# Difficulty-Net: Learning to Predict Difficulty for Long-Tailed Recognition

Saptarshi Sinha
Hitachi Ltd.
Tokyo, Japan
saptarshi.sinha.hx@hitachi.com

Hiroki Ohashi
Hitachi Ltd.
Tokyo, Japan
hiroki.ohashi.uo@hitachi.com

## Abstract

*Long-tailed datasets, where head classes comprise much more training samples than tail classes, cause recognition models to get biased towards the head classes. Weighted loss is one of the most popular ways of mitigating this issue, and a recent work has suggested that class-difficulty might be a better clue than conventionally used class-frequency to decide the distribution of weights. A heuristic formulation was used in the previous work for quantifying the difficulty, but we empirically find that the optimal formulation varies depending on the characteristics of datasets. Therefore, we propose Difficulty-Net, which learns to predict the difficulty of classes using the model's performance in a meta-learning framework. To make it learn reasonable difficulty of a class within the context of other classes, we newly introduce two key concepts, namely the relative difficulty and the driver loss. The former helps Difficulty-Net take other classes into account when calculating difficulty of a class, while the latter is indispensable for guiding the learning to a meaningful direction. Extensive experiments on popular long-tailed datasets demonstrated the effectiveness of the proposed method, and it achieved state-of-the-art performance on multiple long-tailed datasets.*

## 1. Introduction

Despite the outstanding performance of the recent deep learning (DL) models on public datasets, deploying such models in the real world often leads to a performance drop. One of the causes is that the public datasets are usually almost perfectly class-balanced while real-world data are generally long-tailed, where a few classes (called *head classes*) consist of a significantly larger number of training samples than the rest of the classes (called *tail classes*). The 'long-tailed recognition' research domain particularly aims at addressing this issue.

Amongst multiple possible strategies to tackle long-tailed recognition problems, cost-sensitive learning is one of the most popular and promising strategies. Most cost-
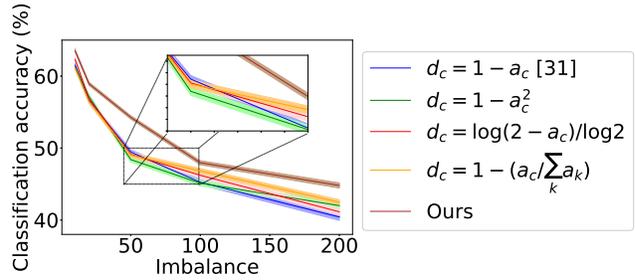


Figure 1. Different quantifications of class-difficulties perform better than others in different situations (imbalance ratios). The imbalance in the data is calculated as the ratio of the frequency of the most frequent class to that of the least frequent class. We compute class-wise difficulty ($d_c$) using four different functions of class-wise accuracy ($a_c$) for the CDB-CE [31] loss function and compare their performance on the CIFAR100-LT. Interestingly, the alternate formulations work better than the originally proposed one ($d_c = 1 - a_c$) in many cases. However, the best performing function changes with the imbalance values. This brings us to the question "Which formulation to choose for my imbalanced dataset?"

sensitive learning techniques modify the cost function to penalize the model differently for different samples. This modification is generally done by scaling the cost value using different weights, and the research direction is mainly aimed at finding an effective weight-assignment strategy. One simple and intuitive way is to assign weights using the inverse of the class-frequencies. Recently, more sophisticated approaches such as class-balanced loss [7] and equalization loss [36] have been proposed. However, most of these approaches give more weights to the tail classes because they assume that tail classes are always the most difficult to learn. Recently Sinha *et al.* [31, 32] empirically showed that the above assumption does not always hold true and further claimed that class-difficulty might be a better clue to decide weights.

While they proposed an intuitive quantification of the class-level difficulties, this quantification is preliminarily determined regardless of the property of a given dataset, and thus may not be optimal in different situations. In

fact, we empirically found that multiple quantifications for class-wise difficulty gave comparable or even better results than [31] as shown in Figure 1. This adds the extra tedious task of selecting the appropriate formulation for a given imbalanced data.

Motivated by recently proposed Meta-Weight-Net (MWN) [29], our research aims to address the above issue by meta-learning a simple model, named Difficulty-Net, to predict class-level difficulty scores and then dynamically distribute the weights based on the scores. Such a strategy removes dependence on any prior formulation for class-wise difficulty and lets the model learn any suitable function to compute it. The key difference with MWN is three folds. First, while MWN is a sample-level weighting method, ours is a class-level weighted approach, whose advantage in long-tailed recognition has been revealed in [31] and also discussed in Sec. 3.2 and Sec. 4.5. Second, we propose to use relative difficulties rather than absolute difficulties that are used in prior works [29, 31] so that the other classes' difficulties are also taken into account when determining the difficulty of a class. Third, we propose a new loss function that drives the learning process of Difficulty-Net in a reasonable direction, without which the performance turned out to degrade.

To summarize, our key contributions are:

- We propose Difficulty-Net, which learns to predict class-difficulty in a meta-learning framework.

- We argue that relative difficulty is more important and effective than absolute difficulty, and provide an empirical evidence for the argument.

- We propose a new loss function, called driver loss, that guides the learning process in a reasonable direction.

- We conducted extensive experiments on multiple long-tail benchmark datasets and achieved state-of-the-art results. In addition, we provide in-depth analysis on the effect and property of the proposed method in comparison to previous works, which revealed the effectiveness of our method.

## 2. Related works

Major strategies to tackle the long-tailed recognition can be broadly categorized as data re-sampling methods [3, 13, 17, 30, 34, 32], metric learning [15, 20, 33, 40], knowledge transfer [24, 43, 44], mixture of experts [41, 45, 50], cost-sensitive learning [2, 7, 25, 29, 32] and decoupled learning [19, 28, 48, 51].

Data re-sampling techniques [3, 13, 17, 30, 34] try to neutralize the long-tail by under-sampling from head classes or over-sampling from tail classes. Under-sampling [34] generally results in poor representation of the head classes, while a straight-forward over-sampling strategy of replicating tail-class samples causes the model to overfit on the repeated samples. Another popular over-sampling technique is synthetic data generation [3, 13] for the tail classes. Certain class-balanced sampling approaches such as class-aware sampling [17, 19, 42] and square-root sampling [19] have been shown to be more effective than using over- or under-sampling. They typically try to increase the sampling rate for the tail classes during training. However, they still result in overfitting due to the repeated sampling of the same samples from tail classes.

Metric learning methods [1, 15, 20, 33, 40, 49] aim to learn a high-quality feature extractor that preserves inter-class and intra-class relationships in the feature space. They achieve this by learning from pairs [20, 33, 49] or triplets [1, 15] of input samples. Metric learning has been used in long-tailed recognition [24, 38] in hope that high-quality feature extractor will mitigate the imbalance between head and tail classes. An effective sampling of the sample groups is the key for efficient training in this scheme. However, such sampling strategies come with the risk of under-representation or overfitting, as explained above.

Knowledge transfer [8, 24, 43, 44, 47] in long-tailed recognition tries to transfer knowledge gained from the head classes to the tail classes. They achieve this either by learning modular transformations from few-shot model parameters to many-shot models [43, 44] or by designing external modules for feature transfer [24, 47]. Designing such modules is usually computationally expensive in real-world usecases [19].

Mixture of experts (MoE) [41, 45, 50, 52] is an ensemble-based technique where the expert models are trained to gain diverse knowledge. The aggregated knowledge of the experts is either used directly to alleviate the long-tail [41, 50] or used to teach a student model for that purpose [45]. Despite the increasing popularity of this domain, our research focuses on the improvement of a single model as it can then easily be combined with any mixture.

Cost-sensitive learning can be achieved by logit-adjustment loss [2, 25, 37] and weighted loss [7, 11, 23, 27, 29, 35, 36] approaches. Most prior methods distribute these adjustment or weight values on the basis of class-frequencies. Recently, Sinha *et al.* [31, 32] showed that class-difficulty is a better metric for the purpose. However, finding an optimal formulation for calculating class-difficulty is not a trivial task as the optimal formulation usually varies depending on datasets as shown in Figure 1. Our research builds on the work of Sinha *et al.* [31] and tries to remove the requirement of any prior formulation by using meta-learning. Meta-learning [12] has previously been used in long-tailed recognition to learn [18, 27] or predict [29] sample weights. The closest to our research is Meta-Weight-Net (MWN) [29], which learns a model to predict sample-level weights from training loss. Different

from them, we use class-level weighting, which is known to be better than sample-level weighting in long-tailed recognition [31].

Recently, Kang *et al.* [19] found that decoupling model learning into representation learning and classifier learning helps long-tailed recognition. Since this finding, most works have tried to improve either the data-representation [28, 46] or the classifier [22, 39, 48, 51]. We show that our method also benefits from this framework and achieves state-of-the-art results based on it.

## 3. Proposed method

### 3.1. Background

Generally, the prior assumption is that the tail classes are the most difficult to learn for the models. However, it has recently been empirically shown that the number of training instances of a class might not be the best clue to determine its difficulty because some classes are well-represented even with fewer training samples. On the basis of this finding, Sinha *et al.* [31] came up with a simple formulation to directly calculate the difficulty of a class from the model's performance. The formulation says that if the model's classification accuracy on a class $c$ is $a_c$, then the difficulty of the class $d_c$ can be computed as $d_c = 1 - a_c$.

However, we found two lacking points in the formulation. First, as stated in Sec. 1, we found multiple decreasing functions of accuracy $a_c$ that outperformed the above formulation in multiple setups. In the meantime, we also found that the best performing formulation varies inconsistently with the data imbalance and thus it is not possible to preliminarily define the best formulation for a given dataset. Second, while the above formulation helps to compute the absolute difficulty of a class, we believe it is more important to compute the difficulty of a class *relative* to the other classes because it is reasonable to assign a high difficulty score to a class with high accuracy (*i.e.* easy class) *if* the other classes have even higher accuracies. For that purpose, all the classes need to be considered when computing the difficulty of a single class, which is not done in [31].

To address these issues, we propose meta-learning the formulation that is most effective for a given dataset, taking relative difficulties into consideration.

### 3.2. Meta-learning via Difficulty-Net

**Difficulty-Net design.** Given a dataset of $C$ classes, we aim to learn a formulation that can compute the relative difficulty for each class. For this purpose, we design the formulation for class-wise difficulty as

$$d_1, d_2, \ldots, d_C = D(\{a_c\}_{c=1}^C; \theta). \qquad (1)$$

Note that both $d_c$ and $a_c$ change as the training progresses, but here we omit the notation of training steps for simplicity.

$D$ is a neural network with parameters $\theta$. In our implementation, we choose $D$ to be a simple MLP model with two hidden layers. The output layer dimension is kept same as the number of classes and a sigmoid activation at the output ensures the difficulty scores to be in the range $(0, 1)$. The input to $D$ are the model's classification accuracies for all $C$ classes. The design of $D$ ensures that while estimating the difficulty for a class, the model's performance on the other classes is also taken into account. We refer to $D$ as 'Difficulty-Net'.

**Meta-learning objective:** Suppose a classification problem in which we are provided a training dataset $S^{train} = \{x_i, y_i\}_{i=1}^N$, where $x_i$ is the $i^{th}$ training sample and $y_i \in \{1, \ldots, C\}$ is its corresponding ground truth label. Given a classifier neural network $f(x; \phi)$ with learnable parameters $\phi$, our primary objective is to learn the optimal parameters $\phi^*$ so that $f(x; \phi^*)$ provides the minimum classification loss on the training set $S^{train}$, *i.e.*

$$\phi^* = \arg \min_\phi \frac{1}{N} \sum_{i=1}^N L(f(x_i; \phi), y_i), \qquad (2)$$

where $L$ computes the loss corresponding to $f$'s prediction for a given sample and is typically the cross-entropy loss.

In long-tailed recognition, the training dataset $S^{train}$ is class-imbalanced. In such cases, optimization using Eq. 2 leads to biased learning of $\phi$. To compensate for the imbalance, we modify the learning objective as most weighted loss approaches do, *i.e.*

$$\phi^* = \arg \min_\phi \frac{1}{N} \sum_{i=1}^N w_i L(f(x_i; \phi), y_i), \qquad (3)$$

where $w_i$ is the weight assigned to the training sample $x_i$. In our proposed approach, $w_i$ is computed by Difficulty-Net $D$ as

$$w_i(A_C, \theta) = D(A_C; \theta)_{y_i}, \qquad (4)$$

where $D(A_C; \theta)_{y_i}$ is the difficulty score for class $y_i$ predicted by Difficulty-Net and $A_C = \{a_c\}_{c=1}^C$ is the set of accuracies of $f(x; \phi)$ for all $C$ classes, evaluated prior to this calculation. Therefore, the learning objective for $\phi^*$ is modified as

$$\phi^*(\theta) = \arg \min_\phi \frac{1}{N} \sum_{i=1}^N w_i(A_C, \theta) L(f(x_i; \phi), y_i). \qquad (5)$$

Since the optimization of our main classifier network depends on the effectiveness of Difficulty-Net, it is important to optimize the parameters $\theta$ of $D$ as well. Inspired by [29], we use a small balanced meta-dataset $S^{meta} =$

$\{x_i^{meta}, y_i^{meta}\}_{i=1}^M$ for optimizing the parameters $\theta$ as follows.

$$\theta^* = \arg\min_\theta \frac{1}{M}\sum_{i=1}^M L_i^{meta}(\phi^*(\theta))$$

$$= \arg\min_\theta \frac{1}{M}\sum_{i=1}^M L(f(x_i^{meta};\phi^*(\theta)), y_i^{meta}). \quad (6)$$

However, we found that $L^{meta}$ alone is not enough for Difficulty-Net to learn to estimate difficulties from accuracy. Even for 2 classes with very different accuracy values, Difficulty-Net learned using Eq. 6 tends to give similar difficulty scores for both classes. To address this issue, we add another loss component to drive the learning of Difficulty-Net in a practically correct direction. We call this loss 'driver loss' and calculate it as

$$L^{dr}(A_C, \theta) = \frac{1}{C}\sum_{c=1}^C \left((1-\hat{a}_c) - D(A_C;\theta)_c\right)^2, \quad (7)$$

where $\hat{a}_c = a_c/\sum_k a_k$ is the normalized accuracy of class $c$. The $L^{dr}$ is built on the motivation that Difficulty-Net should learn to give high difficulty scores to a class, if the accuracy of the class is relatively low. Now the parameters $\theta$ of Difficulty-Net $D$ are optimized as

$$\theta^* = \arg\min_\theta \lambda L^{dr}(A_C, \theta) + \frac{1}{M}\sum_{i=1}^M L_i^{meta}(\phi^*(\theta)), \quad (8)$$

where $\lambda$ is a hyper-parameter controlling the influence of $L^{dr}$. Note that too high value of $\lambda$ will simply cause $D$ to always predict class difficulties as $d_c = 1 - \hat{a}_c$. We ablate over various values of $\lambda$ in our experiments.

**Learning method.** Following [29], our meta-learning method is a 3-step process. Given a classifier network $f(x;\phi_t)$ and Difficulty-Net $D(;\theta_t)$ at time step $t$, the first step aims to learn intermediate classifier parameter $\hat{\phi}_t$ by

$$\hat{\phi}_t(\theta_t) \leftarrow \phi_t - \alpha\frac{1}{b}\sum_{i=1}^b w_i(A_{C,t}, \theta_t)\frac{\partial L(f(x_i;\phi), y_i)}{\partial \phi}\bigg|_{\phi_t}, \quad (9)$$

where $\alpha$ is the step size for gradient descent and $b$ is the number of samples in one mini-batch sampled from the training set $S^{train}$. $A_{C,t}$ is the classification accuracy of $f(x;\phi_t)$ on all the $C$ classes at time step $t$ and is computed on a validation dataset $S^{val}$.

The second step updates the parameters $\theta$ of Difficulty-Net using the obtained intermediate classifier $f(x;\hat{\phi}_t)$ on a mini-batch of size $m$ sampled from the meta-dataset. The

update is done by

$$\theta_{t+1} \leftarrow$$
$$\theta_t - \beta\frac{\partial(\lambda L^{dr}(A_{C,t}, \theta) + \frac{1}{m}\sum_{i=1}^m L_i^{meta}(\hat{\phi}_t(\theta)))}{\partial\theta}\bigg|_{\theta_t}, \quad (10)$$

where $\beta$ is the step size for updating the parameters of Difficulty-Net.

Finally, the third step uses the updated parameters $\theta_{t+1}$ to update the parameters of the classifier network $f(x;\phi_t)$ over the same mini-batch sampled in Eq. 9.

$$\phi_{t+1} \leftarrow \phi_t - \alpha\frac{1}{b}\sum_{i=1}^b w_i(A_{C,t}, \theta_{t+1})\frac{\partial L(f(x_i;\phi), y_i)}{\partial\phi}\bigg|_{\phi_t}. \quad (11)$$

The above three steps are executed iteratively till convergence or the end of the training. The overall algorithm is presented in Algorithm 1 in the supplementary material.

For our experiments, we construct the $S^{meta}$ following exactly the same procedure as [18, 29]. We also found that $S^{meta}$ is reusable as $S^{val}$ for calculating $A_{C,t}$. Therefore, we do not use any extra data compared to previous methods. Also, although it is ideal to calculate $A_{C,t}$ for every time step $t$, we calculate the accuracy only after every epoch in our implementation for saving computational time.

**Difference with MWN and CDB-CE.** Although we share a similar meta-learning framework as MWN [29], our approach is very different from theirs in more than one way. One difference is that MWN is a sample-level weighting strategy while ours is a class-level weighting strategy. The advantages of class-level weighting over sample-level weighting in long-tailed learning is pointed out in [31] and also reflected in our experimental results.

Ours is not the straight-forward combination of MWN [29] and CDB-CE [31]. First, both MWN and CDB-CE use absolute difficulties of a sample or a class to determine the weights. We believe, however, the relative difficulty compared with other samples or classes is more important because it is reasonable to assign a high difficulty score to a class with high accuracy (*i.e.* easy class) *if* the other classes have even higher accuracies. The proposed method estimates relative difficulties of each class amongst all the classes, and it turned out to be more effective as we will show in Sec. 4.5. Second, the straight-forward combination of these prior works without the driver loss turns out to learn almost nothing and predicts almost identical difficulties for all the classes as we will show in Sec. 4.5. The newly proposed driver loss is essential to guide the training in a reasonable direction.

The empirical evidences of these arguments are provided in Sec. 4.5 and 4.6.

# 4. Experiments

## 4.1. Datasets

**CIFAR100-LT.** CIFAR100 [21] is an object-centric balanced classification data-set comprised of tiny images belonging to 100 different classes. Long-tailed versions of the dataset are artificially created by reducing the training samples per class according to an exponential function as given in [7]. Following [7], we use CIFAR100-LT with imbalance varying in 10–200.

**ImageNet-LT.** ImageNet-LT is a long-tail version of ImageNet [9] created by [24]. It contains 1000 object categories with heavy imbalance of 256. We use the same *train, val* and *test* splits as [24].

**Places-LT.** Places-2 [53] is a large-scale scene-centric image dataset, used for scene recognition tasks. Places-LT is a long-tailed subset of Places-2 with 365 classes and imbalance of 996, created by [24]. We use the same splits as [24].

For constructing $S^{meta}$, we followed the setup of previous meta-learning based methods [18, 29] to ensure the fair comparison. Please see the supplementary material for the details. The evaluation results are reported on balanced test sets.

## 4.2. Implementation details

Following previous long-tailed works [7, 24, 36], we use ResNet-32 [14] for CIFAR-100-LT experiments. On ImageNet-LT, we follow [19, 24, 51] and use ResNet-10 [14], ResNet-50 [14]. As in [51], we use pretrained (on ImageNet [9]) ResNet-152 [14] and finetune it on Places-LT. The basic architecture of Difficulty-Net is the same for all the datasets as explained in Sec. 3.2, *i.e.* MLP with two hidden layers, but we change the dimension of hidden and output layers as different datasets have different number of classes. We will explain a simple way to select the dimension of hidden layers in the supplementary material. We evaluate our method both in end-to-end (e2e) learning and decoupled learning [19] settings. For using Difficulty-Net in decoupled learning, we first train the respective model using Difficulty-Net based weighting. Then, following [19], we freeze the feature extractor and re-train the classifier without using Difficulty-Net. We use $\lambda = 0.3$ for all the experiments unless otherwise stated since we find it works reasonably well as we will show in Sec. 4.5. Further details are provided in the supplementary material.

## 4.3. Compared methods

For comparison, we use multiple SOTA methods including (1) data-resampling: class-balanced sampling (CB sampling) [19], (2) cost-sensitive learning: equalization loss (EQL) [36], focal loss [23], class-balanced loss [7], label-distribution-aware-margin (LDAM) loss [2], pre-formulated class-difficulty balanced loss (CDB-CE) [31],

| Method | Imbalance | | | | |
|---|---|---|---|---|---|
| | 200 | 100 | 50 | 20 | 10 |
| *e2e training* | | | | | |
| Focal Loss [23] | 39.64 | 44.03 | 48.91 | 55.57 | 61.10 |
| MWN [29] | 40.25 | 44.81 | 49.68 | 56.53 | 61.44 |
| Class-Balanced [7] | 39.95 | 44.78 | 47.67 | 56.83 | 59.95 |
| CB-DA [18] | 40.89 | 46.24 | 49.80 | 56.67 | 62.16 |
| LDAM [2] | 41.42 | 46.14 | 49.19 | 55.90 | 62.08 |
| EQL [36] | 43.46 | 46.47 | 51.34 | 56.82 | 60.13 |
| CDB-CE [31] | 40.42 | 45.25 | 49.45 | 56.66 | 61.52 |
| PaCo [6] | 43.09 | 47.26 | 52.14 | 58.37 | 63.12 |
|   + Bal. Softmax [26] | 46.72 | 51.47 | 55.88 | 60.32 | 64.10 |
|   + Bal. Softmax [26]† | – | 52.00 | 56.00 | – | 64.20 |
| Ours | 44.80 | 47.96 | 54.27 | 58.93 | 63.52 |
|   + Bal. Softmax | 47.53 | 52.14 | 56.86 | **61.72** | **65.67** |
| *decoupled learning* | | | | | |
| cRT [19] | 44.05 | 48.04 | 53.32 | 58.72 | 63.74 |
| LWS [19] | 44.42 | 48.13 | 53.44 | 59.10 | 63.97 |
| LAS [51] | 44.87 | 48.68 | 53.85 | 59.36 | 64.18 |
| DRO-LT [28] † | – | 47.31 | **57.57** | – | 63.41 |
| BALMS [26] | 46.12 | 50.95 | 54.42 | 59.00 | 63.10 |
| MWN + cRT | 44.56 | 48.34 | 53.62 | 59.05 | 63.99 |
| MWN + LWS | 44.71 | 48.65 | 53.77 | 59.22 | 64.15 |
| MWN + LAS | 45.04 | 49.12 | 53.95 | 59.38 | 64.24 |
| Ours + cRT | 47.45 | 52.01 | 56.34 | 61.08 | 64.80 |
| Ours + LWS | <u>47.91</u> | <u>52.62</u> | 56.61 | 61.38 | 65.08 |
| Ours + LAS | **48.32** | **52.96** | <u>56.90</u> | <u>61.46</u> | <u>65.22</u> |

Table 1. Top-1 classification accuracy (%) on CIFAR-100-LT. † denotes copied results from origin paper [6, 28]. The best results are made bold while the second best results are underlined, which applies for the other tables as well.

(3) metric learning: parametric contrastive learning (PaCo) [6], (4) decoupled learning: classifier normalization ($\tau$-norm) [19], classifier re-training (cRT) [19], learnable weight scaling (LWS) [19], label-aware smoothing (LAS) [51], balanced meta-softmax (BALMS) [26], distribution robustness loss (DRO-LT) [28], (5) meta learning: Meta-Weight-Net (MWN) [29], class-balancing as domain-adaptation (CB-DA) [18]. For the sake of fairness, we do not compare our method directly with MoE methods [41, 45] as they use ensemble of multiple expert models, while we focus on improving the learning for a single expert. However, we verified that the proposed method can exhibit significant performance gains by using simple ensembling techniques and can outperform SOTA MoE methods. The results are found in the supplementary material.

## 4.4. Main results

**CIFAR100-LT.** Following [6, 36], we use AutoAugment [4] and Cutout [10] for all our implementations on CIFAR100-LT. As explained in [36], this achieves a higher baseline than other commonly followed ones. Therefore,

to ensure the fairness of comparison, we re-implemented the compared methods in our training setup using their published codes. Results without using AutoAugment and Cutout are provided in the supplementary material. We achieved better results than originally reported results for all the re-implemented methods except PaCo [6], which uses additional augmentation. Therefore, we list the original results of PaCo for reference in addition to the results in the fair setting. PaCo uses an additional center learning rebalance step, for which they employ Balanced Softmax (Bal. Softmax) [26]. We report the results of PaCo both with and without the use of Bal. Softmax. For the fair comparison with PaCo + Bal. Softmax, we tested Ours + Bal. Softmax in addition to the vanilla variant (Ours).

In e2e learning, our proposed approach without combining any other techniques (Ours) achieved better performance than all previous stand-alone methods as seen in Table 1. The margin of improvement is higher in high-imbalanced situations. End-to-end learning with Ours + Bal. Softmax turned out to be very effective and created new SOTA for low imbalanced cases (*i.e.* 10 and 20).

In decoupled-learning, we find that when we use feature extractors trained using Difficulty-Net, any popular classifier learning method (*e.g.* cRT, LWS, LAS) gives improved performance. This shows that our proposed method learns very powerful data representations. Ours + LAS achieved the best results in high imbalanced situations (*e.g.* 200 and 100), while it achieved the second best in all other cases.

**ImageNet-LT.** Table 2 shows the results on ImageNet-LT. In e2e learning alone, irrespective of the model used, we achieved better overall accuracy than other e2e methods and comparable accuracy with multiple decoupled methods.

Furthermore, Difficulty-Net based representation learning with popular classifier re-training methods achieved state-of-the-art results. Using both ResNet-10 and ResNet-50, Ours + LAS achieved the best overall accuracy, which re-confirms the effectiveness of this method. More ImageNet-LT results with many-/med-/few-shot splits are available in the supplementary material.

**Places-LT.** From Table 1 and Table 2, it is evident that our Difficulty-Net based weighting is consistently effective when used for the representation learning in decoupled training methods. Therefore, for Places-LT, we only report the results of Ours + {cRT, LWS, LAS} and compare them with previous SOTA results in Table 3. The results verify that the representation learned using our method is very powerful and helps us achieve the best overall accuracy by simple classifier re-balancing. Our improvements in overall accuracy is majorly accounted for by significant gains in medium- and few-shot accuracies. Even though our representation learning is effective with any classifier re-training

| Method | ResNet-10 | ResNet-50 |
|---|---|---|
| *e2e training* | | |
| CE | 34.8 | 41.6 |
| Focal loss [23] | 30.5 | – |
| EQL [36] | 36.4 | – |
| CB-DA[18] | 36.7 | 48.0 |
| CDB-CE [31] | 38.5 | – |
| Bal. Softmax [26] | 41.1 | – |
| PaCo [6] * | – | 49.8 |
| + Bal. Softmax [26]* | – | 53.5 |
| Ours | 41.4 | 51.2 |
| + Bal. Softmax | 44.3 | <u>53.7</u> |
| *decoupled learning* | | |
| cRT [19] | 41.8 | 47.3 |
| LWS [19] | 41.4 | 47.7 |
| MiSLAS [51] | – | 52.7 |
| BALMS [26] | 41.8 | – |
| DRO-LT [28] | – | 53.5 |
| Ours + cRT | 43.6 | 53.5 |
| Ours + LWS | <u>44.4</u> | <u>53.7</u> |
| Ours + LAS | **44.6** | **54.0** |

Table 2. Top-1 classification accuracies (%) on ImageNet-LT. * represents reproduced results using author's codes without using RandAugment [5] for fair comparison. Other baseline results are copied from original papers. Results using RandAugment are provided in the supplementary material.

| Method | Many | Med | Few | All |
|---|---|---|---|---|
| CE | <u>45.7</u> | 27.3 | 8.2 | 30.2 |
| CB sampling [19] | – | – | – | 30.3 |
| Focal Loss [23] | 41.1 | 34.8 | 22.4 | 34.6 |
| cRT [19] | 42.0 | 37.6 | 24.9 | 36.7 |
| LWS [19] | 40.6 | 39.1 | 28.6 | 37.6 |
| BALMS [26] | 41.2 | 39.8 | 31.6 | 38.7 |
| LADE [16] | 42.8 | 39.0 | 31.2 | 38.8 |
| DisAlign [48] | 40.4 | 42.4 | 30.1 | 39.3 |
| IEM [54] | **46.8** | 39.2 | 28.0 | 39.7 |
| MiSLAS [51] | 39.6 | 43.3 | 36.1 | 40.4 |
| PaCo [6] | 37.5 | **47.2** | 33.9 | 41.2 |
| Ours + cRT | 43.0 | <u>43.8</u> | 35.0 | **41.7** |
| Ours + LWS | 41.4 | 43.7 | **36.9** | <u>41.5</u> |
| Ours + LAS | 42.4 | 43.7 | <u>36.6</u> | **41.7** |

Table 3. Top-1 classification accuracies (%) for Places-LT.

method, especially Ours + LAS significantly boosts results for the few-shot classes and achieved SOTA in overall accuracy. PaCo achieved the best results for the medium-shot classes, but it sacrificed the performance on the many-shot classes significantly, resulting in lower overall accuracy.

### 4.5. Ablation study

In this sub-section, we first show the ablation study of our key components, namely relative difficulty and the
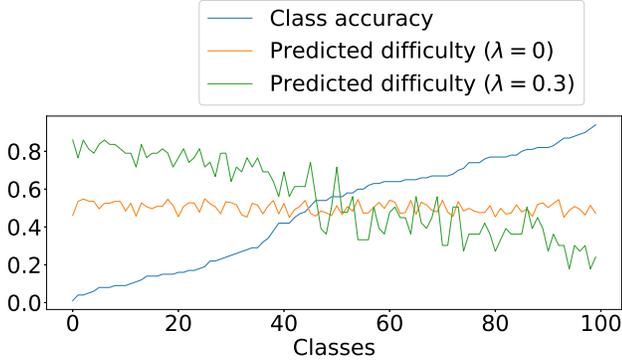
Figure 2. Difficulty scores for CIFAR100-LT (imbalance=100) classes predicted by Difficulty-Net learned with $\lambda = 0$ and $\lambda = 0.3$. The classes are sorted in increasing order of their accuracy.

driver loss. Then we re-verify the effectiveness of the class-level weighting studied in [31] in our meta-learning framework. Further, we verify the effectiveness of using the meta-learning loss in our method. We conclude this sub-section with the effectiveness of the proposed method by comparing it with the straight-forward combination of CDB-CE [31] and MWN [29].

**Absolute difficulty vs. relative difficulty.** For predicting absolute difficulty, we modified Difficulty-Net from Eq. 1 to $d_c = D^{abs}(a_c; \theta)$ and trained it in the same way as before. The comparison is provided in Table 4 (#7 vs. #9). As can be seen, relative difficulty significantly outperforms absolute difficulty for both low and high imbalance. This verifies the effectiveness of relative difficulty.

**Contribution of $L^{dr}$.** The value of $\lambda$ in Eq. 8 controls the impact of $L^{dr}$. Here we analyse the effect of $\lambda$. For that, we evaluate the performance of ResNet-32 trained end-to-end using different values of $\lambda$ and report the results in Table 5. It shows that $\lambda = 0$, which is equivalent to #6 in Table 4, works significantly poor especially in high imbalance case, which asserts the importance of using $L^{dr}$. We find that for higher imbalance, higher $\lambda$ works better. But, too high $\lambda$ leads to significant drop in performance. Irrespective of the imbalance, $\lambda = 0.3$ works consistently well.

To further analyse the usefulness of $L^{dr}$, we visualise the predicted difficulty by Difficulty-Net trained with and without $L^{dr}$. The results are shown in Figure 2. It shows that using $L^{dr}$ with $\lambda = 0.3$ provides a more meaningful learning of Difficulty-Net compared to when not using $L^{dr}$ ($\lambda = 0$). The latter predicts similar difficulty scores for all the classes inspite of the highly biased accuracy. However, using $L^{dr}$ helps Difficulty-Net to predict high difficulty for less accurate classes.

**Sample-level difficulty vs. Class-level difficulty.** We modified the Difficulty-Net to predict sample-level difficulties and compared it with the proposed method. We modified the Difficulty-Net as $d_1, d_2, ..., d_B = D^{sample}(\{l_s\}_{s=1}^{B}; \theta)$ where $B$ is the total number of samples in a single batch, $l_s$ is the model's cross-entropy loss for samples $s$, and $d_s$ is the predicted difficulty for sample $s$. Simply, we meta-learn the $D^{sample}$ to predict difficulty of each sample relative to other samples in the same training batch. Note that this variant uses relative difficulty and the driver loss, and thus is different from MWN.

In Table 4 (#8 vs. #9), it is seen that class-level difficulty significantly outperforms the sample-level difficulty in overall performance. This proves the effectiveness of class-difficulty in our proposed method. We believe that this happens because the head classes have higher *absolute* number of hard samples than the tail classes simply because the head classes have much more training samples. In such case, as pointed out in [31], sample-level weighting gives higher weights to head classes in total, and therefore cause the model to get biased to the head classes. This is verified by the fact that class-level performs much better especially for the tail classes (med and few-shot) as shown in the supplementary material. Another interesting observation is that our sample-level Difficulty-Net even significantly outperforms MWN (#3 vs. #8), which re-verifies the effectiveness of our newly proposed components, namely relative difficulty and the driver loss.

**Contribution of $L^{meta}$.** Here we verify the usefulness of $L^{meta}$ in our Difficulty-Net training. In Table 4 (#5 vs. #6), we see that using meta-learning loss gives a boost of 0.89% (for imb. 100) which confirms the benefit of using ML.

**The straight-forward combination of CDB-CE and MWN does not work.** As stated in Sec. 3.2, ours is not the straight-forward combination of the previous methods. Evidently from Table 4 (#4 vs. #9), such straight-forward combination does not work well, which verifies the contributions of our newly proposed components.

### 4.6. Further analysis

It is evident in Figure 2 that Difficulty-Net successfully learns to predict reasonable difficulty from the class-wise accuracies. Here we further analyse how the predicted difficulties change as the training progresses. For this purpose, we plot the entropy of the difficulty scores with the training steps in Figure 3. We compute the entropy as $E(\{d_c\}_{c=1}^{C}) = -\frac{1}{C}\sum_{c=1}^{C}\log(C\frac{d_c}{\sum_k d_k})$. Figure 3 shows that the entropy decreases with the training steps. This suggests that the predicted difficulty scores gradually become more and more uniform, as the model's class-wise perfor-

| # | Name | S vs. C | A vs. R | ML | $L^{dr}$ | Imb.=100 | Imb.=10 |
|---|------|---------|---------|-----|----------|----------|---------|
| 1 | Focal loss [23] | S | A | | | 44.03 | 61.10 |
| 2 | CDB-CE [31] | C | A | | | 45.25 | 61.52 |
| 3 | MWN [29] | S | A | ✓ | | 44.81 | 61.44 |
| 4 | CDB-CE + MWN | C | A | ✓ | | 45.42 | 61.87 |
| 5 | Ours w/o $L^{dr}$ and ML | C | R | | | 45.51 | 62.44 |
| 6 | Ours w/o $L^{dr}$ | C | R | ✓ | | 46.40 | <u>63.10</u> |
| 7 | Ours w/o relative difficulty | C | A | ✓ | ✓ | <u>46.81</u> | 62.32 |
| 8 | Ours w/o class-level weighting | S | R | ✓ | ✓ | 45.76 | 62.51 |
| 9 | Ours | C | R | ✓ | ✓ | **47.96** | **63.52** |

Table 4. Classification accuracy on CIFAR100-LT with imbalance ratio (Imb.) 100 and 10. "S vs. C" means Sample-level vs. Class-level. "A vs. R" means Absolute difficulty vs. Relative difficulty. ML stands for Meta Learning.

| $\lambda$ | 0 | 0.3 | 0.6 | 0.9 | 1.0 |
|-----------|-----|-----|-----|-----|-----|
| Imbalance=100 | 46.40 | <u>47.96</u> | **48.03** | 47.35 | 46.66 |
| Imbalance=10 | 63.10 | **63.52** | <u>63.44</u> | 62.62 | 62.24 |

Table 5. Accuracy (in e2e learning) for different values of $\lambda$ on CIFAR100-LT.
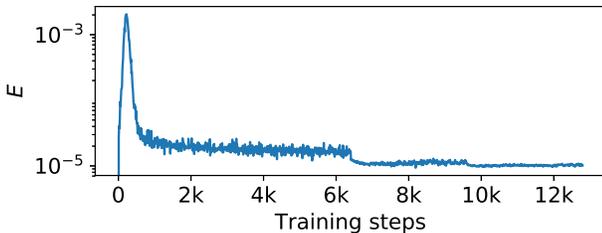


Figure 3. Plotting entropy ($E$) of difficulty scores predicted by Difficulty-Net against number of training steps. We used CIFAR100-LT (imbalance=100) for this plot.
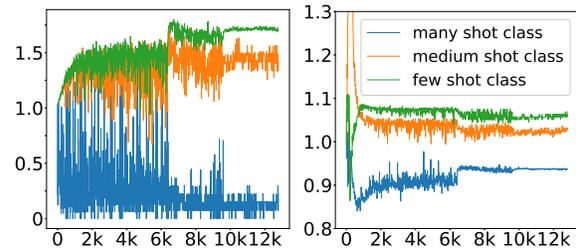


Figure 4. Assigned weights to three different classes during training with CDB-CE [31] (Left) and our Difficulty-Net (Right). The vertical axis represents assigned weights and the horizontal axis represents training steps.

mance gradually gets balanced. this result empirically supports the rationality of Difficulty-Net.

Further, we analyse the characteristics of difficulties estimated by Difficulty-Net in comparison with those by [31]. We pick three classes from the CIFAR100-LT classes (one from each of many-, medium- and few-shot classes) and show how the normalized weights of the classes change as the training progresses. As shown in Figure 4, CDB-CE weighting [31] leads to more fluctuations in the assigned weights, while Difficulty-Net based weighting is more smooth and stable. This suggests that Difficulty-Net has capability of 'remembering' which class is difficult whereas CDB-CE weighting tends to be heavily affected by quick accuracy change at each time step. We believe this characteristic of Difficulty-Net encourages consistent and stable training of the model, ending up in better performance than CDB-CE weighting. Another interesting observation in Figure 4 is that the difficulties of the three classes estimated by Difficulty-Net tend to converge as the training progresses ,which is not observed in the case of CDB-CE. This observation is consistent with Figure 3, which showed that the predicted difficulty scores gradually become more uniform as the model's performance gets balanced.

## 5. Conclusion

This paper has proposed Difficulty-Net, a novel method for long-tailed recognition that learns to predict difficulty of classes in a meta-learning framework. The proposed method has mainly three key features compared to prior works. First, it removes any dependence on heuristic formulations thanks to its ability to learn any suitable difficulty formulation for a given dataset. Second, it estimates relative difficulty of a class compared to the other classes whereas prior works use only absolute difficulty of a class in question. Third, it employs a new driver loss function that helps to drive Difficulty-Net learning in a reasonable direction. We verified the effectiveness of the proposed method by conducting extensive experiments on multiple datasets. Further analysis also demonstrated the usefulness of relative difficulty and the newly proposed driver loss function.

# References

[1] Sumia Abdulhussien Razooqi Al-Obaidi, Davood Zabi-hzadeh, Ali Salim Rasheed, and Reza Monsefi. Robust metric learning based on the rescaled hinge loss. *Int. J. Mach. Learn. Cybern.*, 11:2515–2528, 2020.

[2] Kaidi Cao, Colin Wei, Adrien Gaidon, Nikos Arechiga, and Tengyu Ma. Learning imbalanced datasets with label-distribution-aware margin loss. In *NeurIPS*, 2019.

[3] Nitesh Chawla, Kevin Bowyer, Lawrence Hall, and W. Kegelmeyer. Smote: Synthetic minority over-sampling technique. *J. Artif. Intell. Res. (JAIR)*, 16:321–357, 2002.

[4] Ekin Dogus Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V. Le. Autoaugment: Learning augmentation policies from data. In *CVPR*, 2019.

[5] Ekin D Cubuk, Barret Zoph, Jonathon Shlens, and Quoc V Le. Randaugment: Practical automated data augmentation with a reduced search space. In *CVPR Workshops*, pages 702–703, 2020.

[6] Jiequan Cui, Zhisheng Zhong, Shu Liu, Bei Yu, and Jiaya Jia. Parametric contrastive learning. In *ICCV*, pages 715–724, 2021.

[7] Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge Belongie. Class-balanced loss based on effective number of samples. In *CVPR*, 2019.

[8] Yin Cui, Yang Song, Chen Sun, Andrew Howard, and Serge Belongie. Large scale fine-grained categorization and domain-specific transfer learning. In *CVPR*, pages 4109–4118, 2018.

[9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *CVPR*, 2009.

[10] Terrance DeVries and Graham W Taylor. Improved regularization of convolutional neural networks with cutout. *arXiv preprint arXiv:1708.04552*, 2017.

[11] Qi Dong, Shaogang Gong, and Xiatian Zhu. Class rectification hard mining for imbalanced deep learning. In *ICCV*, 2017.

[12] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *ICML*, page 1126–1135, 2017.

[13] Hui Han, Wen-Yuan Wang, and Bing-Huan Mao. Borderline-smote: A new over-sampling method in imbalanced data sets learning. In *Advances in Intelligent Computing*, volume 3644, pages 878–887, 2005.

[14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016.

[15] Elad Hoffer and Nir Ailon. Deep metric learning using triplet network. In *Similarity-Based Pattern Recognition*, pages 84–92, 2015.

[16] Youngkyu Hong, Seungju Han, Kwanghee Choi, Seokjun Seo, Beomsu Kim, and Buru Chang. Disentangling label distribution for long-tailed visual recognition. In *CVPR*, pages 6626–6636, 2021.

[17] Chen Huang, Yining Li, Chen Change Loy, and Xiaoou Tang. Learning deep representation for imbalanced classification. In *CVPR*, pages 5375–5384, 2016.

[18] Muhammad Abdullah Jamal, Matthew Brown, Ming-Hsuan Yang, Liqiang Wang, and Boqing Gong. Rethinking class-balanced methods for long-tailed visual recognition from a domain adaptation perspective. In *CVPR*, 2020.

[19] Bingyi Kang, Saining Xie, Marcus Rohrbach, Zhicheng Yan, Albert Gordo, Jiashi Feng, and Yannis Kalantidis. Decoupling representation and classifier for long-tailed recognition. In *ICLR*, 2020.

[20] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. In *NeurIPS*, pages 18661–18673, 2020.

[21] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.

[22] Yu Li, Tao Wang, Bingyi Kang, Sheng Tang, Chunfeng Wang, Jintao Li, and Jiashi Feng. Overcoming classifier imbalance for long-tail object detection with balanced group softmax. In *CVPR*, pages 10991–11000, 2020.

[23] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollar. Focal loss for dense object detection. *TPAMI*, 42(02):318–327, 2020.

[24] Ziwei Liu, Zhongqi Miao, Xiaohang Zhan, Jiayun Wang, Boqing Gong, and Stella X. Yu. Large-scale long-tailed recognition in an open world. In *CVPR*, 2019.

[25] Aditya Krishna Menon, Sadeep Jayasumana, Ankit Singh Rawat, Himanshu Jain, Andreas Veit, and Sanjiv Kumar. Long-tail learning via logit adjustment. In *ICLR*, 2021.

[26] Jiawei Ren, Cunjun Yu, Shunan Sheng, Xiao Ma, Haiyu Zhao, Shuai Yi, and Hongsheng Li. Balanced meta-softmax for long-tailed visual recognition. In *NeurIPS*, 2020.

[27] Mengye Ren, Wenyuan Zeng, Bin Yang, and Raquel Urtasun. Learning to reweight examples for robust deep learning. In *ICML*, 2018.

[28] Dvir Samuel and Gal Chechik. Distributional robustness loss for long-tail learning. In *ICCV*, 2021.

[29] Jun Shu, Qi Xie, Lixuan Yi, Qian Zhao, Sanping Zhou, Zongben Xu, and Deyu Meng. Meta-weight-net: Learning an explicit mapping for sample weighting. In *NeurIPS*, 2019.

[30] Naman D. Singh and Abhinav Dhall. Clustering and learning from imbalanced data. *ArXiv*, abs/1811.00972, 2018.

[31] Saptarshi Sinha, Hiroki Ohashi, and Katsuyuki Nakamura. Class-wise difficulty-balanced loss for solving class-imbalance. In *ACCV*, 2020.

[32] Saptarshi Sinha, Hiroki Ohashi, and Katsuyuki Nakamura. Class-Difficulty Based Methods for Long-Tailed Visual Recognition. *International Journal of Computer Vision*, pages 1–15, 2022.

[33] Hyun Song, Yu Xiang, Stefanie Jegelka, and Silvio Savarese. Deep metric learning via lifted structured feature embedding. In *CVPR*, pages 4004–4012, 2016.

[34] Muhammad Atif Tahir, Josef Kittler, and Fei Yan. Inverse random under sampling for class imbalance problem and its application to multi-label classification. *Pattern Recognition*, 45(10):3738–3750, 2012.

[35] Jingru Tan, Xin Lu, Gang Zhang, Changqing Yin, and Quanquan Li. Equalization loss v2: A new gradient balance approach for long-tailed object detection. In *CVPR*, pages 1685–1694, 2021.

[36] Jingru Tan, Changbao Wang, Buyu Li, Quanquan Li, Wanli Ouyang, Changqing Yin, and Junjie Yan. Equalization loss for long-tailed object recognition. In *CVPR*, 2020.

[37] Jiaqi Wang, Wenwei Zhang, Yuhang Zang, Yuhang Cao, Jiangmiao Pang, Tao Gong, Kai Chen, Ziwei Liu, Chen Change Loy, and Dahua Lin. Seesaw loss for long-tailed instance segmentation. In *CVPR*, 2021.

[38] Peng Wang, K. Han, Xiu-Shen Wei, Lei Zhang, and Lei Wang. Contrastive learning based hybrid networks for long-tailed image classification. In *CVPR*, 2021.

[39] Tao Wang, Yu Li, Bingyi Kang, Junnan Li, Junhao Liew, Sheng Tang, Steven Hoi, and Jiashi Feng. The devil is in classification: A simple framework for long-tail instance segmentation. In *ECCV*, pages 728–744, 2020.

[40] Xun Wang, Xintong Han, Weilin Huang, Dengke Dong, and Matthew R Scott. Multi-similarity loss with general pair weighting for deep metric learning. In *CVPR*, pages 5022–5030, 2019.

[41] Xudong Wang, Long Lian, Zhongqi Miao, Ziwei Liu, and Stella Yu. Long-tailed recognition by routing diverse distribution-aware experts. In *ICLR*, 2021.

[42] Yiru Wang, Weihao Gan, Jie Yang, Wei Wu, and Junjie Yan. Dynamic curriculum learning for imbalanced data classification. In *ICCV*, pages 5016–5025, 2019.

[43] Yu-Xiong Wang and Martial Hebert. Learning to learn: Model regression networks for easy small sample learning. In *ECCV*, pages 616–634, 2016.

[44] Yu-Xiong Wang, Deva Ramanan, and Martial Hebert. Learning to model the tail. In *NeurIPS*, 2017.

[45] Liuyu Xiang, Guiguang Ding, and Jungong Han. Learning from multiple experts: Self-paced knowledge distillation for long-tailed classification. In *ECCV*, pages 247–263, 2020.

[46] Yuzhe Yang and Zhi Xu. Rethinking the value of labels for improving class-imbalanced learning. In *NeurIPS*, 2020.

[47] Xi Yin, Xiang Yu, Kihyuk Sohn, Xiaoming Liu, and Manmohan Chandraker. Feature transfer learning for face recognition with under-represented data. In *CVPR*, 2019.

[48] Songyang Zhang, Zeming Li, Shipeng Yan, Xuming He, and Jian Sun. Distribution alignment: A unified framework for long-tail visual recognition. In *CVPR*, 2021.

[49] Xiao Zhang, Zhiyuan Fang, Yandong Wen, Zhifeng Li, and Yu Qiao. Range loss for deep face recognition with long-tailed training data. In *ICCV*, pages 5419–5428, 2017.

[50] Yifan Zhang, Bryan Hooi, Lanqing Hong, and Jiashi Feng. Test-agnostic long-tailed recognition by test-time aggregating diverse experts with self-supervision. *arXiv*, 2021.

[51] Zhisheng Zhong, Jiequan Cui, Shu Liu, and Jiaya Jia. Improving calibration for long-tailed recognition. In *CVPR*, 2021.

[52] Boyan Zhou, Quan Cui, Xiu-Shen Wei, and Zhao-Min Chen. BBN: Bilateral-branch network with cumulative learning for long-tailed visual recognition. In *CVPR*, pages 1–8, 2020.

[53] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva, and Antonio Torralba. Places: A 10 million image database for scene recognition. *TPAMI*, 2017.

[54] Linchao Zhu and Yi Yang. Inflated episodic memory with region self-attention for long-tailed visual recognition. In *CVPR*, 2020.

# Difficulty-Net: Learning to Predict Difficulty for Long-Tailed Recognition
## Supplementary Material

|  | Imb. | Train | Val / Meta | Test |
|---|---|---|---|---|
| CIFAR-LT | 10–200 | 49–2 / 490 | 10 | 100 |
| ImageNet-LT | 256 | 5 / 1280 | 10 | 50 |
| Places-LT | 996 | 5 / 4980 | 10 | 100 |

Table 5: The imbalance ratio (Imb.) and the number of samples per class in each datasets. Since the training splits are imbalanced, we show the number of samples in the least and most frequent classes. Note that exactly the same set of images are used for both the validation sets and meta sets. This indicates that meta-learning based methods including ours do not exploit any extra data.

## 1. More implementation details

### 1.1. More details on datasets

Table 5 shows the number of samples in training, validation, and test splits in each dataset. As shown in the table, meta-learning (ML) based methods including ours and [18, 29] reuse validation images for constructing $S^{meta}$, a dataset to be used for meta learning. Note that (1) our proposed method was compared with other ML based methods in exactly the same conditions, (2) we re-ran the experiments using public codes of previous methods and therefore all the methods compared in this paper are evaluated using exactly the same split as shown above, and (3) all the ML based methods including ours do not use any extra data, and therefore they do not receive any unfair benefit compared to other methods by having $S^{meta}$. All the hyper-parameters were tuned using the validation sets.

### 1.2. Hyperparameter settings

For CIFAR100-LT, following [36, 6], we use AutoAugment [4] and Cutout [10]. Following [36], we train ResNet-32 [14] for 12.8K steps with a batch size of 128 and an initial learning rate of 0.1. The learning rate is linearly warmed up to 0.2 over the first 400 steps. It is also decayed by 0.1 after 6.4K and 9.6K steps. For the classifier learning stage in decoupled learning methods, we fix the feature extractor and re-train the classifier for 50 steps using class-balanced

sampling following [19]. The learning rate used is 0.1 and is decayed by 0.1 after 30 and 40 epochs. For both the stages, we use a weight decay of $1e-4$. CIFAR100-LT experiments are done on a single NVIDIA Tesla V100 GPU.

For ImageNet-LT, we follow [51] and train the models for 180 steps with an initial learning rate of 0.05. The batch size used is 128. We use cosine learning rate decay and weight decay of $5e-4$. In decoupled training, for the second stage we only re-train the classifier for 10 steps using batch size 128 and cosine decayed learning rate with an initial value of 0.05. The models are trained on four NVIDIA Tesla V100 GPUs.

For Places-LT, following [51, 24] we load a ResNet-152 pretrained on ImageNet and then finetune it for 30 steps using an initial learning rate of 0.01 and weight decay of $5e-4$. The learning rate is decayed by 0.1 after 10 and 20 steps. The batch size used is 128. For the classifier learning stage, we retrain the classifier for 20 steps with a batch size of 256 and initial learning rate of 0.1, which is cosine decayed. The training is done on four NVIDIA Tesla V100 GPUs.

For all the datasets, we use SGD optimizer with momentum 0.9. For Difficulty-Net learning, we use ADAM optimizer with a learning rate of 0.001 and a weight decay $1e-4$. All implementations are done on PyTorch.

### 1.3. Designing the Difficulty-Net

As stated in Sec. 3.2, our Difficulty-Net is a MLP with 2 hidden layers. The illustration of our Difficulty-Net is given in Fig. 5. The output layer dimension changes with the number of classes in the dataset. Here we provide a simple way to select the hidden layer dimensions $H$. To come up with the method, we compare the end-to-end training performance using different values for $H$ on 2 different datasets. The results are given in Table 6.

We find that the best working $H$ is different for different datasets. Therefore, based on the results, we decide to select $H = 2^n$ such that $2^{n-1} \leq C < 2^n$, where $C$ is the number of classes and $n$ is a positive integer.

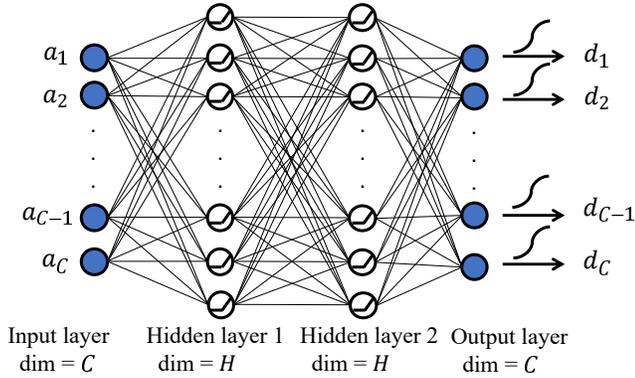The value of $C$ and $H$ for the three different datasets that we used are given in Table 7.

Figure 5: Illustration of our Difficulty-Net

| $H$ | CIFAR100-LT ($C = 100$) | ImageNet-LT ($C = 1000$) |
|---|---|---|
| 128 | <u>47.96</u> | 39.4 |
| 256 | **48.06** | 40.4 |
| 512 | 47.81 | <u>41.2</u> |
| 1024 | 47.34 | **41.4** |
| 2048 | 46.92 | 40.8 |

Table 6: Effect of $H$ on e2e training of ResNet-32 and ResNet-10 on CIFAR100-LT (imbalance=100) and ImageNet-LT respectively.

| Dataset | $C$ | H |
|---|---|---|
| CIFAR100-LT | 100 | 128 |
| ImageNet-LT | 1000 | 1024 |
| Places-LT | 365 | 512 |

Table 7: $C$ and $H$ for datasets used in our experiments.

### 1.4. Algorithm for meta-learning via Difficulty-Net

The algorithm for our Difficulty-Net based learning is provided in Algorithm 1. As stated in Sec. 3.2, our learning method comprises of three main steps (Eq. 9,10 and 11) that are represented by steps 7, 8 and 10 in the algorithm. Note that in our algorithm, $S^{meta}$ is reused as validation set $S^{val}$ for calculating accuracies.

## 2. More Results

### 2.1. CIFAR100-LT results without using extra augmentations

For the CIFAR100-LT results reported in Table 1, we used extra augmentations (AutoAugment [4] and Cutout [10]) to ensure same training setups as recent SOTA meth-

---

**Algorithm 1** Meta-learning using Difficulty-Net

**Require:** Training set $S^{train}$, Meta dataset $S^{meta}$
**Require:** Initial learnable parameters $\theta_1$ and $\phi_1$
**Require:** Max iterations $T$, Value of $\lambda$
**Require:** Learning rates $\alpha$, $\beta$ and batch sizes $b$, $m$
1: **for** $t = 1 \ldots T$ **do**
2:     Compute $A_{C,t}$ using $f(x; \phi_t)$ on $S^{meta}$
3:     Sample mini-batch of size $b$ from $S^{train}$
4:     Sample mini-batch of size $m$ from $S^{meta}$
5:     Compute weights with $A_{C,t}$ and $\theta_t$ using Eq. 4
6:     Compute intermediate $\hat{\phi}_t(\theta_t)$ using Eq. 9
7:     Update $\theta_t$ to $\theta_{t+1}$ using Eq. 10
8:     Re-compute Eq. 4 with $A_{C,t}$ and $\theta_{t+1}$
9:     Update $\phi_t$ to $\phi_{t+1}$ using Eq. 11
10: **end for**
**Output:** $\phi_{T+1}, \theta_{T+1}$

---

ods such as PaCo[6], BALMS [26] and DRO-LT[28] for fair comparison. As expected, these additional augmentation techniques provide a significant boost in the results. To verify that our proposed method is effective independent of these extra augmentations, we compare the results of our method with other SOTA methods without using the augmentation techniques. The results are reported in Table 8. With or without extra augmentations, Ours + LAS proves to be very effective.

### 2.2. ImageNet-LT results on many-, med- and few-shot classes

In Table 2, we saw that our proposed method helps to achieve the best overall accuracy. Here we study the effectiveness of our method for each of many-, medium- and few-shot classes. The comparison results are given in Table 9. We find that in both e2e learning and decoupled learning, Difficulty-Net based weight assignment helps to significantly boost the performance of the few-shot and medium-shot classes. We believe this result indicates the strong capability of Difficulty-Net based weighting in mitigating biased performance caused by the class imbalance. Especially, Ours + LAS is the most effective for the few-shot classes, irrespective of the model used.

### 2.3. ImageNet-LT Results Using RandAugment

In Table 2, we reproduced the results of PaCo [6] without using RandAugment [5] for the sake of fair comparison with all the other methods that do not use RandAugment. However, the originally reported results in [6] use RandAugment as additional augmentation, which are significantly higher than the reproduced results. This suggests that PaCo is greatly benefited by the use of RandAugment. Therefore, we used RandAugment with our method and compared the results with PaCo in Table 10. We only used Ours + LAS for

|  | Imbalance | | | | |
|---|---|---|---|---|---|
| Method | 200 | 100 | 50 | 20 | 10 |
| *e2e training* | | | | | |
| Focal Loss [23] † | 35.62 | 38.41 | 44.32 | 51.95 | 55.78 |
| MWN [29] † | 37.91 | 42.09 | 46.74 | 54.37 | 58.46 |
| Class-Balanced [7] † | 36.23 | 39.60 | 45.32 | 52.99 | 57.99 |
| CB-DA [18] † | 39.31 | 43.35 | 48.53 | 55.62 | 59.58 |
| LDAM [2] † | – | 39.60 | – | – | 56.91 |
| EQL [36] † | 37.34 | 40.54 | 44.70 | 54.12 | 58.32 |
| CDB-CE [31] † | 37.40 | 42.57 | 46.78 | 54.22 | 58.74 |
| PaCo [6] | 36.96 | 40.92 | 46.97 | 53.66 | 59.59 |
| + Bal. Softmax [26] | 39.55 | 44.13 | 48.60 | 55.89 | 60.24 |
| Ours | 39.94 | 43.82 | 49.00 | 55.70 | 60.25 |
| + Bal. Softmax | 41.43 | 45.81 | 51.14 | 56.58 | <u>61.33</u> |
| *decoupled learning* | | | | | |
| cRT [19] | 40.13 | 44.04 | 48.97 | 55.67 | 59.54 |
| LWS [19] | 40.70 | 45.05 | 49.70 | 56.22 | 60.00 |
| LAS [51] | 40.76 | 45.32 | 49.96 | 56.66 | 59.96 |
| BALMS [26] | 39.58 | 44.64 | 48.52 | 54.28 | 58.34 |
| MWN + cRT | 40.57 | 44.00 | 49.47 | 56.05 | 59.64 |
| MWN + LWS | 40.48 | 44.52 | 49.10 | 55.89 | 59.48 |
| MWN + LAS | 40.94 | 44.64 | 49.15 | 55.91 | 59.24 |
| Ours + cRT | 41.12 | 45.41 | 50.50 | 56.30 | 60.86 |
| Ours + LWS | <u>41.67</u> | <u>46.04</u> | <u>51.27</u> | <u>56.66</u> | 61.30 |
| Ours + LAS | **42.19** | **46.42** | **51.60** | **56.82** | **61.47** |

Table 8: Top-1 classification accuracy (%) on CIFAR100-LT without using extra augmentation *i.e.* AutoAugment and Cutout. † denotes copied results from origin paper [31, 18]. The best results are made bold while the second best results are underlined, which applies for the other tables as well.

the comparison because Ours + LAS is the best performing decoupled learning method as seen in Table 1,2 and 3.

From Table 10, we find that using RandAugment benefits our method as well. With or without RandAugment [5], Ours+LAS outperformed PaCo.

## 2.4. Comparison on ImageNet-LT with MoE methods

In Sec. 4.4, we did not compare our proposed method directly to mixture of experts (MoE) methods as the latter uses multiple experts while we focus on improving the learning of a single expert. For the fair comparison with MoE methods, we created an ensemble of Difficulty-Net based trained models. For the ensemble creation, we trained two expert models using Ours + LAS decoupled learning. The backbone architectures of these two models were kept the same. The only difference between these models was that one used a linear classifier and the other used a cosine classifier. During inference, we simply took the mean outputs of these two models. The results of this simple ensemble is provided in

Table 11.

As can be seen, although our ensemble comprises of only two expert models, it performs significantly better than 3-experts and 4-experts RIDE [41]. This shows that our proposed Difficulty-Net is effective in learning expert models for MoE methods. However, the current ensemble is heuristic and a detailed research on contribution of Difficulty-Net in MoE is left for the future.

## 2.5. More results on sample-level v/s class-level difficulty

As empirically verified in Table 4, class-level difficulty is more effective than sample-level difficulty in our Difficulty-Net. We believe that this happens because as stated in [31] and Sec. 4.5, the absolute number of hard samples in head classes is significantly higher than that in tail classes due to the inherent long-tail characteristic of the dataset. Using sample-level difficulty gives high weights to all the hard samples irrespective of their classes, resulting in more weights for the head classes and therefore getting the model biased to the head classes.

We verified this by conducting a simple experiment on CIFAR100-LT. For 2 classes A and B with 376 and 46 training samples respectively, the absolute number of hard samples given high weights by sample-level method was higher for A(50) than B(13). Although higher proportion of samples in B($\approx 28\%$) received high weights compared to A ($\approx 13\%$), A got more weights compared to B due to its higher absolute number of hard samples. As a result, the accuracy for A is improved from 46% to 62% and that for B is decreased from 31% to 21%, hence boosting the bias. In such case, using class-level difficulty gives high weights to all samples of B, resulting in more weights for B. As a result, the accuracy on B was improved from 31% to 40%, while that on A was almost maintained (46% to 44%).

The effectiveness of class-level difficulty in Difficulty-Net for overcoming model bias is further verified in Table 12. Using sample-level difficulty causes the model to get biased towards the many-shot classes while class-level difficulty is particularly useful for improving performance on the med-shot and few-shot classes.

| Backbone Network | ResNet-10 | | | | ResNet-50 | | | |
|---|---|---|---|---|---|---|---|---|
| Method | Many | Medium | Few | Overall | Many | Medium | Few | Overall |
| *e2e training* | | | | | | | | |
| CE | 57.6 | 25.7 | 3.2 | 34.8 | 64.0 | 38.8 | 5.8 | 41.6 |
| Focal Loss [23] | 36.4 | 29.9 | 16.0 | 30.5 | – | – | – | – |
| OLTR [24] | 43.2 | 35.1 | 18.5 | 35.6 | – | – | – | – |
| EQL [36] | – | – | – | 36.4 | – | – | – | – |
| CDB-CE [31] | – | – | – | 38.5 | – | – | – | – |
| Bal. Softmax [26] | 55.8 | 35.7 | 20.9 | 41.1 | – | – | – | – |
| PaCo[6]* | – | – | – | – | **68.4** | 44.8 | 14.7 | 49.8 |
| + Bal. Softmax [26]* | – | – | – | – | 59.9 | **52.6** | 36.1 | 53.5 |
| Ours | **58.8** | 36.4 | 13.9 | 41.4 | 68.1 | 47.2 | 21.5 | 51.2 |
| + Bal. Softmax | 54.6 | 41.6 | 27.8 | 44.3 | 63.6 | 51.4 | 35.8 | 53.7 |
| *decoupled learning* | | | | | | | | |
| cRT [19] | – | – | – | 41.8 | 58.8 | 44.0 | 26.1 | 47.3 |
| LWS [19] | – | – | – | 41.4 | 57.1 | 45.2 | 29.3 | 47.7 |
| MiSLAS [51] | – | – | – | – | 61.7 | 51.3 | 35.8 | 52.7 |
| BALMS [26] | 50.3 | 39.5 | 25.3 | 41.8 | – | – | – | – |
| Ours + cRT | 53.3 | 41.1 | 27.4 | 43.6 | 63.2 | 51.8 | 35.2 | 53.5 |
| Ours + LWS | 51.6 | **43.7** | 29.3 | 44.4 | 62.5 | 52.3 | 36.6 | 53.7 |
| Ours + LAS | 51.8 | 43.6 | **30.2** | **44.6** | 62.9 | 52.6 | **36.8** | **54.0** |

Table 9: Top-1 accuracy (%) on many-, medium- and few-shot classes of ImageNet-LT. * represents results reproduced using author's codes without using RandAugment [5] for fair comparison. Other results are copied from original papers.

| Method | ResNet-10 | ResNet-50 |
|---|---|---|
| PaCo + Bal. Softmax [6] | – | 57.0 |
| Ours + LAS | **46.9** | **57.4** |

Table 10: Top-1 accuracy (%) using RandAugment[5]. Baseline results are copied from the original paper [6].

| Imbalance | 100 | | | | 10 | | | |
|---|---|---|---|---|---|---|---|---|
| Difficulty | Many | Med | Few | All | Many | Med | Few | All |
| Sample-level | **67.00** | 47.46 | 18.99 | 45.76 | **74.50** | 61.01 | 48.02 | 62.51 |
| Class-level (ours) | 64.47 | **51.21** | **24.92** | **47.96** | 70.48 | **64.40** | **53.36** | **63.52** |

Table 12: Comparison of sample-level difficulty and class-level difficulty on CIFAR100-LT.

| Method | ResNet-10 | ResNet-50 |
|---|---|---|
| LFME [45] | 38.8 | – |
| RIDE (2 experts) [41] | – | 54.4 |
| RIDE (3 experts) [41] | – | 54.9 |
| RIDE (4 experts) [41] | – | 55.4 |
| Ours (2 experts) | **47.5** | **56.2** |

Table 11: Comparison with mixture of expert methods. Baseline results are copied from the original papers [45, 41].