# Forest R-CNN: Large-Vocabulary Long-Tailed Object Detection and Instance Segmentation

Jialian Wu[1], Liangchen Song[1], Tiancai Wang[2], Qian Zhang[3] and Junsong Yuan[1]

[1]State University of New York at Buffalo

[2]Tianjin University        [3]Horizon Robotics, Inc.

{jialianw,jsyuan}@buffalo.edu

## ABSTRACT

Despite the previous success of object analysis, detecting and segmenting a large number of object categories with a long-tailed data distribution remains a challenging problem and is less investigated. For a large-vocabulary classifier, the chance of obtaining noisy logits is much higher, which can easily lead to a wrong recognition. In this paper, we exploit prior knowledge of the relations among object categories to cluster fine-grained classes into coarser parent classes, and construct a classification tree that is responsible for parsing an object instance into a fine-grained category via its parent class. In the classification tree, as the number of parent class nodes are significantly less, their logits are less noisy and can be utilized to suppress the wrong/noisy logits existed in the fine-grained class nodes. As the way to construct the parent class is not unique, we further build multiple trees to form a classification forest where each tree contributes its vote to the fine-grained classification. To alleviate the imbalanced learning caused by the long-tail phenomena, we propose a simple yet effective resampling method, NMS Resampling, to re-balance the data distribution. Our method, termed as Forest R-CNN, can serve as a plug-and-play module being applied to most object recognition models for recognizing more than 1000 categories. Extensive experiments are performed on the large vocabulary dataset LVIS. Compared with the Mask R-CNN baseline, the Forest R-CNN significantly boosts the performance with 11.5% and 3.9% AP improvements on the rare categories and overall categories, respectively. Moreover, we achieve state-of-the-art results on the LVIS dataset. *Code is available at* *https://github.com/JialianW/Forest_RCNN*.

## CCS CONCEPTS

• **Computing methodologies** → **Object detection**.

## KEYWORDS

object detection; instance segmentation; large vocabulary; long-tailed data distribution
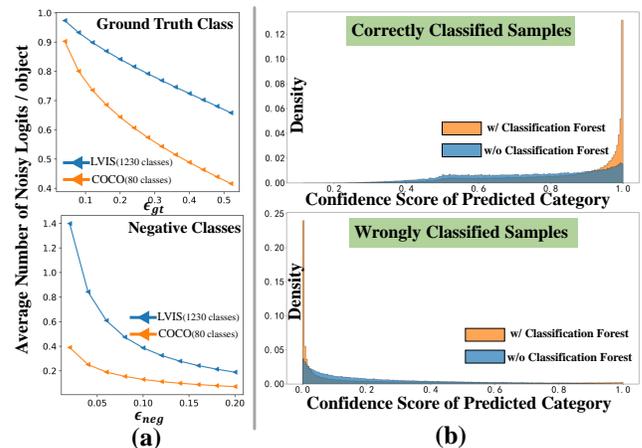
**Figure 1: (a) Statistics of noisy logits of COCO `val` set *vs.* LVIS [15] `val` set tested by Mask R-CNN with ResNet-50-FPN where the noisy logit is defined in Eq. 1. It is seen that the large vocabulary brings a higher chance of noisy logits which can easily lead to a wrong recognition. In the worst case, an object can have $1$ and $N-1$ noisy logits for ground truth class and negative classes, respectively, where $N$ is the total number of classes. (b) Probability density distribution of confidence scores on the LVIS v0.5 `val` set. With our proposed classification forest, the confidence scores of correctly classified samples are improved (top) and those of wrongly classified samples are suppressed (bottom). The top figure is measured by considering the background class.**

## 1 INTRODUCTION

With the renaissance of deep convolutional neural networks (CNNs), recent years have witnessed great progress in object recognition, including object detection [14, 23, 27, 31, 33, 35] and instance segmentation [1, 11, 17, 21, 46]. Object recognition plays a central role in visual learning and has been widely applied to various applications, *e.g.,* person retrieval [45], human-object interaction [40], and visual reasoning [42].

Most great works on object recognition have been thoroughly investigated in the few category regime, *e.g.,* PASCAL VOC [12] (20 classes) and COCO [25] (80 classes). However, practical applications are urging the need for recognizing a large number of categories with a long-tailed data distribution, which is a great challenge for

most existing methods. Generally, the challenge boils down to two aspects: **(i)** As the number of categories grows, the chance of obtaining noisy logits in a classifier becomes higher (*i.e.,* inaccurate classifier predictions on either the ground truth class or other negative classes) as shown in Fig. 1 (a). This increases the difficulty for generating a correct category label or a high confidence score on the ground truth category, therefore easily leading to a wrong recognition. **(ii)** The long-tail phenomena inherently occur in a large vocabulary scenario and cause extreme imbalanced data distribution, where few classes (*a.k.a.* head class) appear very often yet most of other classes (*a.k.a.* tail class) rarely appear. Due to the imbalanced distribution, most tail classes are overwhelmed by head classes in training, making it difficult to well learn effective classifier and feature representations especially for tail classes.

In this paper, we propose a novel classification forest together with a simple yet effective data resampling method, striving to alleviate the above problems **(i)** and **(ii)**, respectively. For **(i)**, we exploit prior knowledge of the relations among categories to cluster thousands of fine-grained classes into tens of parent classes, and construct a classification tree that is responsible for parsing an object into a fine-grained node via its parent node. Since the number of parent classes is significantly less, the parent classifier within the tree obtains fewer noisy logits. We then utilize the parent class probabilities estimated from the parent classifier to suppress the noisy logits produced by the fine-grained classifier. Moreover, in terms of different types of prior knowledge, we construct multiple classification trees to form a classification forest, where each tree will contribute its vote to the fine-grained classification. To illustrate the idea of the classification tree, we show an example in Fig. 2 where the fine-grained classifier wrongly predicts a "toy" as a "sedan" with $p(sedan) = 0.8$ while the parent class probability of "sedan" is $p(vehicle) = 0.05$. In the classification tree, the noisy logit of class "sedan" will be calibrated and suppressed by $p(vehicle)$. By suppressing the noisy logits, our method is able to improve the confidence scores of ground truth class and suppress those of other negative classes as shown in Fig. 1 (b). For **(ii)**, we propose a data resampling method, named as NMS Resampling, to adaptively adjust the Non-maximum Suppression (NMS) threshold for different categories based on their category frequency in training. The NMS Resampling can re-balance the data distribution by preserving more training proposal candidates from tail classes while suppressing those from head classes in the bounding box NMS procedure.

Recent attempts [30, 32, 36–38] also strive to recognize objects under the setting of large vocabulary. For example, [38] studies the influence of loss functions on head classes and tail classes, and propose an equalization loss to reduce the negative influence of loss gradients on tail classes. Note that [38] is complementary to our method and can be employed together for gaining performance. Other works [30, 32, 36] also utilize a hierarchical tree structure for better classification, which however may still easily yield a wrong recognition due to inaccurate parent class nodes generated by a single tree. Our method, instead, exploits different types of prior knowledge to build a forest, in which classification is achieved in the form of plurality vote.

**Contributions:** In this work, we propose a novel classification forest that incorporates relations among fine-grained categories
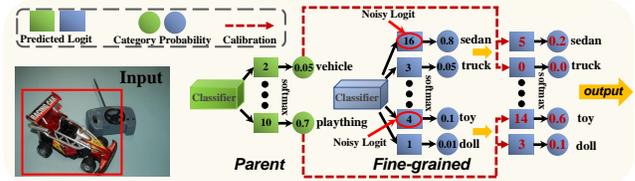


**Figure 2: Brief illustration of the classification tree.**

via different prior knowledge. When dealing with large vocabulary object recognition, our classification forest can better perform via suppressing the noisy logits existed in the fine-grained classifier. In addition, the NMS Resampling method is proposed for re-balancing the data distribution of a long-tailed dataset during training. Our method, termed as Forest R-CNN, is exhaustively evaluated on the large vocabulary object recognition dataset LVIS [15] which contains more than 1000 categories. The Forest R-CNN achieves significant AP gains of 11.5% and 3.9% on the rare categories and overall categories, respectively.

## 2 RELATED WORK

**Object Detection and Instance Segmentation.** Object detection, which plays an important role in various vision applications, has been actively investigated in past decades [5, 7, 8, 24, 33, 35, 43, 44]. With deep learning, two mainstream frameworks, *i.e.,* single-stage detector [6, 24, 27] and two-stage detector [13, 14, 33], have dramatically improved both accuracy and efficiency. Instance segmentation can be viewed as an extension of object detection, where each object instance is bounded by a precise mask instead of a rough bounding box. As a pioneer work, Mask R-CNN [17] achieves instance segmentation by building an extra segmentation head upon the two-stage object detector Faster R-CNN [33]. On the basis of the two-stage fashion, [21, 22, 26] further propose several framework modifications so as to yield better performance. On the other hand, single-stage approaches [1, 4, 11, 46], aiming at a faster inference speed, adopt a more straightforward way to directly generate instance masks from the whole feature maps. Different from the previous works that mostly perform in the few category regime, our method focuses more on object detection and instance segmentation in a large vocabulary with long-tailed data distribution.

**Large-Vocabulary and Long-Tailed Visual Recognition.** To be applied in the natural world, a visual recognition model is expected to deal with a large number of categories with the long-tail phenomena. One prominent method for long-tailed visual recognition is data re-balancing, which is usually divided into two groups: *re-sampling* [2, 16, 19, 34, 49] and *re-weighting* [9, 10, 20]. *Re-sampling* schemes typically oversample the data from minority classes while undersample those from frequent classes in training. *Re-weighting* schemes, instead, assign larger loss weights for the samples from minority classes in training. In addition to data re-balancing, many great works have been made using different ways, *e.g.,* model fine-tuning [9, 30], metric learning [20, 47], meta-learning [28], and knowledge transfer learning [48]. Recently, [30, 32, 36, 38] also strive to solve the problem of object recognition in a large vocabulary. [30, 32, 36] exploit a single hierarchical tree structure to aid the fine-grained classification. Different from the above approaches, we propose a novel classification forest aiming at reducing the noisy

logits existed in the fine-grained classifier to improve classification capability in a large vocabulary. Besides, in contrast to most *re-sampling* methods that resample data on the image level, our proposed NMS Resampling scheme re-balances data on the instance level.

## 3 FOREST R-CNN

### 3.1 Problem Formulation

Given a static image $\mathbf{I} \in \mathbb{R}^{W \times H \times 3}$, an object recognition model is required to perform object detection and instance segmentation simultaneously. Thus, each recognized object instance is associated with four predictions: $\mathbf{b} \in \mathbb{R}^4$, $\mathbf{m} \in [0, 1]^{W \times H}$, $l \in \mathbb{R}^1$, $s \in \mathbb{R}^1$, which are bounding box, segmentation mask, category label, and confidence score, respectively. In this paper, we implement our method based on the baseline model Mask R-CNN [17]. Let us denote by $\mathbf{f}_{\text{img}} = \mathcal{N}_{\text{feat}}(\mathbf{I}) \in \mathbb{R}^{\frac{W}{r} \times \frac{H}{r} \times C}$ the feature maps extracted from input image $\mathbf{I}$, where $\mathcal{N}_{\text{feat}}$, $r$, and $C$ are the backbone network, feature stride, and feature channels, respectively. A region proposal network (RPN) [33] is built upon $\mathbf{f}_{\text{img}}$ to generate a set of proposal candidate boxes $\{\mathbf{b}_p \in \mathbb{R}^4\}$. For each proposal candidate, its corresponding proposal features $\mathbf{f}_{\text{roi}} \in \mathbb{R}^{7 \times 7 \times C}$ are obtained by $\mathbf{f}_{\text{roi}} = \phi(\mathbf{f}_{\text{img}}, \mathbf{b}_p)$, where $\phi$ is the RoI Align operation [17]. On the basis of $\mathbf{f}_{\text{roi}}$, three head networks, $\mathcal{N}_{\text{cls}}$, $\mathcal{N}_{\text{box}}$, and $\mathcal{N}_{\text{mask}}$ are employed to generate $s, l = \mathcal{N}_{\text{cls}}(\mathbf{f}_{\text{roi}})$, $\mathbf{b} = \mathcal{N}_{\text{box}}(\mathbf{f}_{\text{roi}})$, and $\mathbf{m} = \mathcal{N}_{\text{mask}}(\mathbf{f}_{\text{roi}})$, in which $\mathcal{N}_{\text{cls}}$ and $\mathcal{N}_{\text{box}}$ share part of network parameters.

$\mathcal{N}_{\text{cls}}$ is critical for an object recognition model, since classification result is the premise for evaluating boxes and masks. Let $\{x_i\}_{i=1}^N$ be the category set of a given dataset, where $N$ is the number of total classes and $x_i$ is the $i$-th fine-grained class. $p(x_i)$ indicates the classifier inferred probability that a given object belongs to class $x_i$, and it is calculated by the softmax function $p(x_i) = \frac{f_{x_i}}{\sum_{a=1}^N f_{x_a}}$, where $f_{x_i} = e^{z_i}$ and $z_i$ is predicted by the $i$-th neuron in the final layer of $\mathcal{N}_{\text{cls}}$. In this paper, we name $f_{x_i}$ as the logit of class $x_i$ (exponential version) and call $f_{x_i}$ as a noisy logit when:

$$\begin{cases} \frac{f_{x_i}}{\sum_{a=1}^N f_{x_a}} < 1 - \epsilon_{gt}, & \text{if } x_i = x_{gt} \\ \frac{f_{x_i}}{\sum_{a=1}^N f_{x_a}} > \epsilon_{neg}, & \text{if } x_i \neq x_{gt} \end{cases}, \tag{1}$$

where $x_{gt}$ is the ground truth class and $0 < \epsilon_{gt}, \epsilon_{neg} < 1$ are two small constants. In [17], the confidence score of class $x_i$ and the category label $l$ of a given object are obtained by:

$$\begin{aligned} s_i &= p(x_i), \\ l &= \underset{i}{\arg\max}\, s_i, \quad i = 1, 2, ..., N. \end{aligned} \tag{2}$$

Generally speaking, to yield a correct label and a high confidence score on $x_{gt}$, Eq. 2 is satisfactory in the few category regime. However, in a large vocabulary scenario, *e.g.*, $N > 1000$, there exist many more noisy logits within $\{f_{x_i}\}_{i=1}^N$ as shown in Fig. 1 (a), making it difficult to generate a correct label or a high confidence score on $x_{gt}$ using Eq. 2. As shown in Fig. 1 (b), with Eq. 2 the noisy logits result in many samples correctly classified with low confidence scores or misclassified with high confidence scores (blue curve). Moreover, the long-tail phenomena that inherently occur in a large vocabulary make it hard for a model to well learn effective classifier and feature
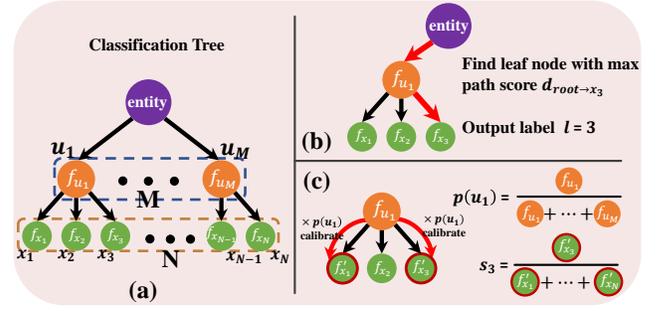


Figure 3: (a) Classification tree where parent and fine-grained class nodes are contained in ⌐⌐⌐⌐ and ⌐ ⌐ ⌐, respectively. 🟣entity and 🟢 indicates the root node and calibrated leaf node, respectively. (b) The category label is estimated by finding the maximum path score from the root node to leaf nodes. (c) The confidence scores of fine-grained classes are produced in terms of the values of calibrated leaf nodes.

representations against tail classes. To alleviate the above problems, we propose a classification forest (§ 3.2 and § 3.3) for enhancing the capability of classifying a large number of categories and an NMS Resampling (§ 3.4) for re-balancing the long-tailed data distribution, respectively.

### 3.2 Classification Tree

**Tree Structure.** In real-world scenes, we are aware of prior knowledge of the relationships among object categories. For example, a "school bus" and a "sedan" have the same parent class "vehicle" when considering lexical relation, while "steering wheel" and "basketball" have the same parent class "circularity" when considering geometrical relation. Accordingly, for each type of prior knowledge, all $N$ fine-grained classes can be clustered into $M$ parent classes, where each parent class is on behalf of a characteristic of its children classes. Note that $M$ may vary when we consider different types of prior knowledge. Based on the hierarchical affiliation obtained by prior knowledge, we can build up a tree structure which is a basic component of the classification forest. As shown in Fig. 3 (a), the proposed classification tree has 3 levels: the first level is a single root node; the second level consists of $M$ parent class nodes $\{u_j\}_{j=1}^M$; the last level comprises $N$ leaf nodes which represent the fine-grained classes. Each node $v$ (excluding the root node) is associated with its corresponding logit $f_v = e^{z_v}$ as the node value. We add additional fully connected layers within $\mathcal{N}_{\text{cls}}$ to generate the parent class logit $f_u$. During training, each level is regarded as an individual classifier and supervised by the softmax function with cross-entropy loss for learning $f_v$. Since $f_v$ is the classifier logit, it essentially reflects the likelihood that an object belongs to class $v$. Namely, the higher $f_v$ is, the more confidently the given object belongs to class $v$ (vice versa). Also, in the proposed tree structure, we define the path score as the product of the values of nodes through which the path passes (excluding the root node). For an instance, the path score from the root node to $x_3$ is $d_{root \to x_3} = f_{x_3} \times f_{u_1}$.

**Preliminary Model.** In Eq. 2, both category label and confidence score are determined by the logits of fine-grained classifier only. Different from Eq. 2, our method takes into account the parent classes
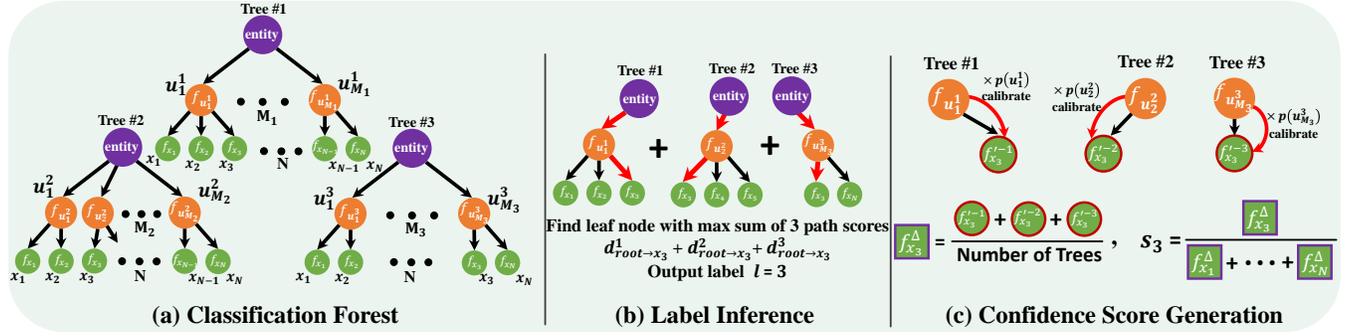
Figure 4: (a) Classification forest with 3 trees. (b) and (c) are brief illustrations of label inference and confidence score generation using (a).

in order to reduce the negative influence of the noisy logits produced by the fine-grained classifier. To this end, a straightforward way is to combine $p(x_i)$ with the corresponding $p(u_j)$ as:

$$s_i = p(x_i) \times p(u_j), \quad x_i \in ch(u_j), \tag{3a}$$

$$l = \operatorname*{argmax}_i s_i, \quad i = 1, 2, ..., N, \tag{3b}$$

where $x_i \in ch(u_j)$ indicates the fine-grained class $x_i$ that is a child of parent class $u_j$. In Eq. 3, to parse an object into a fine-grained category, both the fine-grained classifier and parent classifier need to reach a consensus. As an example, the fine-grained classifier wrongly predicts a "toy" as a "sedan" with $p(x_{sedan}) = 0.8$, while the probability of the parent class of "sedan" is $p(u_{vehicle}) = 0.05$. In this case, the parent classifier has a different "opinion" with the fine-grained classifier, and the final confidence score $s_{toy}$ will be downgraded to 0.04. Therefore, by incorporating the parent class probability, we are able to reduce the confidence scores of those categories misclassified by the fine-grained classifier, enabling a model to be more fault-tolerant with regard to the noisy logits of the fine-grained classifier.

**Confidence Score Generation.** In Eq. 3a, the confidence score is generated by the product of fine-grained class probability and parent class probability, which however fails to improve the confidence score of the ground truth class. For example, given a "sedan", if $p(x_{sedan}) = 0.7$ and $p(u_{vehicle}) = 0.6$, then $s_{sedan}$ becomes 0.42 unexpectedly. To solve this problem, we use the parent class probabilities to directly calibrate the fine-grained class logits instead of scaling the fine-grained class probabilities. Concretely, the logit/node value of $x_i$ is calculated by $f'_{x_i} = f_{x_i} \times p(u_j)$, where $x_i \in ch(u_j)$. Afterwards, similar to the original classifier, the confidence score of the fine-grained class $x_i$ is obtained according to the calibrated fine-grained nodes as:

$$s_i = \frac{p(u_j) \times f_{x_i}}{\sum_{a=1}^{N} p(u_{a^\star}) \times f_{x_a}}, \quad x_i \in ch(u_j), \ x_a \in ch(u_{a^\star}), \tag{4a}$$

$$= \frac{f'_{x_i}}{\sum_{a=1}^{N} f'_{x_a}}. \tag{4b}$$

Compared to the logits of fine-grained classifier $\{f_{x_i}\}_{i=1}^{N}$, $\{f'_{x_i}\}_{i=1}^{N}$ are calibrated by the parent class probabilities which can effectively reduce the noisy logits (as evidenced in Fig. 7). In contrast to Eq. 3a, Eq. 4 suppresses the noisy logits on not only negative classes but

also the ground truth class. That is to say, we not only reduce the confidence scores of negative classes but also further improve those of the ground truth classes (see Fig. 1 (b)). For the above example of $p(x_{sedan}) = 0.7$ and $p(u_{vehicle}) = 0.6$, it indicates that $f_{x_{sedan}}$ account for 70% in $\sum_{a=1}^{N} f_{x_a}$. With Eq. 4, our method will scale $f_{x_{sedan}}$ by $p(u_{vehicle}) = 0.6$ and the other $f_{x \notin u_{vehicle}}$ by a smaller $p(u \neq u_{vehicle}) \leq 0.4$. Accordingly, $f'_{x_{sedan}}$ will be less noisy and make up more than 70% of $\sum_{a=1}^{N} f'_{x_a}$, which leads to a new $s_{sedan}$ higher than 0.7.

**Label Inference.** Based on Eq. 4 and the definition of path score, the category label is inferred by:

$$
\begin{aligned}
l &= \operatorname*{argmax}_i s_i = \operatorname*{argmax}_i f'_{x_i}, \quad i = 1, 2, ..., N, \\
&= \operatorname*{argmax}_i f_{x_i} \times f_{u_j}, \quad x_i \in ch(u_j) \text{ and } i = 1, 2, ..., N, \\
&= \operatorname*{argmax}_i d_{root \to x_i}, \quad i = 1, 2, ..., N.
\end{aligned} \tag{5}
$$

It is seen that the proposed classification tree parses an object into a fine-grained class by finding the maximum path score from the root node to leaf nodes.

## 3.3 Classification Forest

**Forest Structure.** Although a classification tree (§ 3.2) can suppress the noisy logits produced by the fine-grained classifier via incorporating the parent class probabilities, false classification can still easily happen if there exist errors in the parent classifier. To alleviate this problem, we build a classification forest that consists of $T$ different classification trees, and each tree will vote for the final classification decision. As shown in Fig. 4 (a), each tree that is based on one type of prior knowledge of category relations may have a different structure from the others, and the parent class nodes set of the $t$-th tree is denoted by $\{u_j^t\}_{j=1}^{M_t}$. The leaf node set, which represents the fine-grained classes of a given dataset, is the same for all the trees. In this paper, we exploit three types of prior knowledge of the relations among fine-grained classes, $i.e.,$ lexical relation, visual relation, and geometrical relation, in order to take different aspects of object characteristics into consideration. Based on the relations, we then construct three different classification trees, respectively. Note that the proposed classification forest does not constrain the number of classification trees and we observe that
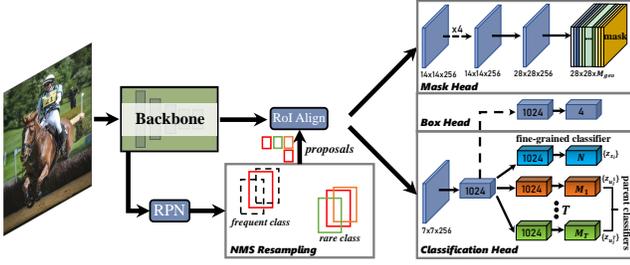
**Figure 5: Network architecture of the Forest R-CNN. The dotted rectangles in the NMS Resampling denote the proposal boxes which are suppressed during the NMS process. The NMS Resampling is only used in the training phase.**

in our experiments three trees are sufficient for achieving improved performance.

**Confidence Score Generation.** In each tree of the classification forest, the value of fine-grained class node $x_i$ is calibrated by its corresponding parent class probability as shown in Fig. 4 (c). Thus, for $x_i$, we have $T$ different calibrated values $\{f_{x_i}'^{-t}\}_{t=1}^{T}$, where $f_{x_i}'^{-t}$ implies the likelihood that the $t$-th classification tree considers a given object to belong to $x_i$. In order to take advantage of each one of $\{f_{x_i}'^{-t}\}_{t=1}^{T}$ and produce a final confidence score $s_i$, we take the average of $f_{x_i}'^{-1}, ..., f_{x_i}'^{-T}$ and denote it as $f_{x_i}^{\triangle}$, which indicates the comprehensive likelihood that all the classification trees consider a given object to belong to $x_i$. Then, similar to Eq. 4, the confidence score of $x_i$ is calculated according to $\{f_{x_i}^{\triangle}\}_{t=1}^{N}$ as:

$$s_i = \frac{f_{x_i}^{\triangle}}{\sum_{a=1}^{N} f_{x_a}^{\triangle}}. \tag{6}$$

Compared to $f_{x_i}'^{-t}$, $f_{x_i}^{\triangle}$ is calibrated by multiple parent classifiers. Since every tree is dedicated to generating confidence score and the likelihood that multiple different parent classifiers yield errors simultaneously is low, the classification forest is robuster than a single classification tree as shown in Tab.3.

**Label Inference.** In § 3.2, the category label is inferred by finding the maximum path score from the root node to leaf nodes. By contrast, to integrate the decisions made by all the trees, we consider a leaf node $x_i$ as the predicted category if the sum of each path score from the root node to $x_i$ of all the trees is maximum. Formally, the category label is inferred as:

$$l = \underset{i}{\arg\max} \sum_{t=1}^{T} d_{root \to x_i}^{t}, \quad i = 1, 2, ..., N, \tag{7}$$

where $d_{root \to x_i}^{t}$ denotes the path score from the root node to $x_i$ in the $t$-th tree. In essence, Eq. 7 is a plurality vote of all the classification trees, which is generally considered to be better than using a single classification tree as in Eq. 5.

## 3.4 NMS Resampling

The long-tail phenomena inherently occur in a large vocabulary dataset and the real visual world, in which few classes appear very often but most other classes rarely appear. Such imbalanced data

distribution introduces great challenges for learning effective classifier and feature representations against tail classes. To re-balance the long-tailed data distribution, we propose a simple yet effective resampling method, termed as NMS Resampling, by adaptively adjusting the NMS threshold during training.

As known, after generating a large amount of proposal boxes from the RPN, the NMS is applied to filter out highly overlapped proposals so as to reduce redundancy. The NMS threshold is class-agnostic and set to a fixed value (*e.g.*, 0.7) for all categories in [17, 33], while our method adaptively adjusts the thresholds for different categories. Specifically, we strive to re-balance the data distribution by utilizing an NMS threshold that is inverse to the data amount of a certain category. That is, we set higher thresholds for tail classes but lower thresholds for head classes. Following this idea, we propose two NMS Resampling schemes to determine the thresholds for specific categories as follows:

**NMS Resampling-Discrete.** The LVIS [15] dataset annotates each category with a category frequency indicating the number of images in which the category appears. Following the descending category frequency, in [15], all 1230 classes are uniformly divided into three groups: *frequent*, *common*, and *rare*. In NMS Resampling-Discrete, we employ three discrete NMS thresholds: $\alpha_f$, $\alpha_c$, and $\alpha_r$ for the *frequent*, *common*, and *rare* classes, respectively, where $\alpha_f < \alpha_c < \alpha_r$. In experiments our method is not sensitive to specific values of $\alpha_f$, $\alpha_c$, and $\alpha_r$ so long as $\alpha_f$, $\alpha_c$, $\alpha_r$ follow ascending order

**NMS Resampling-Linear.** We first uniformly divides three intervals with length $\beta$ for the *frequent*, *common*, and *rare* classes, respectively. Then, in each interval, we linearly assign the NMS threshold to each category as:

$$threshold = \begin{cases} \alpha_r + \beta \times \frac{cf_{max}^{f} - cf_{x_i}}{cf_{max}^{f} - cf_{min}^{f}}, & \text{if } x_i \in frequent \\ \alpha_c + \beta \times \frac{cf_{max}^{c} - cf_{x_i}}{cf_{max}^{c} - cf_{min}^{c}}, & \text{if } x_i \in common \\ \alpha_f + \beta \times \frac{cf_{max}^{r} - cf_{x_i}}{cf_{max}^{r} - cf_{min}^{r}}, & \text{if } x_i \in rare \end{cases} \tag{8}$$

where $\alpha_f < \alpha_c < \alpha_r$ and $cf_{x_i}$ is the category frequency of $x_i$. $cf_{max}^{f}$ and $cf_{min}^{f}$ are the maximum and minimum category frequencies in the *frequent* class group. $\alpha_f$, $\alpha_c$, $\alpha_r$, and $\beta$ are respectively set to 0.65, 0.75, 0.85, and 0.1 by default in our experiments.

Given a foreground proposal box $\mathbf{b}_p$, we first compute its corresponding NMS threshold based on the above schemes. Then, the remaining proposal boxes will be suppressed if they have overlaps with $\mathbf{b}_p$ large then the threshold, otherwise, they will be preserved for the next round of NMS procedure. We use the original NMS threshold of 0.7 for background proposals. The proposed NMS Resampling eases the problem of imbalanced data distribution by preserving more training proposal candidates from the tail classes and suppressing some of those from head classes during training. Compared with the image resampling [15], our method not only is more effective (as shown in Tab.4) but also avoids repeating training images which may cause overfitting and extra training time. In principle, the NMS Resampling is also applicable to single-stage detectors that need the NMS to reduce redundancy during training, and we leave it for future work.

**Table 1: Performance comparison with the baseline Mask R-CNN [17] using different backbone networks on the LVIS v0.5 val set. AP denotes the mask AP and $AP^b$ denotes the box AP. The subscripts "r", "c", and "f" denote performance on the rare, common, and frequent classes, respectively.**

| Backbone | Method | AP | $AP_{50}$ | $AP_{75}$ | $AP_r$ | $AP_c$ | $AP_f$ | $AP^b$ | $AP^b_{50}$ | $AP^b_{75}$ | $AP^b_r$ | $AP^b_c$ | $AP^b_f$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ResNet-50-FPN | Mask R-CNN | 21.7 | 34.7 | 22.8 | 6.8 | 22.6 | 26.4 | 21.8 | 37.1 | 22.5 | 6.5 | 21.6 | 28.0 |
| | Forest R-CNN(Ours) | **25.6** | **40.3** | **27.1** | **18.3** | **26.4** | **27.6** | **25.9** | **42.7** | **27.2** | **16.9** | **26.1** | **29.2** |
| ResNet-101-FPN | Mask R-CNN | 23.6 | 37.1 | 25.0 | 10.0 | 24.8 | 27.6 | 23.5 | 39.9 | 24.4 | 8.7 | 23.1 | 29.8 |
| | Forest R-CNN(Ours) | **26.9** | **42.2** | **28.4** | **20.1** | **27.9** | **28.3** | **27.5** | **44.9** | **29.0** | **20.0** | **27.5** | **30.4** |
| ResNeXt-101-32×4d-FPN | Mask R-CNN | 24.8 | 38.5 | 26.2 | 10.0 | 26.4 | 28.6 | 24.8 | 41.5 | 25.8 | 8.6 | 25.0 | 30.9 |
| | Forest R-CNN(Ours) | **28.5** | **43.8** | **30.9** | **21.6** | **29.7** | **29.7** | **28.8** | **46.3** | **30.9** | **20.6** | **29.2** | **31.7** |

## 3.5 Network Architecture and Loss Function.

**Network Architecture.** The overall network architecture of the proposed Forest R-CNN is shown in Fig. 5. In the classification head $\mathcal{N}_{cls}$, we add $T$ extra fully connected layer (FC) branches for predicting the node values of parent classes compared with [17]. Since the extra FC branches are inserted after the first FC of $\mathcal{N}_{cls}$ whose channel dimensions are 1024, our method only introduces little additional computational overhead. Moreover, we incorporate the prior knowledge of geometrical relation into the mask head $\mathcal{N}_{mask}$. In contrast to the original class-specific mask head, we reduce the number of output channels from $N$ to $M_{geo}$, where $M_{geo}$ is the number of parent classes in the geometrical tree. In inference, the mask of class $x_i$ is fetched from the $j$-th channel, where $x_i \in ch(u_j^{geo})$ and $u_j^{geo}$ is the $j$-th parent class in the geometrical tree.

**Loss Function.** The overall loss function of the Forest R-CNN is defined as:

$$L = L_{MR} + L^1_{cls-p} + ... + L^T_{cls-p}, \qquad (9)$$

where $L_{MR}$ is the original loss of Mask R-CNN and $L^t_{cls-p}$ is the parent classification loss of the $t$-th classification tree.

## 4 EXPERIMENTS

## 4.1 Experiment Setup

**LVIS Dataset.** LVIS [15] is a large vocabulary dataset for object detection and instance segmentation. There are in total 1230 and 1203 categories in LVIS v0.5 and v1.0 datasets, which follow a long-tailed data distribution. All categories are divided into three groups based on the number of images that contains those categories: *rare* (1-10 images), *common* (11-100 images), and *frequent* (>100 images). Our method is trained on the train set and evaluated on the val set. We adopt the evaluation metric AP across IoU threshold from 0.5 to 0.95 for both object detection and instance segmentation results. Our major experiments and ablation studies are performed on LVIS v0.5 dataset. We update the results of Forest R-CNN on LVIS v1.0 dataset in this arxiv V2 version, which is shown in Tab.7.

**Implementation Details.** The Forest R-CNN uses the same basic settings as in [15], *e.g.,* image size, score threshold, and the number of object instances per image. Our method is trained with 24 epochs using the SGD optimizer, and the initial learning rate is set to 0.02 and decreased by a factor of 10 at the 16-th epoch and 22-th epoch, respectively. We use three relations among fine-grained categories, *i.e.,* lexical relation, visual relation, and geometrical relation, respectively, to construct three classification trees before training. For

**Table 2: Ablation study of the proposed NMS Resampling and classification forest. "NR" denotes the NMS Resampling.**

| NR | Classification Forest | AP | $AP^b$ | $AP_r$ | $AP_c$ | $AP_f$ |
|---|---|---|---|---|---|---|
| | | 21.7 | 21.8 | 6.8 | 22.6 | 26.4 |
| ✓ | | 23.5 | 23.5 | 15.6 | 24.1 | 25.9 |
| | ✓ | 23.6 | 24.0 | 10.9 | 24.4 | 27.5 |
| ✓ | ✓ | **25.6** | **25.9** | **18.3** | **26.4** | **27.6** |

the visual and geometrical trees, we employ K-means to cluster the visual features and ground truth binary masks of 1230 fine-grained categories into $M_{vis}$ and $M_{geo}$ parent classes, respectively, where the visual features are obtained by the Mask-RCNN baseline. $M_{vis}$ and $M_{geo}$ are set to 25 and 50, respectively. For the lexical tree, it is constructed according to a subgraph of WordNet [29], which results in $M_{lex}$ = 108 parent classes. In the following experiments, our method is equipped with the NMS Resampling-Discrete and ResNet-50-FPN by default. For the NMS Resampling-Discrete, we empirically set $\alpha_f$, $\alpha_c$, $\alpha_r$ to 0.7, 0.8 and 0.9 by default.

## 4.2 Ablation Studies

**Comparison with Baselines.** As shown in Tab.1, we compare the proposed Forest R-CNN with the baseline Mask R-CNN under different backbone networks. The Forest R-CNN consistently improves AP and $AP^b$ over the baseline with significant gains of 10.1%-12% and 3.3%-4.1% on rare categories and overall categories, respectively. We also separately assess the proposed classification forest (§ 3.3) and NMS Resampling (§ 3.4) in Tab.2. We see from the results that the NMS Resampling improves $AP_r$ from 6.8% to 15.6% with the ResNet-50 [18], demonstrating the strong effectiveness of our re-balancing schemes for tail classes. With the classification forest, our method consistently improves performance on overall categories, which shows that the classification forest is specialized to recognize a large number of categories. Visualized samples can be found in Fig. 8.

**Effectiveness of the Classification Forest.** To study the effectiveness of the classification forest (§ 3.3), we evaluate the average number of noisy logits per object on the LVIS v0.5 val set. As shown in Fig. 7, $\{f^{\triangle}_{x_i}\}^N_{t=1}$ contains fewer noisy logitis than $\{f_{x_i}\}^N_{t=1}$, which demonstrates that the classification forest can effectively suppress the noisy logits produced by the fine-grained classifier. We also investigate the probability density distribution of confidence scores as presented in Fig. 1 (b). We clearly see from the figure that with classification forest our method can effectively improve the
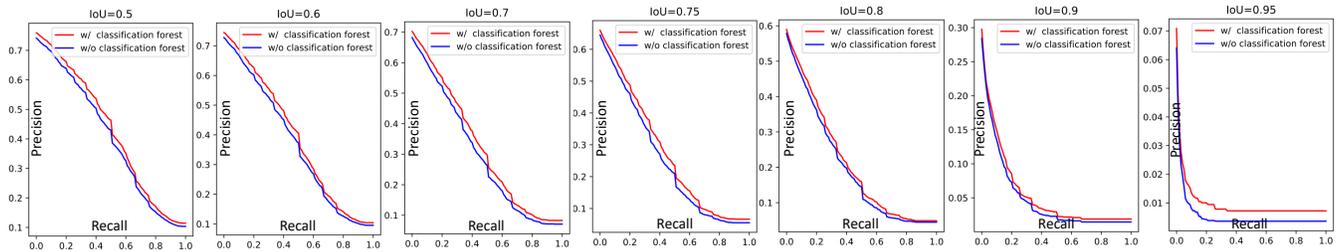
Figure 6: The mask precision-recall (PR) curves w/ and w/o using the classification forest under different IoU thresholds.
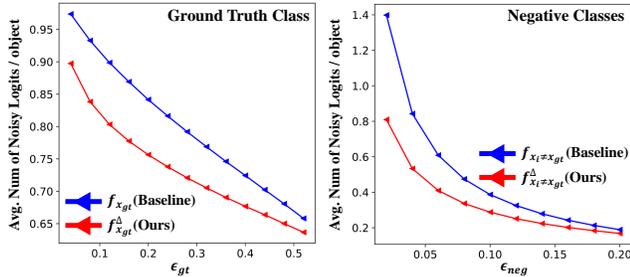


Figure 7: Statistics of the noisy logits w/ and w/o using the classification forest on the LVIS v0.5 `val` set. The results suggest that the proposed classification forest can effectively suppress the noisy logits from the fine-grained classifier. The NMS resampling is enabled for both baseline and ours.

Table 3: Results of the Forest R-CNN using different classification trees.

| Geometrical | Visual | Lexical | AP | $AP^b$ | $AP_r$ | $AP_c$ | $AP_f$ |
|---|---|---|---|---|---|---|---|
| | | | 23.5 | 23.5 | 15.6 | 24.1 | 25.9 |
| ✓ | | | 24.3 | 24.3 | 15.2 | 24.8 | 27.3 |
| | ✓ | | 24.3 | 24.7 | 14.8 | 25.0 | 27.2 |
| | | ✓ | 24.6 | 24.9 | 16.5 | 25.0 | 27.2 |
| | ✓ | ✓ | 24.9 | 25.3 | 16.0 | 25.8 | 27.4 |
| ✓ | ✓ | | 25.2 | 25.4 | 17.5 | 25.8 | 27.5 |
| ✓ | ✓ | ✓ | **25.6** | **25.9** | **18.3** | **26.4** | **27.6** |

confidence scores of correct classified objects and suppress those of wrongly classified objects. This can help to improve the precision/recall of a model under the same recall/precision as evidenced in Fig. 6. Besides, we evaluate the Forest R-CNN with different settings of the classification trees. As shown in Tab.3, the Forest R-CNN with a single classification tree (§ 3.2) improves around 1% AP over the baseline, and the lexical tree achieves slightly better performance. With multiple trees, our method further boosts the AP to 24.9%-25.6%, validating the effectiveness of the classification forest (§ 3.3).

**Effectiveness of the NMS Resampling.** We compare the image resampling method [15] with the two proposed NMS Resampling schemes in Tab.4. We see from the results that both the NMS Resampling-Linear and NMS Resampling-Discrete achieve better performance than the image resampling. The reason for the performance gap is that image resampling may introduce severer overfitting by repeating the same images during training. Also, we assess

Table 4: Performance comparison with the image resampling [15]. "NR-Linear" denotes the NMS Resampling-Linear and "IR" denotes the image resampling.

| IR [15] | NR | AP | $AP^b$ | $AP_r$ | $AP_c$ | $AP_f$ |
|---|---|---|---|---|---|---|
| | | 21.7 | 21.8 | 6.8 | 22.6 | **26.4** |
| ✓ | | 23.0 | 22.7 | 13.8 | 23.4 | 26.1 |
| | ✓ (NR-Linear) | **23.7** | **23.5** | 12.3 | **25.2** | 26.3 |
| | ✓ (NR-Discrete) | 23.5 | **23.5** | **15.6** | 24.1 | 25.9 |

Table 5: Results of the NMS Resampling-Discrete with different thresholds.

| NMS Resampling-Discrete | AP | $AP^b$ | $AP_r$ | $AP_c$ | $AP_f$ |
|---|---|---|---|---|---|
| $\alpha_f = 0.7, \alpha_c = 0.8, \alpha_r = 0.9$ | **23.5** | **23.5** | **15.6** | **24.1** | 25.9 |
| $\alpha_f = 0.6, \alpha_c = 0.7, \alpha_r = 0.8$ | 23.4 | 23.1 | 15.2 | 23.8 | **26.3** |

Table 6: Results of the Forest R-CNN using different number of parent classes. The visual tree is used for experiments.

| Number of Parent Classes | AP | $AP^b$ | $AP_r$ | $AP_c$ | $AP_f$ |
|---|---|---|---|---|---|
| $M_{vis}=25$ | **24.3** | 24.7 | 14.8 | **25.0** | **27.2** |
| $M_{vis}=50$ | 24.1 | **24.8** | **15.0** | 24.8 | 26.8 |
| $M_{vis}=100$ | 23.7 | 24.3 | 12.8 | 24.9 | 26.7 |

Table 7: Performance Comparison on the LVIS v1.0 `val` set. "R" denotes the ResNet with FPN, and "Cascade" denotes the Cascade R-CNN [3].
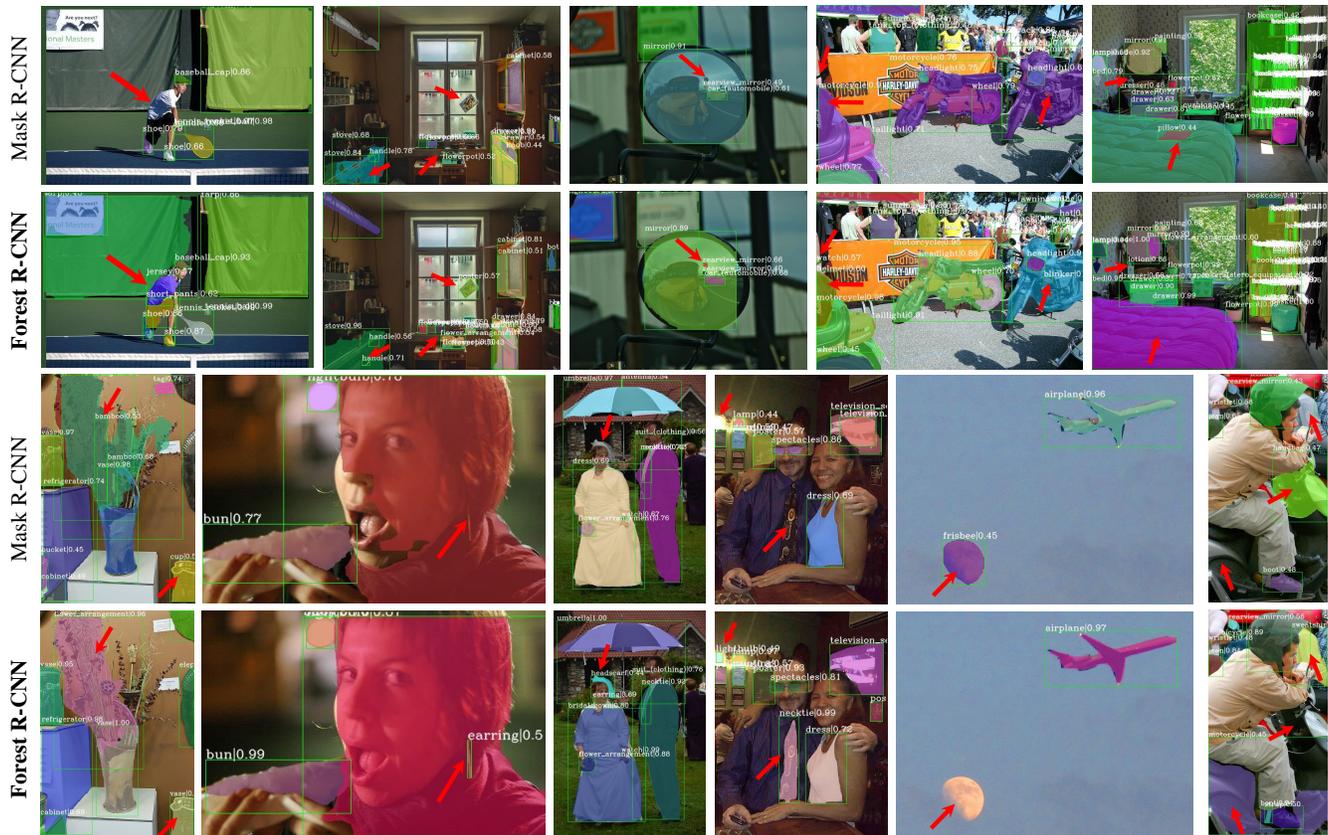
| Method | Setting | AP | $AP^b$ | $AP_r$ | $AP_c$ | $AP_f$ |
|---|---|---|---|---|---|---|
| Mask R-CNN | R-101 | 20.8 | 21.7 | 1.4 | 19.4 | 30.9 |
| EQL [38] | R-101 | 22.9 | 24.2 | 3.7 | 23.6 | 30.7 |
| Cascade Mask R-CNN | Cascade R-101 | 22.6 | 25.2 | 2.0 | 22.0 | 32.5 |
| De-confound [39] | Cascade R-101 | 23.5 | 25.8 | 5.2 | 22.7 | 32.3 |
| Mask R-CNN | R-50 | 19.2 | 20.0 | 0.0 | 17.2 | **29.5** |
| EQL [38] | R-50 | 21.6 | 22.5 | 3.8 | 21.7 | 29.2 |
| Forest R-CNN (Ours) | R-50 | **23.2** | **24.6** | **14.2** | **22.7** | 27.7 |

the proposed NMS Resampling with different threshold settings. As shown in Tab.5, our method is not sensitive to the specific threshold values so long as they are inverse to the data amount of categories.

**Number of the Parent Classes.** To investigate the impact of hyper-parameter $M$ on the classification tree, we experiment the Forest R-CNN with the Visual tree in Tab.6. We see from the table that the performances of $M_{vis} = 25$ and $M_{vis} = 50$ are close and the results of $M_{vis} = 100$ get slightly worse. This suggests that it is better to set $M \leq 50$ when clustering $\sim 1000$ fine-grained class into

**Table 8: Performance comparison with the state-of-the-art methods on the LVIS v0.5 val set. We denote "IR" as the image resampling [15] and "MST" as the multi-scale training.**

| Method | Setting | AP | $AP_{50}$ | $AP_{75}$ | $AP_r$ | $AP_c$ | $AP_f$ | $AP_S$ | $AP_M$ | $AP_L$ | $AP^b$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Class-aware Sampling [34] | ResNet-50-FPN | 18.5 | 31.1 | 18.9 | 7.3 | 19.3 | 21.9 | 13.3 | 24.3 | 30.5 | 18.4 |
| Repeat Factor Sampling [15] | ResNet-50-FPN | 23.2 | - | - | 13.4 | 23.2 | 27.1 | - | - | - | - |
| Class-balanced Loss [10] | ResNet-50-FPN | 20.9 | 33.8 | 22.2 | 8.2 | 21.2 | 25.7 | 15.6 | 28.1 | 35.3 | 21.0 |
| Focal Loss [24] | ResNet-50-FPN | 21.0 | 34.2 | 22.1 | 9.3 | 21.0 | 25.8 | 15.6 | 27.8 | 35.4 | 21.9 |
| LST[19] | ResNet-50-FPN | 23.0 | 36.7 | 24.8 | - | - | - | - | - | - | 22.6 |
| EQL [38] | ResNet-50-FPN | 22.8 | 36.0 | 24.4 | 11.3 | 24.7 | 25.1 | 16.3 | 29.7 | 38.2 | 23.3 |
| Forest R-CNN (Ours) | ResNet-50-FPN | **25.6** | **40.3** | **27.1** | **18.3** | **26.4** | **27.6** | **18.5** | **32.7** | **41.1** | **25.9** |
| EQL [38] | ResNet-101-FPN | 24.8 | 38.4 | 26.8 | 14.6 | 26.7 | 26.4 | - | - | - | 25.2 |
| Forest R-CNN (Ours) | ResNet-101-FPN | **26.9** | **42.2** | **28.4** | **20.1** | **27.9** | **28.3** | **23.6** | **43.5** | **51.8** | **27.5** |
| SOLOv2 [41] | ResNet-50-FPN & MST & IR | 25.5 | - | - | 13.4 | 26.6 | **28.9** | 15.9 | 34.6 | **44.9** | - |
| Forest R-CNN (Ours) | ResNet-50-FPN & MST | **26.7** | **42** | **28.8** | **19.7** | **27.5** | 28.5 | **20.3** | **34.8** | 39.7 | **27** |
| SOLOv2 [41] | ResNet-101-FPN & MST & IR | 26.8 | - | - | 16.3 | 27.6 | **30.1** | 16.8 | 35.8 | **47.0** | - |
| Forest R-CNN (Ours) | ResNet-101-FPN & MST | **28.2** | **44.3** | **29.7** | **20.2** | **29.6** | 29.6 | **21.0** | **36.0** | 42.0 | **28.6** |



**Figure 8: Mask R-CNN [17] *vs.* Forest R-CNN. Mask R-CNN exhibits more wrong classification and miss recognition. For neat visualization, we apply the NMS with threshold of** 0.7 **and filter out the predictions with scores lower than** 0.4.

$M$ parent classes using K-means. We observe in experiments that the results are stable under different K-means initializations.

## 4.3 Comparison with State-of-the-Art

**LVIS v0.5.** In Tab.8, we compare our method with state-of-the-art methods on the LIVS v0.5 dataset. It is worth noting that the

proposed Forest R-CNN achieves state-of-the-art performance under different experimental setups. Moreover, the Forest R-CNN improves 3.9%-7.0% AP over the second-best results of different setups on the rare category. It demonstrates that our method is skilled in recognizing the tail classes as well.

**LVIS v1.0.** We report the result of Forest R-CNN with ResNet-50-FPN on the LVIS v1.0 dataset in Tab.7. Forest R-CNN improves 4% AP and 14.2% AP on overall categories and rare categories, respectively, compared to the baseline Mask R-CNN. Moreover, Forest R-CNN achieves competitive performance compared to EQL [38] and De-confound [39] which are equipped with more complex network settings.

## 5 CONCLUSION

This work presents a novel object recognition model, Forest R-CNN, which is equipped with two key components: (**i**) the classification forest and (**ii**) the NMS Resampling. The classification forest suppresses the noisy logits produced by a fine-grained classifier, enhancing the capability of classifying thousands of categories. The NMS Resampling re-balances the long-tailed data distribution by adaptively adjusting the NMS thresholds for different categories, which aids our method in recognizing more objects from tail classes. The above designs enable strong performance on detecting and segmenting a large number of object instances, outperforming state-of-the-art competitors on the LVIS dataset.

## REFERENCES

[1] Daniel Bolya, Chong Zhou, Fanyi Xiao, and Yong Jae Lee. 2019. YOLACT: Real-time instance segmentation. In *Proceedings of IEEE International Conference on Computer Vision (ICCV)*. 9157–9166.

[2] Jonathon Byrd and Zachary Lipton. 2019. What is the effect of importance weighting in deep learning. In *Proceedings of International Conference on Machine Learning (ICML)*. 872–881.

[3] Zhaowei Cai and Nuno Vasconcelos. 2018. Cascade r-cnn: Delving into high quality object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition (CVPR)*. 6154–6162.

[4] Jiale Cao, Rao Muhammad Anwer, Hisham Cholakkal, Fahad Shahbaz Khan, Yanwei Pang, and Ling Shao. 2020. SipMask: spatial information preservation for fast instance segmentation. In *Proceedings of European conference on computer vision (ECCV)*. 740–755.

[5] Jiale Cao, Hisham Cholakkal, Rao Muhammad Anwer, Fahad Shahbaz Khan, Yanwei Pang, and Ling Shao. 2020. D2Det: Towards high quality object detection and instance segmentation. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 11485–11494.

[6] Jiale Cao, Yanwei Pang, Jungong Han, and XueLong Li. 2019. Hierarchical shot detector. In *Proceedings of IEEE International Conference on Computer Vision (ICCV)*. 9705–9714.

[7] Jiale Cao, Yanwei Pang, and XueLong Li. 2019. Triply supervised decoder networks for joint detection and segmentation. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 7392–7401.

[8] Jiale Cao, Yanwei Pang, Shengjie Zhao, and XueLong Li. 2019. High-level semantic networks for multi-Scale object detection. *IEEE Transactions on Circuits and Systems for Video Technology*.

[9] Kaidi Cao, Colin Wei, Adrien Gaidon, Nikos Arechiga, and Tengyu Ma. 2019. Learning imbalanced datasets with label-distribution-aware margin loss. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*. 1565–1576.

[10] Yin Cui, Menglin Jia, Tsung-Yi Lin, Yang Song, and Serge Belongie. 2019. Class-balanced loss based on effective number of samples. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 9268–9277.

[11] Jifeng Dai, Kaiming He, Yi Li, Shaoqing Ren, and Jian Sun. 2016. Instance-sensitive fully convolutional networks. In *Proceedings of European conference on computer vision (ECCV)*. 534–549.

[12] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. 2010. The pascal visual object classes (voc) challenge. *International Journal of Computer Vision*, 303–338.

[13] Ross Girshick. 2015. Fast r-cnn. In *Proceedings of IEEE International Conference on Computer Vision (ICCV)*. 1440–1448.

[14] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. 2014. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 580–587.

[15] Agrim Gupta, Piotr Dollár, and Ross Girshick. 2019. Lvis: a dataset for large vocabulary instance segmentation. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 5356–5364.

[16] Hui Han, Wen-Yuan Wang, and Bing-Huan Mao. 2005. Borderline-smote: a new over-sampling method in imbalanced data sets learning. In *Proceedings of International Conference on Intelligent Computing (ICIC)*. 878–887.

[17] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. 2017. Mask r-cnn. In *Proceedings of IEEE International Conference on Computer Vision (ICCV)*. 2961–2969.

[18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 770–778.

[19] Xinting Hu, Yi Jiang, Kaihua Tang, Jingyuan Chen, Chunyan Miao, and Hanwang Zhang. 2020. Learning to segment the tail. *arXiv preprint arXiv:2004.00900*.

[20] Chen Huang, Yining Li, Chen Change Loy, and Xiaoou Tang. 2016. Learning deep representation for imbalanced classification. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 5375–5384.

[21] Zhaojin Huang, Lichao Huang, Yongchao Gong, Chang Huang, and Xinggang Wang. 2019. Mask scoring r-cnn. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 6409–6418.

[22] Yi Li, Haozhi Qi, Jifeng Dai, Xiangyang Ji, and Yichen Wei. 2017. Fully convolutional instance-aware semantic segmentation. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2359–2367.

[23] Yu Li, Tao Wang, Bingyi Kang, Sheng Tang, Chunfeng Wang, Jintao Li, and Jiashi Feng. 2020. Overcoming classifier imbalance for long-tail object detection with balanced group softmax. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 10991–11000.

[24] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. 2017. Focal loss for dense object detection. In *Proceedings of IEEE International Conference on Computer Vision (ICCV)*. 2980–2988.

[25] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: common objects in context. In *Proceedings of European conference on computer vision (ECCV)*. 740–755.

[26] Shu Liu, Lu Qi, Haifang Qin, Jianping Shi, and Jiaya Jia. 2018. Path aggregation network for instance segmentation. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 8759–8768.

[27] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C. Berg. 2016. SSD: Single shot multibox detector. In *Proceedings of European conference on computer vision (ECCV)*. 21–37.

[28] Ziwei Liu, Zhongqi Miao, Xiaohang Zhan, Jiayun Wang, Boqing Gong, and Stella X. Yu. 2019. Large-scale long-tailed recognition in an open world. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2537–2546.

[29] George Miller. 1998. WordNet: An electronic lexical database. *MIT press*.

[30] Wanli Ouyang, Xiaogang Wang, Cong Zhang, and Xiaokang Yang. 2016. Factors in finetuning deep model for object detection with long-tail distribution. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 864–873.

[31] Yanwei Pang, Jiale Cao, and XueLong Li. 2016. Learning sampling distributions for efficient object detection. *IEEE Transactions on Cybernetics*, 117–129.

[32] Joseph Redmon and Ali Farhadi. 2017. YOLO9000: Better, faster, stronger. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 7263–7271.

[33] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2015. Faster r-cnn: Towards real-time object detection with region proposal network. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*. 91–99.

[34] Li Shen, Zhouchen Lin, and Qingming Huang. 2016. Relay backpropagation for effective learning of deep convolutional neural networks. In *Proceedings of European conference on computer vision (ECCV)*. 467–482.

[35] Bharat Singh and Larry S. Davis. 2018. An analysis of scale invariance in object detection – SNIP. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 3578–3587.

[36] Bharat Singh, Hengduo Li, Abhishek Sharma, and Larry S. Davis. 2018. R-FCN-3000 at 30fps: Decoupling detection and classification. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 1081–1090.

[37] Jingru Tan, Xin Lu, Gang Zhang, Changqing Yin, and Quanquan Li. 2020. Equalization Loss v2: A New Gradient Balance Approach for Long-tailed Object Detection. *arXiv preprint arXiv:2012.08548*.

[38] Jingru Tan, Changbao Wang, Buyu Li, Quanquan Li, Wanli Ouyang, Changqing Yin, and Junjie Yan. 2020. Equalization Loss for long-tailed object recognition. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 11662–11671.

[39] Kaihua Tang, Jianqiang Huang, and Hanwang Zhang. 2020. Long-tailed classification by keeping the good and removing the bad momentum causal effect. *Advances in Neural Information Processing Systems (NeurIPS)*.

[40] Tiancai Wang, Tong Yang, Martin Danelljan, Fahad Shahbaz Khan, Xiangyu Zhang, and Jian Sun. 2020. Learning human-object interaction detection using interaction points. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 4116–4125.

[41] Xinlong Wang, Rufeng Zhang, Tao Kong, Lei Li, and Chunhua Shen. 2020. Solov2: dynamic, faster and stronger. *arXiv preprint arXiv:2003.10152*.

[42] Jiajun Wu, Erika Lu, Pushmeet Kohli, Bill Freeman, and Josh Tenenbaum. 2017. Learning to see physics via visual deanimation. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*. 153–164.

[43] Jialian Wu, Chunluan Zhou, Ming Yang, Qian Zhang, Yuan Li, and Junsong Yuan. 2020. Temporal-context enhanced detection of heavily occluded pedestrians. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 13430–13439.

[44] Jialian Wu, Chunluan Zhou, Qian Zhang, Ming Yang, and Junsong Yuan. 2020. Self-Mimic learning for small-scale pedestrian detection. In *Proceedings of the 28th ACM International Conference on Multimedia (MM)*.

[45] Tong Xiao, Shuang Li, Bochao Wang, Liang Lin, and Xiaogang Wang. 2017. Joint detection and identification feature learning for person search. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 3415–3424.

[46] Wenqiang Xu, Haiyang Wang, Fubo Qi, and Cewu Lu. 2019. Explicit shape encoding for real-time instance segmentation. In *Proceedings of IEEE International Conference on Computer Vision (ICCV)*. 5168–5177.

[47] Xiao Zhang, Zhiyuan Fang, Yandong Wen, Zhifeng Li, and Yu Qiao. 2017. Range loss for deep face recognition with long-tailed training data. In *Proceedings of IEEE International Conference on Computer Vision (ICCV)*. 5409–5418.

[48] Yaoyao Zhong, Weihong Deng, Mei Wang, Jiani Hu, Jianteng Peng, Xunqiang Tao, and Yaohai Huang. 2019. Unequal-training for deep face recognition with long-tailed noisy data. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 7812–7821.

[49] Boyan Zhou, Quan Cui, Xiu-Shen Wei, and Zhao-Min Chen. 2020. BBN: Bilateral-branch network with cumulative learning for long-tailed visual recognition. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 9719–9728.