: (1) imbalanced data numbers across classes make deep models biased to head classes and performs poorly on tail classes; (2) lack of tail-class samples makes it further challenging to train models for tail-class classification. Such a task is fundamental and may occur in various visual recognition tasks, such as image classification [15], [32], detection [19], [33] and segmentation [26], [34], [35].

TABLE 1
Statistics of long-tailed datasets. "Cls." indicates image classification; "Det." represents object detection; "Seg." means instance segmentation.

| Task | Dataset | # classes | # training data | # test data |
|---|---|---|---|---|
| Cls. | ImageNet-LT [15] | 1,000 | 115,846 | 50,000 |
| | CIFAR100-LT [18] | 100 | 50,000 | 10,000 |
| | Places-LT [15] | 365 | 62,500 | 36,500 |
| | iNaturalist 2018 [23] | 8,142 | 437,513 | 24,426 |
| Det./Seg. | LVIS v0.5 [36] | 1,230 | 57,000 | 20,000 |
| | LVIS v1 [36] | 1,203 | 100,000 | 19,800 |
| Multi-label Cls. | VOC-LT [37] | 20 | 1,142 | 4,952 |
| | COCO-LT [37] | 80 | 1,909 | 5,000 |
| Video Cls. | VideoLT [38] | 1,004 | 179,352 | 51,244 |

## 2.2 Datasets

In recent years, a variety of visual datasets have been released for long-tailed learning, differing in tasks, class numbers and sample numbers. In Table 1, we summarize nine visual datasets that are widely used in the deep long-tailed learning community.

In long-tailed image classification, there are four benchmark datasets: ImageNet-LT [15], CIFAR100-LT [18], Places-LT [15], and iNaturalist 2018 [23]. The previous three are sampled from ImageNet [39], CIFAR100 [40] and Places365 [41] following Pareto distributions, respectively, while iNaturalist is a real-world long-tailed dataset. The imbalance ratio of ImageNet-LT, Places-LT and iNaturalist are 256, 996 and 500, respectively; CIFAR100-LT has three variants with various imbalance ratios $\{10, 50, 100\}$.

In long-tailed object detection and instance segmentation, LVIS [36], providing precise bounding box and mask annotations, is the widely-used benchmark. In multi-label image classification, the benchmarks are VOC-LT [37] and COCO-LT [37], which are sampled from PASCAL VOC 2012 [42] and COCO [43], respectively. Recently, a large-scale "untrimmed" video dataset, namely VideoLT [38], was released for long-tailed video recognition.

## 2.3 Evaluation Metrics

In long-tailed learning, the overall performance on all classes and the performance for head, middle and tail classes are usually reported. The used evaluation metrics differ in various tasks. For example, Top-1 Accuracy (or Error Rate) is the widely-used metric for long-tailed image classification, while mean Average Precision (mAP) [44] is adopted for long-tailed object detection and instance segmentation. Moreover, mAP is also used in long-tailed multi-label image classification as a metric, while video recognition applies both Top-1 Accuracy and mAP for evaluation.

## 2.4 Mainstream Network Backbones

Existing long-tailed learning methods are developed based on generic network backbones, which differ in various datasets. The common practices for ImageNet-LT are ResNet [10] and ResNeXt [45] with different depths, where ResNet-50 and ResNeXt-50 are the most common ones. Moreover, ResNet-32 is generally used for CIFAR100-LT; ResNet-50 is used for iNaturalist 2018; ResNet-152 pre-trained on ImageNet is adopted for Places-LT. For LVIS datasets, the widely-used architectures are Mask R-CNN [46] or Faster R-CNN [7] based on ResNet-50 with Feature Pyramid Networks (FPN) [47]. In multi-label classification, the pre-trained ResNet-50 is the common choice for VOC-LT and COCO-LT, while in video recognition of VideoLT, both the pre-trained ResNet-50 and ResNet-101 are applied. On top of these generic backbones, recent methods also explored multiple network branches (*i.e.,* multi-expert) to improve the backbone [30], [48].

## 2.5 Long-tailed Learning Challenges

The most popular challenge events in long-tailed learning includes iNat [23] and LVIS [36].

**iNat Challenge**. The iNaturalist (iNat) challenge is a large-scale fine-grained species classification competition at CVPR. This challenge seeks to push forward the state of the art in automatic image classification for real-world images with a large number of categories, including plants and animals. In contrast to other classification challenges (*e.g.,* ImageNet Large Scale Visual Recognition Challenge), the iNaturalist dataset [23] in this challenge exhibits a long-tailed class distribution and thus encourages progress in image classification.

**LVIS Challenge**. The Large Vocabulary Instance Segmentation (LVIS) dataset [36] is a high-quality instance segmentation dataset with more than 1,000 object categories. As the categories are long-tailed distributed, LVIS presents a novel instance segmentation and object detection challenge at ICCV/ECCV that is distinct from the famous COCO challenge.

## 2.6 Relationships with Other Tasks

We then briefly discuss the differences of long-tailed learning with class-imbalanced learning, few-shot learning, and out-of-domain generalization. These relationships are consistent between deep learning and non-deep learning.

**Class-imbalanced learning** [5], [49] seeks to train models from class-imbalanced samples. Overall, long-tailed learning can be regarded as a more specific and challenging sub-task within class-imbalanced learning. In comparison, in class-imbalanced learning, the number of classes can be very small (*e.g.,* 2) and the number of minority data is not necessarily small; while in long-tailed learning, there are a large number of classes and the tail-class samples are often very scarce.

**Few-shot learning** [50], [51], [52], [53] aims to train models from a limited number of labeled samples (*e.g.,* 1 or 5). In comparison, few-shot learning can be regarded as a sub-task of long-tailed learning, in which the tail classes generally have a very small number of samples.

**Out-of-domain Generalization** [54], [55] indicates a class of tasks, in which the training distribution is inconsistent with the unknown test distribution. Such inconsistency includes inconsistent data marginal distributions (e.g., domain adaptation [56], [57], [58], [59], [60], [61] and domain generalization [62], [63]), inconsistent class distributions (e.g., long-tailed learning [15], [28], [32], open-set learning [64], [65]), and the combination of the previous two situations. From this perspective, long-tailed learning can be viewed as a specific task within out-of-domain generalization.

## 3 CLASSIC METHODS

As shown in Fig. 2, we divide existing deep long-tailed learning methods into three main categories, including class re-balancing, information augmentation, and module improvement. More specifically, class re-balancing consists of three sub-categories: re-sampling, cost-sensitive learning (CSL), and logit adjustment (LA). Information augmentation comprises transfer learning (TL) and data augmentation (Aug). Module improvement includes representation learning (RL), classifier design (CD), decoupled training (DT) and ensemble learning (Ensemble). According to this taxonomy, we sort out existing deep long-tailed learning methods in Table 2 and will review them in detail as follows.

## 3.1 Class Re-balancing

Class re-balancing, a mainstream paradigm in long-tailed learning, seeks to balance the training sample numbers of different classes during model training. We begin with re-sampling based methods, followed by cost-sensitive learning and logit adjustment.

### 3.1.1 Re-sampling

Re-sampling is one of the most widely-used methods to resolve class imbalance in the last few decades [32], [34], [112], [113], [114], [115], [116]. The common practices of re-sampling are random over-sampling (ROS) and random under-sampling (RUS). To re-balance classes, ROS randomly repeats the samples from tail classes, while RUS randomly discards the samples from head classes. Nevertheless, when the classes are extremely skewed, ROS tends to overfit to tail classes, while RUS tends to degrade the model performance on head classes. Instead of using random re-sampling, recent long-tailed learning studies develop various kinds of sampling methods, including class-balanced re-sampling and scheme-oriented sampling.

**Class-balanced re-sampling**. We begin with Decoupling [32], which empirically evaluated various sampling strategies for representation learning on long-tailed recognition. Specifically, the sampling strategies include instance-balanced sampling, class-balanced sampling, square-root sampling and progressively-balanced sampling. In instance-balanced sampling, each sample has an equal probability of being sampled, while in class-balanced sampling, each class has an equal probability of being selected. In addition, square-root sampling [117] is a variant of instance-balanced sampling, where the sampling probability for each class is related to the square root of sample size in the corresponding class. Progressively-balanced sampling [32] interpolates progressively between instance- and class-balanced sampling.

Simple Calibration (SimCal) [34] proposed a new bi-level class-balanced sampling strategy to handle long-tailed instance segmentation. Specifically, the bi-level sampling strategy combines image-level re-sampling and instance-level re-sampling to alleviate class imbalance in instance segmentation.

Dynamic curriculum learning (DCL) [75] developed a new curriculum strategy to dynamically sample data for class re-balancing. To be specific, the more instances from one class are sampled as training goes by, the lower probabilities of this class would be sampled later. Following this idea, DCL first conducts random sampling to learn general representations, and then samples more tail-class instances based on the curriculum strategy to handle long-tailed class imbalance.

Balanced meta-softmax [86] developed a meta learning based sampling method to estimate the optimal sampling rates of different classes for long-tailed learning. Specifically, the proposed meta learning method, a bi-level optimization strategy, learns the best sample distribution parameter by optimizing the *model classification performance* on a balanced meta validation set.

Feature augmentation and sampling adaptation (FASA) [103] proposed to use the *model classification loss* on a balanced meta validation set (as a metric) to adjust feature sampling rate for different classes, so that the under-represented tail classes can be sampled more.

Long-tailed object detector with classification equilibrium (LOCE) [33] proposed to use the *mean classification prediction score* (*i.e.,* running prediction probability) to monitor model training on different classes, and guide memory-augmented feature sampling for enhancing tail-class performance.

TABLE 2
Summary of existing deep long-tailed learning methods published in the top-tier conferences before mid-2021. There are three main categories: class re-balancing, information augmentation and module improvement. In this table, "CSL" indicates cost-sensitive learning; "LA" indicates logit adjustment; "TL" represents transfer learning; "Aug" indicates data augmentation; "RL" indicates representation learning; "CD" indicates classifier design, which seeks to design new classifiers or prediction schemes for long-tailed recognition; "DT" indicates decoupled training, where the feature extractor and the classifier are trained separately; "Ensemble" indicates ensemble learning based methods. We also make our collected long-tailed learning resources available at https://github.com/Vanint/Awesome-LongTailed-Learning.

| Method | Publication | Year | Class Re-balancing | | | Augmentation | | Module Improvement | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | | Re-sampling | CSL | LA | TL | Aug | RL | CD | DT | Ensemble |
| LMLE [66] | CVPR | 2016 | ✓ | | | | | ✓ | | | |
| HFL [67] | CVPR | 2016 | | | | | | ✓ | | | |
| Focal loss [68] | ICCV | 2017 | | ✓ | | | | | | | |
| Range loss [21] | ICCV | 2017 | | | | | | ✓ | | | |
| CRL [69] | ICCV | 2017 | | | | | | ✓ | | | |
| MetaModelNet [70] | NeurIPS | 2017 | | | | ✓ | | | | | |
| DSTL [71] | CVPR | 2018 | | | | ✓ | | | | | |
| CB [16] | CVPR | 2019 | | ✓ | | | | | | | |
| Bayesian estimate [72] | CVPR | 2019 | | ✓ | | | | | | | |
| FTL [73] | CVPR | 2019 | | | | ✓ | ✓ | | | | |
| Unequal-training [74] | CVPR | 2019 | | | | | | ✓ | | | |
| OLTR [15] | CVPR | 2019 | | | | | | ✓ | | | |
| DCL [75] | ICCV | 2019 | ✓ | | | | | | | | |
| Meta-Weight-Net [76] | NeurIPS | 2019 | | ✓ | | | | | | | |
| LDAM [18] | NeurIPS | 2019 | | ✓ | | | | | | | |
| Decoupling [32] | ICLR | 2020 | ✓ | ✓ | | | | ✓ | ✓ | ✓ | |
| LST [77] | CVPR | 2020 | ✓ | | | ✓ | | | | | |
| BBN [48] | CVPR | 2020 | ✓ | | | | | | | | ✓ |
| BAGS [78] | CVPR | 2020 | ✓ | | | | | | | | ✓ |
| Domain adaptation [28] | CVPR | 2020 | | ✓ | | | | | | | |
| Equalization loss (ESQL) [19] | CVPR | 2020 | | ✓ | | | | | | | |
| DBM [22] | CVPR | 2020 | | ✓ | | | | | | | |
| M2m [79] | CVPR | 2020 | | | | ✓ | ✓ | | | | |
| LEAP [80] | CVPR | 2020 | | | | ✓ | ✓ | ✓ | | | |
| IEM [81] | CVPR | 2020 | | | | | | ✓ | | | |
| SimCal [34] | ECCV | 2020 | ✓ | | | | | | | ✓ | ✓ |
| PRS [82] | ECCV | 2020 | ✓ | | | | | | | | |
| Distribution-balanced loss [37] | ECCV | 2020 | | ✓ | | | | | | | |
| OFA [83] | ECCV | 2020 | | | | ✓ | ✓ | | | ✓ | |
| LFME [84] | ECCV | 2020 | | | | ✓ | | | | | ✓ |
| Deep-RTC [85] | ECCV | 2020 | | | | | | | ✓ | | |
| Balanced Meta-Softmax [86] | NeurIPS | 2020 | ✓ | ✓ | | | | | | | |
| UNO-IC [87] | NeurIPS | 2020 | | | ✓ | | | | | | |
| De-confound-TDE [88] | NeurIPS | 2020 | | | ✓ | | | | ✓ | | |
| SSP [89] | NeurIPS | 2020 | | | | ✓ | | ✓ | | | |
| Logit adjustment [14] | ICLR | 2021 | | | ✓ | | | | | | |
| RIDE [17] | ICLR | 2021 | | | | ✓ | | | | | ✓ |
| KCL [13] | ICLR | 2021 | | | | | | ✓ | | ✓ | |
| LTML [90] | CVPR | 2021 | ✓ | | | | | | | | ✓ |
| Equalization loss v2 [91] | CVPR | 2021 | | ✓ | | | | | | | |
| Seesaw loss [92] | CVPR | 2021 | | ✓ | | | | | | | |
| ACSL [93] | CVPR | 2021 | | ✓ | | | | | | | |
| PML [94] | CVPR | 2021 | | ✓ | | | | | | | |
| LADE [31] | CVPR | 2021 | | ✓ | ✓ | | | | | | |
| RoBal [95] | CVPR | 2021 | | ✓ | ✓ | | | | ✓ | | |
| DisAlign [29] | CVPR | 2021 | | ✓ | ✓ | | | | | ✓ | |
| MiSLAS [96] | CVPR | 2021 | | ✓ | | | ✓ | | | ✓ | |
| CReST [97] | CVPR | 2021 | | | | ✓ | | | | | |
| Conceptual 12M [98] | CVPR | 2021 | | | | ✓ | | | | | |
| RSG [99] | CVPR | 2021 | | | | ✓ | ✓ | | | | |
| MetaSAug [100] | CVPR | 2021 | | | | | ✓ | | | | |
| Hybrid [101] | CVPR | 2021 | | | | | | ✓ | | | |
| Unsupervised discovery [35] | CVPR | 2021 | | | | | | ✓ | | | |
| VideoLT [38] | ICCV | 2021 | ✓ | | | | | | | | |
| LOCE [33] | ICCV | 2021 | ✓ | | ✓ | | | | | | |
| GIST [102] | ICCV | 2021 | ✓ | | | ✓ | | | ✓ | | |
| FASA [103] | ICCV | 2021 | ✓ | | | | ✓ | | | | |
| ACE [104] | ICCV | 2021 | ✓ | | | | | | | | ✓ |
| IB [105] | ICCV | 2021 | | ✓ | | | | | | | |
| DARS [26] | ICCV | 2021 | | | | ✓ | | | | | |
| SSD [106] | ICCV | 2021 | | | | ✓ | | | | | |
| DiVE [107] | ICCV | 2021 | | | | ✓ | | | | | |
| MosaicOS [108] | ICCV | 2021 | | | | ✓ | | | | | |
| PaCo [109] | ICCV | 2021 | | | | | | ✓ | | | |
| DRO-LT [110] | ICCV | 2021 | | | | | | ✓ | | | |
| DT2 [111] | ICCV | 2021 | | | | | | | | ✓ | |

VideoLT [38], seeking to address long-tailed video recognition, introduced a new FrameStack method that conducts frame-level sampling to re-balance class distributions. Specifically, FrameStack dynamically adjusts the sampling rates of different classes based on the *running model performance* during training, so that it can sample more video frames from tail classes (generally with lower running performance) and fewer frames from head classes.

**Scheme-oriented sampling** seeks to facilitate some specific learning scheme for long-tailed learning, such as metric learning and ensemble learning. For example, large margin local embedding (LMLE) [66] developed a new quintuplet sampling scheme for metric learning, to learn high-quality features that maintain both inter-cluster and inter-class margins. Unlike the triplet loss [118] that samples two contrastive pairs, LMLE presented a quintuplet sampler to sample four contrastive pairs, including a positive pair and three negative pairs. The positive pair is the most distant intra-cluster sample, while the negative pairs include two inter-clusters samples from the same class (one is the nearest and one is the most distant within the same cluster) and the most nearest inter-class sample. Following that, LMLE introduced a quintuplet loss to encourage the sampled quintuplet to follow a specific distance order. In this way, the learned representations preserve not only locality across intra-class clusters but also discrimination between classes. Moreover, each data batch in the quintuplet loss contains the same number of samples from different classes for class re-balancing.

Partitioning reservoir sampling (PRS) [82] proposed a replay-based sampling method to handle continual long-tailed learning. One key challenge is that the replay memory is unable to consider the issue of class imbalance because no information about the future input is available. To address this, PRS developed an online memory maintenance algorithm that dynamically maintains the running statistics of samples from different classes. Based on the running statistics, PRS can dynamically adjust the memory size and the scheme of sample-in/out operations for different classes.

Bilateral-branch network (BBN) [48] developed two network branches (*i.e.,* a conventional learning branch and a re-balancing branch) to handle class imbalance based on a new bilateral sampling strategy. To be specific, BBN applies uniform sampling for the conventional branch to simulate the original long-tailed training distributions; meanwhile, BBN applies a reversed sampler for the re-balancing branch to sample more tail-class samples for improving tail-class performance. The final prediction is the weighted sum of two network branches. Afterward, long-tailed multi-label visual recognition (LTML) [90] extended the bilateral branch network to address long-tailed multi-label classification. Geometric structure transfer (GIST) [102] also explored this bilateral sampling strategy for knowledge transfer from the head to tail classes.

Besides sampling for bilateral branches, balanced group softmax (BAGS) [78] proposed to divide classes into several balanced groups based on the number of samples in each class, where each group has classes with a similar number of training data. Following this, BAGS uses different sample groups to train different classification heads so that they perform the softmax operation on classes with a similar number of training data and thus avoid a severely biased classifier due to imbalance. Afterward, learning to segment the tail (LST) [77] also divides the training samples into several balanced subsets, and handles each one based on class-incremental learning. To address catastrophic forget during class-incremental learning, LST developed a class-balanced data reply/sampling strategy, which keeps a relatively balanced sample set for knowledge distillation.

TABLE 3
Summary of losses. In this table, $z$ and $p$ indicate the predicted logits and the softmax probability of the sample $x$, where $z_y$ and $p_y$ correspond to the class $y$. Moreover, $n$ indicates the total number of training data, where $n_y$ is the sample number of the class $y$. In addition, $\pi$ denotes the vector of sample frequencies, where $\pi_y = n_y/n$ represents the label frequency of the class $y$. The class-wise weight is denoted by $\omega$ and the class-wise margin is denoted by $\Delta$, if no more specific value is given. Loss-related parameters include $\gamma$.

| Losses | Formulation |
|---|---|
| Softmax loss | $\mathcal{L}_{ce} = -\log(p_y)$ |
| Weighted Softmax loss | $\mathcal{L}_{wce} = -\frac{1}{\pi_y}\log(p_y)$ |
| Focal loss [68] | $\mathcal{L}_{fl} = -(1-p_y)^\gamma \log(p_y)$ |
| Class-balanced loss [16] | $\mathcal{L}_{cb} = -\frac{1-\gamma}{1-\gamma^{n_y}}\log(p_y)$ |
| Balanced softmax loss [86] | $\mathcal{L}_{bs} = -\log\left(\frac{\pi_y \exp(z_y)}{\sum_j \pi_j \exp(z_j)}\right)$ |
| Equalization loss [19] | $\mathcal{L}_{eq} = -\log\left(\frac{\exp(z_y)}{\sum_j \omega_j \exp(z_j)}\right)$ |
| LDAM loss [18] | $\mathcal{L}_{ldam} = -\log\left(\frac{\exp(z_y-\Delta_y)}{\sum_j \exp(z_j-\Delta_j)}\right)$ |

Instead of division into several balanced groups, ally complementary experts (ACE) [104] divides samples into several skill-diverse subsets, where one subset contains all classes, one contains middle and tail classes, and another one contains only tail classes. Based on these subsets, ACE trains different experts to have specific and complementary skills for ensemble learning.

### 3.1.2 Cost-sensitive Learning

Cost-sensitive learning seeks to re-balance classes by adjusting loss values for different classes during training [119], [120], [121], [122], [123], [124], [125]. Recent studies have developed various cost-sensitive long-tailed learning methods to handle class imbalance, including class-level re-weighting and class-level re-margining.

**Class-level re-weighting.** The most intuitive method is to directly use *label frequencies of training samples* for loss re-weighting, namely weighted softmax loss (c.f. Table 3). Such a loss can be further improved by tuning the influence of label frequencies on loss weights, based on sample influence [105] or distribution alignment between model prediction and a balanced reference distribution [29]. In addition to loss value re-weighting, balanced softmax [86] proposed to use the label frequencies to adjust model predictions during training, so that the bias of class imbalance can be alleviated by the prior knowledge. Following that, LADE [31] introduced a label distribution disentangling loss to disentangle the learned model from the long-tailed training distribution, followed by model adaptation to arbitrary test class distributions if the test label frequencies are available.

Instead of using label frequencies, class-balanced loss (CB) [16] introduced a novel concept of *effective number* to approximate the expected sample number of different classes. Here, the effective number is an exponential function of the training sample number. Following this concept, CB loss enforces a class-balanced re-weighting term, inversely proportional to the effective number of classes, to address class imbalance (c.f. Table 3).

Focal loss [68] explored *class prediction hardness* for re-weighting. To be specific, focal loss is inspired by the observation that *class imbalance usually increases the prediction hardness of tail classes, whose prediction probabilities would be lower than those of head classes*. Following this, Focal loss uses the prediction probabilities to inversely re-weight classes (c.f. Table 3), so that it can assign higher weights to the harder tail classes but lower weights to the easier head classes.

Besides using a pre-defined weighting function, the class weights can also be learned from data. In Meta-Weight-Net [76], guided by a balanced validation set, the weighting function, approximated by a one-layer MLP, is updated for fitting the long-tailed distribution, so that a well-performed model on the uniform test set can be learned. In addition, distribution alignment (DisAlign) [29] developed an adaptive calibration function to calibrate the model classifier. The calibration function is adaptively learned by minimizing the KL-Divergence between the adjusted prediction distribution and a given balanced reference distribution.

Another issue in long-tailed learning is negative gradient over-suppression [19], [126]. That is, each positive sample of one class can be seen as a negative sample for other classes in softmax or sigmoid cross-entropy, leading tail classes to receive more suppressed gradients. To address this, distribution-balanced loss [37] alleviates gradient over-suppression via a new negative-tolerant regularization. Meanwhile, it also evaluates the gap between the expected sampling frequency and the actual sampling frequency of each class, and then uses the division of these two frequencies to re-weight loss values for different classes.

Equalization loss [19] directly down-weights the loss values of tail-class samples when they serve as negative pairs for massive head-class samples. Equalization loss v2 [91] further extended the equalization loss [19] by modeling the multi-class detection problem as a set of independent sub-tasks, where each sub-task focuses on one class. More specifically, equalization loss v2 introduced a novel gradient-guided re-weighting mechanism to dynamically up-weight the positive gradients and down-weight the negative gradients for model training on each sub-task.

Seesaw loss [92] re-balances positive and negative gradients for each class with two re-weighting factors, *i.e.,* a mitigation factor and a compensation factor. To address gradient over-suppression, the mitigation factor alleviates the penalty to tail classes during training based on the dynamic ratios of the cumulative sample number between different classes. Meanwhile, if a false positive sample is observed, the compensation factor up-weights the penalty to the corresponding class for improving model discrimination.

Adaptive class suppression loss (ACSL) [93] uses the *output confidence* to decide whether to suppress the gradient for a negative label. Specifically, if the prediction probability of a negative label is larger than a pre-defined threshold, the model should be confused so the weight for this class is set to 1 to improve model discrimination. Otherwise, the weight is set to 0 to avoid negative over-suppression.

**Class-level re-margining** seeks to handle class imbalance by adjusting the minimal margin (*i.e.,* distance), between the learned features and the model classifier, for different classes. For example, label-distribution-aware margin (LDAM) [18] extended the existing soft margin loss [127], [128] by enforcing class-dependent margins based on *label frequencies* and encouraging tail classes to have larger margins. Nevertheless, simply using LDAM loss is not empirically sufficient to handle class imbalance. Therefore, LDAM further introduced a deferred re-balancing optimization schedule that re-balanced classes by re-weighting LDAM loss in a class-balanced way after learning with LDAM loss for a while.

Bayesian estimate [72] found that *the class prediction uncertainty is inversely proportional to the training label frequency, i.e., tail classes are more uncertain*. Inspired by this, bayesian estimate [72] proposed to use the estimated class-level uncertainty to re-margin losses so that the tail classes with higher class uncertainty would suffer a higher loss value and thus have a larger margin between features and the classifier.

Domain balancing [22] studied a long-tailed domain problem, where a small number of domains (containing multiple classes) frequently appear while other domains exist less. To address this task, this work introduced a novel domain frequency indicator based on the *inter-class compactness of features*, and uses this indicator to re-margin the feature space of tail domains.

LOCE [33] uses the *mean classification prediction score* to monitor the learning status for different classes and apply it to guide class-level margin adjustment for enhancing tail-class performance.

Progressive margin loss (PML) [94] adjusts the class-wise margin for long-tailed learning with two margin terms: the ordinal margin and the variational margin. The ordinal margin seeks to extract discriminative features and maintain the age order relation. The variational margin attempts to progressively suppress head classes to handle class imbalance in long-tailed training samples.

RoBal [95] argued that existing re-margining methods that encourage larger margins for tail classes may degrade the feature learning for head classes. Therefore, RoBal enforces an additional margin term to also enlarge the feature margin for head classes.

### 3.1.3 Logit Adjustment

Logit adjustment, post-hoc shifting the model logits based on label frequencies, is a classic idea to obtain a large relative margin between classes in class-imbalanced problems [14], [129]. Recently, one study [14] comprehensively analyzed logit adjustment methods in long-tailed recognition, and theoretically showed that *logit adjustment is Fisher consistent to minimize the average per-class error*. Following this idea, RoBal [95] applied a post-processing strategy to adjust the cosine classification boundary based on training label frequencies.

Instead of using label frequencies of training data, LADE [31] proposed to use the label frequencies of test data (if available) to post-adjust model outputs, so that the trained model can be calibrated for arbitrary test class distribution. UNO-IC [87] proposed to use a hyper-parameter, tuned on a balanced meta validation set, to calibrate the model classifier for handling class imbalance, leading to better performance on the uniform test set. De-confound [88] introduced a causal classifier (c.f. Section 3.3.2) that records the bias information by computing the exponential moving average of features during training, and then removes the bad causal effect by subtracting the bias information during inference. DisAlign [29] applied an adaptive calibration function for logit adjustment, where the calibration function is learned by matching the calibrated prediction distribution to a relatively balanced class distribution.

### 3.1.4 Discussions

Compared to other long-tailed learning paradigms, class re-balancing methods are relatively simple but can achieve comparably or even better performance. Some of them, especially cost-sensitive learning methods, are theoretically inspired or guaranteed to handle long-tailed problems [16], [18], [31]. These advantages enable class re-balancing to be a good candidate for real-world applications.

However, one drawback of this type of method is that most class re-balancing methods improve tail-class performance at the cost of head-class performance, which is like playing on a performance seesaw. Although the overall performance is improved, it cannot essentially handle the issue of lacking information, particularly on tail classes due to limited data amount. To address this limitation, one feasible solution is to conduct information augmentation for all classes as follows.

## 3.2 Information Augmentation

Information augmentation based methods seek to introduce additional information into model training, so that the model performance can be improved in long-tailed learning. There are two kinds of methods in this method type: transfer learning and data augmentation.

### 3.2.1 Transfer Learning

Transfer learning [70], [83], [99], [130], [131] seeks to transfer the knowledge from a source domain (*e.g.,* datasets, tasks or classes) to enhance model training on a target domain. In deep long-tailed learning, there are four main transfer learning schemes, *i.e.,* head-to-tail knowledge transfer, model pre-training, knowledge distillation, and self-training.

**Head-to-tail knowledge transfer** seeks to transfer the knowledge from head classes to augment model performance on tail classes. For instance, feature transfer learning (FTL) [73] found that *tail-class samples have much smaller intra-class variance than head-class samples, leading to biased feature spaces and decision boundaries*. To address this, FTL exploits the knowledge of intra-class variance from head classes to guide feature augmentation for tail-class samples, so that the tail-class features have higher intra-class variance, leading to better tail-class performance. Following that, LEAP [80] constructs "feature cloud" for each class, and seeks to transfer the knowledge of head-class feature clouds to enhance the intra-class variation of tail-class feature clouds, by augmenting tail-class samples with certain disturbation in the feature space. As a result, the distortion of the intra-class feature variance among classes is alleviated.

Online feature augmentation (OFA) [83] proposed to use class activation maps [132] to decouple sample features into class-specific and class-agnostic ones. Following that, OFA augments tail classes by combining the class-specific features of tail-class samples with class-agnostic features from head-class samples. Afterward, all the augmented and original features would be used to fine-tune the model classifier with a re-balancing sampler, leading to better long-tailed learning performance.

Rare-class sample generator (RSG) [99] also observed that *the feature space of tail classes is much smaller than that of head classes in long-tailed problems*. To address this, RSG proposed to generate new tail-class samples to enlarge the feature space for tail classes and "push away" the decision boundaries. To this end, RSG dynamically estimates a set of feature centers for each class, and uses the feature displacement between head-class sample features and the nearest intra-class feature center to augment each tail sample feature. To further maximize the feature displacement distance and increase the diversity of the generated tail-class sample features, RSG introduced a maximized vector loss to enforce the direction of the feature displacement and the direction of the sample feature to be "co-linear".

In addition to the feature-level head-to-tail transfer, major-to-minor translation (M2m) [79] proposed to augment tail classes by translating head-class samples to tail-class ones via perturbation-based optimization, which is essentially similar to adversarial attack. The translated tail-class samples would be used to construct a more balanced training set for model training. Moreover, GIST [102] proposed to conduct head-to-tail transfer at the classifier level. By enhancing the classifier weights of tail classes with the relatively large classifier geometry information of head classes, GIST is able to obtain better tail-class performance.

MetaModelNet [70] proposed to learn a meta-network that maps few-shot model parameters to many-shot model parameters, where the few-shot model is trained on a small number of samples while the many-shot model is trained on massive samples. To be specific, the meta-network is trained on head classes, where the many-shot model is directly trained on the head-class training set, while the few-shot model is trained on a sample subset from these classes. Following that, the meta-network learned on head classes is applied to map the few-shot model trained on tail classes for obtaining better tail-class performance.

**Model pre-training** is a popular scheme for deep model training [133], [134], [135], [136], [137]. Domain-specific transfer learning (DSTL) [71] first pre-trains the model with all long-tailed samples for representation learning, and then fine-tunes the model on a more class-balanced training subset. In this way, DSTL slowly transfers the learned features to tail classes, obtaining more balanced performance among all classes. In addition, self-supervised pre-training (SSP) [89] proposed to first use self-supervised learning (*e.g.,* contrastive learning [138] or rotation prediction [139]) for model pre-training, followed by standard training on long-tailed data. Empirical results show the effectiveness of SSP, where tail classes exhibit larger performance gains. Such a scheme has also be explored to handle long-tailed data with noisy labels [140]. Recently, a new vision-and-language pre-training dataset (Conceptual 12M [98]) was proposed and has been shown effective for downstream long-tailed recognition.

**Knowledge distillation** seeks to train a student model based on the outputs of a well-trained teacher model [141], [142]. Several recent studies have explored knowledge distillation for long-tailed learning. LST [77] developed a class-incremental learning strategy to handle long-tailed instance segmentation, where knowledge distillation is used to overcome catastrophic forget during incremental learning. Learning from multiple experts (LFME) [84] divides the entire long-tailed dataset into several subsets with smaller degrees of class imbalance, and trains multiple experts with different sample subsets. Based on these experts, LFME trains a unified student model using adaptive knowledge distillation in an easy-to-hard curriculum instance selection manner. Following the multi-expert framework, routing diverse distribution-aware experts (RIDE) [17] introduced a knowledge distillation method to reduce the parameters of the multi-expert model by learning a student network with fewer experts.

Very recently, self-supervision to distillation (SSD) [106] developed a new self-distillation scheme to enhance decoupled training. Specifically, SSD first trains a calibrated model based on supervised and self-supervised information via the decoupled training scheme (c.f. Section 3.3.3), and then uses the calibrated model to generate soft labels for all samples. Following that, both the generated soft labels and original long-tailed hard labels are used to distill a new student model, followed by a new classifier fine-tuning stage. In addition, distill the virtual examples (DiVE) [107] showed the effectiveness of using a class-balanced model as the teacher in knowledge distillation for long-tailed learning.

**Self-training** aims to learn well-performing models from a small number of labeled samples and massive unlabeled samples [143], [144], [145]. To be specific, it firstly uses labeled samples to train a supervised model, which is then applied to generate pseudo labels for unlabeled data. Following that, both the labeled and pseudo-labeled samples are used to re-train models. In this way, self-training exploits the knowledge from massive unlabeled samples to enhance long-tailed learning performance.

Such a paradigm, however, cannot be directly used to handle long-tailed problems, because both labeled and unlabeled datasets may follow long-tailed class distributions. In such cases, the trained model on labeled samples may be biased to head classes and tends to generate more head-class pseudo labels for unlabeled samples, which leads to a more skewed degree of class imbalance. By far, how to enhance self-training to address long-tailed semi-supervised learning is still an under-explored important question.

Class-rebalancing self-training (CReST) [97] explored self-training in long-tailed classification and found that *the precision of the supervised model on tail classes is surprisingly high*. Based on this finding, CReST proposed to select more tail-class samples for online pseudo labeling in each iteration, so that the re-trained model can obtain better performance on tail classes.

Distribution alignment and random sampling (DARS) [26] enhances self-training for handling long-tailed semi-supervised semantic segmentation. To address the potential inconsistency of class imbalance between labeled and unlabeled samples, DARS regards the label frequencies of labeled training data as the true class distribution, and enforces the label frequencies of the generated pseudo labels to be consistent with the labeled ones. Meanwhile, DARS applies a sampling strategy to strictly control the number of pseudo labels in each class. In this way, the generated pseudo labels would be more consistent with the true labeled ones, which leads to model performance improvement.

MosaicOS [108] resorted to additional object-centric images (for image classification) to boost long-tailed object detection. Specifically, it first pre-trains the model with labeled scene-centric images from the original detection dataset, and then uses the pre-trained model to generate pseudo bounding boxes for object-centric images, *e.g.,* ImageNet-1K [39]. After that, MosaicOS fine-tunes the pre-trained model in two stages, *i.e.,* first fine-tuning with the pseudo-labeled object-centric images and then fine-tuning with the original labeled scene-centric images. In this way, MosaicOS alleviates the negative influence of data discrepancies and effectively improves long-tailed learning performance.

### 3.2.2 Data Augmentation

Data Augmentation aims to pack a set of augmentation techniques to enhance the size and quality of datasets for model training [146], [147]. In long-tailed learning, there are two kinds of data augmentation methods having been explored, including transfer-based augmentation (please refer to head-to-tail knowledge transfer in Section 3.2.1) and conventional (non-transfer) augmentation.

Non-transfer augmentation seeks to improve or design conventional data augmentation methods to address long-tailed problems. MiSLAS [96] investigated data mixup in long-tailed learning, and found that (1) *data mixup helps to remedy model over-confidence*; (2) *mixup has a positive effect on representation learning but a negative or negligible effect on classifier learning in the decoupled training scheme* [32]. Following these observations, MiSLAS proposed to use data mixup to enhance representation learning in the decoupled scheme. In addition, Remix [148] also resorted to data mixup for long-tailed learning and introduced a re-balanced mixup method to particularly enhance tail classes.

FASA [103] proposed to generate class-wise features, based on a Gaussian prior with its mean and variance estimated from previously observed samples. Moreover, FASA exploits the model classification loss on a balanced validation set to adjust sampling rates of features for different classes, so that the under-represented tail classes can be augmented more.

Meta semantic augmentation (MetaSAug) [100] proposed to augment tail classes with a variant of implicit semantic data augmentation (ISDA) [149]. To be specific, ISDA estimates the class-conditional statistics (*i.e.,* covariance matrices from sample features) to obtain semantic directions, and generates diversified augmented samples by translating sample features along with diverse semantically meaningful directions. Nevertheless, insufficient tail-class samples make it ineffective to estimate the covariance matrices for tail classes. To address this, MetaSAug explored meta learning to guide the learning of covariance matrices for each class with the class-balanced loss [16]. In this way, the covariance matrices of tail classes are estimated more accurately, and thus the generated tail-class features are more informative.

### 3.2.3 Discussions

Thanks to introducing additional knowledge, transfer learning based methods improve tail-class performance without sacrificing head-class performance. Considering the lack of enough tail-class samples is one of the key problems in long-tailed learning, this type of method is worth further exploring.

Data augmentation is a relatively fundamental technique and can be used for a variety of long-tailed problems, which makes this type of method more practical than other methods in real-world applications. However, simply using existing class-agnostic augmentation techniques for improving long-tailed learning is unfavorable, since they may further increase imbalance considering head classes have more samples and would be augmented more. How to better conduct data augmentation for long-tailed learning is still an open question.

## 3.3 Module Improvement

Besides class re-balancing and information augmentation, researchers also explored methods to improve network modules in long-tailed learning. These methods can be divided into four categories: (1) representation learning improves the feature extractor; (2) classifier design enhances the model classifier; (3) decoupled training boosts the learning of both the feature extractor and the classifier; (4) ensemble learning improves the whole architecture.

### 3.3.1 Representation Learning

Existing representation learning methods for long-tailed learning are based on four main paradigms, *i.e.,* metric learning, sequential training, prototype learning, and transfer learning.

**Metric learning** aims at designing task-specific distance metrics for establishing similarity or dissimilarity between objects; in long-tailed learning, metric learning based methods seek to explore distance-based losses to learn a more discriminative feature space. One example is LMLE [66], which introduced a quintuplet loss to learn representations that maintain both inter-cluster and inter-class margins. Moreover, range loss [21] innovated representation learning by using the overall distances among all sample pairs within one mini-batch. In other words, the range loss uses statistics over the whole batch, rather than instance level, and thus alleviates the bias of data number imbalance overall classes. More specifically, range loss enlarges the inter-class distance by maximizing the distances of any two class centers within mini-batch, and reduces the intra-class variation by minimizing the largest distances between intra-class samples. In this way, the range loss obtains features with better discriminative abilities and less imbalanced bias.

Class rectification loss (CRL) [69] seeks to enhance tail-class sample representations to have a larger degree of intra-class compactness and inter-class distances. To this end, CRL constructs massive hard-pair triplets for tail classes and applies a class rectification loss (similar to the triplet loss [118]) as a class-balanced constraint. In this way, the learned model overcomes the negative influence of class imbalance on representation learning.

Recent studies also explored contrastive learning for long-tailed problems. KCL [13] proposed a $k$-positive contrastive loss to learn a balanced feature space, which helps to alleviate class imbalance and improve model generalization. Following that, Hybrid [101] introduced a prototypical contrastive learning strategy to enhance long-tailed learning. Parametric contrastive learning (PaCo) [109] further innovated supervised contrastive learning by adding a set of parametric learnable class centers, which play the same role as a classifier if regarding the class centers as the classifier weights. DRO-LT [110] extended the prototypical contrastive learning with distribution robust optimization [150], which makes the learned model more robust to data distribution shift.

**Sequential training**. Hierarchical feature learning (HFL) [67] took inspiration from that each class has their individuality in discriminative visual representation. Therefore, HFL hierarchically clusters objects into visually similar class groups, forming a hierarchical cluster tree. In this cluster tree, the model in the original node is pre-trained on ImageNet-1K; the model in each child node inherits the model parameters from its parent node and is then fine-tuned based on samples in the cluster node. In this way, the knowledge from the groups with massive classes is gradually transferred to their sub-groups with fewer classes.

Unequal-training [74] proposed to divide the dataset into head-class and tail-class subsets, and treat them differently in the training process. First, unequal-training uses the head-class samples to train relatively discriminative and noise-resistant features with a new noise-resistant loss. After that, it uses tail-class samples to enhance the inter-class discrimination of representations via hard identities mining and a novel center-dispersed loss. Here, the center-dispersed loss is based on normalized features in each class.

**Prototype learning** based methods seek to learn class-specific feature prototypes to enhance long-tailed learning performance. Open long-tailed recognition (OLTR) [15] innovatively explored the idea of feature prototypes to handle long-tailed recognition in an open world, where the test set includes head, tail and open classes. Here, the open classes indicate the test classes that do not exist in the training set. To address this task, OLTR maintains a visual meta memory containing discriminative feature prototypes, and uses the features induced from the visual memory to augment the original features. In this way, the learned feature space would be more discriminative, and the sample features from novel classes would be far away from the memory and closer to the origin point. Such a feature space enables OLTR to discriminate close-set classes as well as to detect novel classes. Moreover, OLTR also explored a self-attention scheme to enhance feature learning.

Following that, inflated episodic memory (IEM) [81] further innovated the meta-embedding memory by a dynamical update scheme, in which each class has independent and differentiable memory blocks, while each memory block records the most discriminative feature prototype for the corresponding categories. As the dynamic memory banks contain only the most discriminative feature prototypes, they are not influenced by the issue of class number imbalance. In addition, IEM also explored a region self-attention mechanism to further enhance representation learning.

TABLE 4
Summary of classifiers, where $w$, $f$, $b$, $\phi$, and $p$ denote the model classifier, sample features, the bias term, the softmax function, and prediction probabilities, respectively. Moreover, $\hat{d}$ is the unit vector of the exponential moving average features. The temperature factor is denoted by $\tau$, and other classifier-related hyper-parameters include $\gamma$ and $\alpha$.

| Classifier | Formulation |
|---|---|
| Linear classifier | $p = \phi(w^\top f + b)$ |
| Cosine classifier | $p = \phi(\tau \frac{w^\top f}{\|w\|\|f\|} + b)$ |
| $\tau$-norm classifier | $p = \phi(\frac{w^\top f}{\|w\|_2^\tau} + b)$ |
| Causal classifier | $p = \phi(\tau \frac{w^\top f}{(\|w\|+\gamma)\|f\|} - \alpha \frac{cos(x,\hat{d})(w^\top \hat{d})}{\|w\|+\gamma})$ |

**Transfer learning**. In Section 3.2.1, we have introduced some transfer-based long-tailed methods that improve representation learning, including SSP [89] and LEAP [80]. In addition to them, unsupervised discovery (UD) [35] proposed to use self-supervised learning to help discover novel and more fine-grained objects from images with long-tailed objects. To be specific, UD first uses a pre-trained class-agnostic mask proposal network to generate object bounding boxes and segmentation masks for all possible objects. Then, UD applies three new self-supervised triplet losses, based on features of the bounding boxes and semantic masks, to learn a hyperbolic feature space. Based on the learned features, UD lastly conducts unsupervised clustering and exclusive label assignation for clusters to discover novel and more fine-grained objects.

Besides the above learning schemes, decoupling [32] innovated long-tailed representation learning with different sampling strategies, including instance-balanced, class-balanced, square-root, and progressively-balanced sampling. Following the instance-balanced sampling, MiSLAS [96] empirically studied the influence of data mixup on long-tailed representation learning.

### 3.3.2 Classifier Design

In addition to representation learning, researchers also designed different types of classifiers to address long-tailed problems. In generic visual problems [10], [138], the common practice of deep learning is to use linear classifier $p = \phi(w^\top f + b)$, where $\phi$ denotes the softmax function and the bias term $b$ can be discarded.

However, long-tailed class imbalance often results in larger classifier weight norms for head classes than tail classes [73], which makes the linear classifier easily biased to dominant classes. To address this, several studies [80], [95] proposed to use the scale-invariant cosine classifier $p = \phi((\frac{w^\top f}{\|w\|\|f\|})/\tau + b)$, where both the classifier weights and sample features are normalized. Here, the temperature $\tau$ should be chosen reasonably [125], or the classifier performance would be negatively influenced.

The $\tau$-normalized classifier [32] rectifies the imbalance of decision boundaries by adjusting the classifier weight norms through a $\tau$-normalization procedure. Formally, let $\tilde{w} = \frac{w}{\|w\|_2^\tau}$, where $\tau$ is the temperature factor for normalization. When $\tau = 1$, the $\tau$-normalization reduces to $L_2$ normalization, while when $\tau = 0$, no scaling is imposed. Note that, the hyper-parameter $\tau$ can also be trained with class-balanced sampling, and the resulting classifier is named the learnable weight scaling classifier [32]. In addition, the nearest class mean classifier [32] computes the $L_2$-normalized mean features for each class on the training set, and then conducts prediction based on the nearest neighbor algorithm [151] either using the cosine similarity or the Euclidean distance.

Realistic taxonomic classifier (RTC) [85] proposed to address class imbalance with hierarchical classification. Specifically, RTC maps images into a class taxonomic tree structure, where the hierarchy is defined by a set of classification nodes and node relations. Different samples are classified adaptively at different hierarchical levels, where the level at which the prediction is made depends on the sample classification difficulty and the classifier confidence. Such a design favors correct decisions at intermediate levels rather than incorrect decisions at the leaves.

Causal classifier [88] resorted to causal inference for keeping the good and removing the bad momentum causal effects in long-tailed learning. The good causal effect indicates the beneficial factor that stabilizes gradients and accelerates training, while the bad causal effect indicates the accumulated long-tailed bias that leads to poor tail-class performance. To better approximate the bias information, the causal classifier applies a multi-head strategy to divide the channel (or dimensions) of model weights and data features equally into $K$ groups. Formally, the causal classifier calculates the original logits by $p = \phi(\frac{\tau}{K} \sum_{k=1}^{K} \frac{(w^k)^\top f^k}{(\|w^k\|+\gamma)\|f^k\|})$, where $\tau$ is the temperature factor and $\gamma$ is a hyper-parameter. This classifier is essentially the cosine classifier when $\gamma = 0$. In inference, the causal classifier removes the bad causal effect by subtracting the prediction when the input is null, *i.e.,* $p = \phi(\frac{\tau}{K} \sum_{k=1}^{K} \frac{(w^k)^\top f^k}{(\|w^k\|+\gamma)\|f^k\|} - \alpha \frac{cos(x^k, \hat{d}^k)(w^k)^\top \hat{d}^k}{\|w^k\|+\gamma})$, where $\hat{d}$ is the unit vector of the exponential moving average features, and $\alpha$ is a trade-off parameter to control the direct and indirect effects. More intuitively, the classifier records the bias by computing the exponential moving average features during training, and then removes the bad causal effect by subtracting the bias from prediction logits during inference.

GIST classifier [102] seeks to transfer the geometric structure of head classes to tail classes. Specifically, the GIST classifier consists of a class-specific weight center (for encoding the class location) and a set of displacements (for encoding the class geometry). By exploiting the relatively large displacements from head classes to enhance tail-class weight centers, the GIST classifier is able to obtain better performance on tail classes.

### 3.3.3 Decoupled Training

Decoupled training decouples the learning procedure into representation learning and classifier training.

Decoupling [32] was the pioneering work to introduce the two-stage training scheme. It empirically evaluated different sampling strategies for representation learning in the first stage (c.f. Section 3.1.1), and then evaluated different classifier training schemes by fixing the trained feature extractor in the second stage. In the classifier learning stage, there are also four methods, including classifier re-training with class-balanced sampling, the nearest class mean classifier, the $\tau$-normalized classifier, and a learnable weight scaling scheme. The main observations are twofold: (1) *instance-balanced sampling is surprisingly the best strategy for representation learning*; (2) the devil is in classification: *re-adjusting the classifier leads to significant performance improvement* in long-tailed recognition.

Following that, KCL [13] empirically observed that *a balanced feature space is beneficial to long-tailed learning*. Therefore, it innovated the decoupled training scheme by developing a $k$-positive contrastive loss to learn a more class-balanced and class-discriminative feature space, which leads to better long-tailed learning performance.

MiSLAS [96] empirically observed that *data mixup is beneficial to features learning but has a negative/negligible effect on classifier training under the two-stage decoupled training scheme*. Therefore, MiSLAS proposed to enhance the representation learning with data mixup in the first stage. During the second stage, MiSLAS applies a label-aware smoothing strategy for better model generalization, and further addresses the distribution shift between two training stages by keeping updating the running mean and variance in the batch normalization layers.

Several recent studies innovated the decoupled training scheme by enhancing the classifier training stage. OFA [83] innovated the classifier re-training through tail-class feature augmentation by combining the extracted class-specific features of tail classes with the extracted class-generic features from head classes. SimCal [34] proposed to enhance the classifier training stage by calibrating the classification head with a novel bi-level class-balanced sampling strategy for long-tailed instance segmentation. DisAlign [29] innovated the the classifier training with a new adaptive calibration strategy. Specifically, a new adaptive calibration function, learned by minimizing the KL-Divergence between the calibrated prediction distribution and a balanced reference distribution, is used to adjust the output logits of the original classifier. To summarize, DisAlign essentially applies an additional classifier layer to calibrate the original classifier by matching the calibrated prediction distribution to a relatively balanced class distribution.

Very recently, DT2 [111] applied the scheme of decoupled training to the scene graph generation task, which demonstrates the effectiveness of decoupled training in handling long-tailed class imbalance in visual relation learning.

### 3.3.4 Ensemble Learning

Ensemble learning based methods strategically generate and combine multiple network modules (namely, multiple experts) to solve long-tailed visual learning problems. We summarize the main schemes of existing ensemble-based methods in Fig. 3, which will be detailed as follows.

BBN [48] proposed to use two network branches, *i.e.,* a conventional learning branch and a re-balancing branch, to handle long-tailed recognition. To be specific, the conventional learning branch applies uniform sampling to simulate the original long-tailed training distribution; the re-balancing branch applies a reversed sampler to sample more tail-class samples in each mini-batch for improving tail-class performance. The predictions of two branches are dynamically combined during training, so that the learning focus of BBN gradually changes from head classes to tail classes (via the re-balancing branch).

Following BBN, LTML [90] explored the bilateral-branch network scheme to solve long-tailed multi-label classification. To be specific, LTML trains each branch using the sigmoid cross-entropy loss for multi-label classification and enforces a logit consistency loss to improve the consistency of two branches. Moreover, LTML applies label smoothing and logit compensation for improving model generalization.

Similar to BBN, SimCal [34] explored a dual classification head scheme to address long-tail instance segmentation. Specifically, SimCal maintains two classification heads: the original classification head and a calibrated classification head. Based on a new bi-level sampling strategy, the calibrated classification head is able to improve the performance on tail classes, while the original head aims to maintain the performance on head classes.

Fig. 3. Illustrations of existing ensemble-based long-tailed learning methods. The trained experts may have different skills, *e.g.,* being skilled in different class distributions or different class subsets.

Instead of bilateral branches, BAGS [78] explored a multi-head scheme to address long-tailed object detection. Specifically, BAGS took inspiration from an observation that learning a more uniform distribution with fewer samples is sometimes easier than learning a long-tailed distribution with more samples. Therefore, BAGS first divides classes into several sub-groups, where the classes in each sub-group have a similar number of training data. Then, BAGS applies multiple classification heads, upon a shared feature extractor, for prediction, where different classification heads are trained on different data sub-groups. In this way, each classification head performs the softmax operation on classes with a similar number of training data, thus avoiding the negative influence of class imbalance. Moreover, BAGS also introduces "other" classes into each group to alleviate the contradiction among different heads.

Similar to BAGS, LFME [84] divides the long-tailed dataset into several subsets with smaller "class longtailness", and trains multiple experts with different sample subsets. Based on these experts, LFME then learns a unified student model using adaptive knowledge distillation from the multiple teacher experts.

Instead of division into several balanced sub-groups, ACE [104] divides classes into several skill-diverse subsets: one subset contains all classes; one contains middle and tail classes; another one has only tail classes. ACE then trains multiple experts with various class subsets, so that different experts have specific and complementary skills. Moreover, considering various subsets have different sample numbers, ACE also applies a distributed-adaptive optimizer to adjust the learning rate for different experts. A similar idea of ACE was also explored in ResLT [152].

Without data division, RIDE [17] trains each expert independently with softmax loss based on all training samples and enforces a KL-divergence based loss to improve the diversity of different experts. Following that, RIDE applies an expert assignment module to improve computing efficiency. Note that training each expert with the softmax loss independently boosts the ensemble performance on long-tailed learning a lot.

Test-time aggregating diverse experts (TADE) [30] explored the multi-expert scheme to handle test distribution-agnostic long-tailed recognition, where the test class distribution can be either uniform or long-tailed. To be specific, TADE developed a novel spectrum-spanned multi-expert framework, and innovated the expert training scheme by introducing a diversity-promoting expertise-guided loss that trains different experts to handle different class distributions. In this way, the learned experts are more diverse, leading to better ensemble performance, and integratedly span a wide spectrum of possible class distributions. Based on this property, TADE further introduced a novel test-time self-supervised learning method, namely prediction stability maximization, to adaptively aggregate experts for better handling unknown test class distribution, based on only unlabeled test data.

### 3.3.5 Discussions

Representation learning and classifier design are fundamental problems for deep long-tailed learning, being worth further exploring. Decoupled training is attracting increasing attention in recent studies; in this scheme, the second stage of class-balanced classifier learning does not introduce too many computation costs but leads to significant performance gains. One critique [104] is that the accumulated training stages make decoupled training less practical to be integrated with existing well-formulated methods in other long-tailed problems, *e.g.,* object detection and instance segmentation. Despite this, the idea of decoupling training is conceptually simple and thus can be easily used to design new methods for resolving a variety of long-tailed learning problems.

Ensemble-based methods, compared to other types of long-tailed learning methods, generally obtain better performance on both head and tail classes. One concern of these methods is that they generally lead to higher computational costs due to the use of multiple experts. Such a concern, however, can be alleviated by using a shared feature extractor. Moreover, efficiency-oriented expert assignment and knowledge distillation strategies [17], [84] can also reduce computational complexity.

# 4 EMPIRICAL STUDIES

This section empirically analyzes existing long-tailed learning methods. To begin with, we introduce a new evaluation metric.

## 4.1 Novel Evaluation Metric

The key goal of long-tailed learning is to handle class imbalance for better model performance. Therefore, the common evaluation protocol [13], [22] is directly using the top-1 test accuracy (denoted by $A_t$) to judge how well long-tailed methods perform and which method handles class imbalance better. Such a metric, however, cannot accurately reflect the relative superiority among different methods when handling class imbalance, as the top-1 accuracy is also influenced by other factors apart from class imbalance. For example, long-tailed methods like ensemble learning (or data augmentation) also improve the performance of models, trained on a balanced training set. In such cases, it is hard to tell if the performance gain is from the alleviation of class imbalance or from better network architectures (or more data information).

To better evaluate the method's effectiveness in handling class imbalance, we propose a new metric, namely **relative accuracy** $A_r$, to alleviate the influence of unnecessary factors in long-tailed learning. To this end, we first compute an empirically upper reference accuracy $A_u = \max(A_v, A_b)$, which is the maximal value between the vanilla accuracy $A_v$ of the corresponding backbone trained on a balanced training set with cross-entropy and the balanced accuracy $A_b$ of the model trained on a balanced training set with the corresponding long-tailed method. Here, the balanced training set is a variant of the long-tailed training set with a similar total data number but each class has the same data number. This upper reference accuracy, obtained from the balanced training set, is used to alleviate the influence apart from class imbalance; then the relative accuracy is defined by $A_r = \frac{A_t}{A_u}$. In our experiments, all the accuracy, upper reference accuracy and relative accuracy will be used for evaluation.

## 4.2 Experimental Settings

We then introduces the experimental settings.

**Datasets**. We adopt the widely-used ImageNet-LT [15] as the benchmark long-tailed dataset for empirical studies, considering that ImageNet-LT has massive classes of 1,000 and a large imbalance ratio of 256. The corresponding balanced training set variant of ImageNet-LT is sampled based on [13]. The total sample number of ImageNet-LT can be found in Table 1. Besides the performance regarding all classes, we also report performance on three class subsets in ImageNet-LT: Head (more than 100 images), Middle (20∼100 images) and Tail (less than 20 images).

**Baselines**. We select long-tailed learning methods via the following criterion: (1) official source codes are publicly available or easy to re-implement; (2) methods are evaluated on ImageNet-LT in the corresponding papers. As a result, more than 20 methods are empirically evaluated in this paper, including baseline (**Softmax**), cost-sensitive learning (**Weighted Softmax**, **Focal loss** [68], **LDAM** [18], **ESQL** [19], **Balanced Softmax** [86], **LADE** [31]), logit adjustment (**UNO-IC** [87]), transfer learning (**SSP** [89]), data augmentation (**RSG** [99]) representation learning (**OLTR** [15], **PaCo** [109]). classifier design (**De-confound** [88]), decoupled training (**Decouple-IB-CRT** [32], **CB-CRT** [32], **SR-CRT** [32], **PB-CRT** [32], **MiSLAS** [96]), ensemble learning (**BBN** [48], **LFME** [84], **RIDE** [17], **ResLT** [152], **TADE** [30]). More details of these methods can be found in Section 3.

**Implementation details**. We implement all experiments in PyTorch. Following [17], [31], [32], we use ResNeXt-50 as the network backbone for all methods. We conduct model training with the SGD optimizer based on batch size 256, momentum 0.9 and weight decay factor 0.0005, and learning rate 0.1 (linear LR decay). For method-related hyper-parameters, we set the values by either directly following the original papers or manual tuning if the default values perform poorly. Moreover, we use the same basic data augmentation (*i.e.,* random resize and crop to 224, random horizontal flip, color jitter, and normalization) for all methods, while other augmentation techniques proposed in augmentation-based long-tailed methods would be used on top of these basic augmentation operations.

## 4.3 Results on all Classes

Table 5 and Fig. 4 report the average performance over all classes. From these results, we have several observations on overall method progress and different method types.

**Observations on all methods.** As shown in Table 5, almost all long-tailed methods perform better than the Softmax baseline in terms of accuracy, which demonstrates the effectiveness of long-tailed learning. Even so, there are two methods performing slightly worse than Softmax, *i.e.,* Decouple-CB-CRT [32] and BBN [48]. We speculate that the poor performance of Decouple-CB-CRT results from poor representation learning by class-balanced sampling in the first stage of decoupled training (refer to [32] for more empirical observations). The poor results of BBN (based on the official codes) may come from the cumulative learning strategy, which gradually adjusts the learning focus from head classes to tail classes; at the end of the training, however, it may put too much focus on the tail ones. As a result, despite the better tail-class performance, the model accuracy on head classes drops significantly (c.f. Table 6), leading to worse average performance.

In addition to accuracy, we also evaluate long-tailed methods based on upper reference accuracy (UA) and relative accuracy (RA). Table 5 shows that most methods have the same UA as the baseline model, but there are still some methods having higher UA, *e.g.,* SSP, MiSLAS, TADE. For these methods, the performance improvement comes not only from the alleviation of class imbalance, but also from other factors, like data augmentation or better network architectures. Therefore, simply using accuracy for evaluation is not accurate enough, while our proposed RA metric provides a good complement, since it alleviates the influences of factors apart from class imbalance. For example, MiSLAS, based on data mixup, has higher accuracy than Balanced Softmax under 90 training epochs, but it also has higher UA. As a result, the relative accuracy of MiSLAS is lower than Balanced Softmax, which means that Balanced Softmax alleviates class imbalance better than MiSLAS under 90 training epochs. When the training epoch increases to 200, MiSLAS has higher RA than Balanced Softmax. That is, despite having other factors improving performance, MiSLAS with sufficient training also shows a better ability to handle class imbalance than Balanced Softmax. More examples under 200 training epochs can be found in Figs. 4 (a,c).

Although some recent high-accuracy methods have lower RA, the overall development trend of long-tailed learning is still positive, as shown in Fig. 4. Such a performance trend demonstrates that recent studies of long-tailed learning make real progress. Moreover, the RA of the state-of-the-art TADE is 93.0, which implies that there is still room for improvement in the future.

TABLE 5
Results on ImageNet-LT in terms of accuracy (Acc), upper reference accuracy (UA), relative accuracy (RA) under 90 or 200 training epochs. In this table, CR, IA and MI indicate class re-balancing, information augmentation and module improvement, respectively.

| Type | Method | 90 epochs | | | 200 epochs | | |
|---|---|---|---|---|---|---|---|
| | | Acc | UA | RA | Acc | UA | RA |
| Baseline | Softmax | 45.5 | 57.3 | 79.4 | 46.8 | 57.8 | 81.0 |
| CR | Weighted Softmax | 47.9 | 57.3 | 83.6 | 49.1 | 57.8 | 84.9 |
| | Focal loss [68] | 45.8 | 57.3 | 79.9 | 47.2 | 57.8 | 81.7 |
| | LDAM [18] | 51.1 | 57.3 | 89.2 | 51.1 | 57.8 | 88.4 |
| | ESQL [19] | 47.3 | 57.3 | 82.5 | 48.0 | 57.8 | 83.0 |
| | UNO-IC [87] | 45.7 | 57.3 | 81.4 | 46.8 | 58.6 | 79.9 |
| | Balanced Softmax [86] | 50.8 | 57.3 | 88.7 | 51.2 | 57.8 | 88.6 |
| | LADE [31] | 51.5 | 57.8 | 89.1 | 51.6 | 57.8 | 89.3 |
| IA | SSP [89] | 53.1 | 59.6 | 89.1 | 53.3 | 59.9 | 89.0 |
| | RSG [99] | 49.6 | 57.3 | 86.7 | 52.9 | 57.8 | 91.5 |
| MI | OLTR [15] | 46.7 | 57.3 | 81.5 | 48.0 | 58.4 | 82.2 |
| | PaCo [109] | 52.7 | 58.7 | 89.9 | 54.4 | 59.6 | 91.3 |
| | De-confound [88] | 51.8 | 57.7 | 89.8 | 51.3 | 57.8 | 88.8 |
| | Decouple-IB-CRT [32] | 49.9 | 57.3 | 87.1 | 50.3 | 58.1 | 86.6 |
| | Decouple-CB-CRT [32] | 44.9 | 57.3 | 78.4 | 43.0 | 57.8 | 74.4 |
| | Decouple-SR-CRT [32] | 49.3 | 57.3 | 86.0 | 48.5 | 57.8 | 83.9 |
| | Decouple-PB-CRT [32] | 48.4 | 57.3 | 84.5 | 48.1 | 57.8 | 83.2 |
| | MiSLAS [96] | 51.4 | 58.3 | 88.2 | 53.4 | 59.7 | 89.4 |
| | BBN [48] | 41.2 | 57.3 | 71.9 | 44.7 | 57.8 | 77.3 |
| | LFME [84] | 47.0 | 57.3 | 82.0 | 48.0 | 57.8 | 83.0 |
| | ResLT [152] | 51.6 | 57.3 | 90.1 | 53.2 | 58.1 | 91.6 |
| | RIDE [17] | 55.5 | 60.2 | 92.2 | 56.1 | 60.9 | 92.1 |
| | TADE [30] | **57.3** | **61.9** | **92.6** | **58.8** | **63.2** | **93.0** |

TABLE 6
Accuracy results on ImageNet-LT regarding head, middle and tail classes under 90 or 200 training epochs. In this table, WS indicates weighed softmax and BS indicates balanced softmax. The types of methods are the same to Table 5.

| Method | 90 epochs | | | 200 epochs | | |
|---|---|---|---|---|---|---|
| | Head | Middle | Tail | Head | Middle | Tail |
| Softmax | 66.5 | 39.0 | 8.6 | 66.9 | 40.4 | 12.6 |
| WS | 66.3 | 42.2 | 15.6 | 57.9 | 46.2 | 34.0 |
| Focal loss [68] | 66.9 | 39.2 | 9.2 | 67.0 | 41.0 | 13.1 |
| LDAM [18] | 62.3 | 47.4 | 32.5 | 60.0 | 49.2 | 31.9 |
| ESQL [19] | 62.5 | 44.0 | 15.7 | 63.1 | 44.6 | 17.2 |
| UNO-IC [87] | 66.3 | 38.7 | 9.3 | 67.0 | 40.3 | 12.7 |
| BS [86] | 61.7 | 48.0 | 29.9 | 62.4 | 47.7 | 32.1 |
| LADE [31] | 62.2 | 48.6 | 31.8 | 63.1 | 47.7 | 32.7 |
| SSP [89] | 65.6 | 49.6 | 30.3 | 67.3 | 49.1 | 28.3 |
| RSG [99] | **68.7** | 43.7 | 16.2 | 65.0 | 49.4 | 31.1 |
| OLTR [15] | 58.2 | 45.5 | 19.5 | 62.9 | 44.6 | 18.8 |
| PaCo [109] | 59.7 | 51.7 | 36.6 | 63.2 | 51.6 | 39.2 |
| De-confound [88] | 63.0 | 48.5 | 31.4 | 64.9 | 46.9 | 28.1 |
| IB-CRT [32] | 62.6 | 46.2 | 26.7 | 64.2 | 46.1 | 26.0 |
| CB-CRT [32] | 62.4 | 39.3 | 14.9 | 60.9 | 36.9 | 13.5 |
| SR-CRT [32] | 64.1 | 43.9 | 19.5 | 66.0 | 42.3 | 18.0 |
| PB-CRT [32] | 63.9 | 45.0 | 23.2 | 64.9 | 43.1 | 20.6 |
| MiSLAS [96] | 62.1 | 48.9 | 32.6 | 65.3 | 50.6 | 33.0 |
| BBN [48] | 40.0 | 43.3 | 40.8 | 43.3 | 45.9 | **43.7** |
| LFME [84] | 60.6 | 43.5 | 22.0 | 64.1 | 42.3 | 22.8 |
| ResLT [152] | 57.8 | 50.4 | 40.0 | 61.6 | 51.4 | 38.8 |
| RIDE [17] | 66.9 | 52.3 | 34.5 | **67.9** | 52.3 | 36.0 |
| TADE [30] | 65.3 | **55.2** | **42.0** | 67.2 | **55.3** | 40.0 |

In addition, we also evaluate the influence of different training epochs (*i.e.,* 90 and 200) on model training in Table 5. Overall, training with 200 epochs leads to better performance for most long-tailed methods, because sufficient training enables deep models to fit data better and learn better visual representations. However, there are also some methods that perform better when only training 90 epochs, *e.g.,* De-confound and Decouple-CB-CRT. We speculate that, for these methods, 90 epochs are enough to train models well, while training more epochs does not bring additional benefits but increases the training difficulties since it also influences the learning rate decay scheme.

**Observations on different method types.** We also discuss the results in Table 5 from the perspective of different method types. To begin with, almost all class re-balancing (CB) methods is beneficial to long-tailed learning performance, compared to the baseline model. Specifically, LADE, Balanced Softmax and LDAM achieve state-of-the-art in this method type. Moreover, Focal loss was proposed to handle imbalanced object detection [68]; however, when handling an extremely large number of long-tailed classes (*e.g.,* 1,000 in ImageNet-LT), Focal loss cannot perform well and only leads to marginal improvement. In LDAM, there is a deferred re-balancing optimization schedule in addition to the cost-sensitive LDAM loss; note that simply learning with the LDAM loss without the deferred scheme may not achieve promising performance. In addition, as shown in Table 5, the upper reference accuracy of most cost-sensitive learning methods are the same, so the relative accuracy is positively correlated to accuracy (c.f. Fig. 4 (b)). Hence, the accuracy improvement in this method type can accurately reflect the alleviation of class imbalance.

In the method type of information augmentation (IA), both transfer learning (SSP) and data augmentation (RSG) help to handle long-tailed class imbalance. Although SSP also improves upper reference accuracy, the relative accuracy is increased more

significantly, which implies that the performance gain mostly comes from the handling of class imbalance. Considering that lack of enough tail-class samples is one of the key challenges, IA is worth further exploring by bringing more information into training.

In module improvement (MI), all sub-categories of methods contribute to handling class imbalance. By now, the state of the art is ensemble-based long-tailed methods, *i.e.,* TADE [30] and RIDE [17], in terms of both accuracy and relative accuracy. Although ensemble learning also improves upper reference accuracy, the performance gain from handling imbalance is more significant, *i.e.,* higher relative accuracy (c.f. Fig. 4 (d)).

## 4.4 Results on Class Subsets

This section provides the method performance on different class subsets. As shown in Table 6, almost all methods improve tail-class and middle-class performance at the cost of head-class performance. The head classes, however, are also important in long-tailed learning, so it is necessary to improve long-tailed performance without sacrificing the performance on the head. Potential solutions include information augmentation and ensemble learning, *e.g.,* SSP, RIDE and TADE (with sufficient training, *e.g.,* 200 epochs).

By comparing both Tables 5 and 6, one can find that the overall performance gain largely depends on the improvement of middle and tail classes; hence, how to improve their performance is still the most important goal of long-tailed learning in the future.

By now, TADE [30] achieves the best overall performance in terms of accuracy and RA (c.f. Table 5), but TADE does not perform state-of-the-art on all class subsets (c.f. Table 6). For example, when training 200 epochs, the head-class performance of TADE is worse than RIDE and its tail-class performance is worse than BBN. To summarize, the higher average performance of TADE implies that the key to obtaining better long-tailed performance is a better trade-off among all classes.
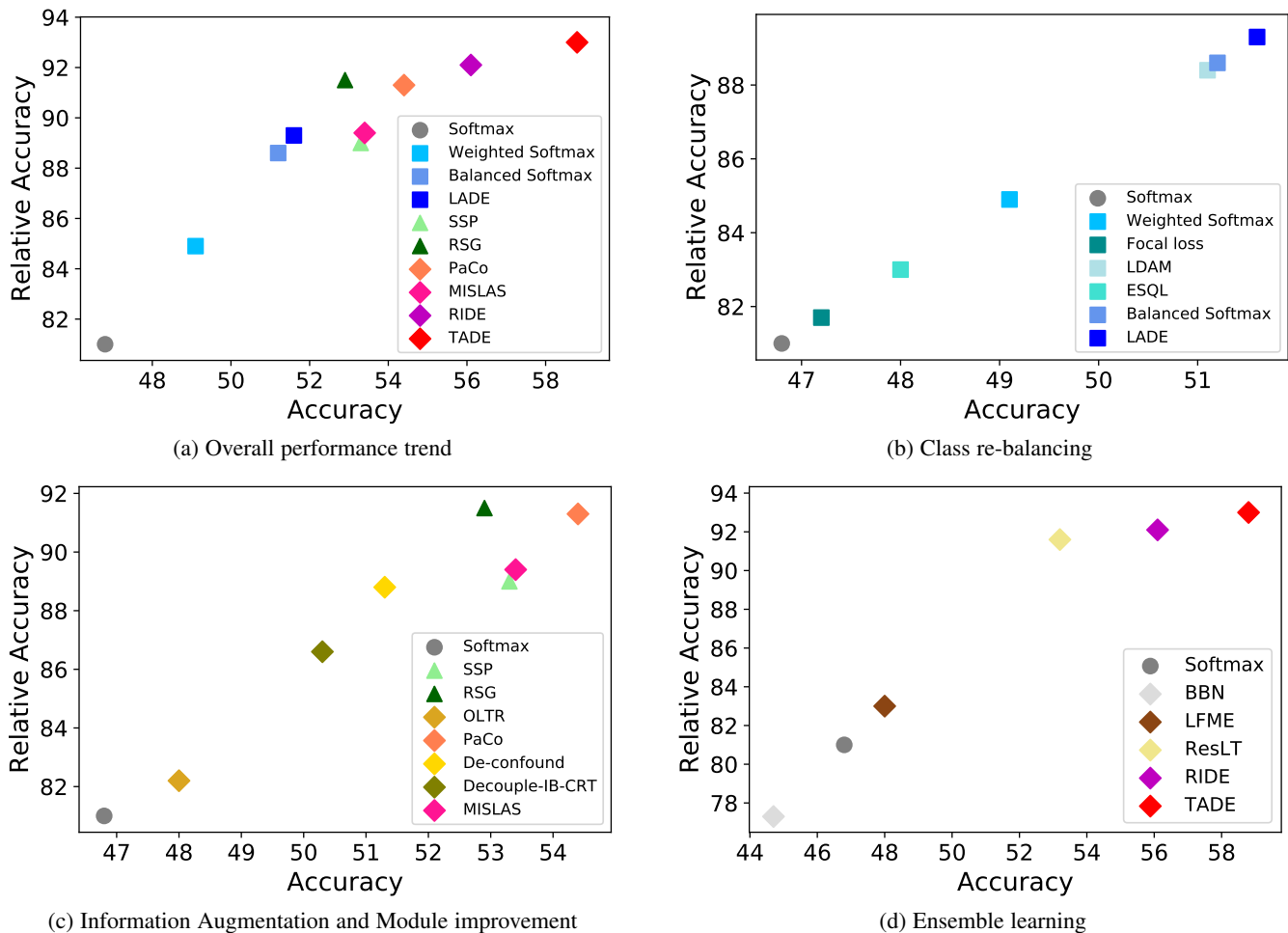
Fig. 4. Performance trends of long-tailed learning methods in terms of accuracy and relative accuracy under 200 epochs. Fig. (a) includes approaches from multiple method types; Figs. (b) refers to class re-balancing; Figs. (c) includes information augmentation and partial module improvement methods; Figs. (d) refers to ensemble learning. Here, the shape of ○ indicates the softmax baseline; □ indicates class re-balancing; △ and ◇ are information augmentation and module improvement methods, respectively. Different colors, consistent in all sub-figures, represent different methods.

In summary, the current best practice for deep long-tailed learning is using ensemble learning and class re-balancing, simultaneously. Note that all these methods, apart from data augmentation based methods, only use basic augmentation operations. If using stronger data augmentation, *e.g.,* RandAugment [153] or Cutmix [154], the model performance can be further improved.

### 4.5 More Discussions on Cost-sensitive Losses

In this section, we further evaluate the performance of different cost-sensitive learning losses based on the decoupled training scheme [32] that decouples representation and classifier learning into two stages. In the first stage, we use different cost-sensitive learning losses to train the model backbone for learning representations, while in the second stage, we use four different strategies for classifier training [32], *i.e.,* joint training without re-training, nearest class mean classifier (NCM), class-balanced classifier re-training (CRT), and learnable weight scaling (LWS).

As shown in Table 7, decoupled training, compared to joint training, can further improve the overall performance of most cost-sensitive methods apart from balanced softmax (BS). These methods under decoupled training can obtain comparable performance to BS that performs the best under joint training. Such results are particularly interesting, as they imply that although these cost-sensitive losses perform differently under joint training, they essentially learn similar quality of feature representations.

TABLE 7
The decoupled training performance of various cost-sensitive losses under 200 training epochs on ImageNet-LT. Here, "Joint" indicates one-stage end-to-end joint training; "NCM" is the nearest class mean classifier [32]; "CRT" represents class-balanced classifier re-training [32]; "LWS" means learnable weight scaling [32]. Moreover, BS indicates the balanced softmax method [86].

| Test Dist. | Accuracy on **all** classes | | | | Accuracy on **head** classes | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Joint | NCM | CRT | LWS | Joint | NCM | CRT | LWS |
| Softmax | 46.8 | 50.2 | 50.2 | 50.8 | 66.9 | 63.5 | 65.0 | 64.6 |
| Focal loss [68] | 47.2 | 50.7 | 50.7 | 51.5 | 67.0 | 62.6 | 64.5 | 64.3 |
| ESQL [19] | 48.0 | 49.8 | 50.6 | 50.5 | 63.1 | 60.2 | 64.0 | 63.3 |
| BS [86] | 51.2 | 50.4 | 50.6 | 51.1 | 62.4 | 62.4 | 64.9 | 64.3 |

| Test Dist. | Accuracy on **middle** classes | | | | Accuracy on **tail** classes | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Joint | NCM | CRT | LWS | Joint | NCM | CRT | LWS |
| Softmax | 40.4 | 45.8 | 45.3 | 46.1 | 12.6 | 28.1 | 25.5 | 28.2 |
| Focal loss [68] | 41.0 | 47.0 | 46.4 | 47.3 | 13.1 | 30.1 | 26.9 | 30.2 |
| ESQL [19] | 44.6 | 46.6 | 46.5 | 46.1 | 17.2 | 31.1 | 27.1 | 29.5 |
| BS [86] | 47.7 | 46.8 | 46.1 | 46.7 | 32.1 | 29.1 | 26.2 | 29.4 |

The worse overall performance of BS under decoupled training than joint training may imply that BS has conducted class re-balancing very well; further using classifier re-training for re-balancing does not bring additional benefits but even degenerate the consistency of network parameters by end-to-end joint training. Similar observations can be found in head-class, middle-class and tail-class performance.

# 5 APPLICATIONS

This section discusses the main visual applications of deep long-tailed learning, including image classification, image detection and segmentation, and visual relation learning.

## 5.1 Image Classification

The most common application of long-tailed learning is multi-class classification [15], [30], [32], [48], [88]. As mentioned in Section 2.2, there are many artificially sampled long-tailed datasets from widely-used image classification datasets, *i.e.,* ImageNet, CIFAR, and Places. Based on these datasets, various long-tailed learning methods have been proposed, as shown in Section 3. Besides these artificial tasks, long-tailed learning is also applied to real-world image classification tasks, including species classification [23], [24], [155], face recognition [21], [22], [74], [80], age classification [94], logo detection [156], rail surface defect detection [157] and medical image diagnosis [25], [158].

In addition to multi-class classification, long-tailed learning is also applied to multi-label classification based on both artificial tasks [37], [90] (*i.e.,* VOC-LT and COCO-LT) and real-world tasks, including web image classification [82], face attribute classification [69] and cloth attribute classification [69].

## 5.2 Image Detection and Segmentation

Object detection and instance segmentation has attracted increasing attention in the long-tailed learning community [19], [35], [68], [78], [103], [126], [159], [160], where most existing studies are conducted based on LVIS and COCO. In addition to these widely-used benchmarks, many other applications have also been explored, including urban scene understanding [26], [161], unmanned aerial vehicle detection [27], point cloud segmentation [162], [163].

## 5.3 Visual Relation Learning

Visual relation learning is important for image understanding and is attracting rising attention in the long-tailed learning community. One important application is long-tailed scene graph generation [111], [164]. In the future, long-tailed visual question answering and image captioning are worth exploring [165], [166].

# 6 FUTURE DIRECTIONS

In this section, we identify several future research directions for deep long-tailed learning from both perspectives of method innovation and task innovation.

## 6.1 New Methodology

We first discuss several potential directions for innovating deep long-tailed learning methods.

**Class re-balancing without label frequencies.** Some real-world long-tailed tasks, *e.g.,* multi-label classification or object detection, may suffer an additional issue besides class imbalance, *i.e.,* label co-occurrence. Specifically, the co-occurrence of labels indicates the situation that head-class labels frequently appear with tail-class labels, which may bias the imbalance degree during model training and make it difficult to obtain exact label frequencies. Considering this issue, existing class re-balancing methods based on label frequencies tend to fail. How to handle this issue in long-tailed learning is an open question.

**Transfer learning with unlabeled data.** One key challenge in long-tailed learning is the lack of enough tail-class samples. Transferring the knowledge from other unlabeled samples is a feasible solution, *e.g.,* self-supervised learning, knowledge distillation and self-training. Existing transfer methods, however, may not handle long-tailed learning very well. For example, CReST [97] found the supervised trained model often has high precision on tail classes in long-tailed image classification and thus proposed to select more tail-class data for pseudo labeling and model training. Such a finding, however, may not hold in long-tailed object detection or multi-label classification. Hence, how to better use unlabeled data for long-tailed learning is worth further exploring.

**Data augmentation for multiple tasks.** Existing long-tailed methods are often designed for a specific task, *e.g.,* image classification or image detection. However, due to the differences among various tasks, existing methods for a specific task may not handle other tasks, leading to poor method generalization. Considering data augmentation is fundamental for all visual tasks, it is valuable to design better augmentation-based long-tailed methods that can resolve multiple long-tailed tasks, simultaneously.

**Ensemble learning for improving all classes.** Most existing long-tailed methods improve tail-class performance at the cost of head-class performance. One solution is ensemble learning, which exploits different expertise of various experts to obtain a better trade-off between head and tail classes, leading to state-of-the-art performance on long-tailed learning [30]. Thanks to potential performance improvement in all classes, ensemble learning would be a promising direction for future research.

## 6.2 New Task Settings

In addition to method innovation, there are several new task settings of long-tailed learning waiting to be resolved.

**Test-agnostic long-tailed learning.** Existing long-tailed learning methods generally hypothesize a balanced test class distribution. The practical test distribution, however, often violates this hypothesis (*e.g.,* being long-tailed or even inversely long-tailed), which may lead existing methods to fail in real-world applications. To overcome this limitation, LADE [31] relaxes this hypothesis by assuming that the test class distribution can be skewed arbitrarily but the prior of test distribution is available. Afterward, TADE [30] further innovates the task, in which the test class distribution is not only arbitrarily skewed but also unknown. Besides class imbalance, this task poses another challenge, *i.e.,* unidentified class distribution shift between the training and test samples.

**Open-set long-tailed learning.** Real-world samples often have a long-tailed and open-ended class distribution. Inspired by this, open-set long-tailed learning [15], [81] seeks to learn from long-tailed data and optimize the classification accuracy over a balanced test set that includes head, tail and open classes. There are two main challenges: (1) how to share visual knowledge between head and tail classes; (2) how to reduce confusion between the tail and open classes.

**Federated long-tailed learning.** Existing long-tailed learning studies generally assume that all the training samples are accessible during model training. However, in real-world applications, long-tailed training data may be distributed on numerous mobile devices or the Internet of Things [167], which requires decentralized training of deep models. Such a task setting is called federated long-tailed learning, which has two main challenges: (1) long-tail class imbalance; (2) unknown class distribution shift among different clients' local data.

**Class-incremental long-tailed learning.** In real-world applications, long-tailed data may come in a continual and class-incremental manner [82], [168]. To deal with this scenario, class-incremental long-tailed learning aims to learn deep models from class-incremental long-tailed data, suffering two key challenges: (1) how to handle long-tailed class imbalance when different classes come sequentially, and the model has no information about the future input regarding classes as well as label frequencies; (2) how to overcome catastrophic forgetting of previous class knowledge when learning new classes. Such a task setting can also be named continual long-tailed learning.

**Multi-domain long-tailed learning.** Current long-tailed methods generally assume that all long-tailed samples come from the same data marginal distribution. However, in practice, long-tailed data may also get from different domains with distinct data distributions [28], [169], *e.g.,* the DomainNet dataset [170]. Motivated by this, multi-domain long-tailed learning seeks to handle both class imbalance and domain distribution shift, simultaneously. One more challenging issue may be the inconsistency of class imbalance among different domains. In other words, various domains may have different class distributions, which further enlarges the domain shift in multi-domain long-tailed learning.

**Robust long-tailed learning.** Real-world long-tailed samples may also suffer image noise [95], [171] or label noise [140], [145]. Most long-tailed methods, however, assume all images and labels are clean, leading to poor model robustness in practical applications. This issue would be particularly severe for tail classes, as they have very limited training samples. Inspired by this, robust long-tailed learning seeks to handle class imbalance and improve model robustness, simultaneously.

**Long-tailed regression.** Most existing studies of long-tailed visual learning focus on classification, detection and segmentation, which have discrete labels with class indices. However, many tasks involve continuous labels, where hard classification boundaries among classes do not exist. Motivated by this, long-tailed regression [172] aims to deal with long-tailed learning with continuous label space. In such a task, how to simultaneously resolve long-tailed class imbalance and handle potential missing data for certain labels remains an open question.

**Long-tailed video learning.** Most existing deep long-tailed learning studies focus on the image level, but ignore that the video domain also suffers the issue of long-tail class imbalance. Considering the additional temporal dimension in video data, long-tailed video learning should be more difficult than long-tailed image learning. Thanks to the recent release of a VideoLT dataset [38], long-tailed video learning can be explored in the near future.

# 7 CONCLUSION

In this survey, we have extensively reviewed classic deep long-tailed learning methods proposed before mid-2021, according to the taxonomy of class re-balancing, information augmentation and module improvement. We have empirically analyzed several state-of-the-art long-tailed methods by evaluating to what extent they address the issue of class imbalance, based on a newly proposed relative accuracy metric. Following that, we discussed the main application scenarios of long-tailed learning, and identified potential innovation directions for methods and task settings. We expect that this timely survey not only provides a better understanding of long-tailed learning for researchers and the community, but also facilitates future research.

## REFERENCES

[1] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.
[2] I. Goodfellow, Y. Bengio, and A. Courville, *Deep learning.* MIT press, 2016.
[3] A. Voulodimos, N. Doulamis, A. Doulamis, and E. Protopapadakis, "Deep learning for computer vision: A brief review," *Computational Intelligence and Neuroscience*, 2018.
[4] C. Dong, C. C. Loy, K. He, and X. Tang, "Image super-resolution using deep convolutional networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 38, no. 2, pp. 295–307, 2015.
[5] Z. Wang, J. Chen, and S. C. Hoi, "Deep learning for image super-resolution: A survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
[6] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," *Advances in Neural Information Processing Systems*, vol. 25, pp. 1097–1105, 2012.
[7] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: towards real-time object detection with region proposal networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 6, pp. 1137–1149, 2016.
[8] E. Shelhamer, J. Long, and T. Darrell, "Fully convolutional networks for semantic segmentation." *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 4, pp. 640–651, 2016.
[9] Y. Bengio, Y. LeCun, and G. Hinton, "Deep learning for ai," *Communications of the ACM*, vol. 64, no. 7, pp. 58–65, 2021.
[10] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Computer Vision and Pattern Recognition*, 2016.
[11] C. Szegedy, A. Toshev, and D. Erhan, "Deep neural networks for object detection," 2013.
[12] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *Computer Vision and Pattern Recognition*, 2014, pp. 580–587.
[13] B. Kang, Y. Li, S. Xie, Z. Yuan, and J. Feng, "Exploring balanced feature spaces for representation learning," in *International Conference on Learning Representations*, 2021.
[14] 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.
[15] Z. Liu, Z. Miao, X. Zhan, J. Wang, B. Gong, and S. X. Yu, "Large-scale long-tailed recognition in an open world," in *Computer Vision and Pattern Recognition*, 2019, pp. 2537–2546.
[16] Y. Cui, M. Jia, T.-Y. Lin, Y. Song, and S. Belongie, "Class-balanced loss based on effective number of samples," in *Computer Vision and Pattern Recognition*, 2019, pp. 9268–9277.
[17] X. Wang, L. Lian, Z. Miao, Z. Liu, and S. X. Yu, "Long-tailed recognition by routing diverse distribution-aware experts," in *International Conference on Learning Representations*, 2021.
[18] 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.
[19] J. Tan, C. Wang, B. Li, Q. Li, W. Ouyang, C. Yin, and J. Yan, "Equalization loss for long-tailed object recognition," in *Computer Vision and Pattern Recognition*, 2020, pp. 11 662–11 671.
[20] V. Vapnik, "Principles of risk minimization for learning theory," in *Advances in Neural Information Processing Systems*, 1992, pp. 831–838.
[21] X. Zhang, Z. Fang, Y. Wen, Z. Li, and Y. Qiao, "Range loss for deep face recognition with long-tailed training data," in *International Conference on Computer Vision*, 2017, pp. 5409–5418.
[22] D. Cao, X. Zhu, X. Huang, J. Guo, and Z. Lei, "Domain balancing: Face recognition on long-tailed domains," in *Computer Vision and Pattern Recognition*, 2020, pp. 5671–5679.
[23] 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 *Computer Vision and Pattern Recognition*, 2018, pp. 8769–8778.
[24] Z. Miao, Z. Liu, K. M. Gaynor, M. S. Palmer, S. X. Yu, and W. M. Getz, "Iterative human and automated identification of wildlife images," *arXiv:2105.02320*, 2021.

[25] L. Ju, X. Wang, L. Wang, T. Liu, X. Zhao, T. Drummond, D. Mahapatra, and Z. Ge, "Relational subsets knowledge distillation for long-tailed retinal diseases recognition," *arXiv:2104.11057*, 2021.

[26] R. He, J. Yang, and X. Qi, "Re-distributing biased pseudo labels for semi-supervised semantic segmentation: A baseline investigation," in *International Conference on Computer Vision*, 2021.

[27] W. Yu, T. Yang, and C. Chen, "Towards resolving the challenge of long-tail distribution in uav images for object detection," in *IEEE Winter Conference on Applications of Computer Vision*, 2021, pp. 3258–3267.

[28] M. A. Jamal, M. Brown, M.-H. Yang, L. Wang, and B. Gong, "Rethinking class-balanced methods for long-tailed visual recognition from a domain adaptation perspective," in *Computer Vision and Pattern Recognition*, 2020, pp. 7610–7619.

[29] S. Zhang, Z. Li, S. Yan, X. He, and J. Sun, "Distribution alignment: A unified framework for long-tail visual recognition," in *Computer Vision and Pattern Recognition*, 2021, pp. 2361–2370.

[30] Y. Zhang, B. Hooi, L. Hong, and J. Feng, "Test-agnostic long-tailed recognition by test-time aggregating diverse experts with self-supervision," *arXiv:2107.09249*, 2021.

[31] Y. Hong, S. Han, K. Choi, S. Seo, B. Kim, and B. Chang, "Disentangling label distribution for long-tailed visual recognition," in *Computer Vision and Pattern Recognition*, 2021.

[32] 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.

[33] C. Feng, Y. Zhong, and W. Huang, "Exploring classification equilibrium in long-tailed object detection," in *International Conference on Computer Vision*, 2021.

[34] T. Wang, Y. Li, B. Kang, J. Li, J. Liew, S. Tang, S. Hoi, and J. Feng, "The devil is in classification: A simple framework for long-tail instance segmentation," in *European Conference on Computer Vision*, 2020.

[35] Z. Weng, M. G. Ogut, S. Limonchik, and S. Yeung, "Unsupervised discovery of the long-tail in instance segmentation using hierarchical self-supervision," in *Computer Vision and Pattern Recognition*, 2021.

[36] A. Gupta, P. Dollar, and R. Girshick, "Lvis: A dataset for large vocabulary instance segmentation," in *Computer Vision and Pattern Recognition*, 2019, pp. 5356–5364.

[37] T. Wu, Q. Huang, Z. Liu, Y. Wang, and D. Lin, "Distribution-balanced loss for multi-label classification in long-tailed datasets," in *European Conference on Computer Vision*, 2020, pp. 162–178.

[38] X. Zhang, Z. Wu, Z. Weng, H. Fu, J. Chen, Y.-G. Jiang, and L. Davis, "Videolt: Large-scale long-tailed video recognition," in *International Conference on Computer Vision*, 2021.

[39] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *Computer Vision and Pattern Recognition*, 2009, pp. 248–255.

[40] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," 2009.

[41] B. Zhou, A. Lapedriza, J. Xiao, A. Torralba, and A. Oliva, "Learning deep features for scene recognition using places database," *Advances in Neural Information Processing Systems*, vol. 27, pp. 487–495, 2014.

[42] M. Everingham, S. A. Eslami, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes challenge: A retrospective," *International Journal of Computer Vision*, vol. 111, no. 1, pp. 98–136, 2015.

[43] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *European Conference on Computer Vision*, 2014.

[44] H. Rezatofighi, N. Tsoi, J. Gwak, A. Sadeghian, I. Reid, and S. Savarese, "Generalized intersection over union: A metric and a loss for bounding box regression," in *Computer Vision and Pattern Recognition*, 2019.

[45] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, "Aggregated residual transformations for deep neural networks," in *Computer Vision and Pattern Recognition*, 2017, pp. 1492–1500.

[46] K. He, G. Gkioxari, P. Dollar, and R. Girshick, "Mask r-cnn," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 42, no. 2, pp. 386–397, 2020.

[47] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," in *Computer Vision and Pattern Recognition*, 2017, pp. 2117–2125.

[48] B. Zhou, Q. Cui, X.-S. Wei, and Z.-M. Chen, "Bbn: Bilateral-branch network with cumulative learning for long-tailed visual recognition," in *Computer Vision and Pattern Recognition*, 2020, pp. 9719–9728.

[49] H. He and E. A. Garcia, "Learning from imbalanced data," *IEEE Transactions on Knowledge and Data Engineering*, vol. 21, no. 9, pp. 1263–1284, 2009.

[50] J. Snell, K. Swersky, and R. Zemel, "Prototypical networks for few-shot learning," *Advances in Neural Information Processing Systems*, 2017.

[51] F. Sung, Y. Yang, L. Zhang, T. Xiang, P. H. Torr, and T. M. Hospedales, "Learning to compare: Relation network for few-shot learning," in *Computer Vision and Pattern Recognition*, 2018, pp. 1199–1208.

[52] Q. Sun, Y. Liu, T.-S. Chua, and B. Schiele, "Meta-transfer learning for few-shot learning," in *Computer Vision and Pattern Recognition*, 2019.

[53] Y. Wang, Q. Yao, J. T. Kwok, and L. M. Ni, "Generalizing from a few examples: A survey on few-shot learning," *ACM Computing Surveys*, vol. 53, no. 3, pp. 1–34, 2020.

[54] D. Krueger, E. Caballero *et al.*, "Out-of-distribution generalization via risk extrapolation," in *International Conference on Machine Learning*, 2021, pp. 5815–5826.

[55] Z. Shen, J. Liu, Y. He, X. Zhang, R. Xu, H. Yu, and P. Cui, "Towards out-of-distribution generalization: A survey," *arXiv:2108.13624*, 2021.

[56] S. J. Pan, I. W. Tsang, J. T. Kwok, and Q. Yang, "Domain adaptation via transfer component analysis," *IEEE Transactions on Neural Networks*, vol. 22, no. 2, pp. 199–210, 2010.

[57] E. Tzeng, J. Hoffman, K. Saenko, and T. Darrell, "Adversarial discriminative domain adaptation," in *Computer Vision and Pattern Recognition*, 2017, pp. 7167–7176.

[58] Y. Zhang, H. Chen, Y. Wei, P. Zhao, J. Cao, X. Fan, X. Lou, H. Liu, J. Hou, X. Han *et al.*, "From whole slide imaging to microscopy: Deep microscopy adaptation network for histopathology cancer image classification," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2019, pp. 360–368.

[59] Y. Zhang, Y. Wei *et al.*, "Collaborative unsupervised domain adaptation for medical image diagnosis," *IEEE Transactions on Image Processing*, 2020.

[60] Z. Qiu, Z. Zhang, H. Lin, S. Niu, Y. Liu, Q. Du, and M. Tan, "Source-free domain adaptation via avatar prototype generation and adaptation," in *International Joint Conference on Artificial Intelligence*, 2021.

[61] H. Wu, H. Zhu, Y. Yan, J. Wu, Y. Zhang, and M. K. Ng, "Heterogeneous domain adaptation by information capturing and distribution matching," *IEEE Transactions on Image Processing*, vol. 30, pp. 6364–6376, 2021.

[62] D. Li, Y. Yang, Y.-Z. Song, and T. M. Hospedales, "Deeper, broader and artier domain generalization," in *International Conference on Computer Vision*, 2017, pp. 5542–5550.

[63] H. Li, S. J. Pan, S. Wang, and A. C. Kot, "Domain generalization with adversarial feature learning," in *Computer Vision and Pattern Recognition*, 2018, pp. 5400–5409.

[64] L. Neal, M. Olson, X. Fern, W.-K. Wong, and F. Li, "Open set learning with counterfactual images," in *European Conference on Computer Vision*, 2018, pp. 613–628.

[65] Y. Fu, X. Wang, H. Dong, Y.-G. Jiang, M. Wang, X. Xue, and L. Sigal, "Vocabulary-informed zero-shot and open-set learning," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 42, no. 12, pp. 3136–3152, 2019.

[66] C. Huang, Y. Li, C. C. Loy, and X. Tang, "Learning deep representation for imbalanced classification," in *Computer Vision and Pattern Recognition*, 2016.

[67] W. Ouyang, X. Wang, C. Zhang, and X. Yang, "Factors in finetuning deep model for object detection with long-tail distribution," in *Computer Vision and Pattern Recognition*, 2016, pp. 864–873.

[68] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *International Conference on Computer Vision*, 2017, pp. 2980–2988.

[69] Q. Dong, S. Gong, and X. Zhu, "Class rectification hard mining for imbalanced deep learning," in *International Conference on Computer Vision*, 2017, pp. 1851–1860.

[70] Y.-X. Wang, D. Ramanan, and M. Hebert, "Learning to model the tail," in *Advances in Neural Information Processing Systems*, 2017.

[71] Y. Cui, Y. Song, C. Sun, A. Howard, and S. Belongie, "Large scale fine-grained categorization and domain-specific transfer learning," in *Computer Vision and Pattern Recognition*, 2018, pp. 4109–4118.

[72] S. Khan, M. Hayat, S. W. Zamir, J. Shen, and L. Shao, "Striking the right balance with uncertainty," in *Computer Vision and Pattern Recognition*, 2019, pp. 103–112.

[73] X. Yin, X. Yu, K. Sohn, X. Liu, and M. Chandraker, "Feature transfer learning for face recognition with under-represented data," in *Computer Vision and Pattern Recognition*, 2019, pp. 5704–5713.

[74] 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 *Computer Vision and Pattern Recognition*, 2019, pp. 7812–7821.

[75] Y. Wang, W. Gan, J. Yang, W. Wu, and J. Yan, "Dynamic curriculum learning for imbalanced data classification," in *International Conference on Computer Vision*, 2019, pp. 5017–5026.
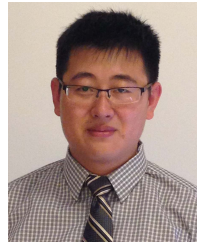
[76] J. Shu, Q. Xie, L. Yi, Q. Zhao, S. Zhou, Z. Xu, and D. Meng, "Meta-weight-net: Learning an explicit mapping for sample weighting," *Advances in Neural Information Processing Systems*, 2019.

[77] X. Hu, Y. Jiang, K. Tang, J. Chen, C. Miao, and H. Zhang, "Learning to segment the tail," in *Computer Vision and Pattern Recognition*, 2020.

[78] Y. Li, T. Wang, B. Kang, S. Tang, C. Wang, J. Li, and J. Feng, "Overcoming classifier imbalance for long-tail object detection with balanced group softmax," in *Computer Vision and Pattern Recognition*, 2020, pp. 10 991–11 000.

[79] J. Kim, J. Jeong, and J. Shin, "M2m: Imbalanced classification via major-to-minor translation," in *Computer Vision and Pattern Recognition*, 2020.

[80] J. Liu, Y. Sun, C. Han, Z. Dou, and W. Li, "Deep representation learning on long-tailed data: A learnable embedding augmentation perspective," in *Computer Vision and Pattern Recognition*, 2020.

[81] L. Zhu and Y. Yang, "Inflated episodic memory with region self-attention for long-tailed visual recognition," in *Computer Vision and Pattern Recognition*, 2020, pp. 4344–4353.

[82] C. D. Kim, J. Jeong, and G. Kim, "Imbalanced continual learning with partitioning reservoir sampling," in *European Conference on Computer Vision*, 2020, pp. 411–428.

[83] P. Chu, X. Bian, S. Liu, and H. Ling, "Feature space augmentation for long-tailed data," in *European Conference on Computer Vision*, 2020.

[84] 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*, 2020, pp. 247–263.

[85] T.-Y. Wu, P. Morgado, P. Wang, C.-H. Ho, and N. Vasconcelos, "Solving long-tailed recognition with deep realistic taxonomic classifier," in *European Conference on Computer Vision*, 2020, pp. 171–189.

[86] 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*, 2020.

[87] J. Tian, Y.-C. Liu, N. Glaser, Y.-C. Hsu, and Z. Kira, "Posterior re-calibration for imbalanced datasets," in *Advances in Neural Information Processing Systems*, 2020.

[88] K. Tang, J. Huang, and H. Zhang, "Long-tailed classification by keeping the good and removing the bad momentum causal effect," *Advances in Neural Information Processing Systems*, vol. 33, 2020.

[89] Y. Yang and Z. Xu, "Rethinking the value of labels for improving class-imbalanced learning," in *Advances in Neural Information Processing Systems*, 2020.

[90] H. Guo and S. Wang, "Long-tailed multi-label visual recognition by collaborative training on uniform and re-balanced samplings," in *Computer Vision and Pattern Recognition*, 2021, pp. 15 089–15 098.

[91] J. Tan, X. Lu, G. Zhang, C. Yin, and Q. Li, "Equalization loss v2: A new gradient balance approach for long-tailed object detection," in *Computer Vision and Pattern Recognition*, 2021, pp. 1685–1694.

[92] J. Wang, W. Zhang, Y. Zang, Y. Cao, J. Pang, T. Gong, K. Chen, Z. Liu, C. C. Loy, and D. Lin, "Seesaw loss for long-tailed instance segmentation," in *Computer Vision and Pattern Recognition*, 2021.

[93] T. Wang, Y. Zhu, C. Zhao, W. Zeng, J. Wang, and M. Tang, "Adaptive class suppression loss for long-tail object detection," in *Computer Vision and Pattern Recognition*, 2021, pp. 3103–3112.

[94] Z. Deng, H. Liu, Y. Wang, C. Wang, Z. Yu, and X. Sun, "Pml: Progressive margin loss for long-tailed age classification," in *Computer Vision and Pattern Recognition*, 2021, pp. 10 503–10 512.

[95] T. Wu, Z. Liu, Q. Huang, Y. Wang, and D. Lin, "Adversarial robustness under long-tailed distribution," in *Computer Vision and Pattern Recognition*, 2021, pp. 8659–8668.

[96] Z. Zhong, J. Cui, S. Liu, and J. Jia, "Improving calibration for long-tailed recognition," in *Computer Vision and Pattern Recognition*, 2021.

[97] C. Wei, K. Sohn, C. Mellina, A. Yuille, and F. Yang, "Crest: A class-rebalancing self-training framework for imbalanced semi-supervised learning," in *Computer Vision and Pattern Recognition*, 2021.

[98] S. Changpinyo, P. Sharma, N. Ding, and R. Soricut, "Conceptual 12m: Pushing web-scale image-text pre-training to recognize long-tail visual concepts," in *Computer Vision and Pattern Recognition*, 2021.

[99] J. Wang, T. Lukasiewicz, X. Hu, J. Cai, and Z. Xu, "Rsg: A simple but effective module for learning imbalanced datasets," in *Computer Vision and Pattern Recognition*, 2021, pp. 3784–3793.

[100] S. Li, K. Gong, C. H. Liu, Y. Wang, F. Qiao, and X. Cheng, "Metasaug: Meta semantic augmentation for long-tailed visual recognition," in *Computer Vision and Pattern Recognition*, 2021, pp. 5212–5221.

[101] P. Wang, K. Han, X.-S. Wei, L. Zhang, and L. Wang, "Contrastive learning based hybrid networks for long-tailed image classification," in *Computer Vision and Pattern Recognition*, 2021, pp. 943–952.

[102] B. Liu, H. Li, H. Kang, G. Hua, and N. Vasconcelos, "Gistnet: a geometric structure transfer network for long-tailed recognition," in *International Conference on Computer Vision*, 2021.

[103] Y. Zang, C. Huang, and C. C. Loy, "Fasa: Feature augmentation and sampling adaptation for long-tailed instance segmentation," in *International Conference on Computer Vision*, 2021.

[104] J. Cai, Y. Wang, and J.-N. Hwang, "Ace: Ally complementary experts for solving long-tailed recognition in one-shot," in *International Conference on Computer Vision*, 2021.

[105] S. Park, J. Lim, Y. Jeon, and J. Y. Choi, "Influence-balanced loss for imbalanced visual classification," in *International Conference on Computer Vision*, 2021.

[106] T. Li, L. Wang, and G. Wu, "Self supervision to distillation for long-tailed visual recognition," in *International Conference on Computer Vision*, 2021.

[107] Y.-Y. He, J. Wu, and X.-S. Wei, "Distilling virtual examples for long-tailed recognition," in *International Conference on Computer Vision*, 2021.

[108] C. Zhang, T.-Y. Pan, Y. Li, H. Hu, D. Xuan, S. Changpinyo, B. Gong, and W.-L. Chao, "Mosaicos: A simple and effective use of object-centric images for long-tailed object detection," in *International Conference on Computer Vision*, 2021.

[109] J. Cui, Z. Zhong, S. Liu, B. Yu, and J. Jia, "Parametric contrastive learning," in *International Conference on Computer Vision*, 2021.

[110] D. Samuel and G. Chechik, "Distributional robustness loss for long-tail learning," in *International Conference on Computer Vision*, 2021.

[111] A. Desai, T.-Y. Wu, S. Tripathi, and N. Vasconcelos, "Learning of visual relations: The devil is in the tails," in *International Conference on Computer Vision*, 2021.

[112] 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.

[113] A. Estabrooks, T. Jo, and N. Japkowicz, "A multiple resampling method for learning from imbalanced data sets," *Computational Intelligence*, vol. 20, no. 1, pp. 18–36, 2004.

[114] 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*, 2005, pp. 878–887.

[115] X.-Y. Liu, J. Wu, and Z.-H. Zhou, "Exploratory undersampling for class-imbalance learning," *IEEE Transactions on Systems, Man, and Cybernetics*, vol. 39, no. 2, pp. 539–550, 2008.

[116] Z. Zhang and T. Pfister, "Learning fast sample re-weighting without reward data," in *International Conference on Computer Vision*, 2021.

[117] D. Mahajan, R. Girshick, V. Ramanathan, K. He, M. Paluri, Y. Li, A. Bharambe, and L. Van Der Maaten, "Exploring the limits of weakly supervised pretraining," in *European conference on computer vision*, 2018, pp. 181–196.

[118] A. Hermans, L. Beyer, and B. Leibe, "In defense of the triplet loss for person re-identification," *arXiv:1703.07737*, 2017.

[119] C. Elkan, "The foundations of cost-sensitive learning," in *International Joint Conference on Artificial Intelligence*, 2001.

[120] Z.-H. Zhou and X.-Y. Liu, "Training cost-sensitive neural networks with methods addressing the class imbalance problem," *IEEE Transactions on Knowledge and Data Engineering*, vol. 18, no. 1, pp. 63–77, 2005.

[121] Y. Sun, M. S. Kamel, A. K. Wong, and Y. Wang, "Cost-sensitive boosting for classification of imbalanced data," *Pattern Recognition*, vol. 40, no. 12, pp. 3358–3378, 2007.

[122] P. Zhao, Y. Zhang, M. Wu, S. C. Hoi, M. Tan, and J. Huang, "Adaptive cost-sensitive online classification," *IEEE Transactions on Knowledge and Data Engineering*, vol. 31, no. 2, pp. 214–228, 2018.

[123] Y. Zhang, P. Zhao, J. Cao, W. Ma, J. Huang, Q. Wu, and M. Tan, "Online adaptive asymmetric active learning for budgeted imbalanced data," in *SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2018, pp. 2768–2777.

[124] Y. Zhang, P. Zhao, S. Niu, Q. Wu, J. Cao, J. Huang, and M. Tan, "Online adaptive asymmetric active learning with limited budgets," *IEEE Transactions on Knowledge and Data Engineering*, 2019.

[125] H.-J. Ye, H.-Y. Chen, D.-C. Zhan, and W.-L. Chao, "Identifying and compensating for feature deviation in imbalanced deep learning," *arXiv:2001.01385*, 2020.

[126] T.-I. Hsieh, E. Robb, H.-T. Chen, and J.-B. Huang, "Droploss for long-tail instance segmentation," in *AAAI Conference on Artificial Intelligence*, vol. 35, no. 2, 2021, pp. 1549–1557.

[127] F. Wang, J. Cheng, W. Liu, and H. Liu, "Additive margin softmax for face verification," *IEEE Signal Processing Letters*, vol. 25, no. 7, pp. 926–930, 2018.

[128] V. Koltchinskii and D. Panchenko, "Empirical margin distributions and bounding the generalization error of combined classifiers," *The Annals of Statistics*, vol. 30, no. 1, pp. 1–50, 2002.

[129] F. Provost, "Machine learning from imbalanced data sets 101," in *AAAI Workshop on Imbalanced Data Sets*, vol. 68, no. 2000, 2000, pp. 1–3.

[130] S. J. Pan and Q. Yang, "A survey on transfer learning," *IEEE Transactions on Knowledge and Data Engineering*, vol. 22, no. 10, pp. 1345–1359, 2009.

[131] C. Tan, F. Sun, T. Kong, W. Zhang, C. Yang, and C. Liu, "A survey on deep transfer learning," in *International Conference on Artificial Neural Networks*, 2018, pp. 270–279.

[132] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba, "Learning deep features for discriminative localization," in *Computer Vision and Pattern Recognition*, 2016, pp. 2921–2929.

[133] D. Erhan, A. Courville, Y. Bengio, and P. Vincent, "Why does unsupervised pre-training help deep learning?" in *International Conference on Artificial Intelligence and Statistics*, 2010, pp. 201–208.

[134] K. He, R. Girshick, and P. Dollár, "Rethinking imagenet pre-training," in *International Conference on Computer Vision*, 2019, pp. 4918–4927.

[135] D. Hendrycks, K. Lee, and M. Mazeika, "Using pre-training can improve model robustness and uncertainty," in *International Conference on Machine Learning*, 2019, pp. 2712–2721.

[136] B. Zoph, G. Ghiasi, T.-Y. Lin, Y. Cui, H. Liu, E. D. Cubuk, and Q. Le, "Rethinking pre-training and self-training," *Advances in Neural Information Processing Systems*.

[137] Y. Zhang, B. Hooi, D. Hu, J. Liang, and J. Feng, "Unleashing the power of contrastive self-supervised visual models via contrast-regularized fine-tuning," in *Advances in Neural Information Processing Systems*, 2021.

[138] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," in *Computer Vision and Pattern Recognition*, 2020.

[139] S. Gidaris, P. Singh, and N. Komodakis, "Unsupervised representation learning by predicting image rotations," in *International Conference on Learning Representations*, 2018.

[140] S. Karthik, J. Revaud, and C. Boris, "Learning from long-tailed data with noisy labels," *arXiv:2108.11096*, 2021.

[141] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," *arXiv:1503.02531*, 2015.

[142] J. Gou, B. Yu, S. J. Maybank, and D. Tao, "Knowledge distillation: A survey," *International Journal of Computer Vision*, vol. 129, no. 6, pp. 1789–1819, 2021.

[143] X. J. Zhu, "Semi-supervised learning literature survey," 2005.

[144] C. Rosenberg, M. Hebert, and H. Schneiderman, "Semi-supervised self-training of object detection models," 2005.

[145] T. Wei, J.-X. Shi, W.-W. Tu, and Y.-F. Li, "Robust long-tailed learning under label noise," *arXiv:2108.11569*, 2021.

[146] L. Perez and J. Wang, "The effectiveness of data augmentation in image classification using deep learning," *arXiv:1712.04621*, 2017.

[147] C. Shorten and T. M. Khoshgoftaar, "A survey on image data augmentation for deep learning," *Journal of Big Data*, vol. 6, no. 1, pp. 1–48, 2019.

[148] H.-P. Chou, S.-C. Chang, J.-Y. Pan, W. Wei, and D.-C. Juan, "Remix: Rebalanced mixup," in *European Conference on Computer Vision Workshop*, 2020, pp. 95–110.

[149] 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*, vol. 32, 2019, pp. 12 635–12 644.

[150] J. Goh and M. Sim, "Distributionally robust optimization and its tractable approximations," *Operations Research*, vol. 58, no. 4-part-1, pp. 902–917, 2010.

[151] T. Cover and P. Hart, "Nearest neighbor pattern classification," *IEEE transactions on information theory*, vol. 13, no. 1, pp. 21–27, 1967.

[152] J. Cui, S. Liu, Z. Tian, Z. Zhong, and J. Jia, "Reslt: Residual learning for long-tailed recognition," *arXiv:2101.10633*, 2021.

[153] E. D. Cubuk, B. Zoph, J. Shlens, and Q. Le, "Randaugment: Practical automated data augmentation with a reduced search space," *Advances in Neural Information Processing Systems*, vol. 33, 2020.

[154] S. Yun, D. Han, S. J. Oh, S. Chun, J. Choe, and Y. Yoo, "Cutmix: Regularization strategy to train strong classifiers with localizable features," in *International Conference on Computer Vision*, 2019.

[155] M. R. Keaton, R. J. Zaveri, M. Kovur, C. Henderson, D. A. Adjeroh, and G. Doretto, "Fine-grained visual classification of plant species in the wild: Object detection as a reinforced means of attention," *arXiv:2106.02141*, 2021.

[156] X. Jia, H. Yan, Y. Wu, X. Wei, X. Cao, and Y. Zhang, "An effective and robust detector for logo detection," *arXiv:2108.00422*, 2021.

[157] Z. Zhang, S. Yu, S. Yang, Y. Zhou, and B. Zhao, "Rail-5k: a real-world dataset for rail surface defects detection," *arXiv:2106.14366*, 2021.

[158] A. Galdran, G. Carneiro, and M. A. G. Ballester, "Balanced-mixup for highly imbalanced medical image classification," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 2021.

[159] T. Weyand, A. Araujo, B. Cao, and J. Sim, "Google landmarks dataset v2-a large-scale benchmark for instance-level recognition and retrieval," in *Computer Vision and Pattern Recognition*, 2020, pp. 2575–2584.

[160] J. Wu, L. Song, T. Wang, Q. Zhang, and J. Yuan, "Forest r-cnn: Large-vocabulary long-tailed object detection and instance segmentation," in *ACM International Conference on Multimedia*, 2020, pp. 1570–1578.

[161] J. Mao, M. Niu, C. Jiang, H. Liang, X. Liang, Y. Li, C. Ye, W. Zhang, Z. Li, J. Yu *et al.*, "One million scenes for autonomous driving: Once dataset," in *NeurIPS 2021 Datasets and Benchmarks Track*, 2021.

[162] Y. Zhang, Z. Zhou, P. David, X. Yue, Z. Xi, B. Gong, and H. Foroosh, "Polarnet: An improved grid representation for online lidar point clouds semantic segmentation," in *Computer Vision and Pattern Recognition*, 2020, pp. 9601–9610.

[163] X. Chen, C. Zhang, G. Lin, and J. Han, "Compositional prototype network with multi-view comparision for few-shot point cloud semantic segmentation," *arXiv:2012.14255*, 2020.

[164] N. Dhingra, F. Ritter, and A. Kunz, "Bgt-net: Bidirectional gru transformer network for scene graph generation," in *Computer Vision and Pattern Recognition*, 2021, pp. 2150–2159.

[165] J. Chen, A. Agarwal, S. Abdelkarim, D. Zhu, and M. Elhoseiny, "Reltransformer: Balancing the visual relationship detection from local context, scene and memory," *arXiv:2104.11934*, 2021.

[166] Z. Li, E. Stengel-Eskin, Y. Zhang, C. Xie, Q. Tran, B. Van Durme, and A. Yuille, "Calibrating concepts and operations: Towards symbolic reasoning on real images," in *International Conference on Computer Vision*, 2021.

[167] M. Luo, F. Chen, D. Hu, Y. Zhang, J. Liang, and J. Feng, "No fear of heterogeneity: Classifier calibration for federated learning with non-iid data," in *Advances in Neural Information Processing Systems*, 2021.

[168] S. Niu, J. Wu, G. Xu, Y. Zhang, Y. Guo, P. Zhao, P. Wang, and M. Tan, "Adaxpert: Adapting neural architecture for growing data," in *International Conference on Machine Learning*, 2021, pp. 8184–8194.

[169] Y. Zhang, S. Niu, Z. Qiu, Y. Wei, P. Zhao, J. Yao, J. Huang, Q. Wu, and M. Tan, "Covid-da: Deep domain adaptation from typical pneumonia to covid-19," *arXiv:2005.01577*, 2020.

[170] X. Peng, Q. Bai, X. Xia, Z. Huang, K. Saenko, and B. Wang, "Moment matching for multi-source domain adaptation," in *International Conference on Computer Vision*, 2019, pp. 1406–1415.

[171] K. Cao, Y. Chen, J. Lu, N. Arechiga, A. Gaidon, and T. Ma, "Heteroskedastic and imbalanced deep learning with adaptive regularization," in *International Conference on Learning Representations*, 2021.

[172] Y. Yang, K. Zha, Y.-C. Chen, H. Wang, and D. Katabi, "Delving into deep imbalanced regression," in *International Conference on Machine Learning*, 2021.

**Yifan Zhang** is working toward the Ph.D. degree in the School of Computing, National University of Singapore. He received the M.E. degree from South China University of Technology in 2020. His research interests are broadly in machine learning, now with high self-motivation to solve generalization problems of deep neural networks. He has published papers in top venues, including NeurIPS, SIGKDD, ICML, IJCAI, TIP, and TKDE. He has been invited as a PC member or reviewer for top international conferences and journals, including CVPR, NeurIPS, ICML, ICLR, AAAI, TIP, IJCV, and TNNLS.

**Jiashi Feng** is currently Principal Scientist at SEA AI Lab and was an Assistant Professor with the department of Electrical and Computer Engineering, National University of Singapore. He was a Postdoctoral Researcher with University of California at Berkeley from 2014 to 2015. He received the Ph.D. degree from National University of Singapore in 2014. His current research interests include machine learning and computer vision techniques for large-scale data analysis. He 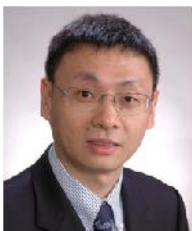received the best student paper award from ACM MM 2018, best paper award from TASK-CV ICCV 2015, and best technical demo award from ACM MM 2012. He was also the recipient of Innovators Under 35 Asia, MIT Technology Review 2018.

**Bingyi Kang** got his Ph.D degree in Electronic and Computer Engineering from National University of Singapore. He received his B.E. degree in automation from Zhejiang University, Hangzhou, Zhejiang in 2016. His current research interest focuses on sample-efficient learning and reinforcement learning.

**Bryan Hooi** is an Assistant Professor in the Computer Science Department and the Institute of Data Science in National University of Singapore. He received his PhD degree in Machine Learning from Carnegie Mellon University, USA in 2019. His research interests include machine learning on graph-structured data, robustness and novelty detection, and spatiotemporal data mining. His work aims to develop efficient and practical data mining algorithms, with applications including fraud detection, online commerce, and automatic monitoring of medical, industrial, weather and environmental sensor data.

**Shuicheng Yan** is currently the director of Sea AI Lab and group chief scientist of Sea. He is an IEEE Fellow, ACM Fellow, IAPR Fellow, and Fellow of Academy of Engineering, Singapore. His research areas include computer vision, machine learning and multimedia analysis. Till now, he has published over 1,000 papers in top international journals and conferences, with Google Scholar Citation over 79,000 times and H-index 121. He had been among "Thomson Reuters Highly Cited Researchers" in 2014, 2015, 2016, 2018, 2019. His team has received winner or honorable-mention prizes for 10 times of two core competitions, Pascal VOC and ImageNet (ILSVRC), which are deemed as "World Cup" in the computer vision community. Also, his team won over 10 best paper or best student paper prizes and especially, a grand slam in ACM MM, the top conference in multimedia, including Best Paper Award, Best Student Paper Award and Best Demo Award.