

Anchor free object detection with mask attention

HE Yang

Shanghai University

Beibei Fan (✉ fanbeibei@shu.edu.cn)

Shanghai University

Ling ling Guo

Shanghai University

Research

Keywords: object detection, anchor-free, convolutional neural network, attention mechanism

Posted Date: June 4th, 2020

DOI: <https://doi.org/10.21203/rs.2.24760/v2>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published on July 14th, 2020. See the published version at <https://doi.org/10.1186/s13640-020-00517-3>.

1 Anchor free object detection with mask attention

2 He Yang , Beibei Fan *, Lingling Guo

3 School of Mechatronic Engineering and Automation, Shanghai University, Shanghai, 200444, China;

4 * Correspondence: fanbeibei@shu.edu.cn

5 **Abstract:** The anchor-free method based on key point detection has made great progress. However, the anchor-free method is
6 too dependent on using a convolutional network to generate a rough heat map. This is difficult to detect for objects with a large
7 size variation and dense and overlapping objects. To solve this problem, first, we propose a mask attention mechanism for object
8 detection methods. And make full use of the advantages of the attention mechanism to improve the accuracy of network detection
9 heat map generation. Then, we designed an optimized fire model to reduce the size of the model. The fire model is an extension
10 of grouped convolution. The fire model allows each group of convolutional network features to learn the same feature through
11 purposeful grouping. In this paper, the mask attention mechanism uses object segmentation images to guide the generation of
12 corner heat maps. Our approach achieved an accuracy of 91.84% and a recall 89.83% in the Tencent-100K dataset. Compared
13 with the popular object detection methods, the proposed method has advantages in model size and accuracy.

14 **Keywords:** object detection; anchor-free; convolutional neural network; attention mechanism

15 1. Introduction

16 In recent years, applications based on object detection technology have become more widespread[1]. Common applications
17 such as pedestrian detection[2], vehicle detection[3], image retrieval[4], and traffic sign detection[5]. This imposes higher
18 requirements on the detection performance and size of the object detection method. As a typical method of object detection,
19 anchor-based object detection can often be seen in common application scenarios. Anchor-based object detection methods have
20 made great progress in location and recognition performance[6,7]. These anchor-based methods regress the boundaries of objects
21 by generating dense coordinates on the feature map. In order to get a higher quality coordinate box, non-maximum suppression
22 method is used to filter to most of the overlap bounding boxes. However, computing anchor boxes and filter boxes requires a
23 large amount of computing resources. At the same time, the pre-designed anchor box is not aligned with the real boundary of the
24 object. As a result, anchor-based has defects in calculation and accuracy.

25 To solve the anchor box problem in the anchor box, the first step is to prioritize the limitations of manually setting the
26 anchor box. Recently, many anchor box-free methods have emerged such as CornerNet and CenterNet[8,9]. The early
27 DenseBox[10] method directly used landmark heatmaps and face score map to predict human face. It used the combination of
28 key points and heatmaps. The effect of box regression using only convolutional networks is not obvious, leading to reduced
29 accuracy and recall. YOLO divides the feature map into 7x7 squares. Object detection is carried out in each square. Because the
30 scale of the object varies too much, the network is difficult to learn. As a result, YOLO[6] mAP is low and recall is low. In order
31 to enhance the detection performance of small objects, YOLO-v2[11] and YOLO-v3[12] join the anchor-boxed side. CornerNet
32 transforms the position detection of an object into the detection of the key points of the top-left corner and the bottom-right
33 corner of the boundary box. This method based on key point detection simplifies the process of object detection and greatly
34 simplifies the task of convolution network.

35 CornerNet simplifies the process of object detection by using key point detection technology. Because CornerNet is based
36 on Hourglass network[13], the Hourglass network is redundantly calculated, which makes it difficult to train in the absence of
37 computing resources. The optimized version of CornerNet, CornerNet-Squeeze[14], solves this problem to some extent.
38 However, the detection performance of CornerNet-Squeeze is unsatisfactory due to the excessive pursuit of compression and
39 simplification and the lack of consideration of global spatial information. CornerNet-Saccade[14] is designed to increase the
40 network global attention to image. Another optimization strategy uses a similar two-stage object detection method. The method
41 guides the cropped image by the attention mechanism, and uses the cropped image to detect to improve efficiency. This method
42 is not contiguous in backpropagation, which is detrimental to the overall weight update.

43 Inspired by CornerNet-Squeeze and CornerNet-Saccade, we propose a Mask-CornerNet. In order to make the convolutional
44 network pay attention to the region information of objects contained in the space, this paper proposes a mask module. The mask
45 module is similar to the method of segmentation model. It is different from other object segmentation methods. The results of
46 the mask module in this paper are fed back to the corner detection network. The mask module can constrain the network to pay
47 more attention to the area containing the object. The recognition results of our method are shown in Figure 1. The SGE
48 module[15] gets inspiration from grouped convolution. It can greatly improve network performance by introducing few
49 parameters, which is very helpful for enhancing channel information. We optimized the fire model with this idea and designed
50 a new fire model.

51 In this paper, we optimized the design of CornerNet-Squeeze. While improving the detection accuracy, the calculation
52 efficiency is maintained. Improve detection accuracy by adding a mask branch. Figure 6 shows some of the results of our method.
53 Our method is compatible with speed and accuracy. The main contributions are as follows:

54 1. We propose a fire model module with fewer parameters and higher efficiency. In the fire model, we take advantage of
55 the group convolution to fully exploit the local and global relationships of the convolution channel.

56 2. We propose to add a mask module. **The mask module has the ability to segment objects. It can use the object segmentation
57 feature map to assist in generating the corner heat map, which effectively enhances the expression ability of the convolution
58 feature.**

59 3. Our approach takes advantage of the features of the Hourglass Network. We use the fusion method of context features
60 to improve the ability to express features.

61 **The rest of this paper will be organized as follows.** Section 2 mainly analyzes the related work of current object detection.
62 The main work description of our work is in Section 3. Section 4 shows experimental results. Finally, we summarize the work
63 and give our conclusion.

64 2. Related Work

65 2.1 Anchor-based Detectors

66 Since the advent of the R-CNN series of object detection methods based on convolutional networks, anchor-based detectors
67 have become very popular. Today, there are many advanced object detection methods such as Faster-rcnn[7], SSD[16], FPN[17],
68 RetinaNet[18], Mask-rcnn[19]. These methods benefit from the clever design of the anchor, and the detection performance is
69 steadily increasing. The R-CNN approach opens the era of anchor-based methods. R-FCN[20], FPN, Mask-RCNN are typical
70 RPN methods. At the same time, there are methods to improve the expression ability of model features by optimizing the shape
71 and size of the convolution kernel. For example Atrous convolution[21], Depthwise separable convolutions[22], Deformable
72 Convolutional[23]. Deformable convolutional networks is also an innovative deep learning optimization direction, breaking
73 through the innovation direction of traditional convolutional network. Recently, it's best to use multiple receptive field branches
74 and take a weight-sharing approach to achieve the best results based on Anchor-based[24].

75 **The anchor-based method constrains the regression results by proposing dense anchor points on the image. This method
76 simplifies the learning task of convolutional networks, but makes the detection method less flexible.** Thus, this results in a
77 convolutional network predicting regression coordinate size limits, fixed shape, fixed aspect ratio, and small plasticity. The
78 anchor box is always not aligned with the regression results, resulting in a small intersection over union (IOU) value of the box.
79 The anchor aspect ratio, size and shape in reality are highly uncertain. A manually set anchor cannot match it. In the training
80 phase, the regression task needs to be performed by comparing each result with the IOU of anchor. The time cost and space
81 consumption are high. **Based on the above reasons, this article uses anchor-free object detection method. This method has no
82 fixed anchor constraints and can more flexibly match objects of various proportions.**

83 2.2 Anchor-free Detectors

84 Anchor-free method is another popular method, YOLO has been considered as the representative of anchor-free object
85 detection. YOLO method was to get rid of the R-CNN series of methods for two-stage inference. YOLO coordinates and
86 classifies objects directly from the image. This method of directly returning to the frame coordinates without using an anchor
87 point greatly improves the speed of object detection. However, this rough detection method results in low accuracy and inability
88 to detect small objects. DenseBox is considered to be the earliest anchor-free method[25]. DenseBox uses five heatmaps to
89 classify faces and two corner coordinates. Because of the defects of DenseBox in the overlapping processing of heatmaps, the
90 recall rate of DenseBox in general object detection is low.

91 CornerNet removes the anchor box setting and uses the top-left corner and the bottom-right corner to represent the
92 boundaries of a box[26]. Main task of CornerNet is to predict the position of two corners, and classify tasks are distributed in
93 corner detection tasks. In order to improve the recall rate, corner pool is used to maximize each other at the two boundaries in
94 order to improve the probability of the occurrence of objects. Inspired by Cornernet, CenterNet uses the midpoint of the object
95 as the detection center[8]. It can detect the object by adding the boundary and size information of the object. FCOS[27] uses
96 Center-ness to suppress excessive and inappropriate box boundaries, and uses multiple and classified objects to classify.
97 ExtremeNet[28] uses standard key points to estimate four poles (top, left, bottom, right) and a center point of the network
98 detection object. **Compared with other anchor-free methods, CornerNet makes full use of the edge information of objects. In
99 particular, the corner embedding vector makes the pair of corners clear.**

100 2.3 Object Segmentation

101 **The performance of object segmentation methods based on convolutional networks is increasingly higher.** FCN[29] directly
102 uses a convolutional network to make binary predictions for each category. **The FCN method uses a convolutional neural network**

103 to extract the semantic features of the image. This segmentation method gets rid of traditional color and shape features, thereby
 104 reducing the difficulty of image segmentation. Mask-rcnn directly adds the segmentation task to the object recognition task,
 105 which provides a new idea for object detection. After that, U-Net[30] and V-Net[31] encoder-decoder networks made full use of
 106 skip structures for feature fusion, and made great progress in the fields of medicine and 3D detection. Our approach is different
 107 from the above. We use the results of the segmentation module to guide the convolution to generate new features.



108
 109 **Fig 1.** Object detection results of some real traffic scenes. Our system works very well in complex scene where objects are small
 110 and highly occluded.

111 3. The Proposed Method

112 Our approach is based on CornerNet-Squeeze. This method uses corner points to represent the bounding box of the object.
 113 Use the heat map to directly return coordinates. In order to eliminate the deviation of the coordinates, use offset to correct. By
 114 analyzing the network structure and overall method of CornerNet, it is found that the heat map is the most critical point of the
 115 whole method, especially the key position and the value of the heat map value. But this feature is not really utilized in the entire
 116 network model. And the detection image of the corner point is completely dependent on the heat map. The accuracy of the corner
 117 heat map depends entirely on the training data and the feature extraction ability of the convolutional neural network. However,
 118 it is more difficult to completely rely on the convolutional network to generate an accurate corner heat map. Object segmentation
 119 can highlight the area containing the target. Inspired by this, this paper uses a masking module to enhance the expression of
 120 object. The main function of the mask module is to generate a complete segmented image based on the original image. The
 121 segmentation result image is fused with the internal feature map of the convolutional network to generate enhanced features.

122 3.1 Base Framework

123 This section describes the overall framework and process of the proposed method. Figure 2 illustrates the overall framework
 124 of our approach. Our method uses two hourglass network structures in tandem, which effectively utilizes the contextual features
 125 of the image. Based on the symmetrical nature of the hourglass network, the backbone network is elegant. The hourglass network
 126 is followed by the mask module. Our mask model is more than just an additional separate module, it is a feature-enhanced
 127 function module. The mask module results are directly fed back to the original feature map to constrain the feature map.

128 The mask module focuses on the critical areas containing objects by constraining the convolutional network, thereby
 129 improving the detection accuracy of the network. In order to fully utilize and mine the contextual features of the feature map, we
 130 combine the outputs of the two Hourglass Network elements by element.

131 The Squeeze method can sufficiently reduce the number of channels in the convolutional network. Because the rough
 132 grouping method does not highlight the relationship between the channels, the relationship between the groups cannot be directly
 133 mined. The Spatial Grouping Enhancement (SGE) module can generate an attention factor for each spatial position in each
 134 convolutional semantic group[32]. Therefore, we adjust the importance of each channel group function through SGE so that each
 135 individual group can autonomously enhance its learned expression and suppress possible noise. With almost no additional
 136 learning parameters added, we optimized the fire module. The advantages of group convolution can reduce the computational
 137 and parametric advantages. By separately performing local and global feature similarity for each set of convolutions, the spatial
 138 distribution of semantic features is enhanced to achieve a clear division of labor for each channel.

139

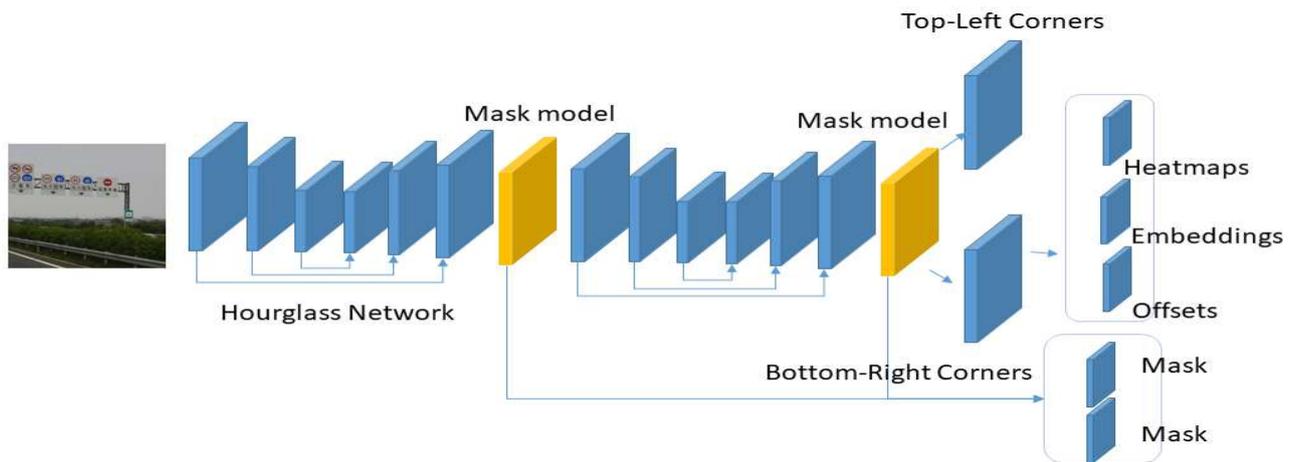


Fig 2. Describe the proposed method. By stacking multiple Hourglass Network and mask modules to enhance the network extraction ability, the mask module can get the image segmentation result. The corner detection module has two parts, and the top left and bottom right corners of the object are obtained

The coordinate regression and classification of the network is realized by detecting the top-left and bottom-right. The final detection of the network includes two branches, Top-left and bottom-right detection modules. Detection modules include heatmaps, offsets, and embeddings.

The heatmap contains C channels (C is the target category, no background category), and each channel is a binary mask indicating the angular position of the corresponding category. For each corner point, there is only one ground-truth, and the other positions are negative samples. During the training process, the model reduces the negative samples and sets a positive sample in the radius r region at each ground-truth corner. This is because the vertices falling within the radius r region can still generate valid boundary of objects.

Offsets is a correction module that is specifically set to correct errors in the heatmap corners. Due to the presence of down-sampling, the heatmap size generated by the model is smaller than the original input image. After n times of down-sampling, the position of the point (i, j) mapped onto the heat map on the original image becomes $(\lfloor x/n \rfloor, \lfloor y/n \rfloor)$. Re-mapping the points on the heatmaps to the original image has a quantization map error, resulting in an offset in the corner position of the map. This offset can seriously affect small object IOU calculations. The correction of the detection result is obtained by adding the heat map prediction result and the offset.

Embedding vectors are the key to corner detection. By encoding the corner points, the relationship of the embedded vectors in the top left and bottom right corners is obtained. The distance of the embedding vector of two corner points is used in the article to determine whether it is a pair of corner points of the same object.

3.2 Fire module

The use of the group convolution method can reduce the use of parameters, but can not improve the feature representation of the convolutional network. In order to enhance the feature representation of the model, we absorbed the inspiration of the SEG module[32]. The fire module is shown in Figure 3. Each SGE module increases the number of parameters by about 2 times the number of groups, and the number of groups is usually 32 or 64.

The fire module groups features on channels. Group convolution divides a channel into multiple sub-functions to represent different semantics. For a particular set of semantics, it is reasonable and beneficial to generate corresponding semantic features at the correct spatial location of the original image. In the fire module of this paper, first, squeeze the internal convolution feature maps of each group to obtain the global semantic parameters of each feature map, such as equation (1). Multiply the semantic parameters and the original features element-by-element, to obtain the feature map c_i of each channel. Then normalize the feature map c_i in each group. The normalization of c_i introduces only two learning parameters, allowing SEG to automatically normalize. The activation function is used to assign weight to feature a_i , and finally a_i is multiplied element-by-element with the original feature.

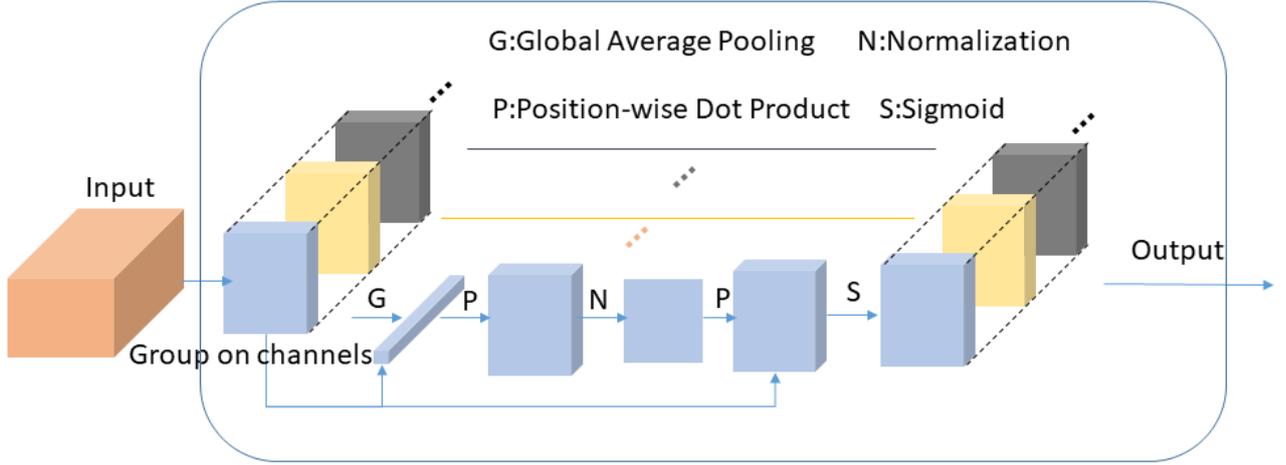


Fig 3. Diagram our improved fire module, which combines global and local features by grouping for incoming feature maps.

First, the features are grouped. Spatial averaging function F_{gp} can get the global semantic feature g . This allows you to quickly get the global semantics of each grouping.

$$g = F_{gp}(x) = \frac{1}{m} \sum_{i=1}^m x_i \quad (1)$$

In the convolution integral group, the global feature of each channel is multiplied by the original feature point to obtain the initial attention feature c_i .

$$c_i = g \cdot x_i \quad (2)$$

The mean μ_c within the group is subtracted from each group and divided by the variance σ_c within the group. The two scaling offset parameters allow the normalize operation to be restored, then the sigmoid to get the final attention mask and scale the feature at each location in the original feature group.

$$\hat{c}_i = \frac{c_i - \mu_c}{\sigma_c}, \quad \mu_c = \frac{1}{m} \sum_j c_j, \quad \sigma_c^2 = \frac{1}{m} \sum_j (c_j - \mu_c)^2 \quad (3)$$

Similar to BN, in order to enhance standardization versatility. Add two additional parameters (λ, β) to represent the standard scale and offset. This is the two sets of non-convolution parameters that need to be learned in the fire model.

$$a_i = \lambda \hat{c}_i + \beta \quad (4)$$

In order to obtain the enhanced feature \hat{x}_i , the weight distribution feature map a_i of each channel is multiplied element-wise by the original feature map.

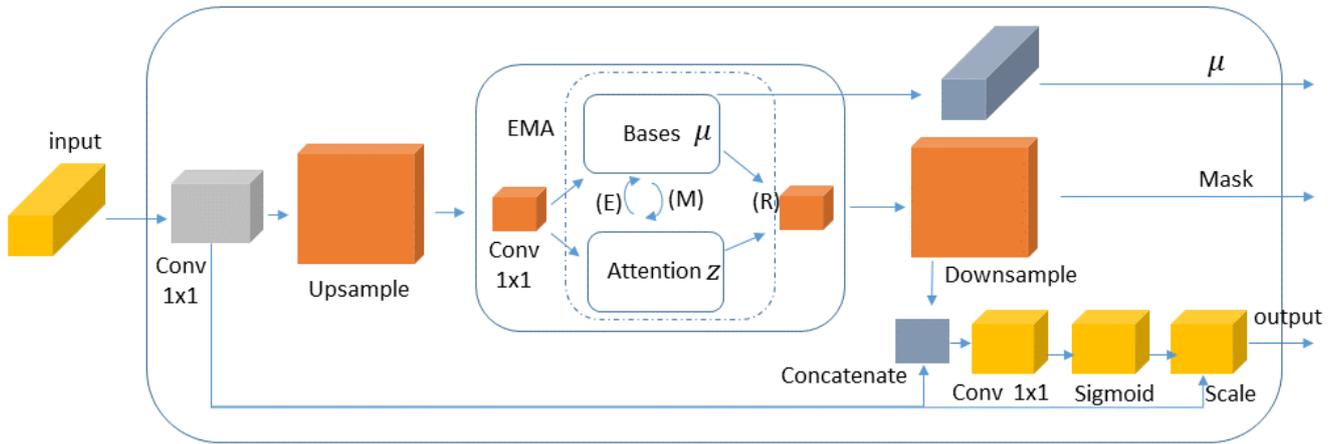
$$\hat{x}_i = x_i \cdot \sigma(a_i) \quad (5)$$

The Spatial averaging function is considered to be a global representation of the channel scale. Each group of feature maps captures a specific semantics during the learning process. We use the global average feature to represent the learning characteristics of each group. Enhanced features can be obtained by using global features and local features for fusion. Our improved fire model just learn two parameters (λ, β) and adjust the mask for feature enhancement.

3.3 Mask module

Our mask is inspired by image segmentation. Since CornerNet relies heavily on the generation of heatmaps, it is difficult to make predictions directly from the heat map. We add a mask module to generate the split image directly. As described in Figure 4. To achieve spatially enhanced features, the results of using the mask module are merged with the feature map. The mask constraint map increases the level of attention of the target object area. To simplify the process of feature extraction, we use a similar expectation-maximization algorithm to get the mask. This method uses the strategy of Maximum Likelihood Estimation to get the mask. The internal features of the mask have the characteristics of low rank. This allows the mask to reduce the internal differences of the categories while maintaining the differences between the categories. The mask model rearranges the original features. The spatial distribution will change to facilitate heat map detection.

In order to ensure the simplicity of the original method, the mask module cannot be complicated. The Maximum Attention Mechanism (EMA) method abandoned the process of calculating the attention map on the full map[33]. The hidden variable is calculated iteratively using the expectation-maximization algorithm. And the attention mechanism is run on the hidden variable, which greatly reduces the complexity of the algorithm.



207
208
209
210
211
212
213
214
215
216
217

Fig 4. Show the details of the mask module. Convolutional features is upsampled by multiple convolution kernels. Using the EMA module to estimate the mask, the mask module can get the segmentation result of the object. Aggregating the mask module can generate spatial constraints on the object and enhance local selection of features.

Inside the mask module, a 1x1 convolution kernel is used to squeeze the input feature channel to reduce the amount of calculation. Squeezing the input image channel can effectively use resources. Then upsample the feature map to the output size. The output size is half of the original image. Use EMA to calculate the feature map to get the output mask and bases parameter μ . The converged bases parameter μ and latent variables Z can be reconstructed output mask. Finally, the out mask and the original feature map are multiplied by the group elements to obtain the final feature. Out mask is the result of instance segmentation. Figure 5 is the result of object segmentation detection on some data.



218
219
220

Fig 5. The data is displayed, the left side is the mark for coordinate positioning, and the right side is the label for image segmentation.

221 The EMA module is composed of three parts: responsibility estimation(E), likelihood maximization (M) and data re-
 222 estimation (R). E and M is the E step and the M step of the EM algorithm. The size of the feature map X is $C \times H \times W$. To simplify
 223 the symbols, reshape X to $N \times C$, where $N = H \times W$. Then, feature map is $X \in R^{N \times C}$, the base initial value is $\mu \in R^{K \times C}$, and
 224 $Z \in R^{N \times K}$ is a hidden variable. A_e estimates Z , and A_m is used to update μ . A_e and A_m alternately perform the T step to obtain
 225 an approximate estimate X of the feature. We set T to the default 3. Ar to get the features \tilde{x} through μ and Z .

226 Step E calculates the posterior distribution of the hidden variable Z .

$$Z^{(t)} = \text{softmax}(\lambda X (\mu^{(t-1)})^T) \quad (6)$$

227 Step M updates μ by maximizing the likelihood function.

$$\mu_k^{(t)} = \frac{z_{nk}^{(t)} \cdot x_n}{\sum_{m=1}^N z_{mk}^{(t)}} \quad (7)$$

228 After E and M alternately perform the T step, μ and Z are used to reconstruct the feature map.

$$\tilde{x} = Z^T \cdot \mu^T \quad (8)$$

230 3.4 Training

231 L_{det} is the loss function for the regression of two diagonal corners of an object. L_{det} indicates that focal losses are used to
 232 constrain the position of the corner points. The category of the object is predicted by a non-standard Gaussian heat maps. Locate
 233 a pair of coordinate points of the object through the top left and bottom right heat maps.

234 L_{mask} is the loss function of the mask model. The mask module is an instance segmentation network, which is trained with
 235 a multi-class cross entropy loss function. The parameters of the mask module are updated using image segmentation.

236 L_{push} and L_{pull} are used to constrain the correlation between the corners of the same object and different objects. L_{pull} is
 237 used to constrain the distance between a pair of corner points of the same object. The distance between the corner points of the
 238 same object is the smallest. L_{push} is used to constrain the corner points of different objects, and the corner points of different
 239 objects are maximized.

240 L_{off} is used to compensate and correct the predicted and true value deviations, constrained by Smooth L1 Loss. Due to the
 241 large amount of down-sampling used in the convolutional network, the final coordinate regression and the coordinates in the
 242 original image are offset. Here we again predict an offset value to adjust the position of the angle. Here we use the smooth L1
 243 Loss function as a penalty function.

244 The loss function for joint training for all tasks is as follows. $\alpha, \beta, \lambda, \gamma$ denote the weight of each task, the set value is 0.1,
 245 0.1, 1,1.

$$L = L_{\text{det}} + \alpha L_{\text{pull}} + \beta L_{\text{push}} + \lambda L_{\text{off}} + \gamma L_{\text{mask}} \quad (9)$$

246 4. Experimental results and discussions

247 4.1 Datasets and Implementation Details

248 In this section, the experimental process and materials are detailed. We validated our method using the Tsinghua-Tencent
 249 100K dataset[34]. The data set contains more than 10,000 images, of which there are 6,000 training images and 3,000
 250 verification images. Three scales are set for the object scale, small, medium and large. The size of each object is [0-30], [30-96],
 251 [96-400].

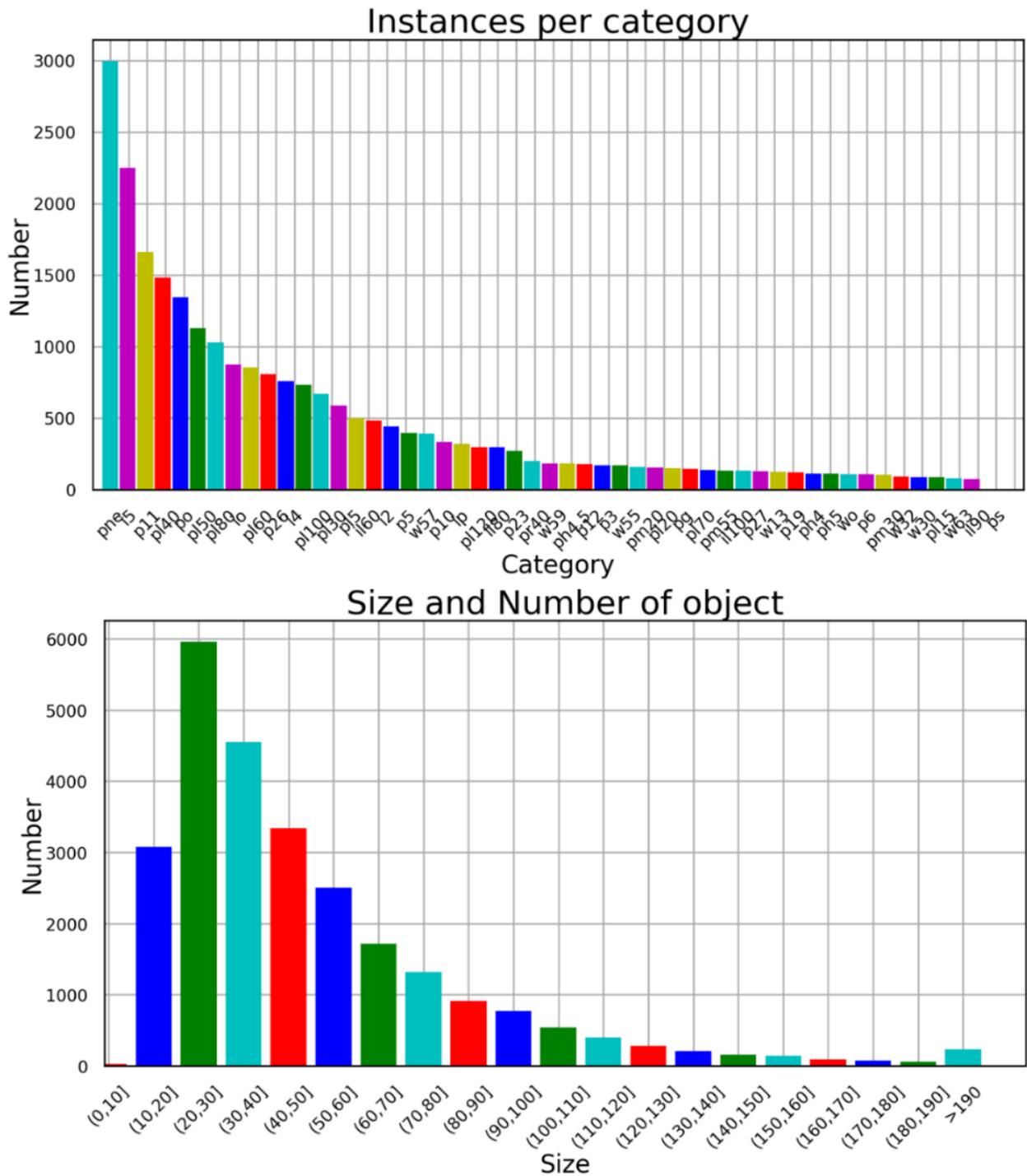


Fig 6. Object size and category statistics

We analyzed the number of labeled objects in multiple small areas and the number of objects in each category. Detailed statistics are described in Figure 6. The statistics of the number of individual categories are shown in the upper part of the figure. There are a total of 182 tag categories in the data set. Only the categories with more than 100 labels are shown in the figure. The lower part of the figure is the statistics of the size of the label object. Objects with size <50 were found to account for 64.3% of all marked objects, while size ≥ 90 was only 11.3%. It can be clearly seen from the figure that small objects account for a large proportion in the data set. This makes it very challenging to detect this data set. This is very consistent with the actual car driving image. To reduce the intensity of training, we used the first 42 categories, which accounted for 89.5% of all tag categories.

Our method is implemented using Pytorch. The entire network parameters are randomly initialized. The size of the network input image is 512x512. The size of the heatmap output is 128x128. The output size of the mask is 256x256. We set both α and β to 0.1 and γ to 1. Training Mask-CornerNet using batch size=13, learning_rate=0.00025, the learning rate per 180,000 iteration is multiplied by 0.1. Use adam optimizes as the optimizer for the network. All experiments were verified on a NVIDIA GTX1080Ti GPU with 32GB memory workstation. In the model evaluation phase, select the top 100 top left corner and the top 100 bottom left corner in the heat map.

252
253
254
255
256
257
258
259
260
261
262
263
264
265
266

4.2 Results

In order to verify our method from multiple aspects, we compared the methods commonly used today, such as traffic sign detection methods based on real environment[34], using pyramidal convolutional networks to detect traffic signs[5], Faster-rcnn[35], yolov3[12], CornerNet[26], CornerNet-Squeeze and CornerNet-Saccade[14]. At the same time, in order to verify the effectiveness of each module from multiple sizes. Set multiple single optimized models based on the baseline method. CornerNet-fusion merges multiple features. CornerNet-fire configures the fire module of this article. CornerNet-mask contains the mask module. Mask-CornerNet is the method proposed in this paper. All reference object detection methods use the same data enhancement method. We study these object methods from three object scales: large, medium and small. Experimental results verify that the method in this paper has great advantages on multiple object scales.

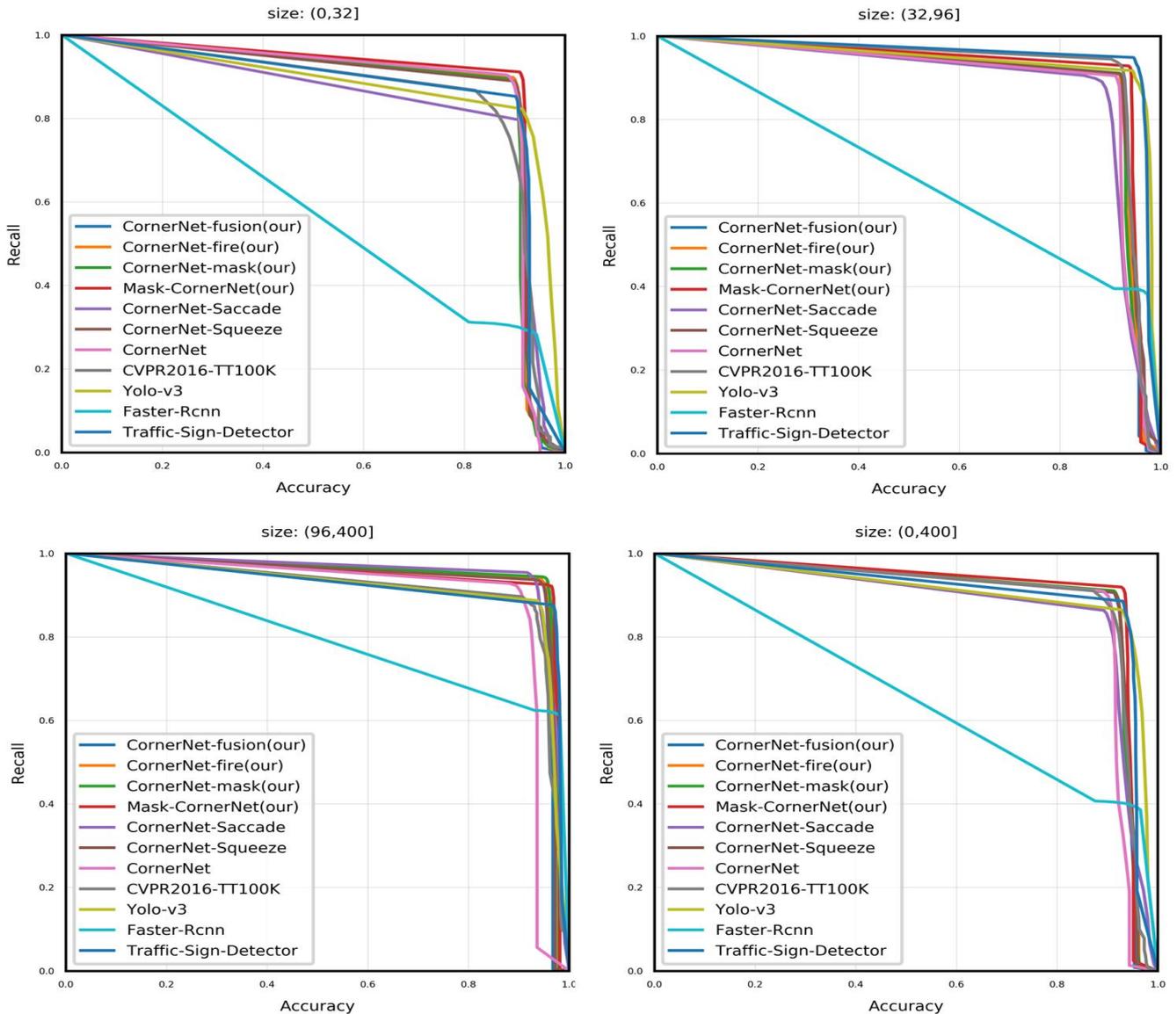


Fig 7. The performance of the proposed method is on the three scales of large, medium and small.

The accuracy-recall curves of our method on multiple scales are shown in Figure 7. Our Mask-CornerNet method performs well on small size and comprehensive size objects. The proposed method has an accuracy rate of 90.11% for small sizes and a recall rate of 88.42%. The accuracy of baseline in small-sized objects is 90.74%, and the recall rate is 86.33%. Pyramid convolutional networks do not have much advantage on small-sized objects. The accuracy and recall rate are 90.23% and 85.25% respectively. Yolov3 detects multiple feature maps in small size with an accuracy rate and recall rate of 91.51% and 82.22%.

From the perspective of full size, the overall size accuracy and recall rate of our method are 91.63% and 89.83%. CornerNet has an accuracy of 90.29% and a recall of 89.22%. Conernet-Saccade has an accuracy of 90.33% at all scales and a recall of 84.14%. Our method improves the accuracy of comprehensive size detection by increasing the detection accuracy of small sizes.

4.3 Ablation study

We compare experiments by adding a single module to explore how much improvement each module can bring to object detection. All experiments were performed under the same conditions, environment and hyper-parameters. This article uses the recall rate (R), accuracy rate (A) and mean average precision (mAP) as the evaluation criteria of the model. In this paper, mAP is the area enclosed by the Precision-recall curve and the X axis. To take full advantage of contextual features, we use the output of a fusion of two Hourglass networks. As can be seen from Table 1, the use of fusion features can significantly improve the network recall rate. The recall rate increased by 0.28%, while the accuracy rate decreased. After adding the fire model, it increased by 0.67% on the basis of the baseline. After adding the mask module, the accuracy of the model has increased by 0.71%. From the experimental results, all the added modules can improve the performance of object classification detection, thereby promoting the accuracy of the model. From mAP, our model's overall performance is optimized. The mAP of the proposed model is 90.01% vs. 89.84% of the baseline.

Table 1. Results of ablation experiments of various components

Baseline	Y	-	-	Y
+ fusion	-	Y	Y	Y
+ fire model	-	-	Y	Y
+ mask	-	-	-	Y
Accuracy	92.13	91.80	92.06	91.63
Recall	88.79	89.07	89.46	89.83
mAP	89.84	89.72	89.91	90.01

Our final model uses all the lifting methods, fuses contextual features, and optimizes the fire module to improve the segmentation network's ability to express features. Finally, on the baseline, we increased the recall rate by 1.04% with less cost.

Table 2. Comparison of the output results of the middle part features

Feature Map	Stacks_1	Stacks_2	Stacks_1+Stacks_2
Accuracy	92.95	91.72	91.63
Recall	82.34	89.76	89.83

In order to explore the specific influence of context on model feature expression. We perform object detection directly on the feature output of each Hourglass networks structure. The first layer is stack-1, the second layer is stack-2, and the fusion feature is stack-1 + stack-2. Table 2 shows our experimental results. The recall rate using the stack-1 model alone was 82.34%. The feature recall rate using fusion was 89.83%. Able to increase by 7.49% on this basis. The tendency bias of recall rate and quasi-curvature was clearly balanced.

4.4 The effect of the number of proposals

Since Mask-Cornernet uses the top left and bottom right to detect objects separately. Generally, the top 50-200 corner points with a score greater than the threshold are used. Use embedding vector similarity to match a pair of corners. We explored the influence of the number of corner points on the performance of object detection. Our experimental results are shown in Figure 8. By adjusting the number of detection points in the experiment, it was found that the number of detection points did not affect the coordinate regression. This eliminates the impact of detection points on our algorithm.

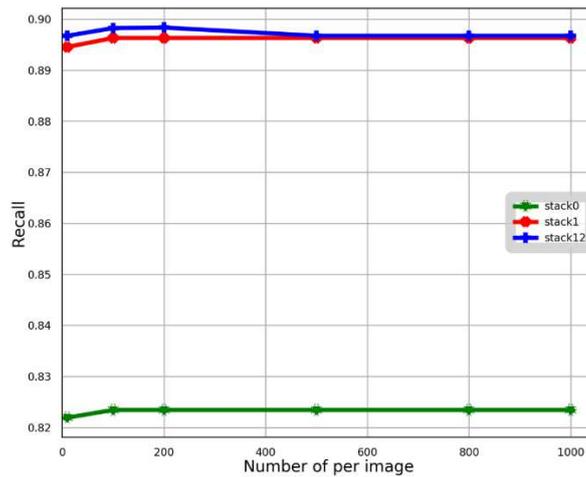


Fig 8. Diagram the effect of the recommended number on the recall

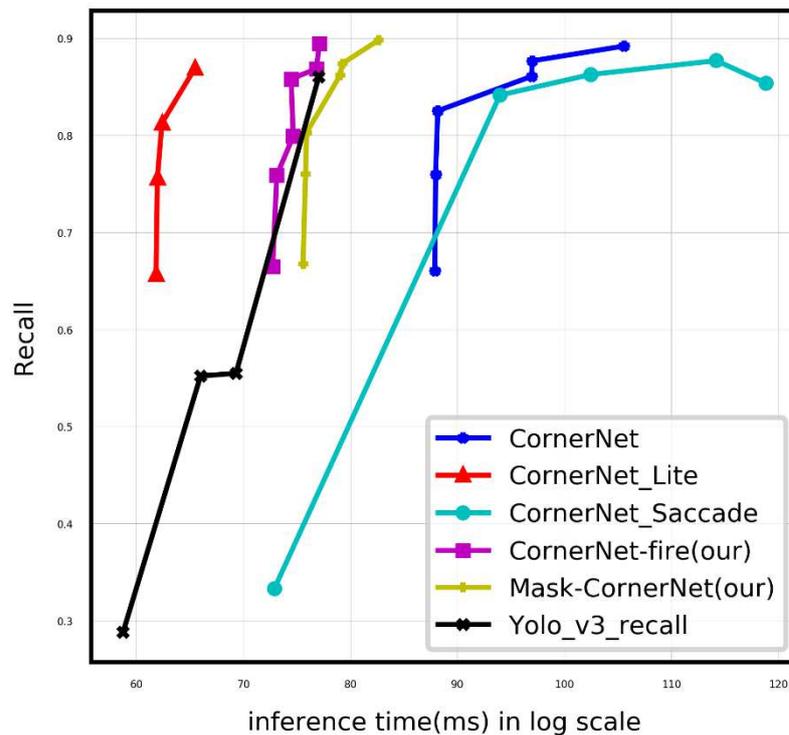
4.5 Accuracy and Efficiency

Our proposed method Mask-CornerNet is compared with ConerNet, YOLOv3 and RetinaNet in terms of time consumption and accuracy. Figure 9 shows the results of our experiment. Explore the performance of all methods on 6 different images of 128, 256, 448, 512, 640. Our model size and schedule relationship are given in Table 3. Our fire-model can directly reduce the size of the entire CornerNet-squeeze by 112Mb, while the accuracy is not reduced too much. After adding the Mask module, our model size has also been reduced to only 112.8Mb. Our anchor-free object detection maintains its advantages in small-size object detection models. In this paper, the fire module and the mask module use a small number of parameters to optimize the baseline model. The increase in inference time is mainly caused by channel fusion and instance segmentation.

Table 3. Time performance comparison result

	Time	Accuracy	Recall	Size
YOLO-v3	47ms	92.90	86.44	247Mb
CornerNet-squeeze	35ms	92.13	88.79	128Mb
CornerNet-fire(our)	41ms	92.06	89.46	112Mb
Mask-cornerNet(our)	62ms	91.63	89.83	112.8Mb

Figure 9 shows a comparison between our model and other models in terms of GPU consumption time. All experiments were performed in a consent environment. Due to the influence of the data set, the accuracy of all models on each input image size fluctuates without reference value. So we just compare the time consumed by the GPU with the recall rate.



332 **Fig 9.** Diagram the effect of the recommended number on the recall
 333

334 It takes 47ms for YOLO-v3 to get a recall of 86.44%. CornerNet-squeeze takes 32ms to the highest recall rate of 88.79% .
 335 After adding the fire module, the size of the model is significantly reduced. It takes only 41ms for the recall rate to reach 85.82%,
 336 and 47.2ms for 89.45%. Mask-CornerNet will take 52ms to reach 87.21%, and it will take 62ms to reach a recall rate of 89.82%.
 337 At the same time, Mask-ConerNet and ConerNet-fire are the same size. The CornerNet-Saccade and CornerNet models do not
 338 perform well on our dataset. CornerNet-Saccade takes 167ms to obtain a recall rate of 86.27%, and it takes 301ms to reach a
 339 recall rate of 87.70%. It takes 127 ms for CornerNet to reach a recall of 86.09%. And it takes 195ms to reach a recall rate of
 340 89.20%.

341 5. Conclusion

342 This paper proposes an anchor free object detection method with mask attention mechanism. In order to make full use of
 343 the feature maps of convolutional networks, this paper fuses feature maps of multiple scales for object detection. The fire module
 344 proposed in this paper is beneficial to the model being transplanted to mobile devices. In this paper, the fire module can
 345 semantically group convolutional network channels to improve the function distribution of the convolution kernel. To directly
 346 use the convolutional network to generate the corner detection heat map is not ideal, this paper proposes to use the Mask
 347 mechanism to guide the corner detection network. The instance segmentation results are fed back into the convolutional network
 348 to increase the diversity of convolutional network feature maps and enhance the network's ability to express. Our method was
 349 evaluated on the Tsinghua-Tencent 100K dataset. And compared with the commonly used object detection method yolov3, and
 350 the latest CornerNet, CornerNet-squeeze method. Experimental results show that our method performs well on the data set.

351 The detection speed and model size of our model have room for further improvement. Next, we will design a more concise
 352 feature extraction network optimization model scale and detection accuracy. The mask module generates instance segmentation
 353 with time consumption. We will further study the instance segmentation module with better time performance. This article mainly
 354 uses a traffic sign dataset containing a large number of small objects, and then we expand to the detection of other natural objects.
 355 This article also provides a feasible idea for improving the detection performance of traffic sign objects.

356 Abbreviations

357 R-CNN: Region with Convolutional Neural Network; YOLO: You only look once; SSD:Single Shot MultiBox Detector;
 358 FPN: Feature Pyramid Networks; R-FCN: Region-based Fully Convolutional Networks; U-Net: Convolutional Networks for
 359 Biomedical Image Segmentation; V-Net: Fully Convolutional Neural Networks for Volumetric Medical Image Segmentation;
 360 SGE: Spatial Group-wise Enhance;

Acknowledgements

The authors are grateful for the comments and reviews from the reviewers and editors.

Funding

Not applicable

Availability of data and materials

Tsinghua-Tencent 100K dataset [25] is available at <http://cg.cs.tsinghua.edu.cn/traffic-sign/>

Authors' contributions

Beibei Fan led the project, provided topics, conducted reviews and provided suggestions. He Yang did experiments and wrote this paper. Lingling Guo improves the paper.

Competing interests

The authors declare that they have no competing interests.

Author details

1. Beibei Fan. She is now a professor at Shanghai University. Her research focuses on data mining, path planning, machine learning, and artificial intelligence.

2. He Yang. He is now a student of School of Mechatronic Engineering and Automation in Shanghai University. His research is mainly on data mining, machine learning and image processing.

3. Lingling Guo. She is now a student of School of Mechatronic Engineering and Automation in Shanghai University. Her research is mainly on path planning and machine learning.

6. References

1. Zou, Z.; Shi, Z.; Guo, Y.; Ye, J. Object Detection in 20 Years: A Survey. *arXiv e-prints* **2019**, arXiv:1905.05055.
2. Liu, W.; Hasan, I.; Liao, S. Center and Scale Prediction: A Box-free Approach for Pedestrian and Face Detection. *arXiv e-prints* **2019**, arXiv:1904.02948.
3. Yang, B.; Tang, M.; Chen, S.; Wang, G.; Tan, Y.; Li, B. A vehicle tracking algorithm combining detector and tracker. *EURASIP Journal on Image and Video Processing* **2020**, *2020*, 17, doi:10.1186/s13640-020-00505-7.
4. Sarwar, A.; Mehmood, Z.; Saba, T.; Qazi, K.A.; Adnan, A.; Jamal, H. A novel method for content-based image retrieval to improve the effectiveness of the bag-of-words model using a support vector machine. *Journal of Information ence* **2019**, *45*, 117-135.
5. Liang, Z.; Shao, J.; Zhang, D.; Gao, L. Traffic sign detection and recognition based on pyramidal convolutional networks. *Neural Computing and Applications*.
6. Redmon, J.; Divvala, S.; Girshick, R.; Farhadi, A. You Only Look Once: Unified, Real-Time Object Detection. *Proc Cvpr Ieee* **2016**, 10.1109/Cvpr.2016.91, 779-788, doi:10.1109/Cvpr.2016.91.
7. Ren, S.Q.; He, K.M.; Girshick, R.; Sun, J. Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. *Adv Neur In* **2015**, *28*.
8. Duan, K.; Bai, S.; Xie, L.; Qi, H.; Huang, Q.; Tian, Q. Centernet: Keypoint triplets for object detection. In Proceedings of Proceedings of the IEEE International Conference on Computer Vision; pp. 6569-6578.
9. Law, H.; Deng, J. CornerNet: Detecting Objects as Paired Keypoints. In *The European Conference on Computer Vision (ECCV)*, 2018.
10. Huang, L.; Yang, Y.; Deng, Y.; Yu, Y. DenseBox: Unifying Landmark Localization with End to End Object Detection. *arXiv e-prints* **2015**, arXiv:1509.04874.

- 403 11. Redmon, J.; Farhadi, A. YOLO9000: Better, Faster, Stronger. *30th Ieee Conference on Computer*
404 *Vision and Pattern Recognition (Cvpr 2017)* **2017**, 10.1109/Cvpr.2017.690, 6517-6525,
405 doi:10.1109/Cvpr.2017.690.
- 406 12. Redmon, J.; Farhadi, A. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*
407 **2018**.
- 408 13. Newell, A.; Yang, K.U.; Deng, J. Stacked Hourglass Networks for Human Pose Estimation. *Lect*
409 *Notes Comput Sc* **2016**, 9912, 483-499, doi:10.1007/978-3-319-46484-8_29.
- 410 14. Law, H.; Teng, Y.; Russakovsky, O.; Deng, J. CornerNet-Lite: Efficient Keypoint Based Object
411 Detection. *arXiv preprint arXiv:1904.08900* **2019**.
- 412 15. Li, X.; Hu, X.; Yang, J. Spatial Group-wise Enhance: Improving Semantic Feature Learning in
413 Convolutional Networks. *CoRR* **2019**, abs/1905.09646.
- 414 16. Liu, W.; Anguelov, D.; Erhan, D.; Szegedy, C.; Reed, S.; Fu, C.-Y.; Berg, A.C. SSD: Single Shot
415 MultiBox Detector. In Proceedings of ECCV.
- 416 17. Lin, T.-Y.; Dollar, P.; Girshick, R.; He, K.; Hariharan, B.; Belongie, S. Feature Pyramid Networks
417 for Object Detection. In *The IEEE Conference on Computer Vision and Pattern Recognition*
418 *(CVPR)*, 2017.
- 419 18. Lin, T.-Y.; Goyal, P.; Girshick, R.; He, K.; Dollar, P. Focal Loss for Dense Object Detection. In
420 *The IEEE International Conference on Computer Vision (ICCV)*, 2017.
- 421 19. He, K.; Gkioxari, G.; Dollar, P.; Girshick, R. Mask R-CNN. In *The IEEE International Conference*
422 *on Computer Vision (ICCV)*, 2017.
- 423 20. Dai, J.; Li, Y.; He, K.; Sun, J. R-FCN: Object Detection via Region-based Fully Convolutional
424 Networks. In *Advances in Neural Information Processing Systems 29*, Curran Associates, Inc.:
425 2016; pp. 379-387.
- 426 21. Yu, F.; Koltun, V. Multi-Scale Context Aggregation by Dilated Convolutions. *arXiv e-prints* **2015**,
427 arXiv:1511.07122.
- 428 22. Howard, A.G.; Zhu, M.; Chen, B.; Kalenichenko, D.; Wang, W.; Weyand, T.; Andreetto, M.; Adam,
429 H. MobileNets: Efficient Convolutional Neural Networks for Mobile Vision Applications. *arXiv e-*
430 *prints* **2017**, arXiv:1704.04861.
- 431 23. Dai, J.; Qi, H.; Xiong, Y.; Li, Y.; Zhang, G.; Hu, H.; Wei, Y. Deformable Convolutional Networks.
432 In *The IEEE International Conference on Computer Vision (ICCV)*, 2017.
- 433 24. Li, Y.; Chen, Y.; Wang, N.; Zhang, Z. Scale-Aware Trident Networks for Object Detection. In *The*
434 *IEEE International Conference on Computer Vision (ICCV)*, 2019.
- 435 25. Huang, L.; Yang, Y.; Deng, Y.; Yu, Y. Densebox: Unifying landmark localization with end to end
436 object detection. *arXiv preprint arXiv:1509.04874* **2015**.
- 437 26. Law, H.; Deng, J. Cornernet: Detecting objects as paired keypoints. In Proceedings of Proceedings
438 of the European Conference on Computer Vision (ECCV); pp. 734-750.
- 439 27. Tian, Z.; Shen, C.; Chen, H.; He, T. FCOS: Fully Convolutional One-Stage Object Detection. In
440 Proceedings of Proc. Int. Conf. Computer Vision (ICCV).
- 441 28. Zhou, X.; Zhuo, J.; Krahenbuhl, P. Bottom-Up Object Detection by Grouping Extreme and Center
442 Points. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- 443 29. Long, J.; Shelhamer, E.; Darrell, T. Fully convolutional networks for semantic segmentation. In
444 Proceedings of Proceedings of the IEEE conference on computer vision and pattern recognition; pp.
445 3431-3440.
- 446 30. Ronneberger, O.; Fischer, P.; Brox, T. U-Net: Convolutional Networks for Biomedical Image
447 Segmentation. In *Medical Image Computing and Computer-Assisted Intervention - MICCAI 2015*,
448 Springer International Publishing: Cham, 2015; pp 234-241.
- 449 31. Milletari, F.; Navab, N.; Ahmadi, S. V-Net: Fully Convolutional Neural Networks for Volumetric
450 Medical Image Segmentation. In *2016 Fourth International Conference on 3D Vision (3DV)*, 2016;
451 pp 565-571.
- 452 32. Li, X.; Hu, X.; Yang, J. Spatial Group-wise Enhance: Improving Semantic Feature Learning in
453 Convolutional Networks. *arXiv preprint arXiv:1905.09646* **2019**.

- 454 33. Li, X.; Zhong, Z.; Wu, J.; Yang, Y.; Lin, Z.; Liu, H. Expectation-Maximization Attention Networks
455 for Semantic Segmentation. **2019**.
- 456 34. Zhu, Z.; Liang, D.; Zhang, S.; Huang, X.; Li, B.; Hu, S. Traffic-sign detection and classification in
457 the wild. In Proceedings of Proceedings of the IEEE Conference on Computer Vision and Pattern
458 Recognition; pp. 2110-2118.
- 459 35. Ren, S.; He, K.; Girshick, R.; Sun, J. Faster r-cnn: Towards real-time object detection with region
460 proposal networks. In Proceedings of Advances in neural information processing systems; pp. 91-
461 99.
- 462
- 463

Figures



Figure 1

Object detection results of some real traffic scenes. Our system works very well in complex scene where objects are small and highly occluded.

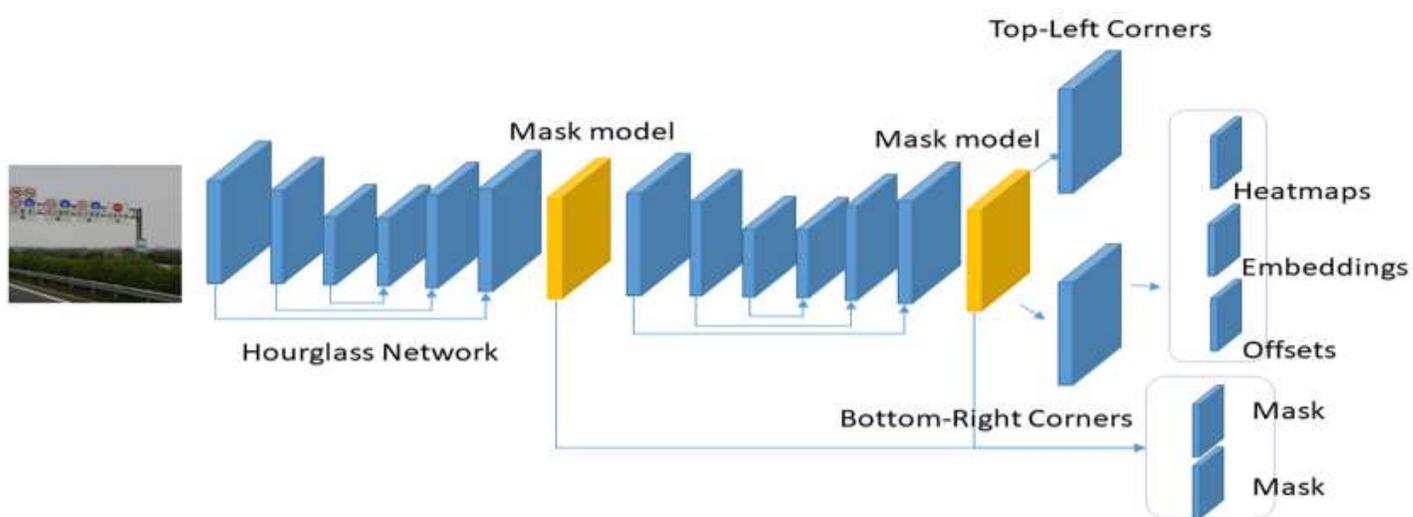


Figure 2

Describe the proposed method. By stacking multiple Hourglass Network and mask modules to enhance the network extraction ability, the mask module can get the image segmentation result. The corner detection module has two parts, and the top left and bottom right corners of the object are obtained

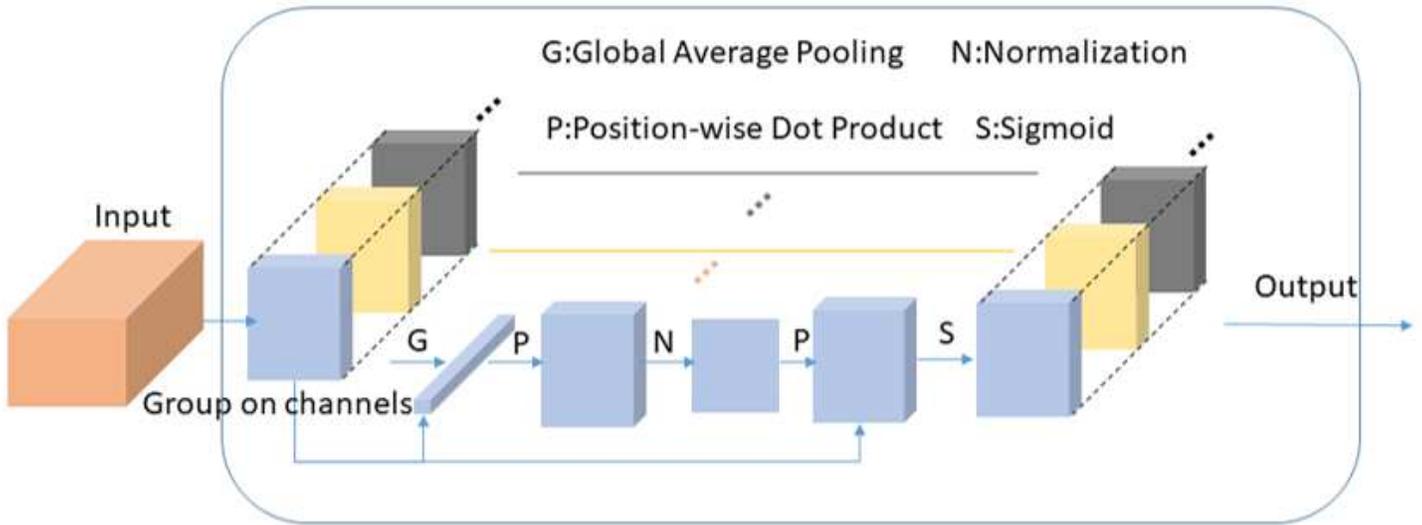


Figure 3

Diagram our improved fire module, which combines global and local features by grouping for incoming feature maps.

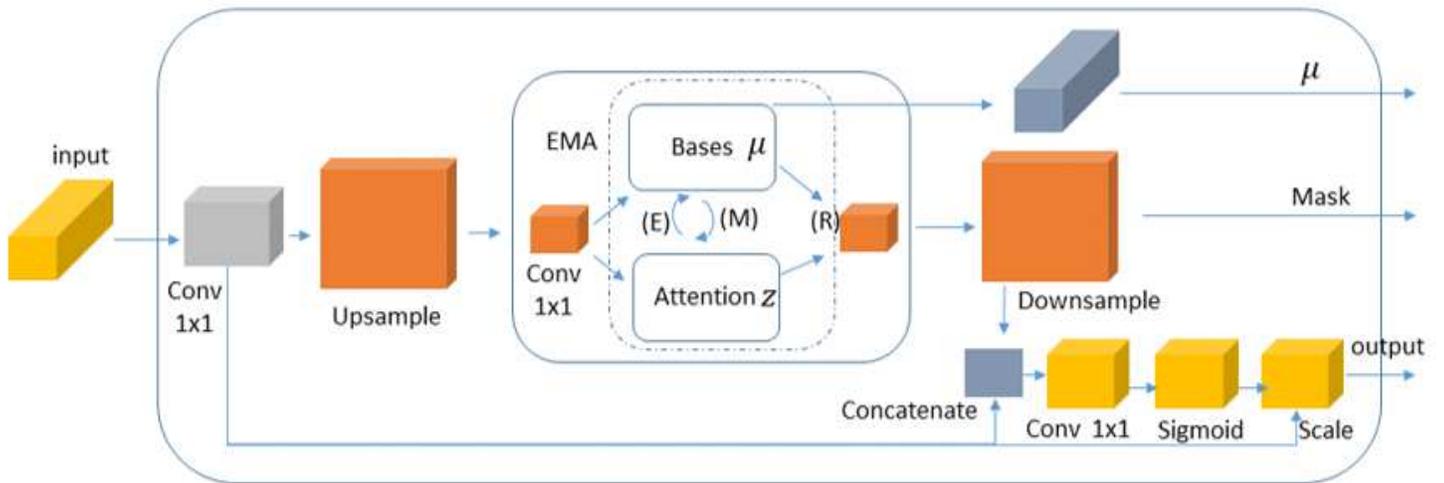


Figure 4

Show the details of the mask module. Convolutional features is upsampled by multiple convolution kernels. Using the EMA module to estimate the mask, the mask module can get the segmentation result of the object. Aggregating the mask module can generate spatial constraints on the object and enhance local selection of features.



Figure 5

The data tag is displayed, the left side is the mark for coordinate positioning, and the right side is the label for image segmentation.

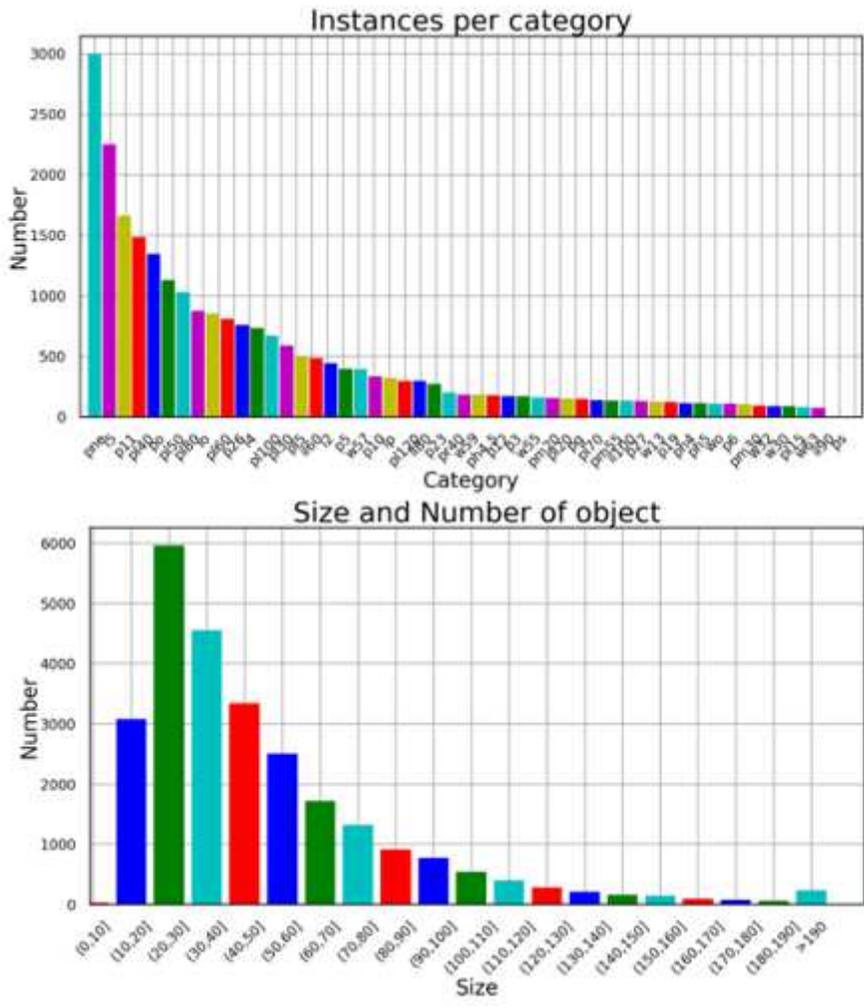


Figure 6

Object size and category statistics

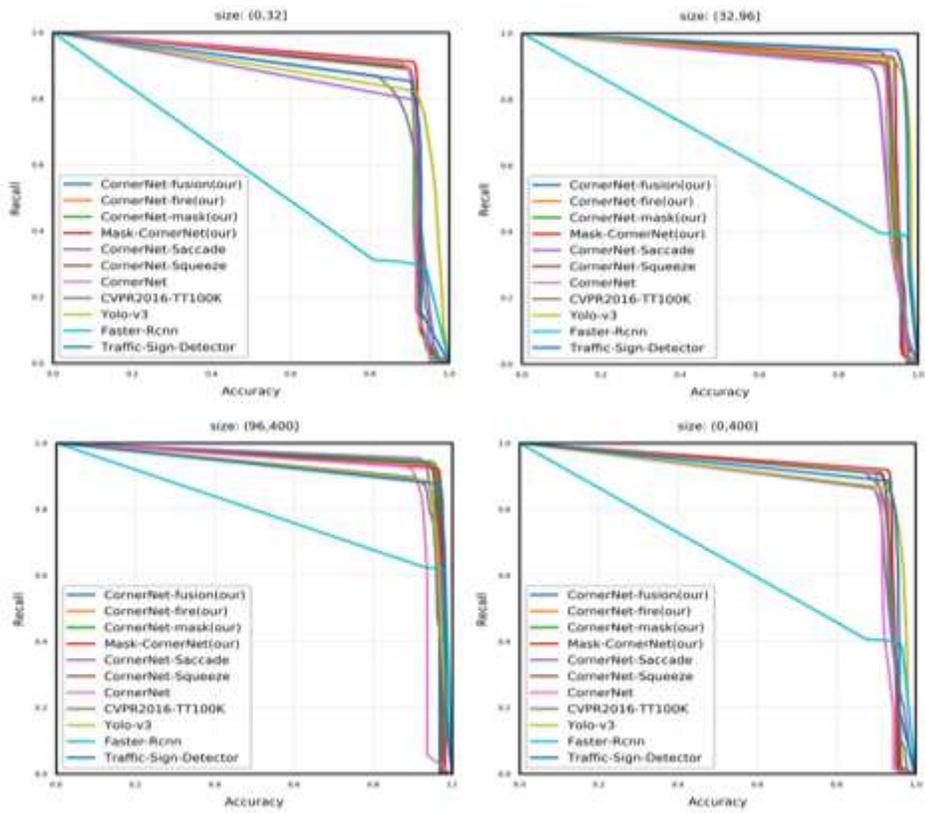


Figure 7

The performance of the proposed method is on the three scales of large, medium and small.

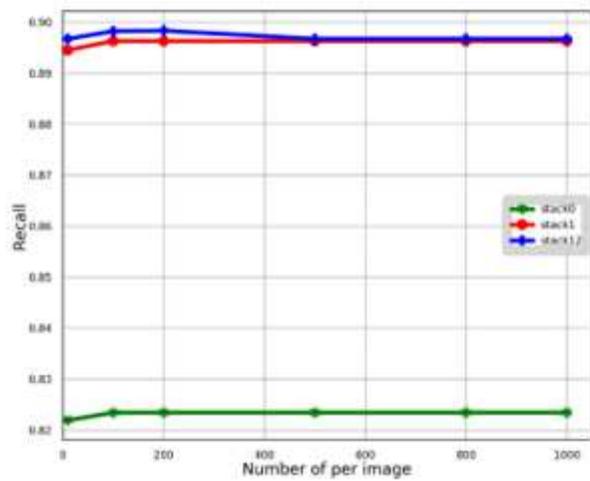


Figure 8

Diagram the effect of the recommended number on the recall

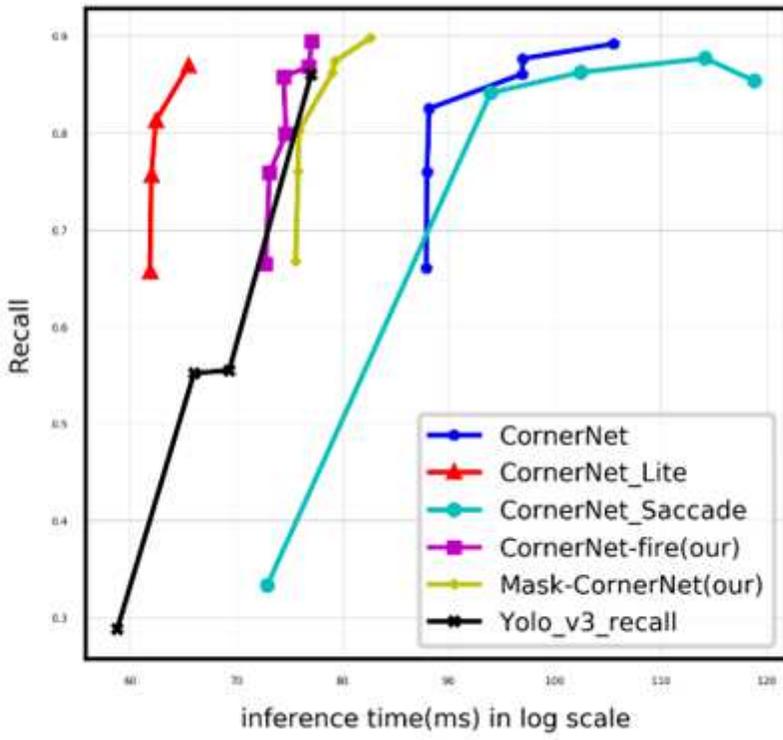


Figure 9

Diagram the effect of the recommended number on the recall