# RetinaMask: Learning to predict masks improves state-of-the-art single-shot detection for free

Cheng-Yang Fu     Mykhailo Shvets     Alexander C. Berg
Computer Science Department of
UNC at Chapel Hill

{cyfu, mshvets, aberg}@cs.unc.edu

## Abstract

*Recently two-stage detectors have surged ahead of single-shot detectors in the accuracy-vs-speed trade-off. Nevertheless single-shot detectors are immensely popular in embedded vision applications. This paper brings single-shot detectors up to the same level as current two-stage techniques. We do this by improving training for the state-of-the-art single-shot detector, RetinaNet, in three ways: integrating instance mask prediction for the first time, making the loss function adaptive and more stable, and including additional hard examples in training. We call the resulting augmented network RetinaMask. The detection component of RetinaMask has the same computational cost as the original RetinaNet, but is more accurate. COCO test-dev results are up to 41.4 mAP for RetinaMask-101 vs 39.1mAP for RetinaNet-101, while the runtime is the same during evaluation. Adding Group Normalization increases the performance of RetinaMask-101 to 41.7 mAP. Code is at:* https://github.com/chengyangfu/retinamask

## 1. Introduction

Single-shot detectors [29, 30, 28, 13, 25] have become extremely popular in applications where speed and computational resources are important design considerations. These include embedded vision applications, self-driving cars, and mobile phone vision. Despite this intense usage, there has been little improvement in state-of-the-art performance for single-shot detectors, e.g. RetinaNet [25]. When published in 2017, RetinaNet effectively cleaned up a range of work on single-shot detection, building on [30] and [28] and introducing innovations in training that resulted in state-of-the-art performance in terms of the speed-versus-accuracy trade-off (sharing the frontier with YOLOv3 [31] on the faster, lower-accuracy, end of the spectrum, and Mask R-CNN on the high end). While there have not been significant improvements in performance on top of Reti-
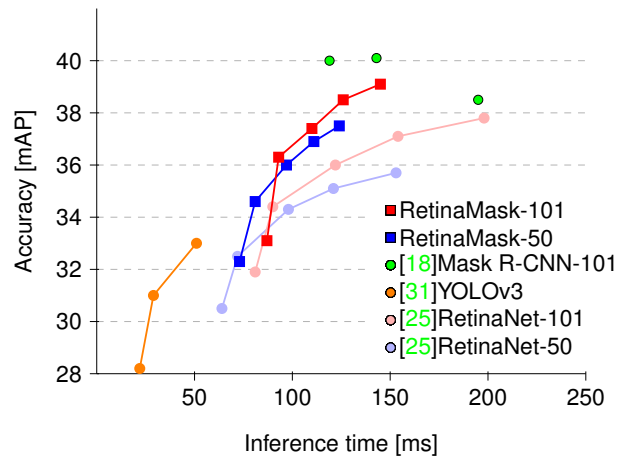


Figure 1: Accuracy versus inference time on COCO test-dev. For fair comparison with RetinaNet, we train our RetinaMask models at {400, 500, 600, 700, 800} resolution without using multi-scale augmentation during training. Our best model, which achieves 42.6 mAP, can be found in Table 5. Our improved versions, RetinaMask-50/101 respectively, are shown with blue/red square markers. The original RetinaNet-50/101 results are shown with blue/red circle markers. The state-of-the-art two-stage detector Mask R-CNN has improved since publication. We show three versions with green circle markers: original paper 38.5/195 detection+mask/M40 (mAP/ms/GPU), Detectron [1] Caffe2 (40.0/119 detection only/P100), and PyTorch [2] (40.1/143 detection only/V100).

naNet, two-stage detectors have advanced over the intervening time, and now outperform RetinaNet on the speed-vs-accuracy trade-off. Part of this improvement has been due to architectures like Mask R-CNN that allow training multiple prediction heads on top of the region proposal and bounding box prediction stages of the detector [18, 17].

In this paper we show how to improve the accuracy of state-of-the-art single-shot detectors (e.g. RetinaNet [25],

SSD [28, 13]) largely by adding the task of **instance mask prediction during training**, but also by introducing an **adaptive loss** that improves robustness to parameter choice during training, and including **more difficult examples** in training. We call the resulting system RetinaMask to make clear that its training included the additional task of instance mask prediction, which in the past had been added to two-stage but not to single-shot detectors. These modifications involve more work during training, but it is possible to evaluate just the detection part of RetinaMask, which has exactly the same computational cost as the original RetinaNet detector—our modifications to training result in detectors with better accuracy at the same computational cost.

Because we keep the structure of the detector at test time unchanged, our approach can be directly applied to the wide range of embedded applications that choose single-shot detectors over two-stage detectors. This is a subtle point, but today (before this paper), two-stage detectors significantly beat single stage detectors on the speed-vs-accuracy trade-off on a standard desktop/workstation + high-end GPU configurations. However, when adapting these implementations to lower-power embedded devices the cost of resampling (e.g. ROI-Align) for the second stage is often exacerbated by communications between the hardware that is used to accelerate convolutions and the CPU. Single-shot approaches avoid this extra implementation challenge, which makes them popular in many implementations.

This allusion to implementation overhead for a variety of special purpose architectures brings up one of the difficulties of keeping track of progress on the speed-vs-accuracy trade-off in detection. In addition to the algorithms, the hardware itself and software libraries are constantly evolving year over year. These complexities are small relative to the variety of software and architectures for embedded systems. The main point of our work is to show a significant improvement in the accuracy of single-shot detectors with keeping the same computational cost as the original network during inference. Note the the plot in Figure 1 shows an improvement in speed as well, but this is partly due to the next generation of libraries and slightly different hardware. We expect the improvements demonstrated in this paper to be applicable to a wide range of single-shot detector implementations.

In summary, this paper shows a significant increase in the accuracy of state-of-the-art single-shot detectors by introducing three new techniques that improve training:

1. Adding a novel instance mask prediction head to the single-shot RetinaNet detector during training.

2. A new self-adjusting loss function that improves robustness during training.

3. Including more of the positive examples in training, even those with low overlap.

Each of these contributions is analyzed with ablation studies, and together they provide a large boost to the accuracy of RetinaNet, bringing it up to state-of-the-art accuracy again.

## 2. Related Work

We review recent developments in object detection using two-stage and single-shot techniques, general techniques for improving detection, and work on integrating instance mask prediction with detection.

**Two-Stage Detectors:** Two-stage detectors follow a long line of reasoning in computer vision about grouping and perception. They first propose potential object locations in an image—region proposals—and then apply a classifier to these regions to score potential detections. Earlier sliding-window approaches ran into scaling problems as the number of potential windows combined with the number of models became unmanageable [8, 12]. Selective Search [37] allowed more expensive and accurate bag-of-visual-words (BoVW) features to be considered by using low-level vision to identify a smaller number of potential locations that needed to be evaluated. A transition to deep learning models was done with R-CNN [15], which used a convolutional neural network to replace the BoVW, resulting in a significant accuracy improvement. SPPNet [23] sped up this process by direct region pooling on the feature layers instead of repetitive image cropping. Then Fast R-CNN [14], and Faster R-CNN [32] sped up detection even further and increased accuracy by replacing Selective Search with a Regional Proposal Network. Faster R-CNN was also the first end-to-end trainable deep learning model for object detection. R-FCN [7] used position-sensitive features and ROI-pooling to make a fully convolutional network design.

**Single-Shot Detectors:** In contrast to two-stage detectors, single-shot detectors avoid image or feature re-sampling. OverFeat [34] and DeepMultiBox [10] were early examples. Then, YOLO [29, 30] and SSD [28] popularized the single-shot approach by demonstrating models that ran in real-time with good accuracy. More recently, RetinaNet [25] proposed Focal Loss to address the extreme class imbalance problem between target and background, and cleaned up a number of design aspects for single-shot detection.

**Techniques for Improving Detectors:** Several techniques for improving detection apply to both Single-Shot and Two-Stage Detectors. First, cleaner training data often helps achieve faster convergence and higher final accuracy. Online Hard Negative Sampling [35] uses non-maximum suppression (nms) during training to provide negative examples diversity. Second, certain model modifications add context information to predictions. SSD [28] and MS-CNN [5] both predict instances across features layers of different resolutions. DSSD [13], Feature Pyramid Network [24]

and TDN [36] combine feature layers in a top-down manner to enrich the context of coarser features, strengthening the feature representation for better detection. Also, additional training information is beneficial for detectors. BlitzNet [9] augments SSD with a semantic segmentation prediction branch, combining these two tasks in a single network, which results in higher detection accuracy.

**Instance Segmentation:** As object detection matured and demonstrated strong accuracy and high speed, the research community started shifting attention to the more detailed task of instance segmentation. In addition to generating a tight bounding box for each object, instance segmentation requires a pixel-level mask for that object. The COCO [26] dataset established a recognized benchmark for this task by holding the Instance Segmentation Challenge starting in 2015. Current state-of-the-art instance segmentation approaches are based on two-stage detectors, adding an instance mask prediction module after detection. MNC [6] breaks down the Instance Segmentation into three stages, namely object detection, class-agnostic mask prediction, and mask categorization. FCIS [40] extends the idea of R-FCN [7] by using position-sensitive score maps for mask prediction. The recent Mask R-CNN [18] identifies the core issue for mask prediction in ROI pooling box misalignment, which arises from pooling box quantization over the coarse feature scale. Bilinear interpolation is introduced in their ROI-Align module to fix this issue. Mask R-CNN has been further improved in the Path Aggregation Network [27], using the ROI-Align operation on multiple feature layers to aggregate better features for instance segmentation.

## 3. Model

We start with the RetinaNet settings in Detectron[1] and rebuild it in PyTorch to form our baseline. Then, we introduce the following modifications to the baseline settings: best matching policy (Sec. 3.1), and modified bounding box regression loss (Sec. 3.2). Finally, we describe how to add the mask prediction module on top of RetinaNet (Sec. 3.3).

### 3.1. Best Matching Policy

In the bounding box matching stage, the RetinaNet policy is as follows. All anchor boxes that have an intersection-over-union (IOU) overlap with a ground truth object greater than 0.5, are considered positive examples. If the overlap is less than 0.4, the anchor boxes are assigned a negative label. All anchors for which the overlap falls between 0.4 and 0.5 are not used in the training. However, there exists an exceptional case for which the assignment can be improved. Specifically, some of the ground truth objects' aspect ratios are outliers, with one side much longer than the other. Thus, no anchor box can be matched to those according to the RetinaNet strategy. For each of these ground

truth boxes we propose to find its best matching anchor box, relaxing the overlapping IOU threshold. We provide an ablation study using different thresholds on the best matching anchors. The results suggests that using best matching anchor with any nonzero overlap gives the best accuracy (notice that such anchors always exists, because single-shot anchors are densely sampled).
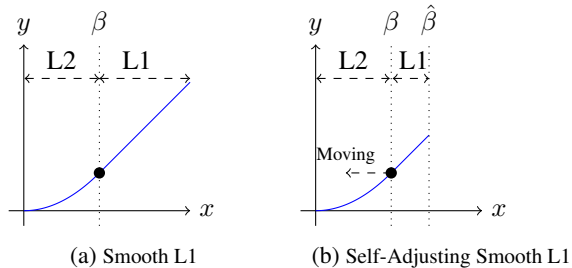


Figure 2: Smooth L1 and Self-Adjusting Smooth L1. In Smooth L1 Loss (a) $\beta$ if a fixed threshold that separates loss into L2 and L1 regions. In the proposed Self-Adjusting Smooth L1 (b), the $\beta$ is calculated as the difference between running mean and running variance of L1 loss and the value is clamped to the $[0, \hat{\beta}]$ range. The $\beta$ is approaching 0 during training.

### 3.2. Self-Adjusting Smooth L1 Loss

Smooth L1 Loss for object detection was originally proposed in Fast R-CNN [14] to make bounding box regression more robust by replacing the excessively strict L2 Loss. Current state-of-the-art methods such as Faster R-CNN [32], R-FCN [7] and SSD [28] still use this loss.

In Smooth L1 Loss described in Equation 1, a point $\beta$ splits the positive axis range into two parts: $L2$ loss is used for targets in range $[0, \beta]$, and $L1$ loss is used beyond $\beta$ to avoid over-penalizing outliers. The overall function is smooth (continuous, together with its derivative), as illustrated in Figure 1. However, the choice of control point($\beta$) is heuristic and is usually done by hyper parameter search.

$$f(x) = \begin{cases} 0.5\frac{x^2}{\beta}, & \text{if } |x| < \beta \\ |x| - 0.5\beta, & \text{otherwise} \end{cases} \quad (1)$$

We propose an improved version of Smooth L1 called Self-Adjusting Smooth L1 Loss. Inside the loss function, the running mean and variance of the absolute loss are recorded. We use the running minibatch mean and variance with `momentum`=0.9 to update these two parameters.

Then, the parameters are used to calculate the control point. Specifically, the control point is chosen to be equal to the difference between the running mean and running variance ($\mu_R - \sigma_R^2$), and the value is clipped to a range $[0, \hat{\beta}]$, as can be seen in Equation 2. Clipping is used because running mean is unstable during training, as the number of positive

examples in each batch is different. Figure 3 shows the running mean of L1 loss for the x offset and for width adjustment prediction in bounding box regression. We observe a decreasing trend for both during training.

$$\mu_B = \frac{1}{n}\sum_{i=1}^{n}|x_i|, \quad \sigma_B^2 = \frac{1}{n}\sum_{i=1}^{n}(|x_i| - \mu_B)^2$$
$$\mu_R = \mu_R * \mathrm{m} + \mu_B * (1 - \mathrm{m}) \qquad (2)$$
$$\sigma_R^2 = \sigma_R^2 * \mathrm{m} + \sigma_B^2 * (1 - \mathrm{m})$$
$$\beta = \max(0, \min(\hat{\beta}, \mu_R - \sigma_R^2))$$
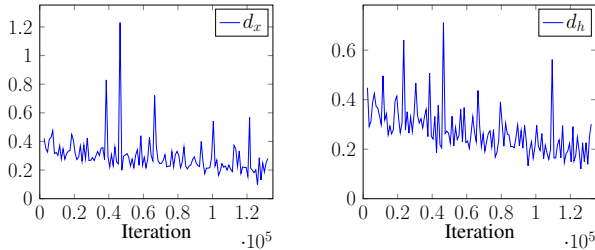


Figure 3: Running mean of the L1 loss applied to bounding box regression variables: $x$ offset prediction $\mathrm{d}_x$ and width prediction $\mathrm{d}_w$. Similar plots for $y$ offset prediction $\mathrm{d}_y$ and height prediction $\mathrm{d}_h$ are omitted for readability.

### 3.3. Mask Prediction Module

In order to add the mask prediction module, single-shot detection predictions are treated as mask proposals. After running RetinaNet for bounding box predictions, we extract the top $N$ scored predictions. Then, we distribute these mask proposals to sample features from the appropriate layers of the FPN according to Equation 3 proposed in FPN [24]. Figure 4 illustrates the assignment process. We use the following equation to determine which feature map, $P_k$ to sample from for predicting the instance mask:

$$k = \lfloor k_0 + \log_2 \sqrt{wh}/224 \rfloor, \qquad (3)$$

where $\mathrm{k}_0 = 4$, and w, h are the width and height of the detection. If the size is smaller than $224^2$, it will be assigned to feature layer $P_3$, between $224^2$ to $448^2$ is assigned to $P_4$, and larger than $448^2$ is assigned to $P_5$.

In our final model, we use the $\{P_3, P_4, P_5, P_6, P_7\}$ [2] feature layers for bounding box predictions and $\{P_3, P_4, P_5\}$ feature layers for mask prediction. In our ablation study, we analyze the impact of using more feature layers for mask proposals assignment, showing that this does not give any performance boost.

**Network:** Figure 4 shows a high-level overview of the model. We use ResNet-50 and ResNet-101 as backbone

---

² We use the same definition as in FPN [24] and RetinaNet [25].

models in our experiments, freezing all of the Batch Normalization layers. Following the Feature Pyramid Network [24] setting, we add extra layers ($P_6$ and $P_7$) and form top-down connections ($P_5$, $P_4$, and $P_3$). The dimensionality of each of the Feature Pyramid layers ($P_3, \ldots, P_7$) is set to 256. The bounding box classification head consists of 4 convolutional layers (conv3x3(256) + ReLU) and uses 1 convolution (conv3x3(number of anchors * number of classes)) with point-wise sigmoid nonlinearities. For bounding box regression, we adopt the class-agnostic setting. We also run 4 convolutional layers (conv3x3(256) + ReLU) and 1 output layer (conv3x3(number of anchors * 4)) to refine the anchors. Once the bounding boxes are predicted, we aggregate them and distribute to the Feature Pyramid layers, as discussed above. The ROI-Align operation is performed at the assigned feature layers, yielding 14x14 resolution features, which are fed into 4 consequent convolutional layers (conv3x3), and a single transposed convolutional layer (convtranspose2d 2x2) that upsamples the map to 28x28 resolution. Finally, a prediction convolutional layer (conv1x1) is applied. We predict class-specific masks.

### 3.4. Training

To train RetinaNet, we follow the settings in the original paper. Images are resized to make the shorter side equal to 800 pixels, while limiting the longer side to 1333 pixels.

We use batch size of 16 images, weight decay $10^{-4}$, momentum 0.9, and train for 90k iterations with the base learning rate of 0.01, dropped to 0.001 and 0.0001 at iterations 60k and 80k. We train our models on servers with 4 1080Ti GPUs. For some models (e.g. a ResNet-101-FPN backbone), there is not enough GPU memory for this batch size. If this is the case, we follow the Linear Scaling policy proposed in [16] and reduce the batch size (increasing the number of training iterations and reducing the learning rate accordingly). In order to train with the Mask Prediction Module, we extend the number of training iterations by a factor of 1.5x, or 2x, during multi-scale training. The multi-scale training is done at scales {640, 800, 1200}.

Thus, for 1.5x training, the number of iterations is set to 135k, and learning rate drops occur at iterations 90k, and 120k. For 2x, we train the network 180k iterations, and drop the learning rate at 120k, 160k. The training time is $\approx 56$ hours when using ResNet-50 with 1.5x train iteration, while for ResNet-101 it goes up to $\approx 75$ hours, for the same number of iterations.

The anchor boxes span 5 scales and 9 combinations (3 aspect ratios [0.5, 1, 2] and 3 sizes [$2^0$, $2^{1/3}$, $2^{2/3}$]), following [25]. The base anchor sizes range from $32^2$ to $512^2$ on Feature Pyramid levels $P_3$ to $P_7$. Each anchor box is matched to no more than one ground truth bounding box. The anchors that have intersection-over-union overlap with a ground truth box larger than 0.5 are considered positive
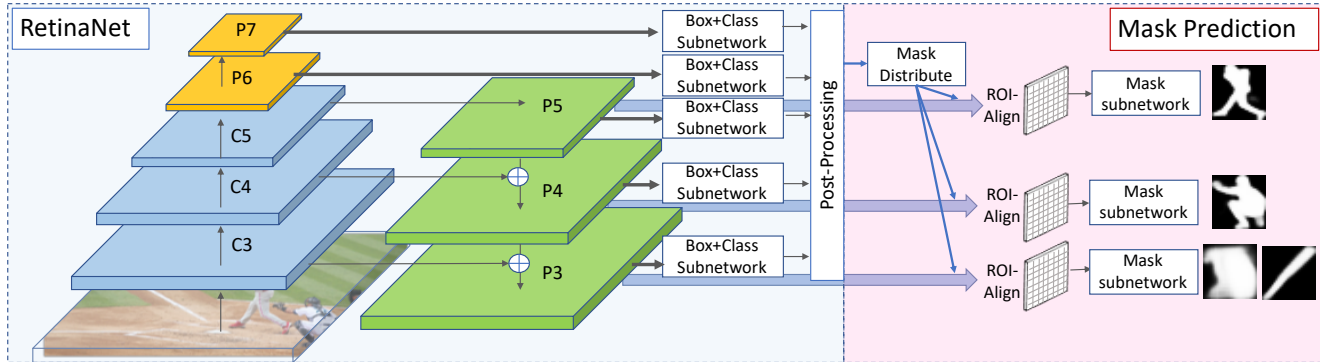
Figure 4: RetinaMask network architecture. This figure demonstrates a single-shot detector extended with the mask prediction module. Left part shows the RetinaNet on the Feature Pyramid Network. The classification and bounding box regression module is applied on the feature layers, $\{P_3, P_4, P_5, P_6, P_7\}$. After the predicted bounding boxes are processed and aggregated, they are distributed to $\{P_3, P_4, P_5\}$ for mask predictions according to the size. This results in $P_5$ predicting masks for larger objects, and $P_3$ predicting smaller objects.

examples. On the other hand, if the overlap is less than 0.4, such anchors are treated as negative examples. Then, we use the proposed best matching policy, as described in Section 3, which can only add positive examples. For the Focal Loss function 4 used in classification, we set $\alpha = 0.25$, $\gamma = 2.0$, and initialize the prediction logits according to $\mathcal{N}(0, 0.01)$ distribution. For the bounding box regression we add the proposed Self-Adjusting Smooth L1 and limit the control point to the range $[0, 0.11]$ ($\hat{\beta} = 0.11$).

$$\mathrm{FL} = -\alpha_t (1 - p_t)^\gamma \log(p_t) \tag{4}$$

For each image during training, we also run suppression and top-100 selection of the predicted boxes (the same processing as single-shot detectors apply during inference). Then, we add ground truth boxes to the proposals set, and run the mask prediction module. Thus, the number of mask proposals is (100+Gt) during training. The final loss function is a sum of the three losses: $Loss_{boxCls} + Loss_{boxReg} + Loss_{mask}$.

### 3.5. Inference

First, during the bounding box inference we use a confidence threshold of $0.05$ to filter out predictions with low confidence. Second, we select the top 1000 scoring boxes from each prediction layer. Third, we apply non-maximum suppression (nms) with threshold $0.4$ for each class separately. Finally, the top-100 scoring predictions are selected for each image. For mask inference, we use the top 50 bounding box predictions as mask proposals. Although there are more intelligent ways to perform post-processing, such as SoftNMS [4] or test-time image augmentations, in order to fairly compare against the baseline models in speed and accuracy, we intentionally do not use those here.

| Threshold | AP | $AP_{50}$ | $AP_{75}$ | $AP_S$ | $AP_M$ | $AP_L$ |
|-----------|------|------|------|------|------|------|
| 0.5 | 35.5 | 53.7 | 38.1 | 19.5 | 39.3 | 47.4 |
| 0.4 | 36.0 | 54.1 | 38.6 | 19.1 | 39.6 | 48.3 |
| 0.3 | 36.1 | 54.5 | **38.9** | 19.8 | 39.6 | **48.6** |
| 0.2 | 36.1 | 54.5 | 38.7 | **20.4** | **39.8** | **48.6** |
| 0.0 | **36.2** | **55.0** | 38.7 | 19.7 | 39.5 | **48.6** |

Table 1: Ablation study of different thresholds used in the best matching case on COCO `minival`. The selected threshold relaxes the regular intersection-over-union threshold of 0.5 for assigning at least one anchor box to each ground truth box. The base threshold is kept at 0.5, so the modification only affects previously unmatched ground truth objects.

## 4. Experiments

**Dataset:** In this paper, we use the COCO [26] dataset, which provides bounding box and segmentation mask annotations. We follow common practice [3], using the COCO `trainval135k` split (union of 2014train 80k and a subset of 35k images from 2014val 40k) for training and the `minival` (remaining 5k images from 2014val 40k) for evaluation[3].

### 4.1. Ablation Study

**Best Matching Policy:** In Table 1, we test the effectiveness of using Best Matching Policy for all ground truth objects, as described in Section 3.1. Threshold 0.5 corresponds to regular matching. We then gradually lower the threshold for the best matching policy, going down all the way to 0 (completely relaxing the threshold). According to Table 1, using best matching anchors with any positive overlap to ground truth gives the best performance.

---

[3]COCO `trainval135k` is also called COCO 2017 train and the `minival` is COCO 2017 val.

| | AP | $AP_{50}$ | $AP_{75}$ | $AP_S$ | $AP_M$ | $AP_L$ |
|---|---|---|---|---|---|---|
| *Fixed* | | | | | | |
| $\beta = 1.0$ | 35.3 | 55.6 | 37.8 | 19.4 | 38.9 | 46.9 |
| $\beta = 0.11$ | 36.2 | 55.0 | 38.7 | 19.7 | 39.5 | 48.6 |
| *Self-Adj* | | | | | | |
| $\beta \leq 1.0$ | 36.4 | 55.4 | **39.0** | 19.9 | 39.9 | 48.1 |
| $\beta \leq 0.11$ | **36.6** | **55.7** | **39.0** | **20.3** | **40.0** | **48.8** |

Table 2: Ablation study of Self-Adjusting Smooth L1 with different $\beta$ on COCO `minival`.

| Method | Train | $AP^{bb}$ | $AP^{bb}_{50}$ | $AP^{bb}_{75}$ | $AP^m$ | $AP^m_{50}$ | $AP^m_{75}$ |
|---|---|---|---|---|---|---|---|
| Base | 1x | 36.2 | 55.0 | 38.7 | — | — | — |
| $P_3$-$P_5$ | 1x | 36.9 | 55.3 | **39.7** | 32.7 | 52.2 | 34.9 |
| $P_3$-$P_5$ | 1.5x | **37.1** | **55.9** | 39.5 | **33.0** | **52.9** | 35.0 |
| $P_2$-$P_5$ | 1x | 36.7 | 54.9 | **39.7** | 32.8 | 52.2 | **35.1** |

Table 3: Ablation study of different settings for adding mask prediction module on COCO `minival`.

Qualitative results suggest that our matching strategy reduces the number of duplicate detections. Indeed, best matching enforces larger changes to anchor boxes (but not too large to destabilize the training process), so different anchors shrink tighter to the ground truth object, and only one survives during non-maximum suppression.

**Self-Adjusting Smooth L1 Loss:** We first ran the Smooth L1 with fixed values (1.0 and 0.11). The choice of $\beta$ is not specified in RetinaNet [25]. According to the released implementation, Detectron, 0.11 is used. The upper part of Table 2 shows that setting $\beta$ is set to 1.0 will favor will favor $AP_{0.5}$ which is widely used in datasets which adopt IOU=0.5 as the evaluation metric such as Pascal VOC [11]. In contrast, the smaller value, 0.11, will favor more restrictive metrics such as IOU=.50:.05:.95 [26].

The bottom part of Table 2 shows the results of using Self-Adjusting Smooth L1 loss. First, we can see that our Self-Adjusting loss with setting 0.11, gives the best results for every metric. It is clear that this method is not dataset dependent. Second, the Self-Adjusting Smooth L1 is robust. When we change the bounding region from 0.11 to 1.0, the decrease of results is minor compared to the original Smooth L1 method. We also tried to share the running mean and variance across channels. The result (36.4 mAP) is slightly worse than the separate channel version.

**Multi-Task Training with Mask Prediction:** Table 3 illustrates bounding box accuracy improvement when running multi-task training with instance segmentation. When training with mask prediction using $\{P_3, P_4, P_5\}$, we see 0.7 mAP improvement. If we train with 1.5x schedule, the improvement is 0.9 mAP. If we add the feature layer with higher resolution, will it be helpful for the prediction? We follow Mask R-CNN to use $\{P_2, P_3, P_4, P_5\}$ for mask prediction. The results are slightly better on mask prediction

but worse on detection. In conclusion, adding mask prediction consistently improves detection results, but requires longer training. It is also worth noting that in the Mask R-CNN ablation study, the authors also show 0.9 mAP improvement on bbox prediction from multi-tasking training.



(a) Tie



(b) Ski



(c) Sports Ball



(d) Toaster

Figure 5: Comparison to RetinaNet baseline (left column). RetinaMask (right column) includes Best Matching policy, Self-Adjusting L1 loss, and Mask Prediction (see Section 4.2). The figure shows all detection results (no confidence threshold applied) only for selected categories (*tie* for (a), *skis* for (b), *sports ball* for (c), and *toaster* for (d)). Prediction scores are also shows, where they do not clutter the image. Our model shows improvement in classes with large aspect ratios (no multiple detections for *tie*, and better recall for *skis*); In (a) our model demonstrates no false negatives (e.g. note false negative sports ball with 0.14 score for the baseline model); (d) shows the failure toaster case, that accounts for the decrease in Figure 7 (only 9 toasters in COCO `minival`). Better viewed electronically, enlarged.

## 4.2. Comparison to RetinaNet

Following [20] we give an explanation of our model's improvement over the RetinaNet baseline. The model evaluated in this section incorporates all three components described in Section 3: Best Matching policy, Self-Adjusting Smooth L1, and Mask Prediction head. ResNet-50 is used as the backbone architecture, and images are resized to a shorter side of 800 pixels. No data augmentation is used.

First, we look at per-class difference of the mean Average Precision in Figure 7, showing improvement in most of the classes. Note that the *toaster* class, whose mAP decreases by 7.9 points (from 28.9 to 21.0), has only 9 ground truth objects in the validation set. On the other hand, hair drier shows a significant increase from 0.9 to 7.1 mAP points. The classes that improve most also include *snowboard*, *sports ball*, *kite*, *refrigerator*, and *scissors* (mAP difference $\geq 5$). See some qualitative examples in Figure 5.
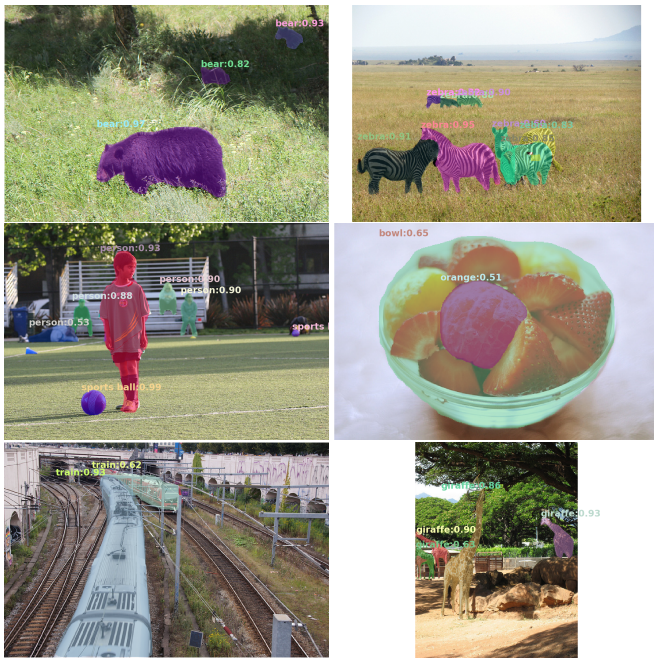


Figure 6: Visualization of RetinaMask with ResNet-101-FPN-GN Model(BBox=41.7 mAP, Mask=36.7 mAP result on COCO `test-dev`).

Figure 8 shows the difference in class-agnostic weak detection at low IoU threshold of 0.1, which ignores localization errors. Moreover, foreground object mis-classification is also ignored, which does not count for errors of attributing an object of one category to a different category. High correlation in these two relative differences for the improved classes (*snowboard*, *sports ball*, etc.) suggest that a large portion of network improvement comes from better localization, rather than better confidence prediction (otherwise class-agnostic weak detection would not improve).

| D | S | M | AP | $AP_{50}$ | $AP_{75}$ | $AP_S$ | $AP_M$ | $AP_L$ | T |
|---|---|---|---|---|---|---|---|---|---|
| 50 | 400 | R | 30.5 | 47.8 | 32.7 | 11.2 | 33.8 | 46.1 | 64 |
| | | O(B) | 32.3 | 49.8 | 34.2 | 11.6 | 34.4 | 47.8 | 73 |
| | | O(M) | 28.7 | 46.9 | 30.1 | 07.0 | 29.1 | 47.4 | 82 |
| 50 | 500 | R | 32.5 | 50.9 | 34.8 | 13.9 | 35.8 | 46.7 | 72 |
| | | O(B) | 34.6 | 52.6 | 36.7 | 14.8 | 36.7 | 48.8 | 81 |
| | | O(M) | 30.7 | 49.8 | 32.3 | 09.6 | 31.7 | 48.2 | 92 |
| 50 | 600 | R | 34.3 | 53.2 | 36.9 | 16.2 | 37.4 | 47.4 | 98 |
| | | O(B) | 36.0 | 54.5 | 38.5 | 17.1 | 38.1 | 49.2 | 97 |
| | | O(M) | 31.9 | 51.6 | 33.8 | 11.5 | 33.3 | 48.4 | 108 |
| 50 | 700 | R | 35.1 | 54.2 | 37.7 | 18.0 | 39.3 | 46.4 | 121 |
| | | O(B) | 36.9 | 55.6 | 39.6 | 18.9 | 39.1 | 48.7 | 111 |
| | | O(M) | 32.8 | 52.9 | 35.0 | 13.0 | 34.5 | 48.4 | 123 |
| 50 | 800 | R | 35.7 | 55.0 | 38.5 | 18.9 | 38.9 | 46.3 | 153 |
| | | O(B) | 37.5 | 56.4 | 40.2 | 19.6 | 39.8 | 48.9 | 124 |
| | | O(M) | 33.4 | 53.7 | 35.4 | 13.6 | 35.1 | 48.7 | 141 |
| 101 | 400 | R | 31.9 | 49.5 | 34.1 | 11.6 | 35.8 | 48.5 | 81 |
| | | O(B) | 33.1 | 49.7 | 35.4 | 11.0 | 35.7 | 49.9 | 87 |
| | | O(M) | 29.4 | 47.3 | 31.2 | 06.8 | 30.4 | 49.5 | 99 |
| 101 | 500 | R | 34.4 | 53.1 | 36.8 | 14.7 | 38.5 | 49.1 | 90 |
| | | O(B) | 36.3 | 54.7 | 38.7 | 15.9 | 39.1 | 50.9 | 93 |
| | | O(M) | 32.0 | 51.8 | 33.8 | 10.2 | 33.8 | 50.2 | 105 |
| 101 | 600 | R | 36.0 | 55.2 | 38.7 | 17.4 | 39.6 | 49.7 | 122 |
| | | O(B) | 37.4 | 56.0 | 39.9 | 17.3 | 39.9 | 51.4 | 110 |
| | | O(M) | 33.2 | 53.2 | 35.2 | 11.6 | 35.0 | 50.7 | 120 |
| 101 | 700 | R | 37.1 | 56.6 | 39.8 | 19.1 | 40.6 | 49.4 | 154 |
| | | O(B) | 38.5 | 57.3 | 41.3 | 19.1 | 41.1 | 51.8 | 126 |
| | | O(M) | 34.1 | 54.5 | 36.3 | 13.0 | 36.2 | 51.0 | 137 |
| 101 | 800 | R | 37.8 | 57.5 | 40.8 | 20.2 | 41.1 | 49.2 | 198 |
| | | O(B) | 39.1 | 58.0 | 41.9 | 20.4 | 41.7 | 51.0 | 145 |
| | | O(M) | 34.7 | 55.4 | 36.9 | 14.3 | 36.7 | 50.5 | 166 |

Table 4: Comparison to RetinaNet with different input resolutions on COCO `test-dev` (Also see Figure 1). For each (D/S) depth/scale, the upper part (R) is the RetinaNet performance from Table 1(e) in RetinaNet [25], our results are in the bottom part. We also report mask prediction accuracies. For Depth, 50:ResNet-50-FPN, 101:ResNet-101-FPN. For Method, R:RetinaNet, O(B): Our Method of BBox prediction, O(M): Our Method of Mask prediction. Our speed number is evaluated on Nvidia 1080 Ti / PyTorch1.0 and FocalLoss results are evaluated on Nvidia M40 / Caffe2.

Table 4 shows comparisons of our model to RetinaNet on different backbone networks and input resolutions. RetinaNet results come from the Table 1(e) of the RetinaNet [25] paper. Our model shows better accuracy for all combinations of backbone network choices and resolutions. We report the speed number evaluated on Nvidia 1080 Ti. We re-implement the network in PyTorch (1080Ti), while RetinaNet is implemented in Caffe2 (M40). Note that speed numbers are reported for different GPU architectures, and thus should not be directly compared. Our network is very similar in inference settings to the original RetinaNet, so most speed performance gains are attributed to better framework implementation.

In Figure 1, we show our results compared with state-of-the-art single-shot and two-stage detectors [31, 25, 18].

| Method | Backbone | $AP^{bb}$ | $AP^{bb}_{50}$ | $AP^{bb}_{75}$ | $AP^{S}_{75}$ | $AP^{M}_{75}$ | $AP^{L}_{75}$ |
|---|---|---|---|---|---|---|---|
| *Two-Stage Detectors* | | | | | | | |
| Faster R-CNN+++ [19] | ResNet-101-C4 | 34.9 | 55.7 | 37.4 | 15.6 | 38.7 | 50.9 |
| Faster R-CNN w FPN [24] | ResNet-101-FPN | 36.2 | 59.1 | 39.0 | 18.2 | 39.0 | 48.2 |
| Faster R-CNN w RoIAlign [18] | ResNet-101-FPN | 37.3 | 59.6 | 40.3 | 19.8 | 40.2 | 48.8 |
| Mask R-CNN [18] | ResNet-101-FPN | 38.2 | 60.3 | 41.7 | 20.1 | 41.1 | 50.2 |
| *single-shot Detectors* | | | | | | | |
| YOLOv2 [30] | Darknet-19 | 21.6 | 44.0 | 19.2 | 5.0 | 22.4 | 35.5 |
| SSD513 [28, 13] | ResNet-101 | 31.2 | 50.4 | 33.3 | 10.2 | 34.5 | 49.8 |
| DSSD513 [13] | ResNet-101-DSSD | 33.2 | 53.3 | 35.2 | 13.0 | 35.4 | 51.1 |
| YOLOv3-608 [31] | Darknet-53 | 33.0 | 57.9 | 34.4 | 18.3 | 35.4 | 41.9 |
| RetinaNet [25] | ResNet-101-FPN | 39.1 | 59.1 | 42.3 | 21.8 | 42.7 | 50.2 |
| RetinaNet [25] | ResNeXt-101-FPN | 40.8 | 61.1 | 44.1 | 24.1 | 44.2 | 51.2 |
| RetinaMask | ResNet-50-FPN | 39.4 | 58.6 | 42.3 | 21.9 | 42.0 | 51.0 |
| RetinaMask | ResNet-101-FPN | 41.4 | 60.8 | 44.6 | 23.0 | 44.5 | 53.5 |
| RetinaMask | ResNet-101-FPN-GN | 41.7 | 61.7 | 45.0 | 23.5 | 44.7 | 52.8 |
| RetinaMask | ResNeXt-101-FPN-GN | **42.6** | **62.5** | **46.0** | **24.8** | **45.6** | **53.8** |

Table 5: Comparison with state-of-the-art methods on COCO `test-dev`. Compared to RetinaNet [25], our model based on ResNet-101-FPN is better by 2.6 mAP. Compared to Mask R-CNN [18], our model shows 3.5 mAP improvement.

Note that YOLOv3 [31] is trained with multi-scale training but ours and ReinaNet [25] are not. Our results show that the detector in RetinaMask has a higher envelop for accuracy-vs-time than RetinaNet when using ResNet-50 and ResNet-101 for the backbone model. All the numbers for Figure 1 can be found in Table 4. Our model shows 1.84 mAP and 1.52 mAP improvement on ResNet-50 and ResNet-101 compared to RetinaNet. Our detection result is better than the original numbers from Mask R-CNN and very close to recent implementation results.

Note that the speed of our implementation on the short side of 400 is surprisingly slow. We think this is an idiosyncrasy of the libraries used, and note that as with all the other resolutions we do see an improvement in accuracy.

### 4.3. Comparisons to the state-of-the-art methods

We use ResNet-50-FPN, ResNet-101-FPN, and ResNeXt32x8d-101-FPN [39] as the backbones in our final models. We train with the multi-scale {640, 800, 1200} and 2x iterations schedule. For the ResNet-101-FPN model, we also train a version using Group Normalization(GN) [38], which is applied only on the extra layers (FPN, localization, and classification). Replacing all the Batch Normalization [22] in ResNet-101 would cause a significant slowdown. The speed of ResNet-101-FPN-GN model is 0.158 s/im (compared to 0.145 s/im without GN). Using ResNeXt32x8d-101-FPN [39] as backbone further improves results by 0.9 mAP and achieves 42.6 mAP on COCO. We provide the quantitative comparison in Table 5 and show some detection examples in Figure 6.

We also acknowledge the recent new architectures for better object detection such as NASNET [42] or efficient networks (MobileNet [21, 33], ShuffleNet [41]), but their evaluation is beyond the scope of this work.

| Method | $AP^{m}$ | $AP^{m}_{50}$ | $AP^{m}_{75}$ | $AP^{bb}$ | $AP^{bb}_{50}$ | $AP^{bb}_{75}$ |
|---|---|---|---|---|---|---|
| *Mask R-CNN* | 36.7 | 59.5 | 38.9 | 39.6 | 61.5 | 43.2 |
| +update baseline | 37.0 | 59.7 | 39.0 | 40.5 | 63.0 | 43.7 |
| +e2e training | 37.6 | 60.4 | 39.9 | 41.7 | 64.1 | 45.2 |
| +ImageNet-5k | 38.6 | 61.7 | 40.9 | 42.7 | 65.1 | 46.6 |
| +train-time augm. | **39.2** | **62.5** | **41.6** | **43.5** | **65.9** | **47.2** |
| RetinaMask | 36.4 | 57.3 | 38.7 | 41.1 | 60.2 | 44.1 |

Table 6: Comparison with Mask R-CNN on mask prediction using ResNet-101 on COCO `minival`. The Mask R-CNN results are from Table 8 in the appendix of Mask R-CNN [18].

### 4.4. Comparison with Mask R-CNN on instance mask prediction

In Table 6, we compare our mask (instance segmentation) results to Mask R-CNN. The Mask R-CNN [18] results are from Table 8 of Mask R-CNN [18]. All the results are using ResNet-101 and Feature Pyramid Network [24] as the backbone model. Our models are trained in a very similar fashion to the *+e2e training* in [18]. Mask R-CNN still shows better accuracy on mask prediction, but the difference is only around 1.2 mAP.

## 5. Conclusion

In this work, we proposed three components to train a more accurate Single-Shot detector, RetinaMask. Our ablations show improvements for each module and our final model shows better accuracy without any network architecture change during inference. The proposed Self-Adjusting Smooth L1 loss can be used beyond the tasks of object detection and instance segmentation.
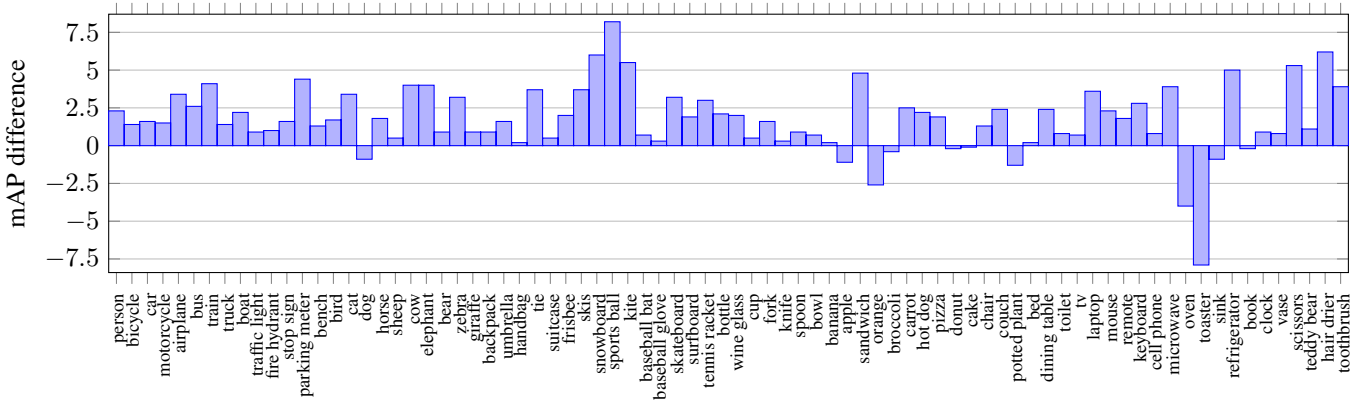
Figure 7: RetinaMask mAP detection improvement over RetinaNet baseline ResNet-50 backbone results are shown. mAP is computed across top 100 detections, and averaged for thresholds in range .50:.05:.95, according to COCO [26].
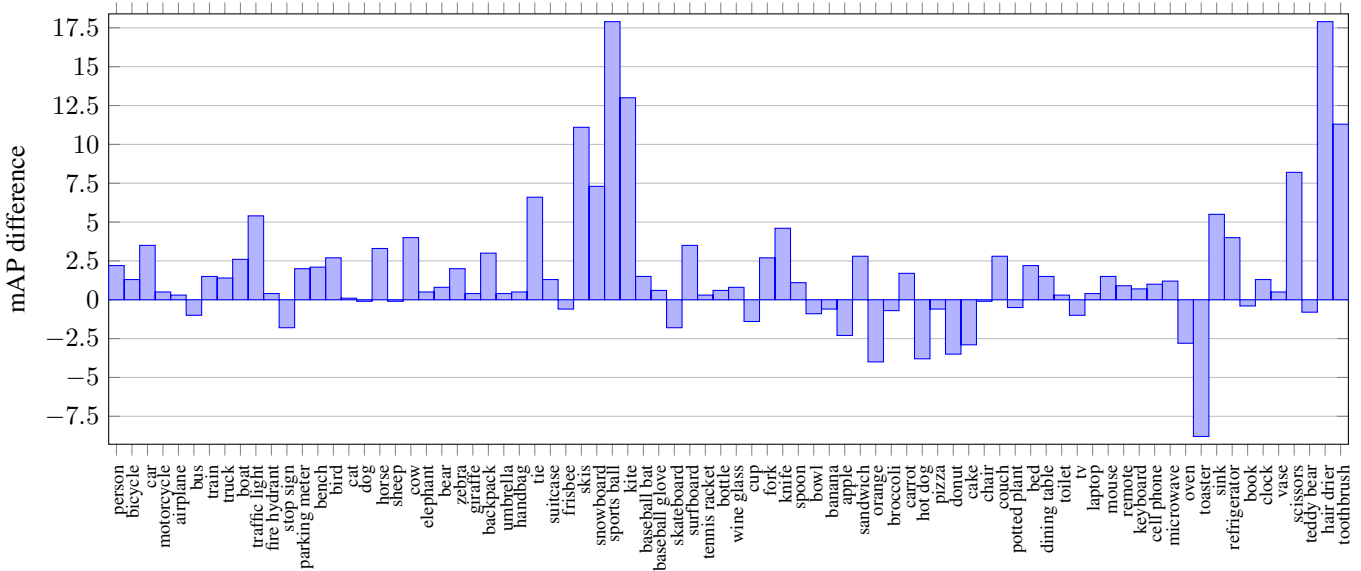


Figure 8: RetinaMask class-agnostic detection improvement over RetinaNet baseline. Localization errors are ignored by setting a low IoU threshold of 0.1, foreground object mis-classification is ignored as well. ResNet-50 backbone results are shown.

# References

[1] Detectron model zoo and baselines. https://github.com/facebookresearch/Detectron/blob/master/MODEL_ZOO.md, 2018. [Online; accessed 2018-11-16]. 1

[2] Faster R-CNN and Mask R-CNN in PyTorch 1.0: Model zoo and baselines. https://github.com/facebookresearch/maskrcnn-benchmark/blob/master/MODEL_ZOO.md, 2018. [Online; accessed 2018-11-16]. 1

[3] S. Bell, C. L. Zitnick, K. Bala, and R. Girshick. Inside-Outside Net: Detecting objects in context with skip pooling and recurrent neural networks. In *CVPR*, 2016. 5

[4] N. Bodla, B. Singh, R. Chellappa, and L. S. Davis. Soft-NMS – improving object detection with one line of code. In *ICCV*, 2017. 5

[5] Z. Cai, Q. Fan, R. Feris, and N. Vasconcelos. A unified multi-scale deep convolutional neural network for fast object detection. In *ECCV*, 2016. 2

[6] J. Dai, K. He, and J. Sun. Instance-aware semantic segmentation via multi-task network cascades. In *CVPR*, 2016. 3

[7] J. Dai, Y. Li, K. He, and J. Sun. R-FCN: Object detection via region-based fully convolutional networks. In *NIPS*, 2016. 2, 3

[8] N. Dalal and B. Triggs. Histograms of oriented gradients for human detection. In *CVPR*, 2005. 2

[9] N. Dvornik, K. Shmelkov, J. Mairal, and C. Schmid. BlitzNet: A real-time deep network for scene understanding. In *ICCV*, 2017. 3

[10] D. Erhan, C. Szegedy, A. Toshev, and D. Anguelov. Scalable object detection using deep neural networks. In *CVPR*, 2014. 2

[11] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL visual object classes (VOC) challenge. *International Journal of Computer Vision*, 2010. 6

[12] P. Felzenszwalb, R. Girshick, D. McAllester, and D. Ramanan. Object detection with discriminatively trained part based models. *Pattern Analysis and Machine Intelligence*, 2010. 2

[13] C.-Y. Fu, W. Liu, A. Ranga, A. Tyagi, and A. C. Berg. DSSD : Deconvolutional single shot detector. *arXiv:1701.06659*, 2017. 1, 2, 8

[14] R. Girshick. Fast R-CNN. In *ICCV*, 2015. 2, 3

[15] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *CVPR*, 2014. 2

[16] P. Goyal, P. Dollár, R. Girshick, P. Noordhuis, L. Wesolowski, A. Kyrola, A. Tulloch, Y. Jia, and K. He. Accurate, large minibatch sgd: Training imagenet in 1 hour. *arXiv:1706.02677*, 2017. 4

[17] R. A. Güler, N. Neverova, and I. Kokkinos. DensePose: Dense human pose estimation in the wild. In *CVPR*, 2018. 1

[18] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask R-CNN. In *ICCV*, 2017. 1, 3, 8

[19] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 8

[20] D. Hoiem, Y. Chodpathumwan, and Q. Dai. Diagnosing error in object detectors. In *ECCV*, 2012. 7

[21] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam. MobileNets: Efficient convolutional neural networks for mobile vision applications. *arXiv:1704.04861*, 2017. 8

[22] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015. 8

[23] S. R. Kaiming He, Xiangyu Zhang and J. Sun. Spatial pyramid pooling in deep convolutional networks for visual recognition. In *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2015. 2

[24] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie. Feature pyramid networks for object detection. In *CVPR*, 2017. 2, 4, 8

[25] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár. Focal loss for dense object detection. In *ICCV*, 2017. 1, 2, 4, 6, 7, 8

[26] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollr, and C. L. Zitnick. Microsoft COCO: Common objects in context. In *ECCV*, 2014. 3, 5, 6, 9

[27] S. Liu, L. Qi, H. Qin, J. Shi, and J. Jia. Path aggregation network for instance segmentation. In *CVPR*, 2018. 3

[28] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg. SSD: Single shot multibox detector. In *ECCV*, 2016. 1, 2, 3, 8

[29] J. Redmon, S. K. Divvala, R. B. Girshick, and A. Farhadi. You only look once: Unified, real-time object detection. In *CVPR*, 2016. 1, 2

[30] J. Redmon and A. Farhadi. YOLO9000: better, faster, stronger. In *CVPR*, 2017. 1, 2, 8

[31] J. Redmon and A. Farhadi. YOLOv3: An incremental improvement. *arXiv:1804.02767*, 2018. 1, 8

[32] S. Ren, K. He, R. Girshick, and J. Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In *NIPS*, 2015. 2, 3

[33] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen. MobileNetV2: Inverted residuals and linear bottlenecks. In *CVPR*, 2018. 8

[34] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and Y. Lecun. OverFeat: Integrated recognition, localization and detection using convolutional networks. In *ICLR*, 2014. 2

[35] A. Shrivastava, A. Gupta, and R. Girshick. Training region-based object detectors with online hard example mining. In *CVPR*, 2016. 2

[36] A. Shrivastava, R. Sukthankar, J. Malik, and A. Gupta. Beyond skip connections: Top-down modulation for object detection. *arXiv:1612.06851*, 2016. 2

[37] J. R. R. Uijlings, K. E. A. van de Sande, T. Gevers, and A. W. M. Smeulders. Selective search for object recognition. *International Journal of Computer Vision*, 2013. 2

[38] Y. Wu and K. He. Group normalization. In *ECCV*, 2018. 8

[39] S. Xie, R. Girshick, P. Dollr, Z. Tu, and K. He. Aggregated residual transformations for deep neural networks. In *CVPR*, 2017. 8

[40] J. D. X. J. Yi Li, Haozhi Qi and Y. Wei. Fully convolutional instance-aware semantic segmentation. In *CVPR*, 2017. 3

[41] X. Zhang, X. Zhou, M. Lin, and J. Sun. ShuffleNet: An extremely efficient convolutional neural network for mobile devices. In *CVPR*, 2018. 8

[42] B. Zoph, V. Vasudevan, J. Shlens, and Q. V. Le. Learning transferable architectures for scalable image recognition. In *CVPR*, 2018. 8