

Relieving Long-tailed Instance Segmentation via Pairwise Class Balance

Yin-Yin He^{1*}, Peizhen Zhang^{2*†}, Xiu-Shen Wei^{3,1}, Xiangyu Zhang², Jian Sun²

¹State Key Laboratory for Novel Software Technology, Nanjing University

²MEGVII Technology

³School of Computer Science and Engineering, Nanjing University of Science and Technology

hey@lamda.nju.edu.cn, weixs.gm@gmail.com

{zhangpeizhen, zhangxiangyu, sunjian}@megvii.com

Abstract

Long-tailed instance segmentation is a challenging task due to the extreme imbalance of training samples among classes. It causes severe biases of the head classes (with majority samples) against the tailed ones. This renders “how to appropriately define and alleviate the bias” one of the most important issues. Prior works mainly use label distribution or mean score information to indicate a coarse-grained bias. In this paper, we explore to excavate the confusion matrix, which carries the fine-grained misclassification details, to relieve the pairwise biases, generalizing the coarse one. To this end, we propose a novel Pairwise Class Balance (PCB) method, built upon a confusion matrix which is updated during training to accumulate the ongoing prediction preferences. PCB generates fightback soft labels for regularization during training. Besides, an iterative learning paradigm is developed to support a progressive and smooth regularization in such debiasing. PCB can be plugged and played to any existing method as a complement. Experimental results on LVIS demonstrate that our method achieves state-of-the-art performance without bells and whistles. Superior results across various architectures show the generalization ability. The code and trained models are available at <https://github.com/megvii-research/PCB>.

1. Introduction

The success of modern object detectors and instance segmentors has been verified on rich and balanced datasets. However, their performance drops dramatically when ap-

*Equal contribution. This paper is supported by the National Key R&D Plan of the Ministry of Science and Technology (Project No. 2020AAA0104400), CAAI-Huawei MindSpore Open Fund, and Beijing Academy of Artificial Intelligence (BAAI). This work is done during Yin-Yin He’s internship at MEGVII Technology.

†Corresponding author.

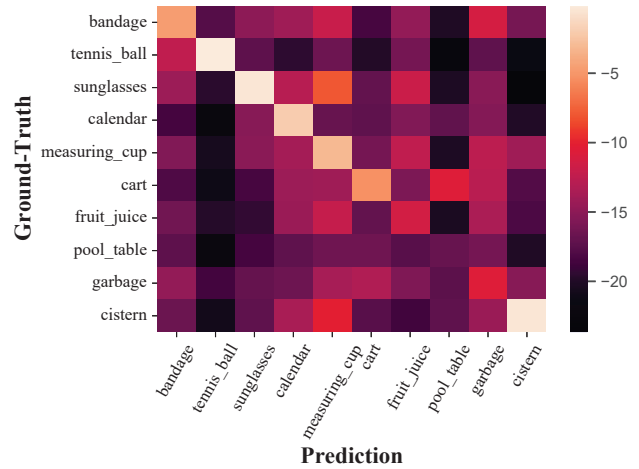


Figure 1. Visualization of a randomly sampled 10-class categorical confusion matrix of Mask R-CNN with ResNet-50-FPN on LVIS v0.5 after taking logarithm with base 2. In current palette, higher probabilities are shown with lighter colors. It reflects the model biases clearly (e.g., the probability of misclassifying *sunglasses* to *measuring cup* is larger than the reverse).

plied to datasets like LVIS [13] which is closer to the long-tailed and large-vocabulary category distribution in the real-world scenario. The devil lies in the classification prediction bias caused by extreme imbalanced training sample volumes among foreground classes [39]. Typically, it is common to leverage instructive indications with positive correlation with the prediction bias into modeling for de-biasing and thus achieve better results. Prior works exploit the class-wise sample frequencies in the training set as an intuitive indication [4, 13, 30, 38]. However, the learning quality of a class is not only about the distribution prior but also factors like optimization hardness [9], relevant to model learning process. Feng *et al.* proposed the *mean classification score* [10] metric. Such train-time model statistics can reflect the learning

Table 1. The performance of models before and after post-hoc calibration using different model statistics on the validation set. Experiments were conducted on LVIS v0.5 using Mask R-CNN with ResNet-50-FPN. MS stands for mean classification score [10] and CM stands for confusion matrix. Calibration using CM could achieve almost fully unbiased performance.

Sampler	Info.	AP	AP _r	AP _c	AP _f	AP ^b
Random	/	21.9	4.6	21.7	28.9	21.9
	MS	25.7	14.3	27.3	28.2	25.4
	CM	26.3	23.7	26.0	27.9	25.9
RFS [13]	/	25.6	16.0	26.4	28.5	25.6
	MS	27.0	20.6	28.2	28.0	27.0
	CM	28.4	28.6	28.8	27.9	28.2

quality of classes beyond mere label distribution. Yet, we notice that it considers only the sample classification statistics within each class, ignoring the inter-class similarities.

To involve the inner-and-inter class relationship, the categorical confusion matrix is already a weapon in hand. It carries the dynamic misclassification conditional probability distribution between pairs of classes (Fig. 1). A proof study is conducted on LVIS v0.5 [13] to verify our conjecture. The upper bound performance of post-hoc calibration using confusion matrix and mean score (both of the confusion matrix and mean score are collected on the validation set) is testified. Detailed results are summarized in Tab. 1. As shown, calibration using mean score and confusion matrix can both lift the performance, especially for the rare classes. More importantly, the upper bound of confusion matrix calibration is much higher than that of mean score calibration (+ 9.4 AP_r on the random sampler and + 8.0 AP_r on RFS [13] sampler with a similar performance on frequent classes). Confusion matrix calibration achieves almost fully unbiased performance. Please refer to Sec. 3.2 for details.

Accordingly, the fine-grained misclassification probabilities between pairwise classes in the confusion matrix, which we conclude as *pairwise bias*, is powerful as an indicator. One vital issue which has to be resolved is that the validation set is unavailable in practical training. One straightforward intuition is to utilize instead the confusion matrix over the train set for calibration. Unfortunately, this malfunctions as shown in Tab. 4. This is probably due to the mismatch between the confusion matrix calculated at train and test time, which is originated from diversified patterns of samples.

Inspired by [17] who conducted train-time disentangling to replace test-time post compensation for long-tailed classification, we develop an online pairwise bias-driven calibration method named PCB (**P**airwise **C**lass **B**alance). It maintains a confusion matrix during training with fightback targets generated in a matrix-transposed posterior manner to balance the pairwise bias for each ongoing proposal learning. However, naive exploitation might limit the efficacy, as

stronger regularization could be detrimental to the discrimination ability while weaker regularization can not relieve the pairwise bias well. To fully absorb the merits of the above PCB regularization and meanwhile facilitate the debiasing, a light prediction-dependent iterative paradigm is equipped. It is accomplished by using the discriminative predictions trained by the original one-hot labels to be embedded back to enhance the features recurrently. PCB regularization is gradually applied to the predictions from each step of the enhanced feature. In this way, a more friendly, progressive pairwise class balancing is achieved which is also proven to be effective in later experimental sections.

To summarize, our contributions are as follows:

- We explore the use of confusion matrix to indicate pairwise model bias in the field of long-tailed instance segmentation, which shows a promising upper bound.
- An Pairwise Class Balance method is proposed to tackle the long-tailed instance segmentation.
- Extensive experiments on LVIS v0.5 and LVIS v1.0 show the effectiveness of our method.

2. Related Work

Object Detection and Instance Segmentation. Object detection has attracted lots of attention in recent years, with remarkable improvements being made. Modern object detection frameworks [2, 11, 24, 31, 36, 50] can be divided into two-stage and one-stage ones. The two-stage detectors [2, 11, 31] first generate a set of proposals in the first stage, then refine the proposals and perform classification in the second stage. While one-stage detectors [24, 36, 50] directly predict bounding boxes. Compared to two-stage detectors, one-stage detectors are faster, yet two-stage detectors can provide better localization.

Mask R-CNN [14] adapts Faster R-CNN [31] to instance segmentation task by adding a mask prediction branch in the second stage. SOLO [41, 43] is another line of instance segmentation that is box-free. Our work is based on Mask R-CNN to stay the same as other long-tailed instance segmentation works.

Long-tail Learning. Long-tailed learning problem is across various domains (e.g., fine-grained recognition [37], multi-label learning [45] and instance segmentation [13]). Two classic solutions are data re-sampling [5, 13] which aims to flatten the data distribution and loss re-weighting [3, 8] which emphasizes more on tail data. Recent works proposed decoupling training [19, 53] which first obtains a good representation through conventional training and then calibrates the classifier. Other techniques like ensemble [42], self-supervised learning [21, 47] and knowledge distillation [16, 21] are verified to be useful in long-tailed learning.

[13] first introduced the long-tailed learning problem to instance segmentation and thus to object detection. They

built a large vocabulary dataset named LVIS and proposed a simple baseline RFS. [39] pointed out that the long-tail property affects classification most. Later on, a series of works tried to alleviate classification bias. One line of works tried to improve the sample strategy [4, 10, 44, 48, 52], while another major line of works focus on loss engineering. Equalization Loss [34] and its improvements [33] down-weight the negative gradients for tail classes from head classes, while droploss [18] further takes the gradients from background into consideration. Similarly, ACSL [40] only penalize negative classes over threshold. Separating the categories into some small groups [22, 44] and simple calibration [28, 51] helps, too. [30, 38] modified the original soft-max function by embedding the distribution prior, achieving success. [10] first introduce model statistics., it utilized mean classification scores in place of the model-agnostic prior. It fails to point out the direction the prediction of one class biased to. So we take one step further, a two-dimensional statistics (i.e., confusion matrix) is utilized to indicate fine-grained pairwise bias.

Confusion Matrix. Confusion matrix is a classical tool for error analysis. In many fields, it has shown powerful abilities. It’s used to estimate the target distribution under label shift [25]. In the field of label noise, rather than hard prediction, [29] uses soft prediction on the cleanest sample of each class to generate a confusion matrix and to assume the noisy ratio. Similarly, [49] actually keeps a confusion matrix using soft prediction for label smoothing. To the best of our knowledge, we are the first to adopt confusion matrix in the field of long-tailed learning.

3. Methodology

In this section, we first discuss the long-tailed phenomenon in a pairwise bias perspective revealed by the confusion matrix (cf. Sec. 3.1). Next, a post-hoc calibration verification shows a promising upper bound of balancing the bias (cf. Sec. 3.2). In realistic training, we propose an on-line iterative regularization paradigm to relieve the pairwise class bias that facilitates long-tailed instance segmentation (cf. Sec. 3.3).

3.1. Confusion matrix on indicating pairwise bias

Most prior works dealing with long-tailed problems aim at relieving the prediction bias between data rich classes (i.e., head classes) and data scarce classes (i.e., tail classes). They mainly convey the spirit of data re-sampling or loss re-weighting. However, these are confined to sample-level, without considering the model learning dynamics, and could be sub-optimal. LOCE [10] proposed to use *mean classification score* for each category across training to reflect the run-time predictive preference. Whereas, it failed to utilize the inter-class relationship that is critical in long-tailed literature. Instead, in this paper, we propose to leverage

the classification confusion matrix to indicate the learning preference which we elaborate on in the following.

For ease of illustration and taking the classical two-stage instance segmentation model Mask R-CNN [14] for instance, we denote by $F(\cdot)$ the R-CNN classification head. Without loss of generality, we mainly discuss the most-adopted cross-entropy loss (CE) (see also the applicability to binary cross-entropy in the last paragraph of Sec. 3.3). Normally, by taking as input a proposal feature map X , F predicts a categorical distribution $\mathbf{z} = F(X) \in \mathbb{R}^{C+1}$ (C foreground classes plus 1 background class). Since misclassification is accustomed to being among foreground classes under long-tailed setting [35], we investigate the foreground classes only. This is achieved by excluding the background logit ($\mathbf{z}^{fg} \in \mathbb{R}^C$) and re-normalizing the multinomial class probability below.

$$\hat{p}_i = \frac{\exp(z_i^{fg})}{\sum_{k=1}^C \exp(z_k^{fg})}. \quad (1)$$

We denote $M \in \mathbb{R}^{C \times C}$ as the confusion matrix which could be calculated given histogram votes in two-dimensional bins of label-to-prediction statistics:

$$M_{i,j} = \frac{\sum_{(\mathbf{x},y)} \mathbb{I}[\arg \max(\mathbf{z}^{fg}) = j, y = i]}{\sum_{(\mathbf{x},y)} \mathbb{I}[y = i]}, \quad (2)$$

where $1 \leq i, j \leq C$ and $\mathbb{I}[\cdot]$ serves as an indicator which evaluates 1 when the inner condition satisfies and 0 otherwise. To keep the finer misclassification distribution details (that are truncated by *argmax* in Eq. 2) and also for more stable training, we opt for a softened version by aggregating the predicted probabilities at the ground-truth indice.

$$M_{i,j} = \frac{\sum_{(\mathbf{x},y)} \hat{p}_j \cdot \mathbb{I}[y = i]}{\sum_{(\mathbf{x},y)} \mathbb{I}[y = i]}. \quad (3)$$

$M_{i,j}$ is the statistical probability, specifying to what extension a sample of class i is classified as j by the model. Unequal values between $M_{i,j}$ and $M_{j,i}$ reflect the asymmetric model preference between the two classes. We term the phenomenon as *pairwise bias*. The pairwise bias between class i and j becomes balanced when $M_{i,j} = M_{j,i}$. Usually, in the long-tailed scenario, pairwise biases imbalance are more likely to happen between head and tail classes, with $M_{i,j} > M_{j,i}, \forall i \in \mathcal{T}, j \in \mathcal{H}$ where sets of tail classes and head classes are denoted as \mathcal{T} and \mathcal{H} respectively. Under extreme cases, there could be $M_{i,j} \gg M_{j,i}, \exists i \in \mathcal{T}, j \in \mathcal{H}$. Officially, the head and tail classes are represented by three class splits, i.e., *frequent* (f), *common* (c) and *rare* (r). Practically, for Mask R-CNN ResNet-50-FPN trained on LVIS v0.5, the probability of a frequent class instance being misclassified as a rare one, i.e., $M_{f,r}$ is 0.01 while the versus $M_{r,f}$ is 0.19.

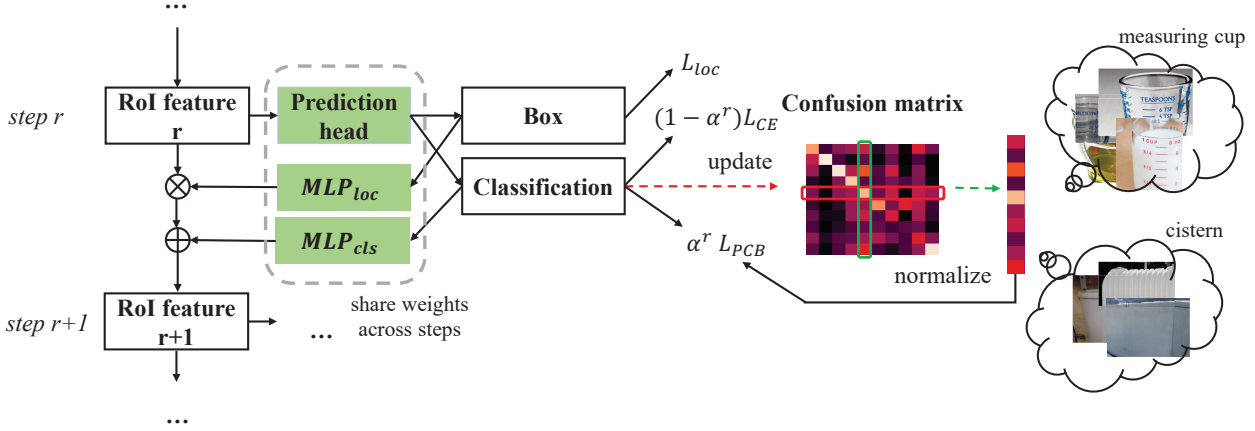


Figure 2. Our proposed PCB framework: At each recurrent step r , the RoI features generated by the previous step are fed to the shared prediction head to obtain predicted boxes and scores. Soft targets are generated according to the confusion matrix for classification regularization L_{PCB} , traded-off by α_r over L_{CE} . Subsequently, the RoI features are updated by the predictions for next step usage, and the confusion matrix is updated by the current iteration of score statistics.

3.2. A post-hoc calibration trial

Intuitive exploitation of the confusion matrix is for post-hoc calibration as mentioned in Sec. 1. Typically, a post-hoc calibration could be conducted by following the spirits of Bayes’ total probability theorem.

$$P(y = i|x) = \sum_{j=1}^C P(y = i|\hat{y} = j, x)P(\hat{y} = j|x), \quad (4)$$

where condition x omit $X = x$, representing a proposal feature. $y = i$ and $\hat{y} = j$ represent the events of “ x belongs to class i indeed” and “ x is predicted as class j ” respectively. Naturally, the $P(\hat{y} = j|x)$ term is instantiated by the predicted probability on j -th position (\hat{p}_j) of the classification output. We instantiate $P(y = i|\hat{y} = j, x)$ by $\hat{M}_{i,j}$ via performing normalization on the confusion matrix such that elements of each *column* (not row) add up to 1.

$$\hat{M}_{i,j} = \frac{M_{i,j}}{\sum_{k=1}^C M_{k,j}} \quad (5)$$

$\hat{M}_{i,j}$ indicates how likely a sample is attributed to class i overall, conditioned on being predicted as class j . It could be viewed as an approximate expectation of $P(y = i|\hat{y} = j, x)$. Hence, a post-hoc calibration of the predicted probability of sample x upon class i , dubbed \tilde{p}_i , is corrected as:

$$\tilde{p}_i = \sum_{j=1}^C \hat{M}_{i,j} \hat{p}_j \quad (6)$$

In another view, the coefficient $\hat{M}_{i,j}$ together with $\hat{M}_{j,i}$ comprise a *pairwise bias* defined in Sec. 3.1.

We also compare the above post-hoc calibration with the one aided by unary *mean classification scores* [10] below:

$$s_i = \frac{\sum_{(x,y)} \hat{p}_i \cdot \mathbb{I}[y = i]}{\sum_{(x,y)} \mathbb{I}[y = i]} \quad (7)$$

Apparently, they are exactly the diagonal elements in the confusion matrix (Eq. 3), *i.e.*, $s_i = M_{i,i}$ (As a complementary, the extra pairwise bias items retain fine-grained inter-class misclassification hints, reflecting model preference in pairs of classes). Following the empirical findings [10] that positive correlation exists between the mean classification score and the number of instances, we conduct below adjustment for calibration.

$$\tilde{p}_i = \frac{\hat{p}_i / s_i}{\sum_{k=1}^C \hat{p}_k / s_k} \quad (8)$$

Beyond above calibration via Eq. 6 or 8 over foreground classes, we keep using the same background probability ($\tilde{p}_{C+1} = \hat{p}_{C+1}$) as suggested in [35, 38]. To ensure the summation-1 property of the predictions, we first divide each foreground probability by their sum, rendering $\sum_{i=1}^C \tilde{p}_i = 1$ (still using the same symbol with slight abuse of the notation). We then re-scale each foreground probability by multiplying coefficient $\delta = 1 - \tilde{p}_{C+1}$.

We testify the post-hoc calibration efficacy *w.r.t.* the two metrics (*confusion matrix* or *mean classification score*). Experiments are conducted upon LVIS v0.5 with confusion matrix computed on the validation set. See Tab. 1 for the details. For post-hoc evaluation which enables partial labels for upper bound examination, actual categories of the proposals are utilized to compute the confusion matrix (through matching using *Intersection-over-Union* (IoU) between proposals and ground-truth boxes). Notably, using confusion

matrix for calibration is much more promising in respect of AP on rare classes ($\approx 10\%$ higher in random mode and similar to AP_c and AP_f). Such advantages show a higher upper bound of confusion matrix against mean classification score by considering the non-diagonal *pairwise biases*. This motivates our method in Sec. 3.3 below.

3.3. Online iterative confusion matrix learning

The above post-hoc calibration relies on validation labels which are unavailable in usual training. Alternatively, one might use the confusion matrix statistics over the training set instead. However, the confusion matrix patterns on the training and validation set can not match perfectly, *e.g.*, due to diversified object appearances even in the same category. Results in Tab. 4 verify this experimentally. Inspired by [17] who conducted train-time disentangling to replace the test-time post compensation for classification, we propose to conduct online pairwise class balancing during training using the train-time confusion matrix information for long-tailed instance segmentation. See Fig. 2 for our framework.

We update the confusion matrix during training through exponential moving average (EMA). Specifically for the t -th minibatch, we update the matrix rows, each corresponding to a foreground class y by the averaged prediction outputs of proposals assigned to it.

$$\mathbf{m}_y^t = \gamma \mathbf{m}_y^{t-1} + (1 - \gamma) \bar{\mathbf{p}}_y^t, \quad (9)$$

where \mathbf{m}_y^t and \mathbf{m}_y^{t-1} denote the y -th row vector of the confusion matrix at $t-1$ -th and t -th iteration, respectively. $\bar{\mathbf{p}}_y^t$ is the averaged predicted probability vector and $\gamma \in (0, 1)$ is the momentum. Updates to categories that do not appear in current iteration are skipped.

As specified in Sec. 3.2, the transformed confusion matrix items ($\hat{M}_{i,j}, 1 \leq i, j \leq C$) by Eq. 5 reflect the conditional posterior probability for rectification in the form of pairwise biases. We thus directly enforce the model to learn such information as a fightback regularization. Specifically, we leverage the transformed confusion matrix items as soft targets. Formally, for each k -th foreground proposal of label y at current iteration, we apply a regularization upon foreground-classes learning via cross entropy below.

$$L_{PCB}(k) = - \sum_{i=1}^C \hat{M}_{i,y}^t \log \hat{p}_i \quad (10)$$

Beyond the rationale elaborated above, the regularization intuitively aims at balancing the pairwise biases. Taking a pair of classes a and b for illustration, the predicted confidence on class a and b is suppressed and raised simultaneously if there exists a preference $M_{a,b} < M_{b,a}$.

The regularization aims at relieving train-time pairwise class balance in the macro model dynamics level. We keep

the plain cross-entropy loss $L_{CE}(\cdot)$ for the micro sample level classification. Denote K the number of proposals in current minibatch, the total classification loss function becomes $L_{cls} = \sum_{k=1}^K L_{cls}(k)$ with each term as:

$$L_{cls}(k) = \begin{cases} \alpha L_{PCB}(k) + (1 - \alpha) L_{CE}(k), & y_k \neq C + 1 \\ L_{CE}(k), & otherwise \end{cases}, \quad (11)$$

where α is a coefficient to trade-off the two loss functions.

As shown in Tab. 6, raw exploitation of PCB regularization (Eq. 10) verifies the profits. Whereas, the performance deteriorates as the regularization term becomes dominant ($\alpha \geq 0.2$). We consider this is because current predictions are less discriminative at the first time to achieve the joint goal of the pairwise bias balancing and fundamental classification. Motivated by this and also the refinement inspiration in [7] (used to distinguish duplicated boxes in crowd detection), we propose a light prediction-dependent iterative paradigm to fully absorb the merits by PCB regularization. The core lies in using predictions stemming from a proposal feature to refine itself (then give a finer prediction next). The predicted boxes and confidences are projected as sparse spatial and channel attention to enhance the proposal feature, rather than direct concatenation as in [7] which could lead to semantic gaps. See the supplementary material for detailed analysis. Overall, Tab. 4 show the magic complementary benefits. The whole process takes place only in the final head part and remains the proposal box unchanged, unlike [2]. Technically, for each step r :

1. Calculating output classification logits $\mathbf{z}^r \in \mathbb{R}^{C+1}$ and predicted box $\mathbf{b}^r \in \mathbb{R}^{4C}$ with regard to classification and box head, given compact proposal feature map $X^r \in \mathbb{R}^{H_p \times W_p \times D}$ (H_p, W_p and D are pooled 7×7 resolution and channel depth). X^0 ($r = 0$) is the regular proposal feature X .
2. Computing the regular localization loss L_{loc}^r together with the modified classification loss L_{cls}^r (Eq. 11).
3. Using separate MLPs (FC-ReLU-FC) to project \mathbf{z}^r and \mathbf{b}^r as D - and $H_p W_p$ -dimensional features. Namely, obtaining $\mathbf{f}_z^r = MLP_{cls}(\mathbf{z}^r)$ and $\mathbf{f}_b^r = MLP_{loc}(\mathbf{b}^r)$ which are then warped into shapes of $\mathbb{R}^{1 \times 1 \times D}$ and $\mathbb{R}^{H \times W \times 1}$ respectively, denoted by X_z^r and X_b^r .
4. Refining proposal feature $X^{r+1} = X_b^r \otimes X^r \oplus X_z^r$. \otimes and \oplus denote broadcast element-wise multiplication and addition, respectively.

Naturally, proposal features in later refinement steps become more discriminative and could bear stronger regularization. We simply impose a linear increase of $\alpha_r = \frac{r-1}{R-1} \alpha$ as steps move on. Additionally, an outer step-wise coefficient

w_r is introduced, incorporating both the classification and localization loss and we derive the overall objective:

$$L = \left[\sum_{r=1}^R w_r (L_{cls}^r + L_{loc}^r) \right] + L_{mask}, \quad (12)$$

The mask prediction loss is calculated once since the segmentation sub-task has its own branch.

Extension to BCE classifier. Beyond cross-entropy, binary cross-entropy (BCE) variants of classifiers with sigmoid activation function are also prevalent in modern instance segmentors. Our proposed method could be seamlessly applied without modification apart from discarding Equation 5 since the PCB regularization is designed for foreground classes, naturally matching the BCE context. Experiments on EQL v2 [33] show the effectiveness.

4. Experiments

In this section, we conduct experiments on the LVIS dataset [13] to validate the effectiveness of our method. We demonstrate the complementarity of our method to other state-of-the-art long-tailed instance segmentation methods.

4.1. Experimental setup

Datasets. Experiments are performed on the large vocabulary instance segmentation dataset LVIS. The LVIS v0.5 contains 1230 categories with both instance mask annotations and bounding box annotations. The latest version LVIS v1.0 includes 1203 categories. We conduct experiments and ablation studies mainly on LVIS v0.5, while the results on the prevalent LVIS v1.0 are also provided. The train set is used for training and the validation set is used for evaluation. All the categories are divided into three splits according to the number of images that each category appears in the train set: *rare* (1-10 images), *common* (11-100 images) and *frequent* (> 100).

Evaluation metrics. Following common protocol, we adopt Average Precision (AP) as the evaluation metric, which is averaged across *IoU* threshold from 0.5 to 0.95. AP for the main mask prediction task are omitted as AP and AP for object detection is denoted as AP^b . AP50 and AP75 are also evaluated for comparison. We report detailed AP results on *rare*, *common* and *frequent* splits, denoted as AP^r , AP^c and AP^f . Besides, we also evaluate methods on a new metric to examine the pairwise class balance of confusion matrix on validation set, named Pairwise Bias (*PwB* in short). In detail, $PwB(M) = \|M - M^T\|_F$, where M^T stands for the transpose of M and $\|\cdot\|_F$ is the Frobenius Normalization.

Implementation details. For the implementation of our PCB, we set the dimension of hidden layers in MLP_{loc} and MLP_{cls} to 512. For EQL v2 [33], we additionally apply a LayerNorm [1] before MLP_{cls} only to make the training stable, which will not increase the performance. For a RoI

feature, it will go through the same prediction head for three iterations and generate classification scores and bounding box predictions three times, while only the last one is used for evaluation. The classification loss weights w_t for each of the three iterations are set to 0.2, 0.2 and 0.6 as to guarantee the performance of the last iteration. The EMA momentum γ is set to 0.99 as default, and we choose $\alpha = 0.4$ for any method without modification of the loss function. For implementation on EQL v2 [33] and Seesaw [38], L_{PCB} and L_{CE} should be modified accordingly. As PCB regularization only replaces one-hot label to soft target, this could be easily achieved by applying this to these methods. Due to the change of loss function, the model is strongly re-balanced, α is to 0.2 and 0.05 respectively, as a complement. The PCB regularization term is applied after the 16th epoch. For Seesaw [38], an RFS sampler is applied to have a fair comparison with LOCE [10] which also resamples the data. For training strategies, we follow the common recipe [38], please refer to our supplemental material for details.

4.2. Benchmark results

Our PCB, which aims to achieve pairwise class balance, can be a complement to any existing long-tailed instance segmentation method. Experiments are conducted on LVIS v0.5 and LVIS v1.0 across various base methods, including two state-of-the-art solutions EQL v2 [33] and Seesaw Loss [38]. Results are shown in Tab. 2. When equipped with PCB, the AP of all base methods gets improved, especially for AP_r . Even on the strong re-balancing method Seesaw, there is still an obvious rise of the rare AP (e.g., + 3.5 AP_r on LVIS v0.5). This is due to the exploration of the pairwise bias in PCB, which has not been considered in the previous works and thus still exists in those SOTA models. Another interesting thing is that PCB hardly ever hurts AP_f to compensate for the AP_r , and the AP_f is even improved sometimes. Additionally, the *PwB* metric value and the model AP are negatively correlated, our PCB can effectively decreases the *PwB*.

In Tab. 3, we report the comparison with those state-of-the-art methods across different backbone networks on LVIS v0.5 and LVIS v1.0. In all experiments, the implemented Seesaw + PCB achieve both the best mask AP and box AP. The gap between AP_r and AP_f is narrowed.

4.3. Ablation study

In this part, A comprehensive ablation study is conducted to analyze various components in our PCB. Experiments are conducted on LVIS v0.5 using Mask R-CNN with ResNet-50-FPN. RFS is applied unless otherwise specified.

Components in PCB. There are two components in PCB, the regularization term and the learning paradigm. Experiments are conducted to verify the choice of each, results of which is summarized in Tab. 4. To evaluate the effectiveness of designed PCB regularization, we conduct a comparison

Table 2. Results on the validation set of LVIS v0.5 and LVIS v1.0 with a ResNet-50 backbone. Our PCB can be complementary to various methods, including the baseline method Softmax, sampling methods (e.g., RFS [13]), one with binary classifier (e.g., EQL v2 [33]), and even strong state-of-the-art methods (e.g., Seesaw [38]). The results of RFS and EQL v2 on LVIS v1.0 are direct copied from [10], and the results of Seesaw on LVIS v1.0 are directly copied from [6].

Dataset		LVIS v0.5						LVIS v1.0					
Method	PCB	AP	AP _r	AP _c	AP _f	AP ^b	PwB ↓	AP	AP _r	AP _c	AP _f	AP ^b	PwB ↓
Softmax	✗	21.9	4.6	21.7	28.9	21.9	13.8	19.0	1.3	16.7	29.3	19.9	14.5
	✓	25.1	12.6	25.5	29.5	25.2	8.4	22.6	7.7	21.8	29.9	24.1	8.8
RFS [13]	✗	25.6	16.0	26.4	28.5	25.6	13.6	23.7	13.5	22.8	29.3	24.7	-
	✓	27.7	21.8	28.0	29.7	28.2	8.8	26.5	18.5	26.5	30.2	28.3	9.0
EQL v2 [33]	✗	26.9	17.8	27.2	29.5	26.5	10.9	25.5	17.7	24.3	30.2	26.1	-
	✓	27.8	20.9	28.4	29.9	28.1	7.4	26.2	18.2	25.9	30.1	27.3	6.5
Seesaw [38]	✗	27.8	19.9	28.9	29.5	27.3	7.9	26.8	19.8	26.3	30.5	27.6	7.5
	✓	28.8	23.4	29.6	30.0	28.6	7.0	27.2	19.0	27.1	30.9	28.1	6.4

Table 3. Comparing with the state-of-the-art on LVIS v0.5 and LVIS v1.0 with various backbones. † indicates results copied from [10]. ‡ indicates results copied from [6].

Dataset	Backbone	Method	AP	AP50	AP75	AP _r	AP _c	AP _f	AP ^b
LVIS v0.5	R-50-FPN	BAGS [22]†	26.3	-	-	18.0	26.9	28.7	25.8
		EQL v2 [33]†	26.9	41.5	28.9	17.8	27.2	29.5	26.5
		LOCE [10]†	28.4	-	-	22.0	29.0	30.2	28.2
		Seesaw [38]	27.8	42.6	29.6	19.9	28.9	29.5	27.3
		Seesaw + PCB	28.8	43.8	30.9	23.4	29.6	30.0	28.6
LVIS v1.0	R-50-FPN	EQL v2 [33]†	25.5	-	-	17.7	24.3	30.2	26.1
		LOCE [10]†	26.6	-	-	18.5	26.2	30.7	27.4
		Seesaw [38]‡	26.8	41.3	28.4	19.8	26.3	30.5	27.6
		Seesaw + PCB	27.2	41.7	29.4	19.0	27.1	30.9	28.1
LVIS v1.0	R-101-FPN	BAGS [22]†	25.6	-	-	17.3	25.0	30.1	26.4
		EQL v2 [33]†	27.2	-	-	20.6	25.9	31.4	27.9
		LOCE [10]†	28.0	-	-	19.5	27.8	32.0	29.0
		Seesaw [10]‡	28.2	42.7	30.2	21.0	27.8	31.8	28.9
		Seesaw + PCB	28.8	43.3	30.9	22.6	28.3	32.0	29.9

with label smoothing [27] and online label smoothing [49]. Apparently, PCB regularization surpasses both. It obtains improvements on all splits, especially for the rare classes (+2.9 AP_r), which differs from traditional schemes doing trade-offs. On the contrary, online label smoothing which does not consider the long-tailed distribution exacerbates the bias. The comparison between deep supervision (DSN) [20, 32] and our proposed iterative learning paradigm is also conducted. The performance of PCB gets improved when equipped with either DSN or iterative learning paradigm even if applying DSN alone hurts the performance, which shows the effectiveness of such a progressive learning manner. The prediction-based self-calibrated iterative learning paradigm still has advantages over DSN, the performance of all splits is better than that of which with DSN. The results of post-hoc confusion matrix calibration (CM) are also provided, of which the confusion matrix is calculated on the train set. The performance of PCB regularization surpasses the post-hoc confusion matrix calibration on all category

splits, which supports our choice to do online learning.

Hyper-parameters. We test the two main hyper-parameters in PCB, the momentum γ for the update of confusion matrix, and the PCB regularization coefficient α . The results are summarized in Tab. 5 and Tab. 6. In Tab. 5, with the increase of γ , the performance of rare classes is gradually improved, which is opposite to common classes. Overall, however, the performance is robust to the choice of γ . We vary α from 0 to 1. With the increase of α , AP_r first gets improved quickly then drops slowly. A similar phenomenon happens on AP_c and AP_f too. When $\alpha = 0$, the model gets no regularization to relieve the bias. And if $\alpha = 1$, too strong regularization harms the basic classification. We thus choose an $\alpha = 0.4$ to trade-off between pairwise balance and discrimination.

Prediction in each recurrent step. PCB outputs predictions at each recurrent step which is embedded to the input of next recurrent step. In Tab. 7, we evaluate the classification predictions of each step. The performance of rare classes is gradually raised through iterations, while the performance

Table 4. Ablation study of each component in PCB. The choice of regularization (Regu.), and the learning paradigm (Paradigm). Label smoothing (LS) [27] and online label smoothing (OLS) [49] is to compare with PCB regularization. Deep supervision (DSN) [20] is to compare with our iterative learning paradigm. The results of post-hoc confusion matrix calibration (CM) are also provided. † indicates using train set confusion matrix, which differs from that in Tab. 1.

Regu.	Paradigm	AP	AP _r	AP _c	AP _f	AP ^b
N/A	N/A	25.6	16.0	26.4	28.5	25.6
CM [†]	N/A	25.4	17.8	25.8	27.8	25.2
LS [27]	N/A	25.9	16.9	26.3	29.1	26.0
OLS [49]	N/A	25.6	15.4	26.4	28.6	25.7
PCB	N/A	26.7	18.9	27.5	28.8	27.1
N/A	DSN [20]	24.8	15.4	25.2	27.9	24.2
N/A	Iterative	26.5	17.4	27.3	29.2	26.8
PCB	DSN [20]	26.9	21.1	27.0	29.1	28.0
PCB	Iterative	27.7	21.8	28.0	29.7	28.2

Table 5. Analysis on the influence of different EMA momentum γ .

γ	AP	AP _r	AP _c	AP _f	AP ^b
0.9	27.6	19.3	28.8	29.3	28.3
0.99	27.7	21.8	28.0	29.7	28.2
0.999	27.5	22.1	27.6	29.5	28.4

Table 6. Analysis on the influence of different PCB regularization coefficient α . Experiments are conducted with RFS on LVIS v0.5.

α	AP	AP _r	AP _c	AP _f	AP ^b
0.0	26.5	17.4	27.3	29.2	26.8
0.2	27.4	20.3	27.8	29.6	28.1
0.4	27.7	21.8	28.0	29.7	28.2
0.6	27.7	21.3	28.2	29.5	28.5
0.8	27.4	21.3	27.9	29.1	28.4
1.0	26.6	20.9	26.9	28.6	27.6

of frequent classes is affected little or even improved. It indicates the mechanism of our PCB that, from recurrent step to recurrent step, the bias of the rare towards the frequent is relieved. The progressive debiasing manner ensures the performance of frequent classes. Interestingly, even the performance of the first recurrent step that is trained only by CE loss surpasses the corresponding base method in Tab. 2 and predictions from the middle step may even be comparable with the final prediction. It might be brought by the weight sharing of different recurrent steps.

Comparison with other complementary methods. Besides our PCB, the recent proposed NORCAL [28] is also a complementary method which does post-hoc calibration according to the training distribution. To show the superiority of PCB, experiments are conducted on various methods. The results are summarized in Tab. 8. It’s not surprising that NORCAL can improve the performance of RFS while fails

Table 7. The performance of each recurrent step’s prediction of a 3-step PCB. Experiments are conducted on LVIS v0.5.

Method	step id	AP	AP _r	AP _c	AP _f	AP ^b
Softmax	1	22.9	8.3	22.6	29.1	23.0
	2	24.6	10.6	25.0	29.6	24.7
	3	25.1	12.6	25.5	29.5	25.2
RFS [13]	1	26.8	20.7	27.1	29.0	27.0
	2	27.8	21.0	28.4	29.7	28.4
	3	27.7	21.8	28.0	29.7	28.2

Table 8. Comparison with NORCAL [28] which is also complementary to other methods. NORCAL is grid searched and reported the results from optimal hyper-parameters for each method.

Method	AP	AP _r	AP _c	AP _f	AP ^b
RFS	25.6	16.0	26.4	28.5	25.6
RFS + NORCAL	27.4	19.6	28.8	28.8	27.4
RFS + PCB	27.7	21.8	28.0	29.7	28.2
EQL v2	26.9	17.8	27.2	29.5	26.5
EQL v2 + NORCAL	26.8	20.5	27.0	29.2	26.4
EQL v2 + PCB	27.8	20.9	28.4	29.9	28.1
Seesaw	27.8	19.9	28.9	29.5	27.3
Seesaw + NORCAL	27.8	20.9	28.9	29.2	27.4
Seesaw + PCB	28.8	23.4	29.6	30.0	28.6

to do so on strong baselines. NORCAL only relies on the label distribution that lacks the inter-class relation, so will get embarrassed when the baseline is strongly re-balanced. It’s not the situation for PCB, that PCB can still relieve the pairwise bias to obtain further improvements even on strong strong EQL v2 and Seesaw as shown in Tab. 8.

5. Conclusions and Limitations

In this paper, we propose to utilize the pairwise biases existing in the confusion matrix statistics as strong and intuitive indicator to facilitate more balanced long-tailed instance segmentation. Such indicator possesses more fine-grained inter-class relationship details which help achieve much higher performance upper bound. In short, an online calibration method aiming at Pairwise Class Balance (PCB) is proposed to relieve the long-tailed instance segmentation by generating fightback soft targets in the form of a simple regularization. Towards more friendly regularization, an iterative learning paradigm is devised to progressively relieve the pairwise bias. Experimentally, our proposed PCB method improves the performance of various existing long-tailed instance segmentation methods, establishing a new state-of-the-art on the very challenging LVIS benchmark.

LIMITATION: There could be room for improving statistical manner of the confusion matrix which is a tradeoff between historical and coming-batch statistics. The former introduces lags in reflecting current model status, while the later is unrepresentative. We will ameliorate PCB in the future.

References

- [1] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016. [6](#)
- [2] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: Delving into high quality object detection. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 6154–6162, 2018. [2](#), [5](#)
- [3] Kaidi Cao, Colin Wei, Adrien Gaidon, Nikos Arechiga, and Tengyu Ma. Learning imbalanced datasets with label-distribution-aware margin loss. *arXiv preprint arXiv:1906.07413*, 2019. [2](#)
- [4] Nadine Chang, Zhiding Yu, Yu-Xiong Wang, Anima Anandkumar, Sanja Fidler, and Jose M Alvarez. Image-level or object-level? a tale of two resampling strategies for long-tailed detection. *arXiv preprint arXiv:2104.05702*, 2021. [1](#), [3](#)
- [5] Nitesh V Chawla, Kevin W Bowyer, Lawrence O Hall, and W Philip Kegelmeyer. Smote: synthetic minority over-sampling technique. *Journal of artificial intelligence research*, 16:321–357, 2002. [2](#)
- [6] Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jiarui Xu, et al. Mmdetection: Open mmlab detection toolbox and benchmark. *arXiv preprint arXiv:1906.07155*, 2019. [7](#), [11](#)
- [7] Xuangeng Chu, Anlin Zheng, Xiangyu Zhang, and Jian Sun. Detection in crowded scenes: One proposal, multiple predictions. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 12214–12223, 2020. [5](#), [12](#)
- [8] Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge Belongie. Class-balanced loss based on effective number of samples. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 9268–9277, 2019. [2](#)
- [9] Rahul Duggal, Scott Freitas, Sunny Dhamnani, Duen Horng Chau, and Jimeng Sun. Elf: An early-exiting framework for long-tailed classification. *arXiv preprint arXiv:2006.11979*, 2020. [1](#)
- [10] Chengjian Feng, Yujie Zhong, and Weilin Huang. Exploring classification equilibrium in long-tailed object detection. In *Int. Conf. Comput. Vis.*, pages 3417–3426, 2021. [1](#), [2](#), [3](#), [4](#), [6](#), [7](#)
- [11] Ross Girshick. Fast r-cnn. In *Int. Conf. Comput. Vis.*, pages 1440–1448, 2015. [2](#)
- [12] Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. Accurate, large mini-batch sgd: Training imagenet in 1 hour. *arXiv preprint arXiv:1706.02677*, 2017. [12](#)
- [13] Agrim Gupta, Piotr Dollar, and Ross Girshick. Lvis: A dataset for large vocabulary instance segmentation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 5356–5364, 2019. [1](#), [2](#), [6](#), [7](#), [8](#), [11](#)
- [14] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Int. Conf. Comput. Vis.*, pages 2961–2969, 2017. [2](#), [3](#), [11](#)
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 770–778, 2016. [11](#)
- [16] Yin-Yin He, Jianxin Wu, and Xiu-Shen Wei. Distilling virtual examples for long-tailed recognition. *arXiv preprint arXiv:2103.15042*, 2021. [2](#)
- [17] Youngkyu Hong, Seungju Han, Kwanghee Choi, Seokjun Seo, Beomsu Kim, and Buru Chang. Disentangling label distribution for long-tailed visual recognition. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 6626–6636, 2021. [2](#), [5](#)
- [18] Ting-I Hsieh, Esther Robb, Hwann-Tzong Chen, and Jia-Bin Huang. Droploss for long-tail instance segmentation. *arXiv preprint arXiv:2104.06402*, 2021. [3](#)
- [19] Bingyi Kang, Saining Xie, Marcus Rohrbach, Zhicheng Yan, Albert Gordo, Jiashi Feng, and Yannis Kalantidis. Decoupling representation and classifier for long-tailed recognition. *arXiv preprint arXiv:1910.09217*, 2019. [2](#), [12](#)
- [20] Chen-Yu Lee, Saining Xie, Patrick Gallagher, Zhengyou Zhang, and Zhuowen Tu. Deeply-supervised nets. In *Artificial intelligence and statistics*, pages 562–570. PMLR, 2015. [7](#), [8](#)
- [21] Tianhao Li, Limin Wang, and Gangshan Wu. Self supervision to distillation for long-tailed visual recognition. In *Int. Conf. Comput. Vis.*, pages 630–639, 2021. [2](#)
- [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 *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 10991–11000, 2020. [3](#), [7](#)
- [23] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 2117–2125, 2017. [11](#)
- [24] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Int. Conf. Comput. Vis.*, pages 2980–2988, 2017. [2](#)
- [25] Zachary Lipton, Yu-Xiang Wang, and Alexander Smola. Detecting and correcting for label shift with black box predictors. pages 3122–3130. PMLR, 2018. [3](#)
- [26] Ziwei Liu, Zhongqi Miao, Xiaohang Zhan, Jiayun Wang, Boqing Gong, and Stella X Yu. Large-scale long-tailed recognition in an open world. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 2537–2546, 2019. [12](#)
- [27] Rafael Müller, Simon Kornblith, and Geoffrey Hinton. When does label smoothing help? *arXiv preprint arXiv:1906.02629*, 2019. [7](#), [8](#)
- [28] Tai-Yu Pan, Cheng Zhang, Yandong Li, Hexiang Hu, Dong Xuan, Soravit Changpinyo, Boqing Gong, and Wei-Lun Chao. On model calibration for long-tailed object detection and instance segmentation. *arXiv preprint arXiv:2107.02170*, 2021. [3](#), [8](#)
- [29] Giorgio Patrini, Alessandro Rozza, Aditya Krishna Menon, Richard Nock, and Lizhen Qu. Making deep neural networks robust to label noise: A loss correction approach. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1944–1952, 2017. [3](#)

- [30] Jiawei Ren, Cunjun Yu, Shunan Sheng, Xiao Ma, Haiyu Zhao, Shuai Yi, and Hongsheng Li. Balanced meta-softmax for long-tailed visual recognition. *arXiv preprint arXiv:2007.10740*, 2020. [1](#), [3](#), [12](#)
- [31] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Adv. Neural Inform. Process. Syst.*, 28:91–99, 2015. [2](#)
- [32] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1–9, 2015. [7](#)
- [33] 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 *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1685–1694, 2021. [3](#), [6](#), [7](#), [11](#)
- [34] Jingru Tan, Changbao Wang, Buyu Li, Quanquan Li, Wanli Ouyang, Changqing Yin, and Junjie Yan. Equalization loss for long-tailed object recognition. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 11662–11671, 2020. [3](#)
- [35] Kaihua Tang, Jianqiang Huang, and Hanwang Zhang. Long-tailed classification by keeping the good and removing the bad momentum causal effect. 2020. [3](#), [4](#)
- [36] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. Fcos: Fully convolutional one-stage object detection. In *Int. Conf. Comput. Vis.*, pages 9627–9636, 2019. [2](#)
- [37] Grant Van Horn, Oisin Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam, Pietro Perona, and Serge Belongie. The inaturalist species classification and detection dataset. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 8769–8778, 2018. [2](#)
- [38] 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 *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 9695–9704, 2021. [1](#), [3](#), [4](#), [6](#), [7](#), [11](#)
- [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 *Eur. Conf. Comput. Vis.*, pages 728–744. Springer, 2020. [1](#), [3](#)
- [40] Tong Wang, Yousong Zhu, Chaoyang Zhao, Wei Zeng, Jinqiao Wang, and Ming Tang. Adaptive class suppression loss for long-tail object detection. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 3103–3112, 2021. [3](#)
- [41] Xinlong Wang, Tao Kong, Chunhua Shen, Yuning Jiang, and Lei Li. Solo: Segmenting objects by locations. In *Eur. Conf. Comput. Vis.*, pages 649–665. Springer, 2020. [2](#)
- [42] Xudong Wang, Long Lian, Zhongqi Miao, Ziwei Liu, and Stella X Yu. Long-tailed recognition by routing diverse distribution-aware experts. *arXiv preprint arXiv:2010.01809*, 2020. [2](#)
- [43] Xinlong Wang, Rufeng Zhang, Tao Kong, Lei Li, and Chunhua Shen. Solov2: Dynamic and fast instance segmentation. *arXiv preprint arXiv:2003.10152*, 2020. [2](#)
- [44] Jialian Wu, Liangchen Song, Tiancai Wang, Qian Zhang, and Junsong Yuan. Forest r-cnn: Large-vocabulary long-tailed object detection and instance segmentation. In *ACM Int. Conf. Multimedia*, pages 1570–1578, 2020. [3](#)
- [45] Tong Wu, Qingqiu Huang, Ziwei Liu, Yu Wang, and Dahua Lin. Distribution-balanced loss for multi-label classification in long-tailed datasets. In *Eur. Conf. Comput. Vis.*, pages 162–178. Springer, 2020. [2](#)
- [46] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 1492–1500, 2017. [12](#)
- [47] Yuzhe Yang and Zhi Xu. Rethinking the value of labels for improving class-imbalanced learning. *arXiv preprint arXiv:2006.07529*, 2020. [2](#)
- [48] Cheng Zhang, Tai-Yu Pan, Yandong Li, Hexiang Hu, Dong Xuan, Soravit Changpinyo, Boqing Gong, and Wei-Lun Chao. Mosaicos: a simple and effective use of object-centric images for long-tailed object detection. In *Int. Conf. Comput. Vis.*, pages 417–427, 2021. [3](#)
- [49] Chang-Bin Zhang, Peng-Tao Jiang, Qibin Hou, Yunchao Wei, Qi Han, Zhen Li, and Ming-Ming Cheng. Delving deep into label smoothing. *IEEE Trans. Image Process.*, 30:5984–5996, 2021. [3](#), [7](#), [8](#)
- [50] Shifeng Zhang, Cheng Chi, Yongqiang Yao, Zhen Lei, and Stan Z Li. Bridging the gap between anchor-based and anchor-free detection via adaptive training sample selection. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 9759–9768, 2020. [2](#)
- [51] Songyang Zhang, Zeming Li, Shipeng Yan, Xuming He, and Jian Sun. Distribution alignment: A unified framework for long-tail visual recognition. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 2361–2370, 2021. [3](#)
- [52] Yongshun Zhang, Xiu-Shen Wei, Boyan Zhou, and Jianxin Wu. Bag of tricks for long-tailed visual recognition with deep convolutional neural networks. In *AAAI*, volume 35, pages 3447–3455, 2021. [3](#)
- [53] Boyan Zhou, Quan Cui, Xiu-Shen Wei, and Zhao-Min Chen. Bbn: Bilateral-branch network with cumulative learning for long-tailed visual recognition. In *IEEE Conf. Comput. Vis. Pattern Recog.*, pages 9719–9728, 2020. [2](#)

A. Details in implementation

Implementation of MS calibration. We found the original MS calibration mentioned in the main paper did not work well. The deep reason is: If instances of rare class i are always predicted as frequent class j with very high confidence (i.e., $M_{i,j} \approx 1$ and $M_{i,i} \approx 0$), then plenty of instances will be miscalibrated into class i ($s_i \approx 0$). To avoid this, we did two slight modifications to the original MS calibration in our experiments. Firstly, instead of using $s_i = M_{i,i}$, we adopted $s_i = \sum_{k=1}^C M_{k,i}$ to soften the distribution. Secondly, if s_i was still close to 0, we did not predict class i in this case. As shown in Table 9, the modifications do help improve the performance of MS calibration on each split.

Table 9. The performance of MS calibration before and after applied two modifications. Experiments are conducted on LVIS v0.5 using Mask R-CNN with ResNet-50-FPN and RFS [13] sampler.

Modification #1	Modification #2	AP	AP _r	AP _c	AP _f	AP ^b
	✓	20.0	13.5	19.1	23.7	19.9
✓		26.5	20.0	27.7	27.5	26.3
✓	✓	27.0	20.6	28.2	28.0	27.0

Applied to EQL v2 [33] and Seesaw [38]. PCB regularization can be easily applied to EQL v2 [33] and Seesaw loss [38] as they only change the loss weight or model activation. For each k -th proposal of label y at current iteration, the original equalization loss v2 can be formulated as follows:

$$L_{EQL\ v2}(k) = - \sum_{i=1}^C w_i [y_i \log \hat{p}_i + (1 - y_i) \log(1 - \hat{p}_i)], \quad (13)$$

where $\hat{p}_i = 1 / (1 + e^{-z_i^{fg}})$ is the post-sigmoid probability for class i , y_i indicates whether the proposal belongs to class i , and w_i is a weight calculated from gradient perspective. When applied to EQL v2, our PCB regularization becomes

$$L_{PCB} = - \sum_{i=1}^C w_i [\hat{M}_{i,y}^t \log \hat{p}_i + (1 - \hat{M}_{i,y}^t) \log(1 - \hat{p}_i)]. \quad (14)$$

And the total classification loss function for the proposal is

$$L_{cls}(k) = \alpha L_{PCB}(k) + (1 - \alpha) L_{EQL\ v2}. \quad (15)$$

Similarly, we define the Seesaw variant of PCB regularization. The original Seesaw loss terms as

$$L_{seesaw}(k) = - \sum_{i=1}^C y_i \log \hat{p}_i, \quad (16)$$

$$\text{with } \hat{p}_i = \frac{\exp(z_i^{fg})}{\sum_{j \neq i}^C \mathcal{S}_{ij} \exp(z_j^{fg}) + \exp(z_i^{fg})}.$$

Table 10. Analysis on the influence of calculating the regression loss in each step or only last step. Experiments are conducted on LVIS v0.5.

Method	Regression	AP	AP _r	AP _c	AP _f	AP ^b
Softmax	each	25.1	11.1	25.9	29.6	25.3
	last	25.1	12.6	25.5	29.5	25.2
RFS	each	27.5	20.4	28.2	29.4	27.8
	last	27.7	21.8	28.0	29.7	28.2

The \mathcal{S}_{ij} is composed by a mitigation factor and a compensation factor. So, the PCB regularization can be written as

$$L_{PCB}(k) = - \sum_{i=1}^C \hat{M}_{i,y}^t \log \hat{p}_i, \quad (17)$$

$$\text{with } \hat{p}_i = \frac{\exp(z_i^{fg})}{\sum_{j \neq i}^C \mathcal{S}_{ij} \exp(z_j^{fg}) + \exp(z_i^{fg})}.$$

We combine them to get

$$L_{cls}(k) = \alpha L_{PCB}(k) + (1 - \alpha) L_{seesaw}. \quad (18)$$

Regression loss calculation. In practice, we calculate the regression loss only in the last recurrent step, rather than in each step. So the overall objective becomes:

$$L = \left[\sum_{r=1}^R w_r L_{cls}^r \right] + L_{loc}^R + L_{mask}. \quad (19)$$

It shows similar performance to the later, while obtains higher AP_r. The results of comparison are shown in Tab. 10. While the overall mask AP and box AP are comparable for the two manners, there is a constant improvement in performance of rare classes (over 1 AP_r) for the former.

Training details. Following [38], we implement our method with mmdetection [6]. Mask R-CNN [14] with ResNet-50-FPN and ResNet-101-FPN [15, 23] is adopted as our baseline model. We utilize the standard 2× schedule for LVIS of both versions. The models are trained using SGD with 0.9 momentum and 0.0001 weight decay for 24 epochs. With batch size of 16 on 8 GPUs, the initial learning rate is set to 0.02 and is decreased by 0.1 after 16 and 22 epochs, respectively. The training data augmentations include scale jittering (640-800) and horizontal flipping. For evaluation, we set the maximum number of detections per image to 300 and the minimum score threshold to 0.0001, as [13].

B. Analysis on α w/o iterative learning paradigm

As discussed in Sec. 3.3 of the main paper, the performance will deteriorate soon with the increase of α if the

Table 11. Analysis on the influence of different PCB regularization coefficient α without iterative learning paradigm. Experiments are conducted with RFS on LVIS v0.5.

α	AP	AP _r	AP _c	AP _f	AP ^b
0.0	25.6	16.0	26.4	28.5	25.6
0.2	26.7	18.9	27.5	28.8	27.1
0.4	26.5	18.3	27.4	28.7	26.7
0.6	26.2	19.7	26.4	28.4	26.6
0.8	26.0	18.5	26.7	28.1	26.5
1.0	25.1	18.7	25.9	26.5	25.3

Table 12. Comparison with refine module (RM in short) in CrowdDet [7]. Experiments are conducted with RFS on LVIS v0.5. PCB regularization is applied.

Paradigm	AP	AP _r	AP _c	AP _f	AP ^b
N/A	26.7	18.9	27.5	28.8	27.1
RM	26.4	20.8	26.5	28.5	26.9
Iterative	27.7	21.8	28.0	29.7	28.2

iterative learning paradigm is not applied. Tab. 11 shows an example. As α increases, AP_r gets improved until $\alpha = 0.6$, while AP_c and AP_f decline soon after $\alpha > 0.2$. So the PCB regularization hurts the fundamental classification, and the flexibility of debiasing is limited. By applying iterative learning paradigm, which guarantees the fundamental classification, such worry gets relieved. The room for debiasing is increased.

C. Comparing with refinement module in CrowdDet [7]

We notice that the refinement module (RM) in [7] is similar to our iterative learning paradigm. There are two main differences, RM concatenates the predictions and features rather than element-wise operation, and it utilizes features of the penultimate layer. We also provide the results of adopting RM as the learning paradigm, which are summarized in Tab. 12. While RM achieves promising AP_r compared to PCB regularization, it hurts the performance of common classes and frequent classes much, so the overall AP drops. Different from RM, our proposed iterative learning paradigm guarantees the performance of common and frequent classes.

D. Extension to long-tailed classification.

We also extend our PCB to long-tailed classification to testify its generalization ability. The commonly used ImageNet-LT [26] dataset is adopted in our experiment, and we use ResNeXt-50 [46] as the backbone network. Models are trained for 90 epochs with batch size 512. The initial learning rate is set to 0.2 and the first 5 epochs are trained with linear warm-up learning rate schedule [12]. The learning rate is decayed at 60th and 80th epoch by 0.1. For the implementa-

Table 13. Accuracy on ImageNet-LT with a ResNeXt-50 backbone.

Method	PCB	Many	Medium	Few	Overall
CE	✗	67.76	38.89	7.44	45.73
	✓	61.68	49.10	21.97	50.26
BSCE [30]	✗	62.65	48.75	25.44	50.94
	✓	61.72	49.66	34.85	52.29

tion of PCB, we ignore the MLP_{loc} and set the dimension of hidden layers in MLP_{cls} to 256 for simplicity. For a feature vector from the backbone, it will go through the same classifier for two times, and the last prediction is used for evaluation. γ is set to 0.9999. We train PCB in a decoupled manner [19], so the PCB regularizer is only applied in the fine-tune phase.

Two methods are utilized as baseline, CE and BSCE [30]. The results are in Tab. 13. Equipped with PCB (α is set to 0.8 and 0.15 respectively), the performance gain is significant and consistent, the accuracy of few split is raised almost 10% even on the strong baseline. The results fully demonstrate the generalization ability of our PCB.