# PointINS: Point-based Instance Segmentation

Lu Qi[1]*, Xiangyu Zhang[2], Yingcong Chen[1], Yukang Chen[3]*, Jian Sun[2‡], and
Jiaya Jia[1,4‡]

[1] The Chinese University of Hong Kong
[2] MEGVII Technology
[3] Institute of Automation, Chinese Academy of Sciences
[4] SmartMore

**Abstract.** A single-point feature has shown its effectiveness in object detection. However, for instance segmentation, it does not lead to satisfactory results. The reasons are two folds. Firstly, it has limited representation capacity. Secondly, it could be misaligned with potential instances. To address the above issues, we propose a new point-based framework, namely PointINS, to segment instances from single points. The core module of our framework is instance-aware convolution, including the instance-agnostic feature and instance-aware weights. Instance-agnostic feature for each Point-of-Interest (PoI) serves as a template for potential instance masks. In this way, instance-aware features are computed by convolving this template with instance-aware weights for following mask prediction. Given the independence of instance-aware convolution, PointINS is general and practical as a one-stage detector for anchor-based and anchor-free frameworks. In our extensive experiments, we show the effectiveness of our framework on RetinaNet and FCOS. With ResNet101 backbone, PointINS achieves 38.3 mask mAP on challenging COCO dataset, outperforming its competitors by a large margin. The code will be made publicly available.
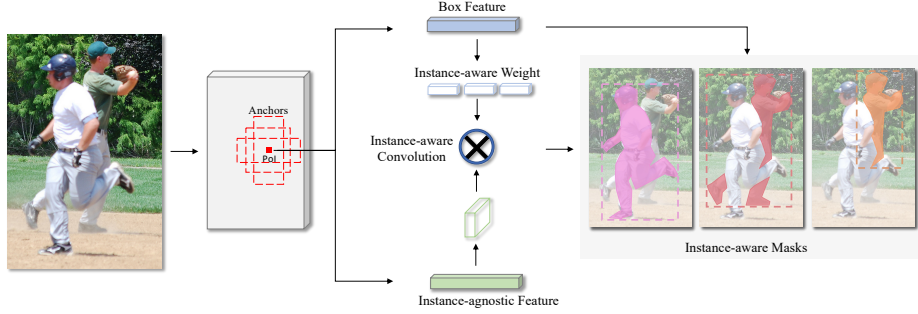
**Keywords:** instance segmentation · single-point feature

## 1 Introduction

Instance segmentation aims to precisely detect all objects in pixel-level in an image. From this point of view, it involves both object detection [10,11,31,19,9] and semantic segmentation [25,39,6]. This task is fundamental in computer vision and benefits a large number of applications, such as autonomous vehicles[28], robotics[26], and video surveillance[32].

---

**Fig. 1.** The illustration of our PointINS framework. The instance-agnostic feature is first obtained from a single-point feature. It serves as a template for various instances. By introducing the unique instance-aware weights, we obtain the masks by instance-aware convolution.

*Existing Methods* Due to the close connection to object detection, most of the SOTA instance segmentation methods [12,23,5,15,16,27] are built upon detectors [10,11,31,19,24,30,20,33] whatever they are one- or two- stage. They first predict Region-of-Interest, and then extract RoI features to segment objects. Although to obtain RoI features is a complicated process, these features are explicit enough for mask prediction thanks to RoI pooling or alignment.

From the view of semantic segmentation, methods [38,34,2,22] are proposed to reduce massive operations for RoIs. An example is to generate masks from the whole feature map. Albeit simple with fully connected structure, this line of methods still costs much computation during inference with feature embedding or the need to producing all masks. Therefore, boxes are still necessary to indicate valid regions for mask prediction.

Following the instance generation of one-stage detectors without RoI feature extraction, Polarmask [36] and Tensormask [7] explored mask prediction only from a single-point feature. However, performance or speed is still unsatisfying. In this paper, we analyze the difficulties to apply the point-of-interest (PoI) feature in instance segmentation, and accordingly propose a new point-based instance segmentation approach, namely PointINS, as a decent solution.

*Difficulties in Point-based Instance Segmentation* The first obstacle is about the *representation capacity* of features. For applying the point features in instance segmentation, without using RoI alignment, they cannot provide sufficient information for mask prediction. One-stage detectors, e.g., RetinaNet [20] and FCOS [33], directly predict bounding boxes of objects from point features. A single-point feature may be sufficient for object detection. But it is not for instance segmentation, since the representation of masks is more complicated. Tensormask [7] uses two representations to predict masks from points. Both of them require a complex tensor bipyramid head to produce masks in various resolutions. Inevitably, the features in Tensormask are heavy, leading to notable computation cost. On the other hand, Polormask [36] represents masks with

**Table 1.** Performance comparison with various RoI sizes in Mask R-CNN. For a fair comparison, we use consecutive de-convolution layers and a bilinear interpolation layer to let the size of final masks reach $32 \times 32$. The mAP is the mean average precision from 0.5 to 0.95 threshold with 0.05 stride.

| RoI feature size | $14 \times 14$ | $7 \times 7$ | $3 \times 3$ | $1 \times 1$ |
|:---:|:---:|:---:|:---:|:---:|
| mAP | 34.5 | 31.5 | 24.7 | 18.4 |

polygons in polar coordinates. However, the representation is coarser than the original ground truth.

In Table 1, we show the results of Mask R-CNN [12] with different RoI feature size. Performance drops dramatically as RoI feature size decreases. When the size reduces to $1 \times 1$, the mask is predicted from a single point, which is the same case of point-based frameworks. In this case, the mAP drops from 34.4 to 18.4. This suggests that single-point features contain insufficient information for mask prediction.

The second difficulty is about the *misalignment* between features and instances. In one-stage detectors, a single-point feature could be related to numerous bounding boxes. However, a mask is for only one box. If the mask is predicted by a single-point feature, it is unclear which bounding box to match. As a result, performance could drop if the segmentation mask mismatches the bounding box. This issue also exists in other point-based instance segmentation frameworks. To reduce this issue, Tensormask [7] and Polarmask [36] fix the number of instances at each point. YOLACT [3] uses mask assembly, and yet complicates the inference process.

*Our Solution* With the analysis above, it is non-trivial to build an instance segmentation model based on the point feature. In this paper, we propose PointINS, a simple and yet effective approach to evolving one-stage object detection frameworks for instance segmentation. The key in PointINS is to predict masks by the more *robust and aligned* single-point feature. Instead of directly interpolating point features in Tensormask, we focus on the solution to enhance the representation power of point features, as well as alignment of features and instances. Instance-aware convolution is the core module in our PointINS, with the two inputs including instance-agnostic feature and instance-aware weights. As shown in Fig. 1, the pyramid features are processed to instance-agnostic feature, serving as a template for potential masks. Meanwhile, instance information from the box feature is used to predict a dynamic weight for instance-aware convolution. The instance-aware feature is generated from this dynamic convolution, which aligns features and instances. Two critical inputs of our instance-aware convolution guarantee the robustness and uniqueness of features, and avoid the mismatch between segmentation masks and bounding boxes. Finally, the instance masks are predicted from the instance-aware features. Our contributions are summarized as follows.

– We propose PointINS, a new pointed-based framework for instance segmentation. The instance-aware convolution including instance-agnostic feature

and instance-aware weight, aiming to address the issues of feature representation and misalignment.

- In our experiments, the proposed PointINS leads to satisfactory results. Under the single-scale 1x training schedule, we obtain 34.5 mask mAP with ResNet101-FPN. Using data augmentation or with longer training time improves performance consistently to 38.3.
- The instance-aware convolution is independent to the box head. PointINS thus can be part of most one-stage detectors, including RetinaNet [20] and FCOS [33]. It serves as a general framework for point-based instance segmentation.

## 2 Related Work

**Instance Segmentation** Current instance segmentation approaches can be roughly classified into two categories of segmentation-based [22,2,34] and detection-based [12,23,5,15,16,27] methods. Segmentation-based methods usually learn a specially designed transformation or instance boundaries and cluster the points with similar embedding into instance masks. The clustering process makes it relatively hard to apply to complicated situations. SOLO [34] directly predicts instances masks from channels of the whole image feature map. Without the box indicator, it costs too much computation to segment invalid regions.
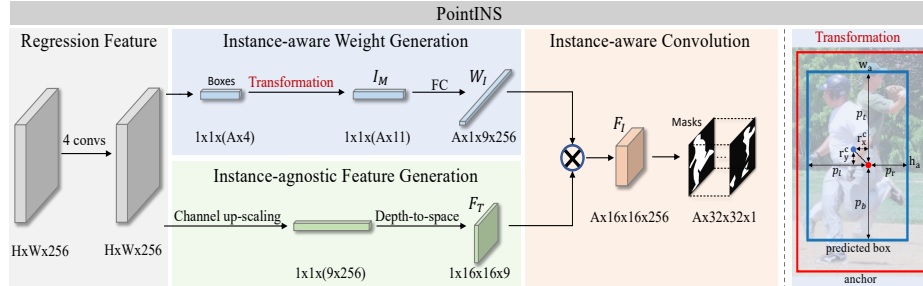
Popular instance segmentation frameworks are mainly detection-based [10,11,31,9,19,27]. They usually predict a series of bounding boxes and then segment masks. Mask R-CNN is the representative framework with two-stage detectors. Several methods, including PANet [23], HTC [5], and Mask Scoring R-CNN [15], refine some details of the structure from different aspects. All of them use the RoI-Align [12] to produce pixel-level features. For another, based on the typical one-stage detector FCOS [33], Centermask [16] and Embedmask [38] use point features to detect boxes but RoI features to segment instances.

To explore this problem, Tensormask [7] and Polarmask [36] force a single-point feature to segment only one or two instances. YOLACT [3] ensembles the point response to several prototypes for final mask prediction. However, the prototypes are pre-defined and involve additional hyper-parameters. Different from these methods [7,36,3], our framework generates a point-wise instance-agnostic feature, serving as a template for numerous potential instance masks.

**Dynamic Weight** As a specific strategy of meta-learning [1,29,35], weight prediction [17] aims to learn parameters by an independent structure. In other words, the network parameters are dynamic with various inputs.

Cai et al. [4] used the predicted parameters of a classifier to discriminate new categories. These parameters are generated by the memory of the support set. In object detection, MetaAnchor [37] is a flexible anchor generator to produce higher quality boxes. The weight of the last regression layer is predicted with properties of customized prior anchors. To balance the instance annotation cost between box and mask, Hu et al. [13] proposed a semi-supervised method to

segment everything by re-weighting box features. Moreover, MetaSR [14] uses a dynamic up-sampling module to super resolve a single image with arbitrary scale factors.

We use dynamic weight convolution to align between an instance-agnostic feature and potential instances. With distinct instance property, we generate unique instance-aware features for mask prediction.



**Fig. 2.** The detailed structure of our PointINS framework. The core module is instance-aware convolution, including the *instance-agnostic* feature and *instance-aware* weight. To commence, we generate the instance-agnostic feature $F_T$ by channel up-scaling (a convolution layer) and depth-to-space (a reshape operation). Meanwhile, the instance-aware weight $W_I$ is obtained by transforming the box feature to explicit instance information $I_M$. Then the two features serve as the input and weight of the *instance-aware* convolution, in order to produce instance-aware feature $F_I$ for final mask prediction. For clear illustration, we only sample a single point and annotate each 4D/3D feature map with N×H×W×C or H×W×C dimension. N, H, W and C denote the batch, height, width and channel numbers respectively. A is the number of anchors covered by a single point, with 9 and 1 in RetinaNet and FCOS. The right-most part is the geometric illustration for our transformation module.

## 3   Our Method

The overall architecture of our PointINS framework is shown in Fig. 2. Different from previous work [7,36], the single point in our PointINS is responsible for several potential masks. The core module of our framework is instance-aware convolution. It has two inputs including *instance-agnostic feature* and the *instance-aware weight*. The former plays a template role for potential instances, while the latter indicates the unique one. Finally, instance-aware features are generated by dynamic instance-aware convolution and then used for final class-agnostic mask prediction.

Similar to the boxes, the masks are also the representations of instance locations. Thus the instance-agnostic feature and instance-aware weight share the last feature map in the box regression branch, which usually has four convolutions.

It is worth noting that our proposed method can be naturally used in the most one-stage detectors, whatever they are anchor-free or anchor-based. For simplicity, we take the RetinaNet as an example to elaborate on our motivation and design of our instance-aware convolution.

### 3.1 Instance-agnostic Feature Generation

This module aims to provide robust and informative features for the subsequent instance segmentation branch. In our framework, a single-point feature $F_T$ could produce several masks. Thus, the single-point feature should have a sufficiently sizeable receptive field to cover all possible masks, and contain enough semantic information for following mask prediction.

To this end, we design the module as containing a channel up-scaling and a depth-to-space operation.

(1) Channel up-scaling operation is a convolution layer with input and output channel numbers 256 and 2304. Since $2304 = 9 \times 256$, the output channel increases 9 times with its $3 \times 3$ convolution kernel size. This operation is essential in practice, as the increased channels allow our tensor to encode more context messages for mask prediction without information loss. We are noting that this operation is relatively cheap, as will be discussed in subsection 3.3.

(2) Depth-to-space operation is used to balance the template spatial resolution (spatial dimension) and feature representation (channel dimension). We transform the point feature at each location into a 3D tensor with size $(9 \times 16 \times 16)$. In this way, we force each channel of our point feature to have explicit relative-location information[5]. Table 4 in our ablation study gives a detailed analysis.

For the FPN design, instances with different sizes are distributed to corresponding feature levels. As such, our point feature is robust enough to capture the whole part of instances, even if it is in the instance boundary.

***Essential Difference from Tensormask*** Note that this module is by nature different from Tensormask [7], because we encode more information in the single-point feature.

At first, we enhance our single-point feature by a channel up-scaling convolution. Compared with using interpolation in Tensormask, our convolution captures obviously more information for a single point. Unlike previous fixed interpolation, the convolution is flexible to fuse its context by learnable weights.

Further, in Tensormask [7], two types of instance representation, i.e., nature and aligned representation, are used for a single point. Both of them have two spatial dimensions, representing the mask resolution at each feature map location. Unlike our solution, there was no channel dimension, making the network

---

[5] It is similar to R-FCN [9] and FCIS [18], whose channels explicitly encode relative location scores of instances.

only focus on the current mask location, and accordingly cannot provide much information outside the template region.

In a nutshell, these two representations in Tensormask only apply to the fixed-size sliding window, while we significantly expand the region to consider for much higher confidence.

### 3.2   Instance-aware Convolution

***Motivation*** The current one-stage detectors usually generate multiple boxes from only a single point, e.g., RetinaNet [20]. However, for instance segmentation, it is even hard to produce a mask from a single point directly. The reasons are as follows.

(1) *Channel explosion:* If we directly generate multiple masks from the single point, as what one-stage detectors do, the output channels of convolution in the mask head would be unacceptably massive. For example, predicting a $28 \times 28$ mask requires 7056 channels, where $7056 = 28 \times 28$.
(2) *Misalignment between features and instances:* As shown in Fig. 2, it is common that different proposals could overlap. Thus, there is a severe issue of misalignment if one predicts multiple masks from single points. The feature of a single point is unique, but it contains too much information of instances.
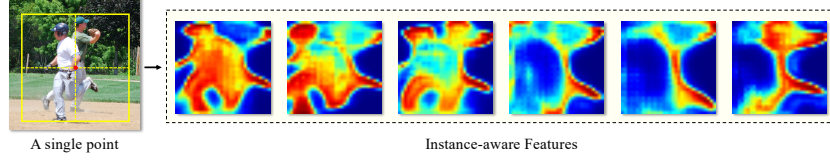
To address the above issues, we propose the instance-aware convolution, which converts the instance-agnostic features to instance-aware ones. The convolution weight is location-sensitive and learned from the specific instance information. This design, no matter how many instances one point should be responsible for, allows us to align the template feature with the specific instance. Therefore, this module is general to most of the one-stage detectors like RetinaNet [20] or FCOS [33].

***Our Design*** We propose a specific instance-aware convolution by dynamic weights. Given a point feature that may be related to several masks, using exact bounding box property can be viewed as a selector that specifies which mask should be predicted.

As shown in Fig 2, the weight prediction module is light-weighted, with only a transformation module and an FC layer. To get the potential instance info as much as possible, a transformation module is first used. For clarity, we divide the input $I_M$ into anchor indicator $I_A$ and instance indicator $I_P$. $I_A$ shows which pre-defined anchors we will use. $I_P$ exhibits the detail offsets between anchors and regressed potential instances. The form is

$$I_M = \{I_A, I_P\}, \quad \text{where}$$
$$I_A = \{\frac{1}{s}, \frac{h_a}{w_a}, \frac{h_a}{s}, \frac{w_a}{s}\}. \tag{1}$$

Here, $s$ is the stride of the global feature level. For example, the strides of pyramid features $P3 - P7$ are $\{8, 16, 32, 64, 128\}$. $h_a$ and $w_a$ are the height and

A single point           Instance-aware Features

**Fig. 3.** The visualization of instance-aware features under the same point feature. With various instance information, instance-aware convolution generates distinct instance-aware features for following mask prediction.

width of anchors. Therefore, the last three components indicate the aspect ratio and scale of the anchor. The definition of $I_A$ contains information about which specific anchor the potential mask is generated from at a feature level.

We define $I_P$ as

$$I_P = \{r_x^c, r_y^c, p_l, p_r, p_b, p_t\}, \tag{2}$$

to represent the instance offset from the original anchors. $r_x^c$ and $r_y^c$ are distances between centers of anchor and proposal. The other elements represent the distances between the box location and the center of corresponding anchor in four directions. Each component in $I_P$ is calculated as

$$
\begin{aligned}
r_x^c &= \frac{x_p^c - x_a^c}{s}, \quad r_y^c = \frac{y_p^c - y_a^c}{s} \\
p_l &= \frac{0.5 * w_p}{s} - r_x^c, \quad p_r = \frac{0.5 * w_p}{s} + r_x^c \\
p_t &= \frac{0.5 * h_p}{s} - r_y^c, \quad p_b = \frac{0.5 * h_p}{s} + r_y^c
\end{aligned}
\tag{3}
$$

where $x_p^c$ and $x_a^c$ are the x-axis center location of proposals and anchors. They are the same as $y_p^c$ and $y_a^c$. $h_p$ and $w_p$ are the height and width of the proposal. The right-most part of Fig. 2 shows the geometric meaning of our transformation module.

For instances with $I_M$, we predict their corresponding convolution weight $W_I$ by a FC layer, which is a fully connected layer followed by a ReLU operation. It is noted that $I_M$ only contains $I_P$ for the anchor-free detectors. Taking FCOS [33] for an example, we regard the points of each feature level as our anchors, in which $x_a^c$ and $y_a^c$ are the sampled-point locations of the images.

Finally, we obtain our instance-aware feature $F_I$ by dynamically convolving with the instance-agnostic template as

$$F_I = conv(F_T, W_I). \tag{4}$$

In this way, our instance-agnostic feature is easily aligned to corresponding boxes. Then via the standard mask head of Mask R-CNN, it is ready for final class-agnostic mask prediction. Fig. 3 shows visualization of the instance-aware feature by dynamically convolving with the same point feature.

### 3.3    Computational Complexity

We elaborately design our instance-aware convolution for high effectiveness. It is achieved with two implementation techniques. First, we only sample positive points to produce the template. In our implementation, a point is defined positive if the Intersection of Union (IoU) between ground truth and both anchors and regressed proposals are larger than 0.5. In the COCO dataset, there are 6.43 positive points for each instance following our sample principle, even only considering the anchor condition. We randomly sampled 128 points in training out of $22,000+$ points on all feature levels. For inference, only 100 detected points are used for final mask prediction.

Thus, the subsequent processing is applied only on these points, which brings remarkable acceleration. Due to the channel up-scaling convolution with stride 3 in the instance-agnostic feature generation module, we only need to sample $3 \times 3$ receptive fields from the feature map to produce our instance-agnostic features. Besides, since these points are independent, they can be processed in parallel for further notable acceleration and memory usage reduction. Our inference achieves 13.4 FPS with ResNet101 backbone on V100.

## 4    Experiments

### 4.1    Experiments on COCO Dataset [21]

We compare our method with other state-of-the-art approaches on the challenging COCO dataset [21]. Following common practice [12,19,23,15], we train our models with 115K train images and report results on the 5K val subset for ablation study. We also report results on 20K test-dev for comparison. Comprehensive ablation studies are conducted on this dataset with RetinaNet-ResNet50-based PointINS. We follow the standard evaluation metrics, i.e., $AP$, $AP_{50}$, $AP_{75}$, $AP_S$, $AP_M$ and $AP_L$. The last three measure performance is with respect to objects in different scales.

***Hyper-Parameters*** We train our network using batch size 16 for 12 epochs. The shorter and longer edges of the images are 800 and 1333. Adam gradient descent with 0.00005 learning rate, $(0.9, 0.999)$ beta and $1e^{-08}$ eps is used as the optimizer. We decay the learning rate with 0.001 at 8 and 11 epochs respectively. Using SGD leads to about 0.3 performance drop. The loss weight for boxes and masks are equal to 1. The edge of the mask is taken into more account with setting weight 4. We initialize our backbone networks with the pre-trained models on ImageNet.

**Comparing with Point-based methods** We have implemented our method based on RatinaNet [20] and FCOS [33], and compare them with three state-of-the-art point-based instance segmentation frameworks, i.e., TensorMask [7], PolarMask [36], and YOLACT [3] in Table 2. It is clear that our PointINS

**Table 2.** Comparison with current point-based instance segmentation frameworks on COCO test-dev split. The 'aug' means data augmentation, including multi-scale and random crop. ✓(or ∘) means that the network is trained with (or without) data augmentation. The 'epoch' indicates the training time. 12 epochs give a baseline, widely adopted by most previous work [12,23,20,33].

| method | backbone | aug | epochs | $AP$ | $AP_{50}$ | $AP_{75}$ | $AP_S$ | $AP_M$ | $AP_L$ |
|---|---|---|---|---|---|---|---|---|---|
| TensorMask [7] | R-50-FPN | ✓ | 72 | 35.4 | 57.2 | 37.3 | 16.3 | 36.8 | 49.3 |
|  | R-101-FPN | ✓ | 72 | 37.1 | 59.3 | 39.4 | 17.4 | 39.1 | 51.6 |
| YoLACT [3] | R-50-FPN | ✓ | 45 | 28.2 | 46.6 | 29.2 | 9.2 | 29.3 | 44.8 |
|  | R-101-FPN | ✓ | 45 | 31.2 | 50.6 | 32.8 | 12.1 | 33.3 | 47.1 |
| PolarMask [36] | R-50-FPN | ∘ | 12 | 29.1 | 49.5 | 29.7 | 12.6 | 31.8 | 42.3 |
|  | R-101-FPN | ∘ | 12 | 30.4 | 51.1 | 31.2 | 13.5 | 33.5 | 43.9 |
|  | R-101-FPN | ✓ | 24 | 32.1 | 53.7 | 33.1 | 14.7 | 33.8 | 45.3 |
| RetinaNet + ours | R-50-FPN | ∘ | 12 | 32.2 | 51.6 | 33.4 | 13.4 | 34.4 | 48.4 |
|  | R-101-FPN | ∘ | 12 | 33.5 | 53.6 | 34.1 | 14.2 | 35.3 | 49.0 |
|  | R-50-FPN | ✓ | 72 | 36.2 | 58.1 | 37.8 | 16.5 | 37.7 | 50.5 |
|  | R-101-FPN | ✓ | 72 | **37.9** | **60.1** | **39.7** | **17.8** | **40.1** | **52.1** |
| FCOS + ours | R-50-FPN | ∘ | 12 | 33.4 | 53.7 | 35.2 | 14.8 | 36.3 | 48.8 |
|  | R-101-FPN | ∘ | 12 | 34.5 | 54.5 | 36.6 | 15.4 | 37.0 | 49.4 |
|  | R-50-FPN | ✓ | 72 | 36.7 | 58.2 | 38.1 | 16.8 | 38.0 | 50.9 |
|  | R-101-FPN | ✓ | 72 | **38.3** | **60.3** | **40.0** | **18.1** | **40.3** | **52.4** |

performs best among all of them. In the following, We explain the comparison with these approaches.

PolarMask [36] uses polar representation as an approximation to the instance mask. With 12 epochs training, PolarMask achieves 29.1 mAP with R-50-FPN backbone, and 30.4 with R-101-FPN. With RatinaNet, our PointINS achieves 32.2, with 3.1 points improvement. The margin is even larger when we use FCOS as detection base – we achieve 33.4 and 34.5 with R-50-FPN and R-101-FPN backbones respectively. Note that polar representation is an approximation to the ground truth. This intrinsically causes precision loss. In contrast, our method uses the pixel-wise mask and does not suffer from this problem. Ours with ResNet50 under only single-scale 1x training still outperforms the one with ResNet101 under multi-scale 2x training.

YOLACT [3] first predicts mask prototypes, then linearly ensembles them for final prediction. Note that even if it takes a longer time for training, and applies data augmentation, the performance is still not as good as ours. It achieves 28.2 and 31.2 with R-50-FPN and R-101-FPN backbones, while ours+FCOS achieves 33.4 and 34.5 with R-50-FPN and R-101-FPN backbone respectively, even if our model is trained with only 12 epochs and without data augmentation. The reason we outperform YOLACT is that one point in our framework only covers a few instances within its receptive field instead of all the prototypes.

It is noted that by applying data augmentation and training for a longer time, our performance further improves significantly. By training as long as 72

**Table 3.** Comparison among state-of-the-arts instance segmentation methods. All methods reported are conducted with ResNet101 backbone. 'type' means the source of mask features like semantic map (SEG), region feature (RoI) or point feature (PoI).

| method | type | FPS | AP | GPU | method | type | FPS | AP | GPU |
|---|---|---|---|---|---|---|---|---|---|
| Mask R-CNN [12] | RoI | 10.6 | **38.3** | V100 | Tensormask [7] | PoI | 2.7 | 37.1 | V100 |
| SOLO [34] | SEG | 10.4 | 37.8 | V100 | PolarMask [36] | PoI | 12.3 | 31.9 | V100 |
| EmbedMask [38] | RoI | 13.7 | 37.7 | V100 | YOLACT-700 [3] | PoI | **23.4** | 31.2 | XP |
| **RetinaNet+ours** | PoI | 13.1 | 37.9 | V100 | **FCOS+ours** | PoI | 13.4 | **38.3** | V100 |

**Table 4.** Ablation study for instance-agnostic feature representation. The size of the tensor is $(C, H, W)$ for the three dimensions. The column "modification" means the ways to change this tensor.

| $(C, H, W)$ | modification | $AP$ | $AP_{50}$ | $AP_{75}$ |
|---|---|---|---|---|
| $(9, 16, 16)$ | default | **32.5** | **51.9** | **33.8** |
| $(1, 16, 16)$ | Decrease Channels | 30.4 | 49.2 | 31.6 |
| $(4, 16, 16)$ | | 31.6 | 50.1 | 32.4 |
| $(36, 8, 8)$ | Reshape Styles | 31.8 | 51.0 | 32.4 |
| $(1, 48, 48)$ | | 29.6 | 48.5 | 30.1 |
| $(16, 16, 16)$ | Increase Channels | 32.3 | 51.7 | 33.7 |
| $(25, 16, 16)$ | | **32.5** | 51.8 | 33.6 |

epochs and using data augmentation, our model with RetinaNet achieves 37.9 mAP with R-101-FPN backbone, while with FCOS it achieves 38.3 mAP. For a fair comparison with previous champion, TensorMask [7], we use the same hyper-parameter setting.

As shown, we achieve 0.8 (37.9-37.1) and 1.2 (38.3-37.1) improvement with RatinaNet and FCOS respectively. Our success stems from twofold. First, as mentioned in Section 3.1, the instance-agnostic feature we use is more informative and robust. Second, we propose the instance-aware module, which allows our framework to capture accurate masks for the predicted instances.

**Comparing with State-of-the-arts** Table. 3 reports the performance and frame-per-second (FPS) of current state-of-the-art instance segmentation methods with ResNet101 backbone. By only using a single-point feature, PointINS is still comparable to the RoI- or segmentation- based methods (EmbedMask [38], SOLO [34]). It validates PointINS is promising without any RoI feature extraction, and the point feature can be further improved with more delicate designs.

**Ablation Study** All ablation studies are conducted by RetinaNet-based PointINS with ResNet50. Note that this baseline yields a slight performance increase in COCO validation split (results 32.2 for test-dev and 32.5 for validation split).

*Instance-agnostic Feature Generation* The representation of the template tensor is first ablated. As illustrated in Table 4, we explore the performance influence on different shapes. The default tensor shape is $(9, 16, 16)$ with 9 denoting the channel dimension and 16 for two spatial dimensions. It is evident that the

**Table 5.** Ablation study for the input to generate instance-aware weights. As described in Section 3.2, we split the input into three parts and explore them. ∘ denotes "not used" for this part. And it is vice versa to ✓.

| $I_A$ | $\{r_x^c, r_y^c\}$ | $\{p_l, p_r, p_b, p_t\}$ | $AP$ | $AP_{50}$ | $AP_{75}$ |
|---|---|---|---|---|---|
| ∘ | ∘ | ∘ | 27.9 | 46.5 | 28.5 |
| ✓ | ∘ | ∘ | 25.5 | 44.9 | 26.1 |
| ∘ | ✓ | ∘ | 26.3 | 45.0 | 26.6 |
| ∘ | ∘ | ✓ | 28.9 | 47.1 | 29.3 |
| ✓ | ✓ | ∘ | 30.4 | 49.1 | 31.2 |
| ∘ | ✓ | ✓ | 31.7 | 50.6 | 32.8 |
| ✓ | ∘ | ✓ | 31.3 | 50.5 | 32.6 |
| ✓ | ✓ | ✓ | **32.5** | **51.9** | **33.8** |

**Table 6.** The performance comparison with other heuristic $I_M$ designs. The first column describes the detailed way of design.

| the heuristic $I_M$ design | $AP$ | $AP_{50}$ | $AP_{75}$ |
|---|---|---|---|
| Generating $I_M$ with an embed layer | 27.9 | 46.5 | 28.5 |
| $I_P$ are with FPN's box regression style | 31.8 | 50.3 | 32.7 |
| $\mathrm{Log}(I_A)$ | 32.2 | 51.6 | 33.6 |
| $\mathrm{Sigmoid}(I_M)$ | 32.3 | 51.6 | 33.6 |
| $\mathrm{Tanh}(I_M)$ | 32.3 | 51.8 | 33.4 |
| default Design | **32.5** | **51.9** | **33.8** |

performance drops when decreasing the channel numbers. This indicates that the channel information is vitally important.

Intriguingly, as shown in the last two rows, further increasing the channel numbers does not improve results. The reason lies in the origin of our template tensor. No matter how we transform the information, the template is reshaped from a single-point feature. In our framework, this single-point feature is generated by a $3 \times 3$ convolution with input channel number 256. Therefore, the information capacity is $3 \times 3 \times 256$. More output convolution channels do not increase valid messages.

This analysis also explains the performance drop using different reshape styles. The $3 \times 3$ kernel encodes the specific spatial sequence inherently. Therefore, reshaping to $(9, 16, 16)$ tensor is more intuitive than the way that each channel represents the relative location response. We note this is also the core idea of R-FCN [9] and FCIS [18]. For the tensor with size $(1, 16, 16)$, we regard it as the nature representation [7] proposed by the TensorMask work. It causes about 3 mAP decrease.

***Instance-aware Weight Generation*** Table 5 shows the ablation study on the input used for generating instance-aware weight. We first explore the two extreme cases, shown in the first and last rows of Table 5. Our framework accomplishes the best performance when involving all three parts. In the first row, we directly use the features from the detection branch instead of our explicit box information. The performance decrease indicates that the abstract information of bounding boxes is more helpful in generating instance-aware weights.

As alternatives to use this piece of information, only using the $I_A$ or $\{r_x^c, r_y^c\}$ leads to lower mAPs compared with directly using features of detection head. It means that the offset between anchors and boxes is vitally important in our module. Our design makes it possible to know the detailed shift to guide polishing the instance-agnostic template. It is verified in the last five rows of Table 5. In our module, $\{p_l, p_r, p_b, p_t\}$ play vital roles, and $I_A$ (and $\{r_x^c, r_y^c\}$) is peripheral.

Table 6 shows the performance with other heuristic designed for $I_M$. We conclude that the unique region indicators are critical to our performance. The design of the first row is the same as the one in Table 5. For another, no matter how the hand-crafted features are designed, the final performance is relatively robust. We suppose such robustness is led by our proposed weight prediction structure (FC Layer). This structure could handle the variations of our designed features and then produce the instance-aware weights. The performance change among various element-wise operations such as Log, Sigmoid or Tanh is marginal for $I_A$ or $I_P$. Using FPN's box regression style in $I_P$ obtains lower mAP about 0.7 than using FCOS's style.

Table 7 shows the influence using different weight prediction structures. It is clear that there is no big difference when using or not using a single fc layer. The performance by two layers is just 0.1 higher than a single layer. The reason is that our input is abstract enough with only 11 elements. Using more layers does not enhance information much in this setting. Moreover, the ReLU operation guarantees the non-linearity in our module.

**Table 7.** Ablation study for instance-aware weight prediction structure. The fc and ReLU mean fully connected layer and ReLU operation respectively. "+" means cascading two operations in sequence. $2\times$ means repeating the operation module by two times.

| structure | $AP$ | $AP_{50}$ | $AP_{75}$ |
|---|---|---|---|
| fc+ReLU | 32.5 | 51.9 | **33.8** |
| fc | 31.7 | 51.0 | 33.4 |
| $2\times$(fc+ReLU) | **32.6** | **52.0** | 33.6 |
| fc+ReLU+fc | 32.3 | 51.7 | 33.6 |

**Table 8.** Ablation study for information fusion between the instance-agnostic template and instance-aware weight. For feature addition and concatenation, we reshaped both our instance-agnostic feature and instance-aware weight to $16 \times 16 \times 9$.

| combination | $AP$ | $AP_{50}$ | $AP_{75}$ |
|---|---|---|---|
| addition | 31.1 | 50.9 | 33.5 |
| concatenation | 30.9 | 50.6 | 33.6 |
| $3\times3$ convolution | **32.5** | 51.8 | 33.7 |
| $1\times1$ convolution | **32.5** | **51.9** | **33.8** |

***Instance-aware Convolution*** Table 8 ablates different fusion ways between the instance-agnostic template and instance-aware weight. It is obvious that directly adding or concatenating them leads to a performance drop. The reason is that these two features come from different sources, and thus, have a distinct meaning. Specifically, the instance-agnostic template has strong semantic information, while instance-aware weight indicates shifting information between template and proposals. The difference makes using dynamic convolution more appropriate to fuse them – by the intuition of interpolating semantic information by geometric offsets to an extent. With this consideration, our design is effective in aligning instance-agnostic mask features given different box information. There is no performance difference whether $1 \times 1$ or $3 \times 3$ the convolution ker-

**Table 9.** Resuts on cityscapes dataset. All experiments are conducted with ResNet50.

| Method | $AP_{val}$ | $AP$ | $AP_{50}$ | person | rider | car | truck | bus | train | mcycle | bicycle |
|---|---|---|---|---|---|---|---|---|---|---|---|
| SGN [22] | 29.2 | 25.0 | 44.9 | 21.8 | 20.1 | 39.4 | **24.8** | **33.2** | **30.8** | 17.1 | 12.4 |
| Mask R-CNN [12] | 31.5 | 26.2 | 49.9 | 30.5 | 23.7 | 46.9 | 22.8 | 32.2 | 18.6 | 19.1 | 16.0 |
| Ours | **32.4** | **27.3** | **50.5** | **34.1** | **26.1** | **51.9** | 20.9 | 30.9 | 15.9 | **20.5** | **19.2** |

nel size is, which further decreases the computation cost of our instance-aware convolution.

### 4.2 Visualization

Figure 4 shows the visualization of our predicted instances. Compared with segmentation-based methods [34,22], our approach could easily detect the small objects with excellent pixel-level prediction. For large objects, the boundary is a little worse than that of small ones. It means a single-point feature in our PointINS could have compressed a level of detail information for large regions. It is caused by the pyramid features in object detection and instance segmentation frameworks. Like RetinaNet and FCOS, large objects are usually detected with P6 or P7 features. These features are highly abstract to estimate the location of large objects. However, they are not suitable for pixel-level instance segmentation because they lose many details by downsampling. Longer training time and data augmentation could relieve this difficulty as shown in Table 2.

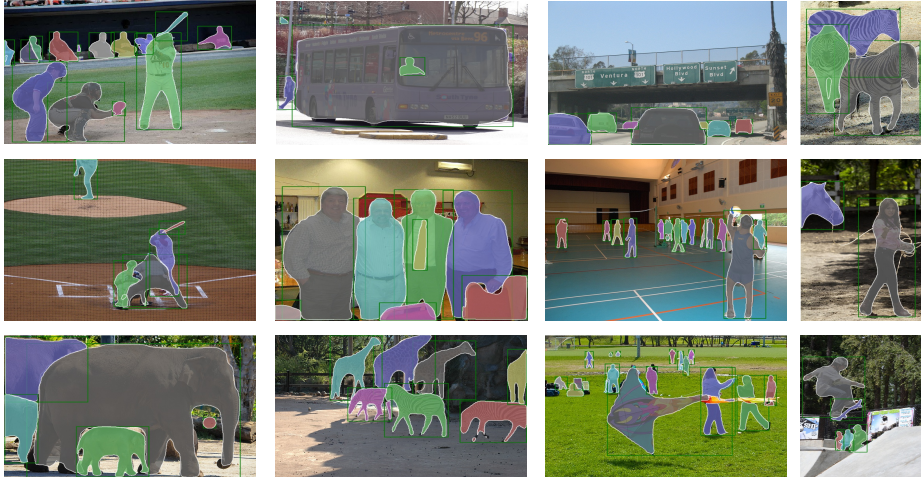### 4.3 Experiments on Cityscapes

***Dataset and Metrics*** Cityscapes [8] contains street scenes captured by carmounted cameras. There are 2,975 training, 500 validation and 1,525 testing images with fine annotation. We report our results on both validation and secret test subset. Eight semantic classes are annotated with instance masks. Each image is with size $1024 \times 2048$. We evaluate results based on AP and AP50.

***Hyper-parameters*** We use images with shorter edges randomly sampled from $\{800, 1024\}$ for training and those with shorter edge length 1024 for inference. We train our model with learning rate 0.0005 for 18k iterations and with rate 0.00005 in another 6k iterations. Eight images (1 image per GPU) are in one image batch.

***Results*** Table 9 shows the results of our PointINS with FCOS-Res50 backbone. Compared with RoI or segmentation based methods, PointINS also outperforms them with more simple but efficient structure.

## 5 Conclusion

Following current one-stage detectors, PointINS aims to segment masks with a single-point feature efficiently. The key idea is the alignment between the mask feature and various boxes. A core module is proposed in our framework – that is, instance-aware convolution. It has two inputs including instance-agnostic feature and instance-aware weight. The former produces a novel 3D tensor, providing sufficient information for potential masks. It is noted that the masks

**Fig. 4.** Visualization of RetinaNet-based PointINS with ResNet101 on COCO dataset under 1× training. For all images, we set mask confidence threshold 0.5.

generated from the same point share this template feature. As such, the latter aligns the feature of a specific region by instance-aware dynamic convolution with the template. The module of instance-aware convolution helps our PointINS achieve state-of-the-art performance on various datasets among point-based instance segmentation frameworks. PointINS can serve as a fundamental baseline for point-based instance segmentation for its high generality and practicality.

# References

1. Andrychowicz, M., Denil, M., Gomez, S., Hoffman, M.W., Pfau, D., Schaul, T., Shillingford, B., De Freitas, N.: Learning to learn by gradient descent by gradient descent. In: NeurIPS (2016)
2. Bai, M., Urtasun, R.: Deep watershed transform for instance segmentation. In: CVPR (2017)
3. Bolya, D., Zhou, C., Xiao, F., Lee, Y.J.: Yolact: Real-time instance segmentation. In: ICCV (2019)
4. Cai, Q., Pan, Y., Yao, T., Yan, C., Mei, T.: Memory matching networks for one-shot image recognition. In: CVPR (2018)
5. Chen, K., Pang, J., Wang, J., Xiong, Y., Li, X., Sun, S., Feng, W., Liu, Z., Shi, J., Ouyang, W., et al.: Hybrid task cascade for instance segmentation. In: CVPR (2019)
6. Chen, L.C., Papandreou, G., Kokkinos, I., Murphy, K., Yuille, A.L.: Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. In: PAMI (2016)
7. Chen, X., Girshick, R., He, K., Dollár, P.: Tensormask: A foundation for dense object segmentation. In: ICCV (2019)
8. Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R., Franke, U., Roth, S., Schiele, B.: The cityscapes dataset for semantic urban scene understanding. In: CVPR (2016)
9. Dai, J., Li, Y., He, K., Sun, J.: R-FCN: object detection via region-based fully convolutional networks. In: NeurIPS (2016)
10. Girshick, R., Donahue, J., Darrell, T., Malik, J.: Rich feature hierarchies for accurate object detection and semantic segmentation. In: CVPR (2014)
11. Girshick, R.B.: Fast R-CNN. In: ICCV (2015)
12. He, K., Gkioxari, G., Dollár, P., Girshick, R.B.: Mask R-CNN. In: ICCV (2017)
13. Hu, R., Dollár, P., He, K., Darrell, T., Girshick, R.: Learning to segment every thing. In: CVPR (2018)
14. Hu, X., Mu, H., Zhang, X., Wang, Z., Tan, T., Sun, J.: Meta-sr: A magnification-arbitrary network for super-resolution. In: CVPR (2019)
15. Huang, Z., Huang, L., Gong, Y., Huang, C., Wang, X.: Mask scoring r-cnn. In: CVPR (2019)
16. Lee, Y., Park, J.: Centermask: Real-time anchor-free instance segmentation. In: CVPR (2020)
17. Lemke, C., Budka, M., Gabrys, B.: Metalearning: a survey of trends and technologies. In: Artificial intelligence review (2015)
18. Li, Y., Qi, H., Dai, J., Ji, X., Wei, Y.: Fully convolutional instance-aware semantic segmentation. In: CVPR (2017)
19. Lin, T., Dollár, P., Girshick, R.B., He, K., Hariharan, B., Belongie, S.J.: Feature pyramid networks for object detection. In: CVPR (2017)
20. Lin, T.Y., Goyal, P., Girshick, R., He, K., Dollár, P.: Focal loss for dense object detection. In: ICCV (2017)
21. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: ECCV (2014)
22. Liu, S., Jia, J., Fidler, S., Urtasun, R.: Sgn: Sequential grouping networks for instance segmentation. In: ICCV (2017)
23. Liu, S., Qi, L., Qin, H., Shi, J., Jia, J.: Path aggregation network for instance segmentation. In: CVPR (2018)

24. Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C.Y., Berg, A.C.: Ssd: Single shot multibox detector. In: ECCV (2016)
25. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: CVPR (2015)
26. Morrison, D., Tow, A.W., McTaggart, M., Smith, R., Kelly-Boxall, N., Wade-McCue, S., Erskine, J., Grinover, R., Gurman, A., Hunn, T., et al.: Cartman: The low-cost cartesian manipulator that won the amazon robotics challenge. In: ICRA (2018)
27. Peng, S., Jiang, W., Pi, H., Bao, H., Zhou, X.: Deep snake for real-time instance segmentation. In: CVPR (2020)
28. Qi, L., Jiang, L., Liu, S., Shen, X., Jia, J.: Amodal instance segmentation with kins dataset. In: CVPR (2019)
29. Ravi, S., Larochelle, H.: Optimization as a model for few-shot learning. In: ICLR (2016)
30. Redmon, J., Divvala, S., Girshick, R., Farhadi, A.: You only look once: Unified, real-time object detection. In: CVPR (2016)
31. Ren, S., He, K., Girshick, R.B., Sun, J.: Faster R-CNN: towards real-time object detection with region proposal networks. In: NeurIPS (2015)
32. Shu, G.: Human detection, tracking and segmentation in surveillance video (2014)
33. Tian, Z., Shen, C., Chen, H., He, T.: Fcos: Fully convolutional one-stage object detection. In: ICCV (2019)
34. Wang, X., Kong, T., Shen, C., Jiang, Y., Li, L.: Solo: Segmenting objects by locations. In: arXiv (2019)
35. Wang, Y.X., Hebert, M.: Learning to learn: Model regression networks for easy small sample learning. In: ECCV (2016)
36. Xie, E., Sun, P., Song, X., Wang, W., Liu, X., Liang, D., Shen, C., Luo, P.: Polarmask: Single shot instance segmentation with polar representation. In: CVPR (2020)
37. Yang, T., Zhang, X., Li, Z., Zhang, W., Sun, J.: Metaanchor: Learning to detect objects with customized anchors. In: NeurIPS (2018)
38. Ying, H., Huang, Z., Liu, S., Shao, T., Zhou, K.: Embedmask: Embedding coupling for one-stage instance segmentation. In: arXiv (2019)
39. Zhao, H., Shi, J., Qi, X., Wang, X., Jia, J.: Pyramid scene parsing network. In: CVPR (2017)