

Overcoming Classifier Imbalance for Long-tail Object Detection with Balanced Group Softmax

Yu Li^{1,2,3*}, Tao Wang^{3,4}, Bingyi Kang³, [†]Sheng Tang^{1,2}, Chunfeng Wang², Jintao Li^{1,2}, Jiashi Feng³

¹Key Laboratory of Intelligent Information Processing, Institute of Computing Technology, Chinese Academy of Sciences, Beijing, China

²University of Chinese Academy of Sciences, Beijing, China

³Department of Electrical and Computer Engineering, National University of Singapore, Singapore

⁴Institute of Data Science, National University of Singapore, Singapore

{liyu, ts, jtlj}@ict.ac.cn, twangnh@gmail.com, kang@u.nus.edu,

wangchunfeng14@mails.ucas.ac.cn, elefjia@nus.edu.sg

Abstract

Solving long-tail large vocabulary object detection with deep learning based models is a challenging and demanding task, which is however under-explored. In this work, we provide the first systematic analysis on the underperformance of state-of-the-art models in front of long-tail distribution. We find existing detection methods are unable to model few-shot classes when the dataset is extremely skewed, which can result in classifier imbalance in terms of parameter magnitude. Directly adapting long-tail classification models to detection frameworks can not solve this problem due to the intrinsic difference between detection and classification. In this work, we propose a novel balanced group softmax (BAGS) module for balancing the classifiers within the detection frameworks through group-wise training. It implicitly modulates the training process for the head and tail classes and ensures they are both sufficiently trained, without requiring any extra sampling for the instances from the tail classes. Extensive experiments on the very recent long-tail large vocabulary object recognition benchmark LVIS show that our proposed BAGS significantly improves the performance of detectors with various backbones and frameworks on both object detection and instance segmentation. It beats all state-of-the-art methods transferred from long-tail image classification and establishes new state-of-the-art. Code is available at <https://github.com/FishYuLi/BalancedGroupSoftmax>.

*This work was supported by the National Key Research and Development Program of China (2017YFC0820605), and the National Natural Science Foundation of China (61525206, 61572472, U1703261, 61871004), and 242 project (2019A010). Jiashi Feng was partially supported by A1SG R-263-000-D97-490, NUS ECRA R-263-000-C87-133 and MOE Tier-II R-263-000-D17-112. Yu Li was partially supported by the program of China Scholarships Council (No.201904910801).[†]: corresponding author.

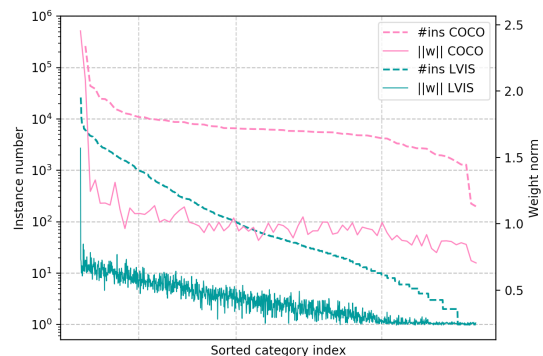


Figure 1. Sorted number of training instances (#ins) for categories in COCO and LVIS training set, and the corresponding classifier weight norm $\|w\|$ from Faster R-CNN model trained on COCO and LVIS. The x-axis represents the sorted category index of COCO and LVIS. We align 80 classes of COCO with 1230 classes of LVIS for better visualization. Category 0 indicates background.

1. Introduction

Object detection [31, 29, 25, 23, 21, 1] is one of the most fundamental and challenging tasks in computer vision. Recent advances are mainly driven by large-scale datasets that are manually balanced, such as PASCAL VOC [9] and COCO [24]. However in reality, the distribution of object categories is typically long-tailed [30]. Effective solutions that adapt state-of-the-art detection models to such class-imbalanced distribution are highly desired yet still absent. Recently, a long-tail large vocabulary object recognition dataset LVIS [14] is released, which greatly facilitates object detection research in much more realistic scenarios.

A straightforward solution to long-tail object detection is to train a well-established detection model (e.g., Faster R-CNN [31]) on the long-tail training data directly. However, big performance drop would be observed when adapting detectors designed for fairly balanced datasets (e.g., COCO) to a long-tail one (e.g., LVIS), for which the rea-

sons still remain unclear due to multiple entangled factors. Inspired by [20], we decouple the representation and classification modules within the detection framework, and find the weight norms of the proposal classifier corresponding to different categories are severely imbalanced, since low-shot categories get few chances to be activated. Through our analysis, this is one direct cause of the poor long-tail detection performance, which is intrinsically induced by data imbalance. As shown in Figure 1, we sort the category-wise classifier weight norms of models trained on COCO and LVIS respectively by the number of instances in the training set. For COCO, the relatively balanced data distribution leads to relatively balanced weight norms for all categories, except for *background* class (CID=0, CID for Category ID). For LVIS, it is obvious that the category weight norms are imbalanced and positively correlated with the number of training instances. Such imbalanced classifiers (w.r.t. their parameter norm) would make the classification scores for low-shot categories (tail classes) much smaller than those of many-shot categories (head classes). After standard softmax, such imbalance would be further magnified thus the classifier wrongly suppresses the proposals predicted as low-shot categories.

The classifier imbalance roots in data distribution imbalance—classifiers for the many-shot categories would see more and diverse training instances, leading to dominating magnitude. One may consider using solutions to long-tail classification to overcome such an issue, including re-sampling training instances to balance the distribution [15, 7, 32, 26] and re-weighting classification loss at category level [5, 2, 18] or instance level [23, 33]. The re-sampling based solutions are applicable to detection frameworks, but may lead to increased training time and overfitting risk to the tail classes. Re-weighting based methods are unfortunately very sensitive to hyper-parameter choices and not well applicable to detection frameworks due to difficulty in dealing with the special *background* class, an extremely many-shot category. We empirically find none of these methods works well on long-tail detection problem.

In this work, to address the classifier imbalance, we introduce a simple yet effective balanced group softmax (BAGS) module into the classification head of a detection framework. We propose to put object categories with similar numbers of training instances into the same group and compute group-wise softmax cross-entropy loss separately. Treating categories with different instance numbers separately can effectively alleviate the domination of the head classes over tail classes. However, due to the lack of diverse negative examples for each group training, the resultant model suffers too many false positives. Thus, BAGS further adds a category *others* into each group and introduces the *background* category as an individual group, which can alleviate the suppression from head classes over tail classes

by keeping their classifiers balanced while preventing false positives by categories *background* and *others*.

We experimentally find BAGS works very well. It improves by 9% – 19% the performance on tail classes of various frameworks including Faster R-CNN [31], Cascade R-CNN [1], Mask R-CNN [16] and HTC [4] with ResNet-50-FPN [17, 22] and ResNeXt-101-x64x4d-FPN [37] backbones consistently on the long-tail object recognition benchmark LVIS [14], with the overall mAP lifted by around 3% – 6%.

To sum up, this work makes following contributions:

- Through comprehensive analysis, we reveal the reason why existing models perform not well for long-tail detection, *i.e.* their classifiers are imbalanced and not trained equally well, reflected by the observed imbalanced classifier weight norms.
- We propose a simple yet effective balanced group softmax module to address the problem. It can be easily combined with object detection and instance segmentation frameworks to improve their long-tail recognition performance.
- We conduct extensive evaluations with state-of-the-art long-tail classification methods for object detection. Such benchmarking not only deepens our understandings of these methods as well as the unique challenges of long-tail detection, but also provides reliable and strong baselines for future research in this direction.

2. Related Works

Compared with balanced distribution targeted object detection [11, 31, 1], and few-shot object detection [19, 3, 38, 10], the challenging and practical long-tail object detection problem is still underexplored. Though Ouyang *et al.* [28] proposes the concept of long-tail object detection, their work focuses on the imbalanced training data distribution on ILSVRC DET dataset [6] without few-shot setting for tail classes like LVIS [14]. [14] proposes repeat factor sampling (RFS) serving as a baseline. Classification calibration [36] enhances RFS by calibrating classification scores of tail classes with another head trained with ROI level class-balanced sampling strategy. Below we first review general object detection methods, and then long-tail classification methods.

General object detection Deep learning based object detection frameworks are divided into anchor-based and anchor-free ones. Anchor-based approaches [12, 11, 31, 29, 23] explicitly or implicitly extract features for individual regions thus convert object detection into proposal-level classification which have been largely explored. In contrast, anchor-free approaches focus on detecting key points of objects and construct final detection boxes by properly combining detected key points [21, 8, 40] or expanding the rep-

resentation of key points [39, 35]. For such detectors, proposal classification is achieved by classifying the key points.

These popular object detection frameworks all employ a softmax classifier for either proposal classification or key-point classification. Our proposed balanced group softmax module can be easily plugged into such mainstream detectors by simply replacing the original softmax classifier. For simplicity, we mainly experiment with anchor-based detectors Faster R-CNN [31] and Cascade R-CNN [1] as well as their corresponding instance segmentation approaches Mask R-CNN [16] and HTC [4].

Long-tail classification Long-tail classification is attracting increasing attention due to its realistic applications. Current works leverage data re-sampling, cost-sensitive learning, or other techniques. For data re-sampling methods, training samples are either over-sampled (adding copies of training samples for tail classes) [15], under-sampled (deleting training samples for head classes) [7], or class-balanced sampled [32, 26], which motivates RFS [14]. For cost-sensitive learning, the network losses are re-weighted at category level by multiplying different weights on different categories to enlarge the influence of tail-class training samples [5, 2, 18] or at instance level by multiplying different weights on different training samples for more fine-grained control [23, 33]. Some other approaches optimize the classifier trained with long-tail data such as Nearest Class Mean classifier (NCM) [27, 13], and τ -normalized classifier [20]. These methods are usually sensitive to hyper-parameters and do not perform well when transferred to detection frameworks due to the inherent difference between classification and detection as stated in Sec. 1.

Therefore, an approach specifically designed for long-tail object detection is desirable, and our work is the first successful attempt to overcome classifier imbalance through group-wise training without extra sampling from tail classes.

3. Preliminary and Analysis

3.1. Preliminary

We first revisit the popular two-stage object detection framework, by taking Faster R-CNN [31] as an example. We adopt such a two-stage framework to develop and implement our idea.

The backbone network f_{back} takes an image I as input, and generates a feature map $F = f_{\text{back}}(I)$. The feature map is then passed to ROI-align [16] or ROI-pooling [11] to produce K proposals with their own feature $F_k = \text{ROIAlign}(F, b_k)$. Here b_k denotes proposal k . The classification head f_{head} then extracts a d -dimensional feature $h = f_{\text{head}}(F_k)$ for each of the proposals. Finally, one FC (fully connected) layer is used to transfer h to the $(C + 1)$ -category prediction (C object classes plus *background*) by

$z = Wh + b$, where $W \in \mathbb{R}^{d \times (C+1)}$ is the classifier weights with each column $w_j \in \mathbb{R}^d$ related to one specific category j , and b is the bias term.

During training, with ground truth label $y \in \{0, 1\}^{C+1}$, softmax cross entropy is applied to compute loss for a specific proposal:

$$\mathcal{L}_k(p, y) = - \sum_{j=0}^C y_j \log(p_j), \quad (1)$$

$$p_j = \text{softmax}(z_j) = \frac{e^{z_j}}{\sum_{i=0}^C e^{z_i}}. \quad (2)$$

Here z_j denotes the j -th element of z and p_j is the predicted probability of the proposal being an instance of category j .

3.2. Analysis

Current well-performing detection models often fail to recognize tail classes when the training set follows a long-tailed distribution. In this section, we try to investigate the underlying mechanism behind such performance drop from balanced dataset to long-tailed dataset, by conducting contrast experiments on their representative examples, *i.e.*, COCO and LVIS.

We adopt a Faster R-CNN [11] model with R50-FPN backbone. By directly comparing the mAP on the two datasets, the performance drops notably from 36.4%(COCO) to 20.9%(LVIS). Despite the unfairness as LVIS contains much more classes than COCO (1230 v.s. 80), we can still draw some interesting observations. On head classes, the LVIS model achieves comparable results with COCO. However, when it comes to tail classes, the performance decreases to 0 rapidly. Such a phenomenon implies current detection models are indeed challenged by data imbalance. To further investigate how the performance degradation is induced by data imbalance, we decouple the detection framework into proposal feature extraction stage and proposal classification stage, following [20].

Specifically, following the notations in Sec. 3.1, we deem the operations used to generate h as proposal feature extraction, and the last FC layer and softmax in Eqn. (2) as a softmax classifier. Then, we investigate the correlation between the number of training instances and the weight norm $\|w_j\|$ in the classifier for each category. The results are visualized in Figure 1. We can see for COCO dataset, most categories contain $10^3 - 10^4$ training instances (at least 10^2); and classifier weight norms are also relatively balanced (0.75-1.25) for all foreground categories¹. In contrast, for the LVIS dataset, a weight norm $\|w_j\|$ is highly related to the number of training instances in the corresponding category j ; the more training examples there are, the larger weight magnitude it will be. For the extreme few-shot categories (tail classes), their corresponding weight norms are extremely

¹Note that the first class is *background*(CID=0).

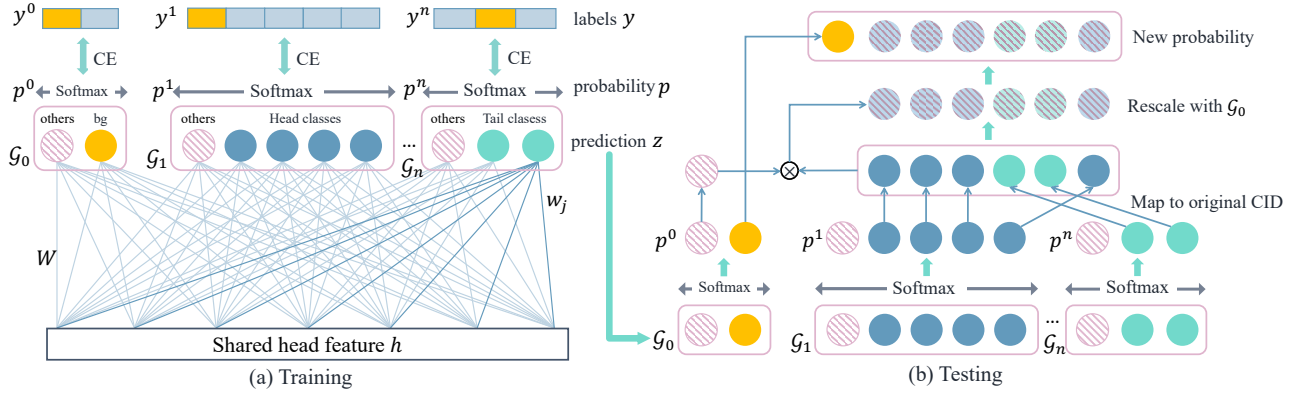


Figure 2. Framework of our balanced group softmax module. (a) Training: Classes containing similar training instances are grouped together. Class *others* is added to each group. \mathcal{G}_0 denotes *background* group. Softmax cross entropy (CE) loss is applied to each group individually. (b) Testing: With new prediction z , softmax is applied in each group, and probabilities are ordered by their original category id (CID) and re-scaled with foreground probability, generating new probability vectors for following post process.

small, even close to zero. Based on such observations, one can foresee that prediction scores for tail classes will be congenitally lower than head classes, and proposals of tail classes will be less likely to be selected after competing with those of head categories within the softmax computation. This explains why current detection models often fails on tail classes.

Why would the classifier weights be correlated to the number of training instances per-class? To answer this question, let us further inspect the training procedure of Faster R-CNN. When proposals from a head class j are selected as training samples, z_j should be activated, while the predictions for other categories should be suppressed. As the training instances for head classes are much more than those of tail classes (e.g., 10,000 vs. 1 in some extreme cases), classifier weights of tail classes are much more likely (frequent) to be suppressed by head class ones, resulting in imbalanced weight norms after training.

Therefore, one may see why re-sampling method [14, 36] is able to benefit tail classes on long-tail instance classification and segmentation. It simply increases the sampling frequency of tail class proposals during training so that the weights of different classes can be equally activated or suppressed, thus balance the tail and head classes to some degree. Also, loss re-weighting methods [5, 2, 18, 23, 33] can take effect in a similar way. Though the resampling strategy is able to alleviate data imbalance, it actually introduces new risks like overfitting to tail classes and extra computation overhead. Meanwhile, loss re-weighting is sensitive to per-class loss weight design, which usually varies across different frameworks, backbones and datasets, making it hardly deployable in real-world applications. Moreover, re-weighting based methods cannot handle the *background* class well in detection problems. Therefore, we propose a simple yet effective solution to balance the classifier weight norms without heavy hyper-parameter engineering.

4. Balanced Group Softmax

Our novel balanced group softmax module is illustrated in Figure 2. We first elaborate on its formulation and then explain the design details.

4.1. Group softmax

As aforementioned, detector performance is harmed by the positive correlation between weight norms and number of training examples. To solve this problem, we propose to divide classes into several disjoint groups and perform the softmax operation separately, such that only classes with similar numbers of training instances are competing with each other within each group. In this way, classes containing significantly different numbers of instances can be isolated from each other during training. The classifier weights of tail classes would not be substantially suppressed by head classes.

Concretely, we divide all the C categories into N groups according to their training instance numbers. We assign category j to group \mathcal{G}_n if

$$s_n^l \leq \mathcal{N}(j) < s_n^h, \quad n > 0 \quad (3)$$

where $\mathcal{N}(j)$ is the number of ground-truth bounding boxes for category j in the training set, and s_n^l and s_n^h are hyper-parameters that determine minimal and maximal instance numbers for group n . In this work, we set $s_{n+1}^l = s_n^h$ to ensure there is no overlap between groups, and each category can only be assigned to one group. N and s_n^l are set empirically to make sure that categories in each group contain similar total numbers of training instances. Throughout this paper, we set $N = 4$, $s_1^l = 0$, $s_2^l = 10$, $s_3^l = 10^2$, $s_4^l = 10^3$, $s_4^h = +\infty$.

Besides, we manually set the \mathcal{G}_0 to contain only the *background* category, because it owns the most training instances (typically 10-100 times more than object categories). We adopt sigmoid cross entropy loss for \mathcal{G}_0 here because it only contains one prediction, while for the other

groups we use softmax cross entropy loss. The reason for choosing softmax is that the softmax function inherently owns the ability to suppress each class from another, and less likely produce large numbers of false positives. During training, for a proposal b_k with ground-truth label c , two groups will be activated, which are background group \mathcal{G}_0 and foreground group \mathcal{G}_n where $c \in \mathcal{G}_n$.

4.2. Calibration via category “others”

However, we find the above group softmax design suffers from the following issue. During testing, for a proposal, all groups will be used to predict since its category is unknown. Thus, at least one category per group will receive a high prediction score, and it will be hard to decide which group-wise prediction we should take, leading to a large number of false positives. To address this issue, we add a category *others* into every group to calibrate predictions among groups and suppress false positives. This category *others* contains categories not included in the current group, which can be either *background* or foreground categories in other groups. For \mathcal{G}_0 , category *others* also represents foreground classes. To be specific, for a proposal b_k with ground-truth label c , the new prediction z should be $z \in \mathbb{R}^{(C+1)+(N+1)}$. The probability of class j is calculated by

$$p_j = \frac{e^{z_j}}{\sum_{i \in \mathcal{G}_n} e^{z_i}}, \{n \mid j \in \mathcal{G}_n\}. \quad (4)$$

The ground-truth labels should be re-mapped in each group. In groups where c is not included, class *others* will be defined as the ground-truth class. Then the final loss function is

$$\mathcal{L}_k = - \sum_{n=0}^N \sum_{i \in \mathcal{G}_n} y_i^n \log(p_i^n), \quad (5)$$

where y^n and p^n represent the label and probability in \mathcal{G}_n .

4.3. Balancing training samples in groups

In the above treatment, the newly added category *others* will again become a dominating outlier with overwhelming many instances. To balance training sample number per group, we only sample a certain number of *others* proposals for training, which is controlled by a sampling ratio β . For \mathcal{G}_0 , all training samples of *others* will be used since the number of *background* proposals is very large. For $\{\mathcal{G}_n \mid n \in \mathbb{R}, 1 \leq n \leq N\}$, m_n *others* instances will be randomly sampled from all *others* instances, where $m_n = \beta \sum_{i \in \mathcal{G}_n} \mathcal{N}_{batch}(i)$. $\beta \in [0, +\infty)$ is a hyperparameter and we conduct an ablation study in Sec. 5.4 to show the impact of β . Normally, we set $\beta = 8$. $\mathcal{N}_{batch}(i)$ indicates the instances for category i in current batch.

Namely, within the groups that contain the ground-truth categories, *others* instances will be sampled proportionally based on a mini-batch of K proposals. If there is no normal

categories activated in one group, all the *others* instances will not be activated. This group is ignored. In this way, each group can keep balanced with a low ratio of false positives. Adding category *others* brings 2.7% improvement over the baseline.

4.4. Inference

During inference, we first generate z with the trained model, and apply softmax in each group using Eqn. (4). Except for \mathcal{G}_0 , all nodes of *others* are ignored, and probabilities of all categories are ordered by the original category IDs. p_0^0 in \mathcal{G}_0 can be regarded as the probability of foreground proposals. Finally, we rescale all probabilities of normal categories with $\tilde{p}_j = p_0^0 \times p_j$. This new probability vector is fed to the following post-processing steps like NMS to produce final detection results. It should be noticed that the \tilde{p} is not a real probability vector technically since the summation of it does not equal to 1. It plays the role of the original probability vector which guides the model through selecting final boxes.

5. Experiments

5.1. Dataset and setup

We conduct experiments on the recent Large Vocabulary Instance Segmentation (LVIS) dataset [14], which contains 1,230 categories with both bounding box and instance mask annotations. For object detection experiments, we use only bounding box annotation for training and evaluation. When exploring BAGS’s generalization to instance segmentation, we use mask annotations. Please refer to our supplementary materials for implementation details.

Following [36], we split the categories in the validation set of LVIS into 4 bins according to their training instance numbers to evaluate the model performance on the head and tail classes more clearly. Bin _{i} contains categories that have 10^{i-1} to 10^i instances. We refer categories in the first two bins as “tail classes”, and categories in the other two bins as “head classes”. Besides the official metrics mAP, AP_r (AP for rare classes), AP_c (AP for common classes), and AP_f (AP for frequent classes) that are provided by LVIS-api², we also report AP on different bins. AP _{i} denotes the mean AP over the categories from Bin _{i} .

5.2. Main results on LVIS

We transfer multiple state-of-the-art methods for long-tail classification to the Faster R-CNN framework, including fine-tuning on tail classes, repeat factor sampling (RFS) [26], category loss re-weighting, Focal Loss [23], NCM [20, 34], and τ -normalization [20]. We carefully adjust the hyperparameter settings to make them suitable for

²<https://github.com/lvis-dataset/lvis-api>

ID	Models	mAP	AP ₁	AP ₂	AP ₃	AP ₄	AP _r	AP _c	AP _f	ACC	ACC ₁	ACC ₂	ACC ₃	ACC ₄	ACC _{bg}
(1)	R50-FPN	20.98	0.00	17.34	24.00	29.99	4.13	19.70	29.30	92.78	0.00	2.47	25.30	45.87	95.91
(2)	x2	21.93	0.64	20.94	23.54	28.92	5.79	22.02	28.26	92.62	0.00	5.60	26.51	45.71	95.69
(3)	Finetune tail	22.28	0.27	22.58	23.89	27.43	5.67	23.54	27.34	94.81	0.00	5.04	5.58	5.86	99.85
(4)	RFS [14]	23.41	7.80	24.18	23.14	28.33	14.59	22.74	27.77	92.71	0.60	7.50	25.62	44.39	95.84
(5)	RFS-finetune	22.66	8.06	23.07	22.43	27.73	13.44	22.06	27.09	92.77	0.60	7.14	25.08	43.79	95.91
(6)	Re-weight	23.48	6.34	22.91	23.88	30.12	11.47	22.41	29.61	94.84	0.00	0.82	9.57	17.40	99.53
(7)	Re-weight-cls	24.66	10.04	24.12	24.57	31.07	14.16	23.51	30.28	94.76	0.00	0.34	7.72	16.02	99.64
(8)	Focal loss [23]	11.12	0.00	10.24	13.36	13.17	2.74	11.13	14.46	3.87	0.00	17.45	40.11	49.31	1.35
(9)	Focal loss-cls	19.29	1.64	19.30	20.64	23.70	6.60	19.81	23.71	2.90	0.00	27.67	48.53	48.89	0.16
(10)	NCM-fc [20]	16.02	5.87	14.13	16.97	21.40	10.31	13.92	20.92	94.29	0.00	0.02	0.23	0.15	100.00
(11)	NCM-conv [20]	12.56	4.20	9.71	13.75	18.46	6.11	10.39	17.85	94.29	0.00	0.00	0.20	0.10	100.00
(12)	τ -norm [20]	11.01	0.00	11.71	12.01	12.36	2.07	12.30	12.97	5.91	0.00	30.32	39.49	49.14	3.42
(13)	τ -norm-select	21.61	0.35	20.07	23.43	29.16	6.18	20.99	28.54	92.43	0.00	13.19	20.62	38.98	95.91
(14)	Ours	25.96	11.33	27.64	25.14	29.90	17.65	25.75	29.54	93.71	2.06	7.50	22.07	35.88	97.41

Table 1. Comparison with state-of-the-art methods transferred from long-tail image classification on LVIS *val* set. **Bold** numbers denote the best results among all models. Model (1) and (4) are initialized with model pre-trained on COCO dataset. All the others are initialized with model (1). “-cls” denotes only train the classification FC layer W and b , and the other parameters are frozen. Model (10) and (11) represent NCM model using classification FC features and ROI-pooled Conv feature to calculate category centers respectively. Model (13) means using τ -norm results only on foreground proposals. Refer to our supplementary materials for more implementation details.

object detection. Implementation details are provided in our supplementary material. We report their detection performance and proposal classification accuracy in Table 1.

How well does naive baseline perform? We take Faster R-CNN with ResNet-50-FPN backbone as the baseline (model (1) in the table) that achieves 20.98% mAP but 0 AP₁. The baseline model misses most tail categories due to domination of other classes. Consider other models are initialized by model (1) and further fine-tuned by another 12 epochs. To make sure the improvement is not from longer training schedule, we train model (1) with another 12 epochs for fair comparison. This gives model (2). Comparing model (2) with model (1), we find longer training mainly improves on AP₂, but AP₁ remains around 0. Namely, longer training hardly helps improve the performance for low-shot categories with instances less than 10. Fine-tuning model (1) on tail-class training samples (model (3)) only increases AP₂ notably while decreases AP₄ by 2.5%, and AP₁ remains 0. This indicates the original softmax classifier cannot perform well when the number of training instances is too small.

Do long-tail classification methods help? We observe sampling-based method RFS (model (4)) improves the overall mAP by 2.5%. The AP for tail classes is improved while maintaining AP for head classes. However, RFS increases the training time cost by 1.7 \times . We also try to initialize the model with model (1), obtaining model (5). But mAP drops by 0.8% due to over-fitting.

For cost sensitive learning methods, both model (6) and (7) improve the performance, while model (7) works bet-

ter. This confirms the observation in [20] that decoupling feature learning and classifier benefits long-tail recognition still applies for object detection. For focal loss, we directly apply sigmoid focal loss at proposal level. It is worth noting that in terms of proposal classification, the accuracy over all the object classes (ACC_{1,2,3,4}) increases notably. However, for *background* proposals, ACC_{bg} drops from 95.8% to 0.16%, leading to a large number of false positives and low AP. This phenomenon again highlights the difference between long-tail detection and classification—the very special *background* class should be carefully treated.

For NCM, we try to use both FC feature just before classier (model (10)), and Conv feature extracted by ROI-align (model (11)). However, our observation is NCM works well for extremely low-shot classes, but is not good for head classes. Moreover, NCM can provide a good 1-nearest-neighbour classification label. But for detection, we also need the whole probability vector to be meaningful so that scores of different proposals on the same categories can be used to evaluate the quality of proposals.

The τ -normalization model (12) suffers similar challenge as Focal Loss model (8). The many-shot *background* class is extremely dominating. Though foreground proposal accuracy is greatly increased, ACC_{bg} drops hugely. Consequently, for model (13), the proposals categorized to *background* inherit prediction of the original model while the others take τ -norm results. However, the improvement is limited. We should notice that AP₁ and ACC₁ are still 0 after τ -norm, but AP₂ and ACC₂ are improved.

How well does our method perform? For our model, ex-

Models	mAP	AP_1	AP_2	AP_3	AP_4	AP_r	AP_c	AP_f
Faster R50	20.98	0.00	17.34	24.00	29.99	4.13	19.70	29.30
Ours	25.96	11.33	27.64	25.14	29.90	17.65	25.75	29.54
Faster X101	24.63	0.79	22.37	27.45	32.73	5.80	24.54	32.25
Ours	27.83	14.99	28.07	27.93	32.02	18.78	27.32	32.07
Cascade X101	27.16	0.00	24.06	31.09	36.17	4.84	27.22	36.00
Ours	32.77	19.03	36.10	31.13	34.96	28.24	32.11	35.41

Table 2. Results with stronger backbone ResNeXt-101-64x4d and stronger framework Cascade R-CNN. All *Ours* models are initialized with their plain counterparts.

cept for \mathcal{G}_0 , we split the normal categories into 4 groups for group softmax computation, with s^l and s^h to be (0, 10), (10, 10^2), (10^2 , 10^3), (10^3 , $+\infty$) respectively, and $\beta = 8$. Our model is initialized with model (1), and the classification FC layer is randomly initialized since the output shape changed. Only this FC layer is trained for another 12 epochs, and all other parameters are frozen. Our results surpass all the other methods by a large margin. AP_1 increases 11.3%, AP_2 increases 10.3%, with AP_3 and AP_4 almost unchanged. This result verifies the effectiveness of our designed balanced group softmax module.

Extension of our method to stronger models. To further verify the generalization of our method, we change Faster R-CNN backbone to ResNeXt-101-64x4d [37]. Results are shown in Table 2. On this much stronger backbone, our method still gains 3.2% improvement. Then, we apply our method to state-of-the-art Cascade R-CNN [1] framework with changing all 3 softmax classifiers in 3 stages to our BAGS module. Overall mAP is increased significantly by 5.6%. Our method brings persistent gain with 3 heads.

5.3. Results for instance segmentation

We further evaluate our method benefits for instance segmentation models including Mask R-CNN [16] and state-of-the-art HTC [4] on LVIS. Here HTC models are trained with COCO stuff annotations for a segmentation branch. Results are shown in Table 3. First, comparing our models (8)(10)(12) with their corresponding baseline models (7)(9)(11), mAPs of both bounding box and mask increase largely. Our models fit the tail classes much better while APs on head classes drop slightly. Second, we compare our results with state-of-the-art results [36, 14] on LVIS instance segmentation task. With both Mask R-CNN framework and ResNet-50-FPN backbone, our model (8) surpass RFS (1) and Calib (4) by at least 1.8%. With both HTC framework and ResNeXt-101-FPN backbone, our model (10) is 1.4% better than Calib (5). With ResNeXt-101-FPN-DCN backbone and multiscale training, our model (12) is 2.3% better than Calib (6). Our method establishes new state-of-the-arts in terms of both bounding box and mask criterion.

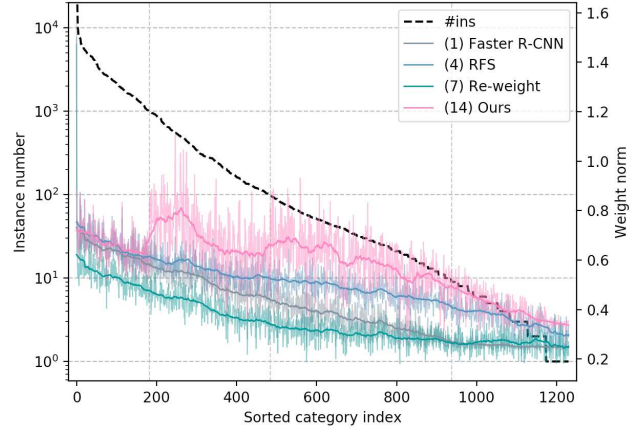


Figure 3. Comparison of weight norms from model (1)(4)(7)(14) in Table 1. The vertical dashed lines split all categories into $\text{Bin}_{1,2,3,4}$.

5.4. Model analysis

Does our method balance classifiers well? We visualize the classifier weight norm W of model (1)(4)(7) and our model (14) of Table 1 in Figure 3. Weights of RFS on tail classes are obviously enlarged. Re-weighting method suppresses the weights of head classes and lifts weights of the tail classes. For ours, since we decouple the relationships of different group of categories, weights of $\mathcal{G}_1, \mathcal{G}_2$ and \mathcal{G}_3 are almost at the same level. Though weights of \mathcal{G}_4 are still smaller, they have been better balanced than the original model. Noting that the weights norm of our model are less related to the training instance numbers in each group, implying such decoupling benefits network training.

How much background and others contribute? See Table 4. With baseline model (0), directly grouping normal categories to 4 sets without adding *background* \mathcal{G}_0 and *others* in each group, we get results of (1). For model (1), during inference, scores of each group are fed to softmax respectively, and concatenated directly for NMS. Though AP_1 improves 5.7%, performance on all the other Bins drops significantly. This is because we do not have any constraints for FPs. For a single proposal, at least one category will be activated in each group, leading to many FPs. When we add \mathcal{G}_0 (model (2)), and use p_0^0 to rescale scores of normal categories, we get 1.9% improvement over model (1), but still worse than model (0). For model (3), we add category *others* into each groups, and not using \mathcal{G}_0 , we obtain 2.7% performance gain.

How many groups to use in BAGS? With rescaling with \mathcal{G}_0 , another 2.2% improvement is obtained (model (5)). If we reduce the group number from 4 to 2, as shown in model (4), the overall mAP drops 0.6. However, specifically, it should be noticed that AP_1 becomes much worse, while AP_4 increases a little. Using more groups does not help as well (model(6)). Since #ins for Bin_1 is too small for $N = 4$, dividing Bin_1 into 2 bins further decreases #ins of per group, leading to highly insufficient training for tails.

ID	Models	Backbone	mAP	AP_1	AP_2	AP_3	AP_4	AP_r	AP_c	AP_f	mAP^m	AP_1^m	AP_2^m	AP_3^m	AP_4^m	AP_r^m	AP_c^m	AP_f^m
(1)	Mask-RFS* [14]	R50	—	—	—	—	—	—	—	—	24.40	—	—	—	—	14.50	24.30	28.40
(2)	Mask-RFS* [14]	R101	—	—	—	—	—	—	—	—	26.00	—	—	—	—	15.80	26.10	29.80
(3)	Mask-RFS* [14]	X101-32x8d	—	—	—	—	—	—	—	—	27.10	—	—	—	—	15.60	27.50	31.40
(4)	Mask-Calib* [36]	R50	—	—	—	—	—	—	—	—	21.10	8.60	22.00	19.60	26.60	—	—	—
(5)	HTC-Calib* [36]	X101	—	—	—	—	—	—	—	—	29.85	16.05	30.60	29.80	33.50	—	—	—
(6)	HTC-Calib* [36]	X101-MS-DCN	—	—	—	—	—	—	—	—	32.10	12.70	32.10	33.60	37.00	—	—	—
(7)	Mask R-CNN	R50	20.78	0.00	15.88	24.61	30.51	3.28	18.99	30.00	20.68	0.00	17.06	23.66	29.62	3.73	19.95	28.37
(8)	Ours	R50	25.76	9.65	26.20	26.09	30.45	15.03	25.45	30.42	26.25	12.81	28.28	25.15	29.61	17.97	26.91	28.74
(9)	HTC	X101	31.28	5.02	31.71	33.24	37.21	12.39	32.58	37.18	29.28	5.11	30.34	30.62	34.37	12.11	31.32	33.58
(10)	Ours	X101	33.68	19.95	36.14	32.82	36.06	25.43	34.12	36.42	31.20	17.33	33.87	30.34	33.29	23.40	32.34	32.89
(11)	HTC	X101-MS-DCN	34.61	5.80	35.36	36.87	40.50	14.24	35.98	41.03	31.94	5.56	33.07	33.75	37.02	13.67	34.04	36.62
(12)	Ours	X101-MS-DCN	37.71	24.40	40.30	36.67	40.00	29.43	37.78	40.92	34.39	21.07	36.69	33.71	36.61	26.79	35.04	36.61

Table 3. Results of bounding box and mask AP when extending our method to instance segmentation on LVIS *val* set. AP^m denotes AP of instance segmentation mask. All backbones are with FPN. X101 denotes X101-64x4d. * results are from the corresponding cited paper. **Bold** numbers indicate the best results among all models. Model (8)(10)(12) are initialized with model (7)(9)(11).

ID	b	o	N	mAP	AP_1	AP_2	AP_3	AP_4	AP_r	AP_c	AP_f
(0)				20.98	0.00	17.34	24.00	29.99	4.13	19.70	29.30
(1)		4		17.82	5.71	17.07	18.09	23.13	8.52	17.44	22.01
(2)	✓	4		19.73	7.18	19.66	18.80	25.95	9.89	19.32	24.19
(3)	✓	4		23.74	9.90	24.06	23.38	28.88	15.46	22.58	28.49
(4)	✓	✓	2	25.31	6.53	27.55	24.19	30.35	15.30	25.14	29.53
(5)	✓	✓	4	25.96	11.33	27.64	25.14	29.90	17.65	25.75	29.54
(6)	✓	✓	8	24.85	7.79	26.05	24.59	29.58	14.11	24.79	29.21

Table 4. Effect of adding different components to our module. *b* for *background* group G_0 . *o* for adding category *others* to all bins. *N* is number of normal category groups.

To sum up, adding category *others* into each group matters a lot, and using specially trained p_0^0 to suppress *background* proposals works better than just *others*. Finally, grouping categories into bins and decoupling the relationship between tail and head classes benefits a lot for learning of tail classes.

Impact of β in BAGS. After adding category *others* to all groups, we need to sample training instances for *others*. Using all *others* proposals will lead to imbalance problem in each group. Thus, our strategy is to sample *others* with ratio β , so that $\#ins\ others:\#ins\ normal = \beta$. As shown in Fig. 4, mAP continuously increases when we increase β until $\beta = 8$. If we use all *others* proposal in activated group (indicated as *n* in x-axis), the performance for head classes keep increasing, but that for tail classes drops a lot. If we train all *others* proposals no matter whether there are normal categories being activated (indicated as *all* in x-axis), mAP gets worse. This verifies our opinion that another imbalance problem could worsen the results.

5.5. Results on COCO-LT

To further verify the generalization ability of our method, we construct a long-tail distribution COCO-LT dataset by

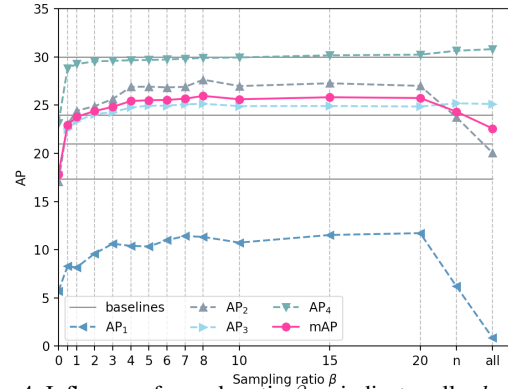


Figure 4. Influence of sample ratio β . *n* indicates all *others* in the activated groups, and *all* indicates all *others* in all the groups.

sampling images and annotations from COCO [24]. We get similar results on COCO-LT as on LVIS. Our model still introduces over 2% improvement on mAP (+2.2% for Faster R-CNN, +2.4% for Mask R-CNN bounding box, +2.3% for Mask R-CNN mask), especially gaining large improvement on tail classes (from 0.1% to 13.0% for bounding box) with both Faster R-CNN and Mask R-CNN frameworks. Please refer to our supplementary materials for dataset construction, data details, and full results.

6. Conclusion

In this work, we first reveal a reason for poor detection performance on long-tail data is that the classifier becomes imbalanced due to insufficiently training on low-shot classes, by analyzing their classifier weight norms. Then, we investigate multiple solid baseline methods transferred from long-tail classification, but we found they are limited in addressing challenges for the detection task. We thus propose a balanced group softmax module to undertake the imbalance problem of classifiers, which achieves notably better results on different strong backbones for long-tail detection as well as instance segmentation.

References

- [1] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: Delving into high quality object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 6154–6162, 2018. 1, 2, 3, 7
- [2] Kaidi Cao, Colin Wei, Adrien Gaidon, Nikos Arachiga, and Tengyu Ma. Learning imbalanced datasets with label-distribution-aware margin loss. *arXiv preprint arXiv:1906.07413*, 2019. 2, 3, 4
- [3] Hao Chen, Yali Wang, Guoyou Wang, and Yu Qiao. Lstd: A low-shot transfer detector for object detection. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018. 2
- [4] Kai Chen, Jiangmiao Pang, Jiaqi Wang, Yu Xiong, Xiao-xiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jianping Shi, Wanli Ouyang, et al. Hybrid task cascade for instance segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4974–4983, 2019. 2, 3, 7
- [5] Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge Belongie. Class-balanced loss based on effective number of samples. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9268–9277, 2019. 2, 3, 4
- [6] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 2
- [7] Chris Drummond, Robert C Holte, et al. C4. 5, class imbalance, and cost sensitivity: why under-sampling beats over-sampling. In *Workshop on learning from imbalanced datasets II*, volume 11, pages 1–8. Citeseer, 2003. 2, 3
- [8] Kaiwen Duan, Song Bai, Lingxi Xie, Honggang Qi, Qingming Huang, and Qi Tian. Centernet: Keypoint triplets for object detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 6569–6578, 2019. 2
- [9] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88(2):303–338, 2010. 1
- [10] Qi Fan, Wei Zhuo, and Yu-Wing Tai. Few-shot object detection with attention-rpn and multi-relation detector. *arXiv preprint arXiv:1908.01998*, 2019. 2
- [11] Ross Girshick. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 1440–1448, 2015. 2, 3
- [12] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Region-based convolutional networks for accurate object detection and segmentation. *IEEE transactions on pattern analysis and machine intelligence*, 38(1):142–158, 2015. 2
- [13] Samantha Guerriero, Barbara Caputo, and Thomas Mensink. Deepncm: Deep nearest class mean classifiers. 2018. 3
- [14] Agrim Gupta, Piotr Dollar, and Ross Girshick. Lvis: A dataset for large vocabulary instance segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5356–5364, 2019. 1, 2, 3, 4, 5, 6, 7, 8
- [15] Hui Han, Wen-Yuan Wang, and Bing-Huan Mao. Borderline-smote: a new over-sampling method in imbalanced data sets learning. In *International conference on intelligent computing*, pages 878–887. Springer, 2005. 2, 3
- [16] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pages 2961–2969, 2017. 2, 3, 7
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 2
- [18] Chen Huang, Yining Li, Change Loy Chen, and Xiaoou Tang. Deep imbalanced learning for face recognition and attribute prediction. *IEEE transactions on pattern analysis and machine intelligence*, 2019. 2, 3, 4
- [19] Bingyi Kang, Zhuang Liu, Xin Wang, Fisher Yu, Jiashi Feng, and Trevor Darrell. Few-shot object detection via feature reweighting. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 8420–8429, 2019. 2
- [20] 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, 3, 5, 6
- [21] Hei Law and Jia Deng. Cornernet: Detecting objects as paired keypoints. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 734–750, 2018. 1, 2
- [22] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2117–2125, 2017. 2
- [23] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017. 1, 2, 3, 4, 5, 6
- [24] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *European conference on computer vision*, pages 740–755. Springer, 2014. 1, 8
- [25] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *European conference on computer vision*, pages 21–37. Springer, 2016. 1
- [26] Dhruv Mahajan, Ross Girshick, Vignesh Ramanathan, Kaiming He, Manohar Paluri, Yixuan Li, Ashwin Bharambe, and Laurens van der Maaten. Exploring the limits of weakly supervised pretraining. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 181–196, 2018. 2, 3, 5
- [27] Thomas Mensink, Jakob Verbeek, Florent Perronnin, and Gabriela Csurka. Distance-based image classification: Generalizing to new classes at near-zero cost. *IEEE transactions on pattern analysis and machine intelligence*, 35(11):2624–2637, 2013. 3

- [28] Wanli Ouyang, Xiaogang Wang, Cong Zhang, and Xiaokang Yang. Factors in finetuning deep model for object detection with long-tail distribution. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 864–873, 2016. [2](#)
- [29] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018. [1](#), [2](#)
- [30] William J Reed. The pareto, zipf and other power laws. *Economics letters*, 74(1):15–19, 2001. [1](#)
- [31] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *Advances in neural information processing systems*, pages 91–99, 2015. [1](#), [2](#), [3](#)
- [32] Li Shen, Zhouchen Lin, and Qingming Huang. Relay back-propagation for effective learning of deep convolutional neural networks. In *European conference on computer vision*, pages 467–482. Springer, 2016. [2](#), [3](#)
- [33] 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. *arXiv preprint arXiv:1902.07379*, 2019. [2](#), [3](#), [4](#)
- [34] Jake Snell, Kevin Swersky, and Richard Zemel. Prototypical networks for few-shot learning. In *Advances in Neural Information Processing Systems*, pages 4077–4087, 2017. [5](#)
- [35] Zhi Tian, Chunhua Shen, Hao Chen, and Tong He. Fcos: Fully convolutional one-stage object detection. *arXiv preprint arXiv:1904.01355*, 2019. [3](#)
- [36] Tao Wang, Yu Li, Bingyi Kang, Junnan Li, Jun Hao Liew, Sheng Tang, Steven Hoi, and Jiashi Feng. Classification calibration for long-tail instance segmentation. *arXiv preprint arXiv:1910.13081*, 2019. [2](#), [4](#), [5](#), [7](#), [8](#)
- [37] Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. Aggregated residual transformations for deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1492–1500, 2017. [2](#), [7](#)
- [38] Chi Zhang, Guosheng Lin, Fayao Liu, Rui Yao, and Chunhua Shen. Canet: Class-agnostic segmentation networks with iterative refinement and attentive few-shot learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5217–5226, 2019. [2](#)
- [39] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points. *arXiv preprint arXiv:1904.07850*, 2019. [3](#)
- [40] Xingyi Zhou, Jiacheng Zhuo, and Philipp Krahenbuhl. Bottom-up object detection by grouping extreme and center points. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 850–859, 2019. [2](#)