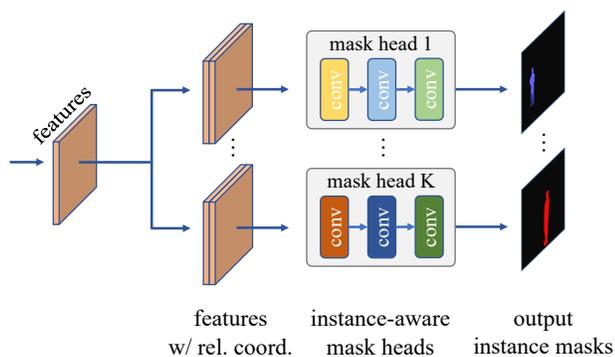# Instance and Panoptic Segmentation Using Conditional Convolutions

Zhi Tian,     Bowen Zhang,     Hao Chen,     Chunhua Shen

**Abstract**—We propose a simple yet effective framework for instance and panoptic segmentation, termed CondInst (conditional convolutions for instance and panoptic segmentation). In the literature, top-performing instance segmentation methods typically follow the paradigm of Mask R-CNN and rely on ROI operations (typically ROIAlign) to attend to each instance. In contrast, we propose to attend to the instances with dynamic conditional convolutions. Instead of using instance-wise ROIs as inputs to the instance mask head of fixed weights, we design dynamic instance-aware mask heads, conditioned on the instances to be predicted. CondInst enjoys three advantages: 1) Instance and panoptic segmentation are unified into a fully convolutional network, eliminating the need for ROI cropping and feature alignment. 2) The elimination of the ROI cropping also significantly improves the output instance mask resolution. 3) Due to the much improved capacity of dynamically-generated conditional convolutions, the mask head can be very compact (*e.g.*, 3 conv. layers, each having only 8 channels), leading to significantly faster inference time per instance and making the overall inference time less relevant to the number of instances. We demonstrate a simpler method that can achieve improved accuracy and inference speed on both instance and panoptic segmentation tasks. On the COCO dataset, we outperform a few state-of-the-art methods. We hope that CondInst can be a strong baseline for instance and panoptic segmentation. Code is available at: https://git.io/AdelaiDet

**Index Terms**—Fully convolutional networks, conditional convolutions, instance segmentation, panoptic segmentation

✦



**Fig. 1** – CondInst uses instance-aware mask heads to predict the mask for each instance. $K$ is the number of instances to be predicted. Note that each output map only contains the mask of one instance. The filters in the mask head vary with different instances, which are dynamically-generated and conditioned on the target instance. ReLU is used as the activation function (excluding the last conv. layer).

features
w/ rel. coord.

instance-aware
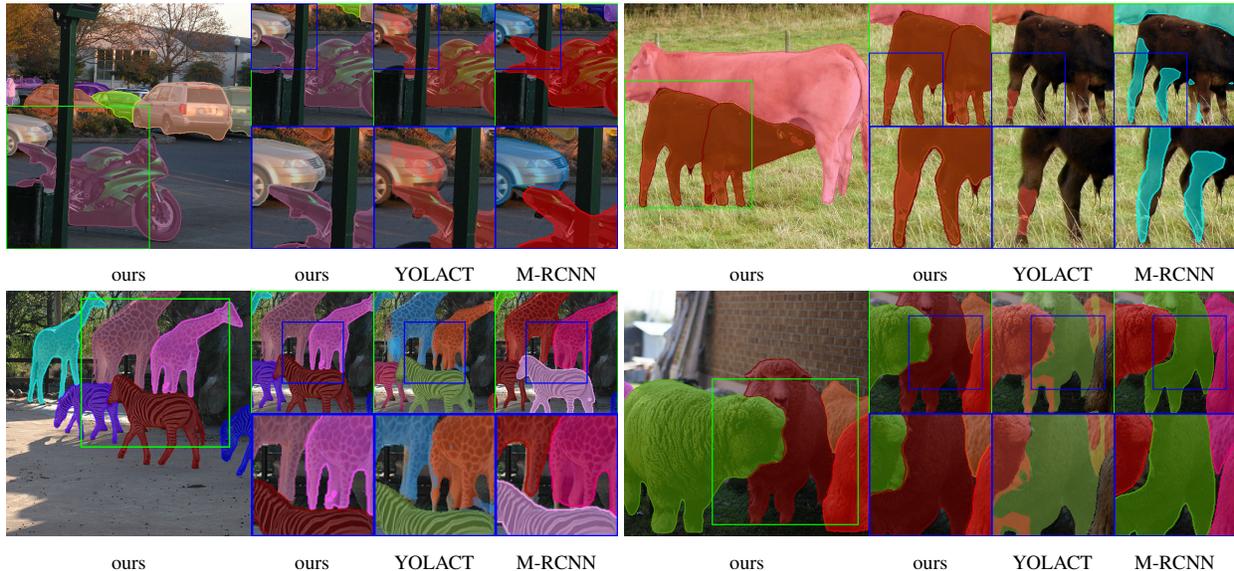mask heads

output
instance masks

## 1 INTRODUCTION

Instance segmentation is a fundamental yet challenging task in computer vision, which requires an algorithm to predict a per-pixel mask with a category label for each instance of interest in an image. Panoptic segmentation further requires the algorithm to segment the stuff (*e.g.*, sky and grass), assigning every pixel in the image a semantic label. Panoptic segmentation is often built on an instance segmentation

framework with an extra semantic segmentation branch. Therefore, both instance and panoptic segmentation share the same key challenge——how to efficiently and effectively distinguish individual instances.

Despite a few works being proposed recently, the dominant method tackling this challenge is still the two-stage method such as Mask R-CNN [2], which casts instance segmentation into a two-stage detection-and-segmentation task. To be specific, Mask R-CNN first employs an object detector Faster R-CNN to predict a bounding-box for each instance. Then for each instance, regions-of-interest (ROIs) are cropped from the networks' feature maps using the ROIAlign operation. To predict the final masks for each instance, a compact fully convolutional network (FCN) (*i.e.*, mask head) is applied to these ROIs to perform foreground/background segmentation. However, this ROI-based method may have the following drawbacks. 1) Since ROIs are often axis-aligned bounding-boxes, for objects with irregular shapes, they may contain an excessive amount of irrelevant image content including background and other instances. This issue may be mitigated by using rotated RoIs [3], but with the price of a more complex pipeline. 2) In order to distinguish between the foreground instance and the background stuff or instance(s) with the fixed mask head, the mask head needs a strong capacity and a relatively larger receptive field to encode sufficiently large context information. As a result, a stack of $3 \times 3$ convolutions is used in the mask head (*e.g.*, four $3 \times 3$ convolutions with 256 channels in Mask R-CNN). It considerably increases computational complexity of the mask head, resulting that the inference time significantly varies in the number of instances. 3) ROIs are typically of different sizes. In order to use effective batched computation in modern deep learning frameworks [4], [5], a resizing operation is often required

**Fig. 2** – Qualitative comparisons with other methods. We compare the proposed CondInst against YOLACT [1] and Mask R-CNN [2]. Our masks are generally of higher quality (*e.g.*, preserving finer details). Best viewed on screen.

to resize the cropped regions into patches of the same size. For instance, Mask R-CNN resizes all the cropped regions to $14 \times 14$ (upsampled to $28 \times 28$ using a deconvolution), which restricts the output resolution of instance segmentation, as large instances would require higher resolutions to retain details at the boundary.

In computer vision, the closest task to instance segmentation is semantic segmentation, for which fully convolutional networks (FCNs) have shown dramatic success [6]–[10]. FCNs also have shown excellent performance on many other per-pixel prediction tasks ranging from low-level image processing such as denoising, super-resolution; to mid-level tasks such as optical flow estimation and contour detection; and high-level tasks including recent single-shot object detection [11], monocular depth estimation [12]–[14] and counting [15]. However, almost all the instance segmentation methods based on FCNs[1] lag behind state-of-the-art ROI-based methods. Why do the versatile FCNs perform unsatisfactorily on instance segmentation? This is due to the fact that the FCNs tend to yield similar predictions for similar image appearance. As a result, the vanilla FCNs are incapable of distinguishing individual instances. For example, if two persons A and B with the similar appearance are in an input image, when predicting the instance mask of A, the FCN needs to predict B as background *w.r.t.* A, which can be difficult as they look similar in appearance. Therefore, an ROI operation is used to crop the person of interest, *i.e.*, A; and filter out B. Essentially, this is the core operation making the model attend to an instance.

In this work, we advocate a new solution for instance segmentation, termed CondInst. Instead of using ROIs, CondInst attends to each instance by using *instance-sensitive convolution filters* as well as relative coordinates that are appended to the feature maps. Specifically, unlike Mask R-CNN, which uses a standard convolution network with

a set of fixed convolutional filters as the mask head for predicting all instances, the network parameters in our mask head are adapted according to the instance to be predicted. Inspired by dynamic filtering networks [16] and CondConv [17], for each instance, a controller sub-network (see Fig. 3) dynamically generates the mask head's filters (conditioned on the center area of the instance), which is then used to predict the mask of this instance. It is expected that the network parameters can encode the characteristics (*e.g.*, relative position, shape and appearance) of this instance, and only fires on the pixels of this instance, which thus bypasses the difficulty in the standard FCNs. These conditional mask heads are applied to the whole high-resolution feature maps, *thus eliminating the need for ROI operations*. At the first glance, the idea may not work well as instance-wise mask heads may incur a large number of network parameters provided that some images contain as many as dozens of instances. However, as the mask head filters are only asked to predict the mask of only one instance, it largely eases the learning requirement and thus reduces the load of the filters. As a result, the mask head can be extremely light-weight. We will show that a very compact mask head with dynamically-generated filters can already outperform previous ROI-based Mask R-CNN. This compact mask head also results in much reduced computational complexity per instance than that of the mask head in Mask R-CNN.

We summarize our main contributions as follow.

- We attempt to solve instance segmentation from a new perspective that uses dynamic mask heads. This novel solution achieves improved instance segmentation performance than existing methods such as Mask R-CNN. To our knowledge, this is the first time that a new instance segmentation framework outperforms recent state-of-the-art both in accuracy and speed.
- CondInst is fully convolutional and avoids the aforementioned resizing operation used in many existing

---

1. By FCNs, we mean the vanilla FCNs in [6] that only involve convolutions and pooling.

methods, as CondInst does not rely on ROI operations. Without having to resize feature maps leads to high-resolution instance masks with more accurate edges, as shown in Fig. 2.

- Since the mask head in CondInst is very compact and light-weight, compared with the box detector FCOS, CondInst needs only ∼10% more computational time (less than 5 milliseconds) to obtain the mask results of all the instances, even when processing the maximum number of instances per image (*i.e.*, 100 instances). As a result, the overall inference time is stable as it almost does not depend on the number of instances in the image.
- With an extra semantic segmentation branch, CondInst can be easily extended to panoptic segmentation [18], resulting a unified fully convolutional network for both instance and panoptic segmentation tasks.
- CondInst achieves state-of-the-art performance on both instance and panoptic segmentation tasks while being fast and simple. We hope that CondInst can be a new strong alternative for instance and panoptic segmentation tasks, as well as other instance-level recognition tasks such as keypoint detection.

## 2 RELATED WORK

Here we review some work that is most relevant to ours.

**Conditional Convolutions/Dynamic filters.** Unlike traditional convolutional layers, which have fixed filters once trained, the filters of conditional convolutions are conditioned on the input and are dynamically generated by another network (*i.e.*, a controller). This idea has been explored previously in dynamic filter networks [16] and CondConv [17] mainly for the purpose of increasing the capacity of a classification network. DGMN [19] also employs dynamic filters to generate the node-specific filters for message calculation, which improves the capacity of the networks and thus results in better performance. In this work, we extend this idea to generate the mask head's filters conditioned on each instance, and present a high-performance instance segmentation method without the need for ROIs.
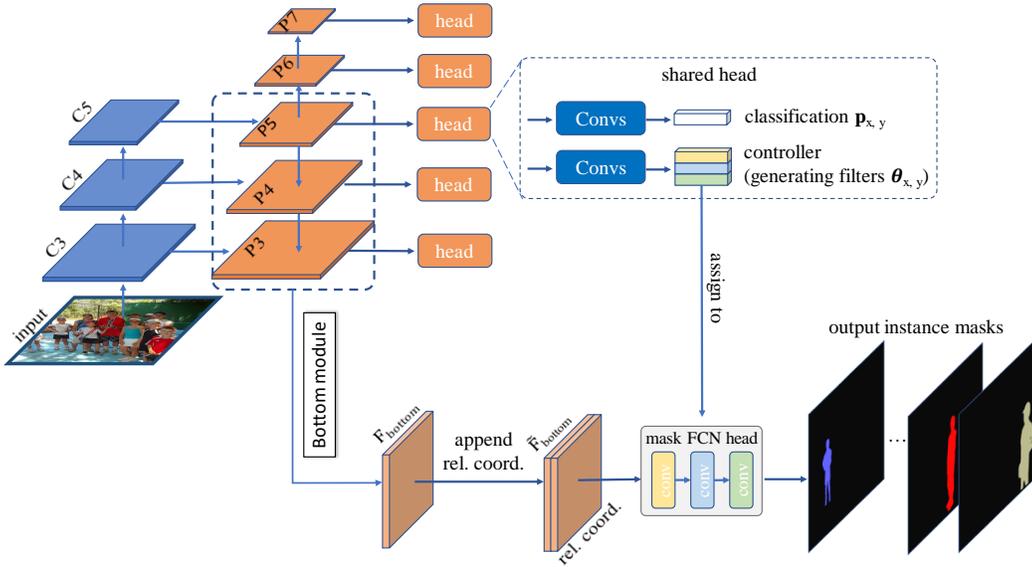
**Instance Segmentation.** To date, the dominant framework for instance segmentation is still Mask R-CNN. Mask R-CNN first employs an object detector to detect the bounding-boxes of instances (*e.g.*, ROIs). With these bounding-boxes, an ROI operation is used to crop the features of the instance from the feature maps. Finally, an FCN head is used to obtain the desired instance masks. Many works [20]–[22] with top performance are built on Mask R-CNN. Moreover, some works have explored to apply the standard FCNs [6] to instance segmentation. InstanceFCN [23] may be the first instance segmentation method that is fully convolutional. InstanceFCN proposes to predict position-sensitive score maps with vanilla FCNs. Afterwards, these score maps are assembled to obtain the desired instance masks. Note that InstanceFCN does not work well with overlapping instances. Others [24]–[26] attempt to first perform image segmentation and then the desired instance masks are formed by assembling the pixels of the same instance. Deep Watershed [27] models instance

segmentation with classical watershed transform, and object instances can be viewed as the energy basins in the energy map of the watershed transform of an image. SGN [28] uses a sequence of networks to gradually group the raw pixels to line segments, connected components, and finally object instances, achieving impressive performance. The single-shot Box2Pix [29] solves instance segmentation in the bottom-up fashion. Novotny *et al.* [30] propose semi-convolutional operators to make FCNs applicable to instance segmentation. Arnab *et al.* [31] propose dynamically instantiated CRF (Conditional Random Field) for instance segmentation, which is able to produce a variable number of instances per image. To our knowledge, thus far none of these methods can outperform Mask R-CNN both in accuracy and speed on the COCO benchmark dataset.

The recent YOLACT [1] and BlendMask [32] may be viewed as a reformulation of Mask RCNN, which decouples ROI detection and feature maps used for mask prediction. Wang *et al.* developed a simple FCN based instance segmentation method, which segments the instances by the their locations, showing competitive performance [33], [34]. PolarMask [35] developed a new simple mask representation for instance segmentation, which extends the bounding box detector FCOS [11].

**Panoptic segmentation.** There are two main approaches for solving this task. The first one is the bottom-up approach. It tackles the task as a semantic segmentation at first and then uses clustering/grouping methods to assemble the pixels into individual instances or stuff [36], [37]. The authors of [37] also explore the weakly- or semi-supervised panoptic segmentation. The second approach is the top-down approach, which is often built on top-down instance segmentation methods. For example, Panoptic-FPN [38] extends an additional semantic segmentation branch from Mask R-CNN and combines the results with the instance segmentation results generated by Mask R-CNN [18]. Moreover, attention based methods recently gain much popularity in many computer vision tasks, which provide a new approach to panoptic segmentation. Axial DeepLab [39] used a carefully designed module to enable attention to be applied to large-size images for panoptic segmentation. CondInst can easily be applied to panoptic segmentation following the top-down approaches. We empirically observe that the quality of the instance segmentation results may be the dominant factor to the final performance. Thus in CondInst, without bells and whistles, by simply applying the same method used by Panoptic-FPN, the panoptic segmentation performance of CondInst is already competitive compared to the state-of-the-art panoptic segmentation methods.

Additionally, AdaptIS [40] recently proposes to solve panoptic segmentation with FiLM [41]. The idea shares some similarity with CondInst in which information about an instance is encoded in the coefficients generated by FiLM. Since only the batch normalization coefficients are dynamically generated, AdaptIS needs a large mask head to achieve good performance. In contrast, CondInst directly encodes them into the conv. filters of the mask head, which is much more straightforward and efficient. Also, as shown in experiments, CondInst can achieve much better panoptic segmentation accuracy than AdaptIS, which suggests that CondInst is much more effective.

**Fig. 3** – The overall architecture of **CondInst**. $C_3$, $C_4$ and $C_5$ are the feature maps of the backbone network (*e.g.*, ResNet-50). $P_3$ to $P_7$ are the FPN feature maps as in [11], [42]. $\mathbf{F}_{bottom}$ is the bottom branch's output, whose resolution is the same as that of $P_3$. Following [32], the bottom branch aggregates the feature maps $P_3$, $P_4$ and $P_5$. $\tilde{\mathbf{F}}_{bottom}$ is obtained by concatenating the relative coordinates to $\mathbf{F}_{bottom}$. The classification head predicts the class probability $\boldsymbol{p}_{x,y}$ of the target instance at location $(x, y)$, same as in FCOS. The controller generates the filter parameters $\boldsymbol{\theta}_{x,y}$ of the mask head for the instance. Similar to FCOS, there are also center-ness and box heads in parallel with the controller (not shown in the figure for simplicity). Note that the heads in the dashed box are repeatedly applied to $P_3 \cdots P_7$. The mask head is instance-aware, and is applied to $\tilde{\mathbf{F}}_{bottom}$ as many times as the number of instances in the image (refer to Fig. 1).

# 3 OUR METHODS: INSTANCE AND PANOPTIC SEGMENTATION WITH CONDINST

We first present CondInst for instance segmentation, and then we show how the instance segmentation framework can be easily extended to panoptic segmentation by using a new semantic branch.

## 3.1 Overall Architecture for Instance Segmentation

Given an input image $I \in \mathbb{R}^{H \times W \times 3}$, the goal of instance segmentation is to predict the pixel-level mask and the category of each instance of interest in the image. The ground-truths are defined as $\{(\mathbf{M}_i, c_i)\}$, where $\mathbf{M}_i \in \{0, 1\}^{H \times W}$ is the mask for the $i$-th instance and $c_i \in \{1, 2, ..., C\}$ is the category. $C$ is 80 on MS-COCO [43]. In semantic segmentation, the prediction target of each pixel are well-defined, which is the semantic category of the pixel. In addition, the number of categories is known and fixed. Thus, the outputs of semantic segmentation can be easily represented with the output feature maps of the FCNs, and each channel of the output feature maps corresponds to a class. However, in instance segmentation, the prediction target of each pixel is hard to define because instance segmentation also requires to distinguish individual instances, but the number of instances changes in different images. This poses a major challenge when applying traditional FCNs [6] to instance segmentation.

In this work, our core idea is that for an image with $K$ instances, $K$ different mask heads will be dynamically generated, and each mask head will contain the characteristics of its target instance in their filters. As a result, when the mask is applied to an input, it will only fire on the

pixels of the instance, thus producing the mask prediction of the instance and distinguishing individual instances. We illustrate the process in Fig. 1. The instance-aware filters are generated by modifying an object detector. Specifically, we add a new controller branch to generate the filters for the target instance of each box predicted by the detector, as shown in Fig. 3. Therefore, the number of the dynamic mask heads is the same as the number of the predicted boxes, which should be the number of the instances in the image if the detector works well. In this work, we build CondInst on the popular object detector FCOS [11] due to its simplicity and flexibility. Also, the elimination of anchor-boxes in FCOS can also save the number of parameters and the amount of computation.

As shown in Fig. 3, following FCOS [11], we make use of the feature maps $\{P_3, P_4, P_5, P_6, P_7\}$ of feature pyramid networks (FPNs) [42], whose down-sampling ratios are 8, 16, 32, 64 and 128, respectively. As shown in Fig. 3, on each feature level of the FPN, some functional layers (in the dash box) are applied to make instance-aware predictions. For example, the class of the target instance and the dynamically-generated filters for the instance. In this sense, CondInst can be viewed as the same as Mask R-CNN, both of which first attend to instances in an image and then predict the pixel-level masks of the instances (*i.e.*, instance-first).

Moreover, recall that Mask R-CNN employs an object detector to predict the bounding-boxes of the instances in the input image. The bounding-boxes are actually the way that Mask R-CNN represents instances. Similarly, CondInst employs the instance-aware filters to represent the instances. In other words, instead of encoding the instance information

with the bounding-boxes, CondInst implicitly encodes it with the parameters of the generated dynamic filters, which is much more flexible. For example, the dynamic filters can easily represent the irregular shapes that are hard to be tightly enclosed by a bounding-box (elaborated in Sec. 4.4). This is one of CondInst's advantages over the previous ROI-based methods.

Besides the detector, as shown in Fig. 3, there is also a bottom branch, which provides the feature maps (denoted by $\mathbf{F}_{bottom}$) that our generated mask heads take as inputs to predict the desired instance mask. The bottom branch aggregates the FPN feature maps $P_3$, $P_4$ and $P_5$. To be specific, $P_4$ and $P_5$ are upsampled to the resolution of $P_3$ with bilinear interpolation and added to $P_3$. After that, four $3 \times 3$ convolutions with 128 channels are applied. The resolution of the resulting feature maps is the same as $P_3$ (*i.e.*, $\frac{1}{8}$ of the input image resolution). Finally, another convolutional layer is used to reduce the number of the output channels $C_{bottom}$ from 128 to 8, resulting in the bottom feature $\mathbf{F}_{bottom}$. The small output channel reduces the number of the generated parameters. We empirically found that using $C_{bottom} = 8$ can already achieve good performance, and as shown in our experiments, a larger $C_{bottom}$ here (*e.g.*, 16) cannot improve the performance. Even more aggressively, using $C_{bottom} = 1$ only degrades the performance by $\sim 1\%$ in mask AP. It is probably because our mask heads only predict relatively simple class-agnostic instance masks and most of the information of an instance has been encoded in the dynamically generated filters.

As mentioned before, the generated filters can also encode the shape and position of the target instance. Since the CNN feature maps do not generally convey the position information, a map of the coordinates needs to be appended to $F_{bottom}$ such that the generated filters are aware of positions. As the filters are generated with the location-agnostic convolutions, they can only (implicitly) encode the shape and position with the coordinates relative to the location where the filters are generated (*i.e.*, using the coordinate system with the location as the origin). Thus, as shown in Fig. 3, $\mathbf{F}_{bottom}$ is combined with a map of the relative coordinates, which are obtained by transforming all the locations on $\mathbf{F}_{bottom}$ to the coordinate system with the location generating the filters as the origin. Then, the combination is sent to the mask head to predict the instance mask in the fully convolutional fashion. The relative coordinates provide a strong cue for predicting the instance mask, as shown in our experiments. It is also interesting to note that even if the generated mask heads only take as input the map of the relative coordinates, a modest performance can be obtained as shown in the experiments. This empirically proves that the generated filters indeed encode the shape and position of the target instance. Finally, sigmoid is used as the last layer of the mask head and obtains the mask scores. The mask head only classifies the pixels as the foreground or background. The class of the instance is predicted by the classification head of the detector, as shown in Fig. 3.

The resolution of the original mask prediction is same as the resolution of $\mathbf{F}_{bottom}$, which is $^1/_8$ of the input image resolution. In order to improve the resolution of instance masks, we use bilinear interpolation to upsample the mask prediction by 2, resulting in $200 \times 256$ instance masks (if the input image size is $800 \times 1024$). The mask's resolution is much higher than that of Mask R-CNN (only $28 \times 28$ as mentioned before).
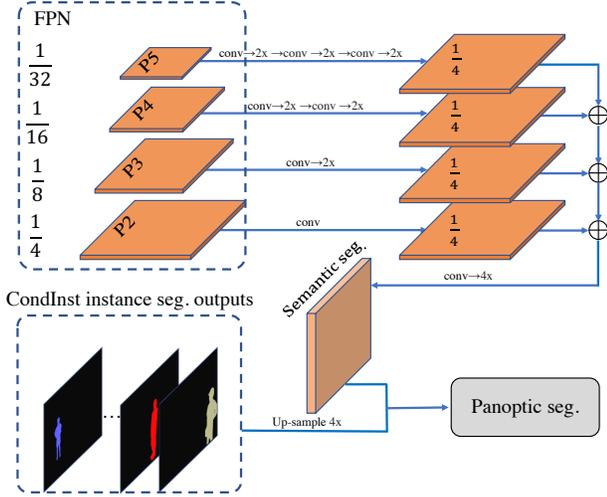
## 3.2 Network Outputs and Training Targets

Similar to FCOS, each location on the FPN's feature maps $P_i$ either is associated with an instance, thus being a positive sample, or is considered as a negative sample. The associated instance and label for each location are determined as follows.

Let us consider the feature maps $P_i \in \mathbb{R}^{H \times W \times C}$ and let $s$ be its down-sampling ratio. As shown in previous works [11], [44], [45], a location $(x, y)$ on the feature maps can be mapped back onto the input image as $(\lfloor s/2 \rfloor + xs, \lfloor s/2 \rfloor + ys)$. If the mapped location falls in the center region of an instance, the location is considered to be responsible for the instance. Any locations outside the center regions are labeled as negative samples. The center region is defined as the box $(c_x - rs, c_y - rs, c_x + rs, c_y + rs)$, where $(c_x, c_y)$ denotes the mass center of the instance mask, $s$ is the down-sampling ratio of $P_i$ and $r$ is a constant scalar being 1.5 as in FCOS [11]. As shown in Fig. 3, at a location $(x, y)$ on $P_i$, CondInst has the following output heads.

**Classification Head.** The classification head predicts the class of the instance associated with the location. The ground-truth target is the instance's class $c_i$ or 0 (*i.e.*, background). As in FCOS, the network predicts a $C$-D vector $\boldsymbol{p}_{x,y}$ for the classification and each dimension of $\boldsymbol{p}_{x,y}$ corresponds to a binary classifier, where $C$ is the number of categories.

**Controller Head.** The controller head, which has the same architecture as the classification head, is used to predict the parameters of the conv. filters of the mask head for the instance at the location. The mask head predicts the mask for this particular instance. This is the core contribution of our work. To predict the parameters, we concatenate all the parameters of the filters (*i.e.*, weights and biases) together as an $N$-D vector $\boldsymbol{\theta}_{x,y}$, where $N$ is the total number of the parameters. Accordingly, the controller head has $N$ output channels. The mask head is a very compact FCN architecture, which has three $1 \times 1$ convolutions, each having 8 channels and using ReLU as the activation function except for the last one. No normalization layer such as batch normalization is used here. The last layer has 1 output channel and uses sigmoid to predict the probability of being foreground. The mask head has 169 parameters in total ($\#weights = (8+2) \times 8(conv1) + 8 \times 8(conv2) + 8 \times 1(conv3)$ and $\#biases = 8(conv1) + 8(conv2) + 1(conv3)$). The masks predicted by the mask heads are supervised with the ground-truth instance masks, which pushes the controller to generate the correct filters.

**Box Head.** The box head is the same as that in FCOS, which predicts a 4-D vector encoding the four distances from the location to the four boundaries of the bounding-box of the target instance. Conceptually, CondInst can eliminate the box head since CondInst needs no ROIs. However, we note that if we make use of box-based NMS, the inference time will be much reduced since we only need to compute the masks for the instances kept after box NMS. Thus, we still predict boxes in CondInst. We would like to highlight that the predicted boxes are *only* used in NMS and do not involve

**Fig. 4** – Illustration of CondInst for panoptic segmentation by attaching a semantic segmentation branch. The semantic segmentation branch follows [38]. Results from the instance segmentation and segmentation segmentation branches are combined together using the same post-processing as in [18].

any ROI operations. Moreover, as shown in Table 6, the box prediction can be removed if other kinds of NMS are used (*e.g.*, mask NMS [34]). This is fundamentally different from previous ROI-based methods, in which the box prediction is mandatory.

**Center-ness Head.** Like FCOS [11], at each location, we also predict a center-ness score. The center-ness score depicts how the location deviates from the center of the target instance. In inference, it is used to down-weight the boxes predicted by the locations far from the center as these boxes might be unreliable. The ground-truth center-ness score is defined as

$$\text{centerness}^* = \sqrt{\frac{\min(l^*, r^*)}{\max(l^*, r^*)} \cdot \frac{\min(t^*, b^*)}{\max(t^*, b^*)}}, \quad (1)$$

where $l^*, r^*, t^*$ and $b^*$ denote the distances from the location to the four boundaries of the ground-truth bounding box. We used the binary cross entropy (BCE) loss to supervise center-ness score as in FCOS.

**Semantic Branch for Panoptic Segmentation.** As mentioned before, we can extend CondInst to panoptic segmentation by adding a new semantic segmentation branch. For the semantic segmentation branch, we use the structure from Panoptic-FPN [38]. To be specific, as shown in Fig. 4, the semantic segmentation branch takes as inputs the feature maps $\{P_2, P_3, P_4, P_5\}$ of FPNs. $\{P_3, P_4, P_5\}$ are up-sampled to the same resolution as $P_2$ and the four feature maps are concatenated together. The resolution of $P_2$ is $1/4$ of the input image, which is also the same as the instance masks predicted by CondInst. Then, it is followed by a $1 \times 1$ convolution and softmax to obtain the semantic segmentation classification scores. The classification scores are trained with the cross-entropy loss. In inference, the semantic segmentation results are merged with the above instance masks to generate the final panoptic segmentation results. The details can be found in Sec. 3.4.

## 3.3 Loss Functions

Formally, the overall loss function of CondInst can be formulated as,

$$L_{overall} = L_{fcos} + \lambda L_{mask} + \mu L_{pano}, \quad (2)$$

where $L_{fcos}$ and $L_{mask}$ denote the original loss of FCOS and the loss for instance masks, respectively. $L_{pano}$ (only available in panoptic segmentation) is the loss for the semantic branch of panoptic segmentation. $\lambda$ and $\mu$ being $1$ and $0.5$, respectively, is used to balance these losses. $L_{fcos}$ is the same as in FCOS. Specifically, $L_{fcos}$ includes the classification head, the box regression head and the center-ness head, which are trained with the focal loss [46], the GIoU loss, and the binary cross-entropy (BCE) loss, respectively. $L_{mask}$ is defined as,

$$L_{mask}(\{\boldsymbol{\theta}_{x,y}\}) = \frac{1}{N_{pos}} \sum_{x,y} \mathbb{1}_{\{c^*_{x,y}>0\}} L_{dice}\Big(\mathbf{M}_{x,y}, \mathbf{M}^*_{x,y}\Big), \quad (3)$$

where $c^*_{x,y}$ is the classification label of location $(x, y)$, which is the class of the instance associated with the location or $0$ (*i.e.*, background) if the location is not associated with any instance. $N_{pos}$ is the number of locations where $c^*_{x,y} > 0$. $\mathbb{1}_{\{c^*_{x,y}>0\}}$ is the indicator function, being $1$ if $c^*_{x,y} > 0$ and $0$ otherwise. $\mathbf{M}^*_{x,y} \in \{0, 1\}^{H \times W}$ is the ground-truth mask of the instance associated with location $(x, y)$, and $\mathbf{M}_{x,y}$ is the mask predicted by the dynamic mask head of location $(x, y)$. Formally,

$$\mathbf{M}_{x,y} = MaskHead(\tilde{\mathbf{F}}_{x,y}; \boldsymbol{\theta}_{x,y}), \quad (4)$$

where $\boldsymbol{\theta}_{x,y}$ is the generated filters' parameters at location $(x, y)$. $\tilde{\mathbf{F}}_{x,y} \in \mathbb{R}^{H_{bottom} \times W_{bottom} \times (C_{bottom}+2)}$ is the combination of $\mathbf{F}_{bottom}$ and a map of coordinates $\mathbf{O}_{x,y} \in \mathbb{R}^{H_{bottom} \times W_{bottom} \times 2}$. As described before, $\mathbf{O}_{x,y}$ is the relative coordinates from all the locations on $\mathbf{F}_{bottom}$ to $(x, y)$ (*i.e.*, the location where the filters are generated). $MaskHead$ consists of a stack of convolutions with dynamic parameters $\boldsymbol{\theta}_{x,y}$.

Moreover, $L_{dice}$ is the Dice loss as in [47], which is used to overcome the foreground-background sample imbalance. We do not employ focal loss here as it requires to initialize the biases with a prior probability [46], which is not trivial if the parameters are dynamically generated. Formally, $L_{dice}$ is defined as

$$L_{dice}(\mathbf{M}, \mathbf{M}^*) = 1 - \frac{2 \sum_{i,j} \mathbf{M}_{i,j} \mathbf{M}^*_{i,j}}{\sum_{i,j} (\mathbf{M}_{i,j})^2 + \sum_{i,j} (\mathbf{M}^*_{i,j})^2}, \quad (5)$$

where $\mathbf{M}_{i,j}$ or $\mathbf{M}^*_{i,j}$ denotes the elements of $\mathbf{M}_{x,y}$ or $\mathbf{M}^*_{x,y}$, and the subscript $(x, y)$ is omitted for clarification. Note that, in order to compute the loss between the predicted mask $\mathbf{M}_{x,y}$ and the ground-truth mask $\mathbf{M}^*_{x,y}$, they need to have the same sizes. As mentioned before, the resolution of the predicted mask $\mathbf{M}_{x,y}$ is $1/4$ of the ground-truth mask $\mathbf{M}^*_{x,y}$. Thus, we down-sample $\mathbf{M}^*_{x,y}$ by 4 to make the sizes equal. The operation is omitted in Eq. (5) for clarification.

By design, all the positive locations on the feature maps should be used to compute the mask loss. For the images having hundreds of positive locations, the model would consume a large amount of memory. Therefore, in our preliminary version [48], the positive locations used in

computing the mask loss are limited up to 500 per GPU (*i.e.*, 250 per image and we have two images on one GPU). If there are more than 500 positive locations, 500 locations will be randomly chosen. In this version, instead of randomly choosing the 500 locations, we first rank the locations by the scores predicted by the FCOS detector, and then choose the locations with top scores for each instance. As a result, the number of locations per image can be reduced to $64$. This strategy works equally well and further reduces the memory footprint. For instance, using this strategy, the ResNet-50 based CondInst can be trained with 4 1080Ti GPUs.

Moreover, as shown in YOLACT [1] and BlendMask [32], during training, the instance segmentation task can benefit from a joint semantic segmentation task (*i.e.*, using instance masks as semantic labels). Thus, we also conduct experiments with the joint semantic segmentation task, showing improved performance. However, unless explicitly specified, all the experiments in the paper are *without* the semantic segmentation task. If used, the semantic segmentation loss is added to $L_{overall}$.

## 3.4 Inference

**Instance Segmentation.** Given an input image, we forward it through the network to obtain the outputs including classification confidence $\boldsymbol{p}_{x,y}$, center-ness scores, box prediction $\boldsymbol{t}_{x,y}$ and the generated parameters $\boldsymbol{\theta}_{x,y}$. We first follow the steps in FCOS to obtain the box detections. Afterwards, box-based NMS with the threshold being $0.6$ is used to remove duplicated detections and then the top $100$ boxes are used to compute masks. Note that each box is also associated with a group of filters generated by the controller. Let us assume that $K$ boxes remain after the NMS, and thus we have $K$ groups of generated filters. The $K$ groups of filters are used to produce $K$ instance-specific mask heads. These instance-specific mask heads are applied, in the fashion of FCNs, to $\tilde{\mathbf{F}}_{x,y}$ (*i.e.*, the combination of $\mathbf{F}_{bottom}$ and $\mathbf{O}_{x,y}$) to predict the masks of the instances. Since the mask head is a very compact network (having three $1 \times 1$ convolutions with $8$ channels and $169$ parameters in total), the overhead of computing masks is extremely small. For example, even with $100$ detections (*i.e.*, the maximum number of detections per image on MS-COCO), only less than $5$ milliseconds in total are spent on the mask heads, which only adds $\sim 10\%$ computational time to the base detector FCOS. In contrast, the mask head of Mask R-CNN has four $3 \times 3$ convolutions with $256$ channels, thus having more than 2.3M parameters and taking longer computational time.

**Panoptic Segmentation.** For panoptic segmentation, we follow [38] to combine instance and semantic results to obtain the panoptic results. We first rank the instance results from CondInst by their confidence scores generated by FCOS. The results with their scores less than $0.45$ are discarded. When overlaps occur between the instance masks, the overlap areas are attributed to the instance with higher score. Moreover, the instance that loses more than 40% of its total area due to the overlap with other higher-score-instances is discarded. Finally, the semantic results are filled to the areas that are not occupied by any instance.

# 4 EXPERIMENTS

We evaluate CondInst on the large-scale benchmark MS-COCO [43]. Following the common practice [2], [11], [46], our models are trained with split `train2017` (115K images) and all the ablation experiments are evaluated on split `val2017` (5K images). Our main results are reported on the `test-dev` split (20K images).

## 4.1 Implementation Details

Unless specified, we make use of the following implementation details. Following FCOS [11], ResNet-50 is used as our backbone network and the weights pre-trained on ImageNet [49] are used to initialize it. For the newly added layers, we initialize them as in [11]. Our models are trained with stochastic gradient descent (SGD) over $8$ V100 GPUs for 90K iterations with the initial learning rate being $0.01$ and a mini-batch of $16$ images. The learning rate is reduced by a factor of 10 at iteration $60K$ and $80K$, respectively. Weight decay and momentum are set as $0.0001$ and $0.9$, respectively. Following `Detectron2` [3], the input images are resized to have their shorter sides in $[640, 800]$ and their longer sides less or equal to 1333 during training. Left-right flipping data augmentation is also used during training. When testing, we do not use any data augmentation and only the scale of the shorter side being $800$ is used. The inference time in this work is measured on a single V100 GPU with 1 image per batch.

## 4.2 Architectures of the Mask Head

In this section, we discuss the design choices of the mask head in CondInst. We show that the performance is not sensitive to the architectures of the mask head. Our baseline is the mask head of three $1 \times 1$ convolutions with $8$ channels (*i.e.*, width $= 8$). As shown in Table 1 (3rd row), it achieves $35.6\%$ in mask AP. Next, we first conduct experiments by varying the depth of the mask head. As shown in Table 1a, apart from the mask head with depth being 1, all other mask heads (*i.e.*, depth $= 2, 3$ and $4$) attain similar performance. The mask head with depth being 1 achieves inferior performance as in this case the mask head is actually a linear mapping, which has overly weak capacity and cannot encode the complex shapes of the instances. Moreover, as shown in Table 1b, varying the width (*i.e.*, the number of the channels) does not result in a remarkable performance change either as long as the width is in a reasonable range. We also note that our mask head is extremely light-weight as the filters in our mask head are dynamically generated. As shown in Table 1, our baseline mask head only takes $4.5$ ms per 100 instances (the maximum number of instances on MS-COCO), which suggests that our mask head only adds small computational overhead to the base detector. Moreover, our baseline mask head only has 169 parameters in total. In sharp contrast, the mask head of Mask R-CNN [2] has more than 2.3M parameters and takes $\sim 2.5 \times$ computational time ($11.4$ ms per 100 instances).

## 4.3 Design Choices of the Bottom Module

We further investigate the impact of the bottom module. We first change $C_{bottom}$, which is the number of channels

| depth | time | AP | AP$_{50}$ | AP$_{75}$ | AP$_S$ | AP$_M$ | AP$_L$ |
|---|---|---|---|---|---|---|---|
| 1 | **2.2** | 30.5 | 52.7 | 30.7 | 13.7 | 32.8 | 44.9 |
| 2 | 3.3 | 35.5 | 56.2 | **37.9** | 17.1 | 38.8 | **51.2** |
| 3 | 4.5 | **35.6** | **56.4** | 37.9 | **18.0** | 38.9 | 50.8 |
| 4 | 5.6 | **35.6** | 56.3 | 37.8 | 17.3 | 38.9 | 51.0 |

(a) Varying the depth (width = 8).

| width | time | AP | AP$_{50}$ | AP$_{75}$ | AP$_S$ | AP$_M$ | AP$_L$ |
|---|---|---|---|---|---|---|---|
| 2 | **2.5** | 33.9 | 55.3 | 35.8 | 15.8 | 37.0 | 48.6 |
| 4 | 2.6 | 35.4 | 56.3 | 37.4 | 16.9 | 38.7 | **51.2** |
| 8 | 4.5 | 35.6 | **56.4** | 37.9 | **18.0** | 39.1 | 50.8 |
| 16 | 4.7 | **35.7** | 56.1 | **38.1** | 16.9 | 39.0 | 50.8 |

(b) Varying the width (depth = 3).

**TABLE 1** – Instance segmentation results with different architectures of the mask head on the MS-COCO `val2017` split. "depth": the number of layers in the mask head. "width": the number of channels of these layers. "time": the milliseconds that the mask head takes for processing 100 instances.

| $C_{bottom}$ | AP | AP$_{50}$ | AP$_{75}$ | AP$_S$ | AP$_M$ | AP$_L$ |
|---|---|---|---|---|---|---|
| 1 | 34.7 | 56.0 | 36.8 | 16.5 | 37.9 | 50.1 |
| 2 | 34.9 | 55.7 | 37.2 | 16.5 | 38.3 | 50.6 |
| 4 | 35.5 | 56.3 | 37.5 | 17.8 | 38.7 | 50.7 |
| 8 | **35.6** | **56.4** | **37.9** | **18.0** | **39.1** | 50.8 |
| 16 | 35.4 | 56.0 | 37.5 | 16.9 | 38.7 | **50.9** |

**TABLE 2** – Instance segmentation results by varying the number of channels of the bottom branch's output (*i.e.*, $C_{bottom}$) on the MS-COCO `val2017` split. The performance keeps almost the same if $C_{bottom}$ is in a reasonable range, which suggests that CondInst is robust to the design choice.

of the mask branch's output feature maps (*i.e.*, $\mathbf{F}_{bottom}$). As shown in Table 2, as long as $C_{bottom}$ is in a reasonable range (*i.e.*, from 2 to 16), the performance keeps almost the same. $C_{bottom} = 8$ is optimal and thus we use $C_{bottom} = 8$ in all other experiments by default.

We conduct experiments by varying the input FPN features of the bottom module. Specifically, we change the FPN feature level from $P_3$ (stride being 8) to $P_2$ (stride being 4) for the bottom module. As shown in Table 4, this can improve the mask AP from 35.6% to 36.0% with 20% more inference time. Moreover, as mentioned before, before taken as the input of the mask heads, the bottom module's output $\mathbf{F}_{bottom}$ is concatenated with a map of relative coordinates, which provides a strong cue for the mask prediction. As shown in Table 3 (2nd row), the performance drops significantly if the relative coordinates are removed (35.6% vs. 31.5%). We also experiment with the absolute coordinates, but it cannot largely boost the performance as shown in Table 3 (32.0%). This is understandable because an instance segmentation model should be translation-equivalence. Besides, as shown in Table 3 (2rd row), only using the relative coordinates can also obtain decent performance (31.3% in mask AP). The qualitative results are shown in 5.

### 4.4 What the Generated Filters Encode?

It is not straightforward to see what the generated filters encode. However, this can be analyzed by varying the inputs of the dynamic filters and visualizing the changes of the results. As shown in 5, it can be noted that if the mask heads only take the relative coordinates as inputs, our model is able to obtain the coarse contour of the instance. This suggests that the generated dynamic filters can attend to the target instance according to the relative coordinates, and it encodes the contour of the target instance. The generated dynamic filters can be also viewed as a representation of a contour. This is different from Mask R-CNN, which attends to a target instance by an axis-aligned RoI produced by Faster R-CNN, CondInst encodes the instance's contour into



**Fig. 5** – Qualitative results without relative coordinates or bottom features as inputs to the dynamic mask heads. From top to bottom: only with relative coordinates, only with bottom features and with both. We can see that the bottom features are crucial to the details of the instance masks, and relative coordinates can help the model distinguish between different instances.

the generated filters. Thus, CondInst can easily represent any shapes including irregular ones, being much more flexible. Moreover, if the bottom features are added, the dynamic filters can produce the details of instance masks. This suggests the generated filters look at the bottom features to obtain the details of the instance masks.

### 4.5 How Important to Upsample Mask Predictions?

As mentioned before, the original mask prediction is upsampled and the upsampling is of great importance to the final performance. We confirm this in the experiment. As shown in Table 5, without using the upsampling (1st row in the table), in this case CondInst can produce the mask prediction with $1/8$ of the input image resolution, which merely achieves 34.6% in mask AP because most of the details (*e.g.*, the boundary) are lost. If the mask prediction is upsampled by factor = 2, the performance can be significantly improved by 1% in mask AP (from 34.6% to 35.6%). In particular, the improvement on small objects is large (from 15.6% to 18.0), which suggests that the upsampling can greatly retain the details of objects. Increasing the upsampling factor to 4 slightly worsens the performance in some metrics, probably due to the relatively low-quality annotations of MS-COCO. Therefore, we use factor = 2 in all other models.

| w/ abs. coord. | w/ rel. coord. | w/ $\mathbf{F}_{bottom}$ | AP | AP$_{50}$ | AP$_{75}$ | AP$_S$ | AP$_M$ | AP$_L$ | AR$_1$ | AR$_{10}$ | AR$_{100}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | ✓ | 31.5 | 53.5 | 32.0 | 14.8 | 34.6 | 44.8 | 28.0 | 43.6 | 45.6 |
| | ✓ | | 31.3 | 55.0 | 31.9 | 15.6 | 34.1 | 44.3 | 27.1 | 43.3 | 45.6 |
| ✓ | | ✓ | 32.0 | 53.4 | 32.7 | 14.6 | 34.1 | 47.0 | 28.7 | 44.6 | 46.6 |
| | ✓ | ✓ | **35.6** | **56.4** | **37.9** | **18.0** | **39.1** | **50.8** | **30.3** | **48.7** | **51.3** |

**TABLE 3** – Ablation study of the input to the mask head on MS-COCO `val2017` split. As shown in the table, without the relative coordinates, the performance drops significantly from 35.6% to 31.5% in mask AP. Using the absolute coordinates *cannot* improve the performance remarkably. In addition, it is worth noting that if the mask head only takes as inputs the relative coordinates (*i.e.*, no appearance features in this case), CondInst also achieves modest performance.

| | AP | AP$_{50}$ | AP$_{75}$ | AP$_S$ | AP$_M$ | AP$_L$ |
|---|---|---|---|---|---|---|
| $P_3$ | 35.6 | 56.4 | 37.9 | **18.0** | **39.1** | 50.8 |
| $P_2$ | **36.0** | **56.6** | **38.4** | 17.6 | 38.9 | **51.7** |

**TABLE 4** – Instance segmentation results on MS-COCO `val2017` split by varying the FPN feature level for the bottom module. Using $P_2$ has better performance but it increases the inference latency by about 20%.

| factor | resolution | AP | AP$_{50}$ | AP$_{75}$ | AP$_S$ | AP$_M$ | AP$_L$ |
|---|---|---|---|---|---|---|---|
| 1 | 1/8 | 34.6 | 55.6 | 36.4 | 15.6 | 38.7 | 51.7 |
| 2 | 1/4 | **35.6** | **56.4** | **37.9** | **18.0** | **39.1** | 50.8 |
| 4 | 1/2 | **35.6** | 56.2 | 37.7 | 16.9 | 38.8 | 50.8 |

**TABLE 5** – The instance segmentation results on MS-COCO `val2017` split by changing the factor used to upsample the mask predictions. "resolution" denotes the resolution ratio of the mask prediction to the input image. Without the upsampling (*i.e.*, factor = 1), the performance drops significantly. The similar results are obtained with ratio 2 or 4.

| NMS | AP | AP$_{50}$ | AP$_{75}$ | AP$_S$ | AP$_M$ | AP$_L$ |
|---|---|---|---|---|---|---|
| box | 35.6 | 56.4 | **37.9** | **18.0** | **39.1** | **50.8** |
| mask | 35.6 | **56.5** | 37.7 | **18.0** | **39.1** | 50.7 |

**TABLE 6** – Instance segmentation results with different NMS algorithms. Mask-based NMS can obtain the same overall performance as box-based NMS, which suggests that CondInst can eliminate the box detection.

### 4.6 CondInst without Bounding-box Detection

Although we still keep the bounding-box detection branch in CondInst, it is conceptually feasible to eliminate it if we make use of the NMS using no bounding-boxes. In this case, all the foreground samples (predicted by the classification head) will be used to compute instance masks, and the duplicated masks will be removed by mask-based NMS. This is confirmed in Table 6. As shown in the table, by removing the box branch in inference and using the mask-based NMS, similar performance can be obtained to box-based NMS (35.6% vs. 35.6% in mask AP). The similar performance of mask and box NMS is probably due to the fact that the instances of MS-COCO are often less dense. Also, although with highly-optimized implementation on GPUs, mask and box NMS can have similar latency, it is worth noting that we need to compute the masks for all the foreground instances before the mask NMS can be applied. The underlying detector FCOS often predicts thousands of foreground instances, and thus it will take significantly longer time to obtain the masks of all the foreground instances. This makes the model with mask NMS significantly slower than the one with box NMS (often more than 2 times slower).

### 4.7 Comparisons with State-of-the-art Methods

We compare CondInst against previous state-of-the-art methods on MS-COCO `test-dev` split. As shown in Table 7, with $1\times$ learning rate schedule (*i.e.*, $90K$ iterations), CondInst outperforms the original Mask R-CNN by 0.7% (35.3% vs. 34.6%). CondInst also achieves a much faster speed than the original Mask R-CNN (49ms vs. 65ms per image on a single V100 GPU). To our knowledge, it is the

first time that a new and simpler instance segmentation method, without any bells and whistles outperforms Mask R-CNN both in accuracy and speed. CondInst also obtains better performance (35.9% vs. 35.5%) and on-par speed (49ms vs. 49ms) than the well-engineered Mask R-CNN in `Detectron2` (*i.e.*, Mask R-CNN* in Table 7). Furthermore, with a longer training schedule (*e.g.*, $3\times$) or a stronger backbone (*e.g.*, ResNet-101), a consistent improvement is achieved as well (37.7% vs. 37.5% with ResNet-50 $3\times$ and 39.1% vs. 38.8% with ResNet-101 $3\times$). Moreover, as shown in Table 7, with the auxiliary semantic segmentation task, the performance can be boosted from 37.7% to 38.6% (ResNet-50) or from 39.1% to 40.0% (ResNet-101), without increasing the inference time. For fair comparisons, all the inference time here is measured by ourselves on the same hardware with the official code.

We also compare CondInst with the recently-proposed instance segmentation methods. Only with half training iterations, CondInst surpasses TensorMask [50] by a large margin (37.7% vs. 35.4% for ResNet-50 and 39.1% vs. 37.1% for ResNet-101). CondInst is also $\sim 8\times$ faster than Tensor-Mask (49ms vs. 380ms per image on the same GPU) with similar performance (37.7% vs. 37.1%). Moreover, CondInst outperforms YOLACT-700 [1] by a large margin with the same backbone ResNet-101 (40.0% vs. 31.2% and both with the auxiliary semantic segmentation task). Moreover, as shown in Fig. 2, compared with YOLACT-700 and Mask R-CNN, CondInst can preserve more details and produce higher-quality instance segmentation results.

### 4.8 Real-time Instance Segmentation with CondInst

We also present a real-time version of CondInst. Following FCOS [53], the $4\times$ conv. layers in the classification and box regression towers in FCOS are shared in the real-time models (denoted by "shtw." in Table 8). Moreover, we reduce the input image from a scale of 800 to 512 during testing, and the FPN levels $P_6$ and $P_7$ are removed since there are not many larger objects with the small input images. In order to compensate for the performance loss due to the smaller input size, we use a more aggressive training strategy here. Specifically, the real-time models are trained for $360K$ iterations (*i.e.*, $4\times$) and the shorter side of the

| method | backbone | aug. | sched. | AP (%) | $AP_{50}$ | $AP_{75}$ | $AP_S$ | $AP_M$ | $AP_L$ |
|---|---|---|---|---|---|---|---|---|---|
| Mask R-CNN [2] | R-50-FPN | | 1× | 34.6 | **56.5** | 36.6 | 15.4 | 36.3 | **49.7** |
| **CondInst** | R-50-FPN | | 1× | **35.3** | 56.4 | **37.4** | **18.2** | **37.8** | 46.7 |
| Mask R-CNN* | R-50-FPN | ✓ | 1× | 35.5 | 57.0 | 37.8 | 19.5 | 37.6 | 46.0 |
| Mask R-CNN* | R-50-FPN | ✓ | 3× | 37.5 | 59.3 | 40.2 | **21.1** | 39.6 | 48.3 |
| TensorMask [50] | R-50-FPN | ✓ | 6× | 35.4 | 57.2 | 37.3 | 16.3 | 36.8 | 49.3 |
| BlendMask w/ sem. [32] | R-50-FPN | ✓ | 3× | 37.0 | 58.9 | 39.7 | 17.3 | 39.4 | 52.5 |
| **CondInst** | R-50-FPN | ✓ | 1× | 35.9 | 57.0 | 38.2 | 19.0 | 38.6 | 46.7 |
| **CondInst** | R-50-FPN | ✓ | 3× | 37.7 | 58.9 | 40.3 | 20.4 | 40.2 | 48.9 |
| **CondInst** w/ sem. | R-50-FPN | ✓ | 3× | **38.6** | **60.2** | **41.4** | 20.6 | **41.0** | **51.1** |
| Mask R-CNN | R-101-FPN | ✓ | 6× | 38.3 | 61.2 | 40.8 | 18.2 | 40.6 | 54.1 |
| Mask R-CNN* | R-101-FPN | ✓ | 3× | 38.8 | 60.9 | 41.9 | 21.8 | 41.4 | 50.5 |
| YOLACT-700 [1] | R-101-FPN | ✓ | 4.5× | 31.2 | 50.6 | 32.8 | 12.1 | 33.3 | 47.1 |
| PolarMask [35] | R-101-FPN | ✓ | 2× | 32.1 | 53.7 | 33.1 | 14.7 | 33.8 | 45.3 |
| TensorMask | R-101-FPN | ✓ | 6× | 37.1 | 59.3 | 39.4 | 17.4 | 39.1 | 51.6 |
| SOLO [33] | R-101-FPN | ✓ | 6× | 37.8 | 59.5 | 40.4 | 16.4 | 40.6 | 54.2 |
| BlendMask* w/ sem. | R-101-FPN | ✓ | 3× | 39.6 | 61.6 | 42.6 | **22.4** | 42.2 | 51.4 |
| SOLOv2 [34] | R-101-FPN | ✓ | 6× | 39.7 | 60.7 | **42.9** | 17.3 | **42.9** | **57.4** |
| **CondInst** | R-101-FPN | ✓ | 3× | 39.1 | 60.8 | 41.9 | 21.0 | 41.9 | 50.9 |
| **CondInst** w/ sem. | R-101-FPN | ✓ | 3× | **40.0** | **62.0** | **42.9** | 21.4 | 42.6 | 53.0 |
| **CondInst** w/ sem. | R-101-BiFPN | ✓ | 3× | 40.5 | 62.4 | 43.4 | 21.8 | 43.3 | 53.3 |
| **CondInst** w/ sem. | DCN-101-BiFPN | ✓ | 3× | **41.3** | **63.3** | **44.4** | **22.5** | **43.9** | **55.2** |

**TABLE 7** – Instance segmentation comparisons with state-of-the-art methods on MS-COCO `test-dev`. "Mask R-CNN" is the original Mask R-CNN [2]. "Mask R-CNN*" and "BlendMask*" mean that the models are improved by Detectron2 [3]. "aug.": using multi-scale data augmentation during training. "sched.": the learning rate schedule. 1× is $90K$ iterations, 2× is $180K$ iterations and so on. The learning rate is changed as in [51]. "w/ sem": using the auxiliary semantic segmentation task.

| method | backbone | sched. | FPS | AP | $AP_{50}$ | $AP_{75}$ |
|---|---|---|---|---|---|---|
| YOLACT-550++ [52] | R-50 | 4.5× | 44 | 34.1 | 53.3 | 36.2 |
| YOLACT-550++ | R-101 | 4.5× | 36 | 34.6 | 53.8 | 36.9 |
| CondInst-RT shtw. | R-50 | 4× | 43 | 36.0 | 57.0 | 38.0 |
| CondInst-RT shtw. | DLA-34 | 4× | **47** | 35.8 | 56.5 | 38.0 |
| CondInst-RT | DLA-34 | 4× | 41 | **36.3** | **57.3** | **38.5** |

**TABLE 8** – The mask AP and inference speed of the real-time CondInst models on the COCO `test-dev` data. "shtw.": sharing the conv. towers between the classification and box regression branches in FCOS. Both YOLACT++ and CondInst use the auxiliary semantic segmentation loss here. As you can see, with the same backbone R-50, CondInst-RT outperforms YOLACT++ by $1.9\%$ AP with almost the same speed. All inference time is measured with a single V100 GPU.

input image is randomly chosen from the range 256 to 608 with step 32. Synchronized BatchNorm (SyncBN) is also used during training. In the real-time models, following YOLACT, we enable the extra semantic segmentation loss by default.

The performance and inference speed of these real-time models are shown in Table 8. As shown in the table, the R-50 based CondInst-RT outperforms the R-50 based YOLACT++ [52] by about $2\%$ AP ($36.0\%$ vs. $34.1\%$) and has almost the same speed (43 FPS vs. 44 FPS). By further using a strong backbone DLA-34 [54], CondInst-RT can achieve 47 FPS with similar performance. Furthermore, if we do not share the classification and box regression towers in FCOS, the performance can be improved to $36.3\%$ AP with slightly longer inference time (41 FPS).

### 4.9 Instance Segmentation on Cityscapes

We also conduct the instance segmentation experiments on Cityscapes [55]. The Cityscapes dataset is designed for the understanding of urban street scenes. For instance segmentation, it has 8 categories, which are person, rider, car, truck, bus, train, motorcycle, and bicycle. It includes 2975, 500 and 1525 images with `fine` annotations for training,

validation and testing, respectively. It also has 20K training images with `coarse` annotations. Following Mask R-CNN [2], we only use the images with `fine` annotations to train our models. All images in Cityscapes have the same resolution 2048×1024. The performance on Cityscapes is also measured with the COCO-style mask AP, which are the averaged mask AP over ten IoU thresholds from $0.5$ to $0.95$.

We follow the training details in Detectron2 [3] to train CondInst on Cityscapes. Specifically, the models are trained for 24K iterations with batch size 8 (1 image per GPU). The initial learning rate is $0.01$, which is reduced by a factor of 10 at step 18K. Since Cityscapes has relatively fewer images, following Mask R-CNN, we may initialize the models with the weights pre-trained on the COCO dataset if specified. Moreover, we use multi-scale data augmentation during training and the shorter side of the images is sampled in the range from 800 to 1024 with step 32. In inference, we only use the original image scale 2048×1024. Additionally, in order to preserve more details on Cityscapes, we increase the mask output resolution of CondInst from $1/4$ to $1/2$ resolution of the input image.

The results are reported in Table 9. As shown in the table, with the same settings, CondInst generally outperforms the previous strong baseline Mask R-CNN by more than $1\%$ mask AP in all the experiments. On Cityscapes, the auxiliary semantic segmentation loss can also improve the instance segmentation performance. The results with the loss are denoted by "w/ sem." in Table 9. By further using the complementary techniques such as deformable convolutions and BiFPN, the performance can be further boosted as expected.

### 4.10 Experiments on Panoptic Segmentation

As mentioned before, CondInst can be easily extended to panoptic segmentation [18] by attaching a new semantic segmentation branch depicted in Fig. 4. Here, we conduct

| method | backbone | training data | AP [val] | AP | AP$_{50}$ | person | rider | car | truck | bus | train | mcycle | bicycle |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Mask R-CNN | ResNet-50-FPN | train | 31.5 | 26.2 | 49.9 | 30.5 | 23.7 | 46.9 | 22.8 | 32.2 | 18.6 | 19.1 | 16.0 |
| CondInst | ResNet-50-FPN | train | 33.3 | **28.6** | **53.5** | **31.3** | 23.4 | 51.7 | **23.4** | **36.0** | **27.3** | 19.1 | 16.6 |
| CondInst w/ sem. | ResNet-50-FPN | train | **33.9** | **28.6** | 53.1 | **31.3** | 24.2 | **51.9** | 21.2 | 35.9 | 26.5 | **20.9** | **17.0** |
| Mask R-CNN | ResNet-50-FPN | train+COCO | 36.4 | 32.0 | 58.1 | 34.8 | 27.0 | 49.1 | 30.1 | 40.9 | 30.9 | 24.1 | 18.7 |
| CondInst | ResNet-50-FPN | train+COCO | 37.5 | 33.2 | 57.2 | 35.1 | 27.7 | 54.5 | 29.5 | 42.3 | 33.8 | 23.9 | **18.9** |
| CondInst w/sem. | ResNet-50-FPN | train+COCO | 37.7 | 33.7 | 57.7 | **35.7** | 28.0 | 54.8 | 29.6 | 41.4 | **36.3** | **24.8** | 18.9 |
| CondInst w/sem. | DCN-101-BiFPN | train+COCO | **39.3** | 33.9 | 58.2 | 35.6 | 28.1 | 55.0 | 32.1 | 44.2 | 33.6 | 24.5 | 18.6 |
| CondInst w/sem. | ResNet-50-FPN | train+val+COCO | - | 34.4 | **59.6** | 36.4 | 28.4 | 55.3 | 32.6 | 43.3 | 33.9 | 24.8 | **20.1** |
| CondInst w/sem. | DCN-101-BiFPN | train+val+COCO | - | **35.1** | 59.0 | 35.9 | **28.7** | 55.4 | 34.4 | 45.7 | 35.5 | 25.5 | 19.6 |

**TABLE 9** – Instance segmentation results on Cityscapes `val` ("AP [val]" column) and `test` (remaining columns) splits. "DCN": using deformable convolutions in the backbones. "+COCO": fine-tuning from the models pre-trained on COCO. "train+val+COCO": using both `train` and `val` splits to train the models evaluated on the `test` split. "w/ sem.": using the auxiliary semantic segmentation loss during training as in COCO.



**Fig. 6** – Panoptic segmentation results on the COCO dataset (better viewed on screen). Color encodes categories and instances.

| method | backbone | sched. | PQ | PQ$_{th}$ | PQ$_{st}$ |
|---|---|---|---|---|---|
| CondInst | R-50-FPN | 1× | 42.1 | 50.4 | 29.7 |
| Unifying [56] | R-50-FPN | - | 43.6 | 48.9 | **35.6** |
| CondInst | R-50-FPN | 3× | **44.6** | **53.0** | 31.8 |
| DeeperLab [36] | Xception-71 [57] | - | 34.3 | 37.5 | 29.6 |
| Panoptic-DeepLab [58] | Xception-71 | - | 39.7 | 43.9 | 33.2 |
| Panoptic-FPN [38] | R-101-FPN | 3× | 40.9 | 48.3 | 29.7 |
| AdaptIS [40] | ResNeXt-101 | 1.7× | 42.8 | 50.1 | 31.8 |
| Aixal DeepLab [39] | Axial-ResNet-L | - | 43.6 | 48.9 | 35.6 |
| Panoptic-FCN [59] | R-101-FPN | 3× | 45.5 | 51.4 | **36.4** |
| **CondInst** | R-101-FPN | 3× | **46.1** | **54.7** | 33.2 |
| UPSNet [60] | DCN-101-FPN | 3× | 46.6 | 53.2 | 36.7 |
| Panoptic-FCN | DCN-101-FPN | 3× | 47.1 | 53.2 | **37.8** |
| Unifying | DCN-101-FPN | - | 47.2 | 53.5 | 37.7 |
| **CondInst** | DCN-101-FPN | 3× | **47.8** | **55.8** | 35.8 |

**TABLE 10** – Panoptic segmentation on the COCO `test-dev` data. All results are with single-model and single-scale testing. Here we report comparisons with state-of-the-art methods using various backbones and training schedules (1× means 90K iterations). CondInst achieves the best results among the compared methods.

| method | backbone | PQ | PQ$_{th}$ | PQ$_{st}$ |
|---|---|---|---|---|
| Li *et al.* [37] | - | 53.8 | 42.5 | 62.1 |
| DeeperLab [36] | Xception-71 | 56.5 | - | - |
| Panoptic-FPN [38] | R-101-FPN | 58.1 | 52.0 | 62.5 |
| AdaptIS [40] | R-50 | 59.0 | 55.8 | 61.3 |
| UPSNet [60] | R-50-FPN | 59.3 | 54.6 | 62.7 |
| Panoptic-DeepLab [58] | R-50 | 59.7 | - | - |
| Unifying [56] | R-50-FPN | 61.4 | 54.7 | 66.3 |
| Panoptic-FCN [59] | R-50-FPN | 61.4 | 54.8 | **66.6** |
| **CondInst** | R-50-FPN | **61.7** | **59.0** | 63.7 |

**TABLE 11** – Panoptic segmentation on the Cityscapes `val` set. All results are with single-model and single-scale with no flipping. Here we report comparisons with state-of-the-art methods.

the panoptic segmentation experiments on the COCO 2018 dataset. Unless specified, the training and testing details (*e.g.*, image sizes, the number of iterations and etc.) are the same as in the instance segmentation task on COCO.

Although panoptic segmentation can be viewed as a combination of instance segmentation and semantic segmentation, there is a discrepancy between the ground-truth annotations of the original instance segmentation and the instance segmentation task in panoptic segmentation. Panoptic segmentation requires that a pixel in the resulting mask has only one label. Therefore if two instances overlap, the pixels in the overlapped region will only be assigned to the front instance. However, in the original instance segmen-

tation, the pixels in the overlapped region belong to both instances, and the ground-truth masks are labeled in such a way. Therefore, when we use the instance segmentation framework for panoptic segmentation, the training targets of the instance segmentation need to be changed to the instance annotations in panoptic segmentation accordingly.

We compare our method with a few state-of-the-art panoptic segmentation methods in Table 10. On the challenging COCO `test-dev` benchmark, we outperform the previous strong baseline Panoptic-FPN [38] by a large margin with the same backbone and training schedule (*i.e.*, from $40.9\%$ to $46.1\%$ in PQ with ResNet-101). Moreover, compared to AdaptIS [40], which shares some similarity with us, the ResNet-101 based CondInst achieves dramatically better performance than ResNeXt-101 based AdaptIS ($46.1\%$ vs. $42.8\%$ PQ). This suggests that using the dynamic filters here might be more effective than using FiLM [41]. In addition, compared to the recent methods such as [56]

and Panoptic-FCN [59], CondInst also outperforms them considerably. Some qualitative results are in Fig. 6. We also conduct experiments on the panoptic segmentation task of Cityscapes [55], and we follow the training strategy of Panoptic-FPN [38] on this benchmark. Similar to previous works [38], [56], [59], we report the results on the Cityscapes `val` set. As shown in Table 11, we outperform previous methods on this benchmark as well.

## 5 CONCLUSION

We have proposed a new and simple instance segmentation framework, termed CondInst. Unlike previous method such as Mask R-CNN, which employs the mask head with fixed weights, CondInst conditions the mask head on instances and dynamically generates the filters of the mask head. This not only reduces the parameters and computational complexity of the mask head, but also eliminates the ROI operations, resulting in a faster and simpler instance segmentation framework. To our knowledge, CondInst is the first framework that can outperform Mask R-CNN both in accuracy and speed, without longer training schedules needed. With simple modifications, CondInst can be extended to solve panoptic segmentation and achieve state-of-the-art performance on the challenging COCO dataset. We believe that CondInst can be a strong alternative for both instance and panoptic segmentation.

## REFERENCES

[1] D. Bolya, C. Zhou, F. Xiao, and Y. J. Lee, "YOLACT: real-time instance segmentation," in *Proc. IEEE Int. Conf. Comp. Vis.*, pp. 9157–9166, 2019.

[2] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask R-CNN," in *Proc. IEEE Int. Conf. Comp. Vis.*, pp. 2961–2969, 2017.

[3] Y. Wu, A. Kirillov, F. Massa, W.-Y. Lo, and R. Girshick, "Detectron2." https://github.com/facebookresearch/detectron2, 2019.

[4] A. Paszke et al., "PyTorch: An imperative style, high-performance deep learning library," in *Proc. Advances in Neural Inf. Process. Syst.*, pp. 8024–8035, 2019.

[5] M. Abadi *et al.*, "TensorFlow: A system for large-scale machine learning," in *USENIX Symp. Operating Systems Design & Implementation*, pp. 265–283, 2016.

[6] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 3431–3440, 2015.

[7] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. Yuille, "Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected CRFs," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 4, pp. 834–848, 2017.

[8] Z. Tian, T. He, C. Shen, and Y. Yan, "Decoders matter for semantic segmentation: Data-dependent decoding enables flexible feature aggregation," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 3126–3135, 2019.

[9] T. He, C. Shen, Z. Tian, D. Gong, C. Sun, and Y. Yan, "Knowledge adaptation for efficient semantic segmentation," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 578–587, 2019.

[10] Y. Liu, C. Shu, J. Wang, and C. Shen, "Structured knowledge distillation for dense prediction," *IEEE Trans. Pattern Anal. Mach. Intell.*, 2020.

[11] Z. Tian, C. Shen, H. Chen, and T. He, "FCOS: Fully convolutional one-stage object detection," in *Proc. IEEE Int. Conf. Comp. Vis.*, pp. 9627–9636, 2019.

[12] F. Liu, C. Shen, G. Lin, and I. Reid, "Learning depth from single monocular images using deep convolutional neural fields," *IEEE Trans. Pattern Anal. Mach. Intell.*, 2016.

[13] W. Yin, Y. Liu, C. Shen, and Y. Yan, "Enforcing geometric constraints of virtual normal for depth prediction," in *Proc. IEEE Int. Conf. Comp. Vis.*, 2019.

[14] J. Bian, Z. Li, N. Wang, H. Zhan, C. Shen, M.-M. Cheng, and I. Reid, "Unsupervised scale-consistent depth and ego-motion learning from monocular video," in *Proc. Advances in Neural Inf. Process. Syst.*, pp. 35–45, 2019.

[15] L. Boominathan, S. Kruthiventi, and R. V. Babu, "Crowdnet: A deep convolutional network for dense crowd counting," in *Proc. ACM Int. Conf. Multimedia*, pp. 640–644, ACM, 2016.

[16] X. Jia, B. De Brabandere, T. Tuytelaars, and L. V. Gool, "Dynamic filter networks," in *Proc. Advances in Neural Inf. Process. Syst.*, pp. 667–675, 2016.

[17] B. Yang, G. Bender, Q. V. Le, and J. Ngiam, "Condconv: Conditionally parameterized convolutions for efficient inference," in *Proc. Advances in Neural Inf. Process. Syst.*, pp. 1305–1316, 2019.

[18] A. Kirillov, K. He, R. Girshick, C. Rother, and P. Dollár, "Panoptic segmentation," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 9404–9413, 2019.

[19] L. Zhang, D. Xu, A. Arnab, and P. H. Torr, "Dynamic graph message passing networks," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, June 2020.

[20] K. Chen, J. Pang, J. Wang, Y. Xiong, X. Li, S. Sun, W. Feng, Z. Liu, J. Shi, W. Ouyang, *et al.*, "Hybrid task cascade for instance segmentation," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 4974–4983, 2019.

[21] S. Liu, L. Qi, H. Qin, J. Shi, and J. Jia, "Path aggregation network for instance segmentation," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 8759–8768, 2018.

[22] Z. Huang, L. Huang, Y. Gong, C. Huang, and X. Wang, "Mask scoring R-CNN," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 6409–6418, 2019.

[23] J. Dai, K. He, Y. Li, S. Ren, and J. Sun, "Instance-sensitive fully convolutional networks," in *Proc. Eur. Conf. Comp. Vis.*, pp. 534–549, Springer, 2016.

[24] D. Neven, B. D. Brabandere, M. Proesmans, and L. V. Gool, "Instance segmentation by jointly optimizing spatial embeddings and clustering bandwidth," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 8837–8845, 2019.

[25] A. Newell, Z. Huang, and J. Deng, "Associative embedding: End-to-end learning for joint detection and grouping," in *Proc. Advances in Neural Inf. Process. Syst.*, pp. 2277–2287, 2017.

[26] A. Fathi, Z. Wojna, V. Rathod, P. Wang, H. O. Song, S. Guadarrama, and K. P. Murphy, "Semantic instance segmentation via deep metric learning," *arXiv: Comp. Res. Repository*, 2017.

[27] M. Bai and R. Urtasun, "Deep watershed transform for instance segmentation," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 5221–5229, 2017.

[28] S. Liu, J. Jia, S. Fidler, and R. Urtasun, "Sgn: Sequential grouping networks for instance segmentation," in *Proc. IEEE Int. Conf. Comp. Vis.*, pp. 3496–3504, 2017.

[29] J. Uhrig, E. Rehder, B. Fröhlich, U. Franke, and T. Brox, "Box2pix: Single-shot instance segmentation by assigning pixels to object boxes," in *Proc. IEEE Intelligent Vehicles Symp.*, pp. 292–299, IEEE, 2018.

[30] D. Novotny, S. Albanie, D. Larlus, and A. Vedaldi, "Semi-convolutional operators for instance segmentation," in *Proc. Eur. Conf. Comp. Vis.*, pp. 86–102, 2018.

[31] A. Arnab and P. Torr, "Pixelwise instance segmentation with a dynamically instantiated network," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 441–450, 2017.

[32] H. Chen, K. Sun, Z. Tian, C. Shen, Y. Huang, and Y. Yan, "Blendmask: Top-down meets bottom-up for instance segmentation," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2020.

[33] X. Wang, T. Kong, C. Shen, Y. Jiang, and L. Li, "SOLO: Segmenting objects by locations," in *Proc. Eur. Conf. Comp. Vis.*, 2020.

[34] X. Wang, R. Zhang, T. Kong, L. Li, and C. Shen, "SOLOv2: Dynamic and fast instance segmentation," in *Proc. Advances in Neural Inf. Process. Syst.*, 2020.

[35] E. Xie, P. Sun, X. Song, W. Wang, D. Liang, C. Shen, and P. Luo, "PolarMask: Single shot instance segmentation with polar representation," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2020.

[36] T.-J. Yang, M. D. Collins, Y. Zhu, J.-J. Hwang, T. Liu, X. Zhang, V. Sze, G. Papandreou, and L.-C. Chen, "Deeperlab: Single-shot image parser," *arXiv preprint arXiv:1902.05093*, 2019.

[37] Q. Li, A. Arnab, and P. H. Torr, "Weakly-and semi-supervised panoptic segmentation," in *Proc. Eur. Conf. Comp. Vis.*, pp. 102–118, 2018.

[38] A. Kirillov, R. Girshick, K. He, and P. Dollár, "Panoptic feature pyramid networks," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 6399–6408, 2019.

[39] H. Wang, Y. Zhu, B. Green, H. Adam, A. Yuille, and L.-C. Chen, "Axial-DeepLab: Stand-alone axial-attention for panoptic segmentation," in *Proc. Eur. Conf. Comp. Vis.*, 2020.

[40] K. Sofiiuk, O. Barinova, and A. Konushin, "Adaptis: Adaptive instance selection network," in *Proc. IEEE Int. Conf. Comp. Vis.*, pp. 7355–7363, 2019.

[41] E. Perez, F. Strub, H. De Vries, V. Dumoulin, and A. Courville, "FiLM: Visual reasoning with a general conditioning layer," in *Proc. AAAI Conf. Artificial Intell.*, 2018.

[42] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 2117–2125, 2017.

[43] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *Proc. Eur. Conf. Comp. Vis.*, pp. 740–755, Springer, 2014.

[44] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," in *Proc. Advances in Neural Inf. Process. Syst.*, pp. 91–99, 2015.

[45] K. He, X. Zhang, S. Ren, and J. Sun, "Spatial pyramid pooling in deep convolutional networks for visual recognition," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 37, no. 9, pp. 1904–1916, 2015.

[46] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 2980–2988, 2017.

[47] F. Milletari, N. Navab, and S.-A. Ahmadi, "V-net: Fully convolutional neural networks for volumetric medical image segmentation," in *Proc. Int. Conf. 3D Vision*, pp. 565–571, IEEE, 2016.

[48] Z. Tian, C. Shen, and H. Chen, "Conditional convolutions for instance segmentation," in *Proc. Eur. Conf. Comp. Vis.*, 2020.

[49] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 248–255, Ieee, 2009.

[50] X. Chen, R. Girshick, K. He, and P. Dollár, "Tensormask: A foundation for dense object segmentation," in *Proc. IEEE Int. Conf. Comp. Vis.*, pp. 2061–2069, 2019.

[51] K. He, R. Girshick, and P. Dollár, "Rethinking imagenet pretraining," in *Proc. IEEE Int. Conf. Comp. Vis.*, pp. 4918–4927, 2019.

[52] D. Bolya, C. Zhou, F. Xiao, and Y. J. Lee, "YOLACT++: Better real-time instance segmentation," *IEEE Trans. Pattern Anal. Mach. Intell.*, 2019.

[53] Z. Tian, C. Shen, H. Chen, and T. He, "FCOS: A simple and strong anchor-free object detector," *IEEE Trans. Pattern Anal. Mach. Intell.*, 2021.

[54] F. Yu, D. Wang, E. Shelhamer, and T. Darrell, "Deep layer aggregation," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 2403–2412, 2018.

[55] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, "The Cityscapes dataset for semantic urban scene understanding," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 3213–3223, 2016.

[56] Q. Li, X. Qi, and P. H. S. Torr, "Unifying training and inference for panoptic segmentation," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 13320–13328, 2020.

[57] F. Chollet, "Xception: Deep learning with depthwise separable convolutions," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 1251–1258, 2017.

[58] B. Cheng, M. D. Collins, Y. Zhu, T. Liu, T. S. Huang, H. Adam, and L.-C. Chen, "Panoptic-DeepLab: A simple, strong, and fast baseline for bottom-up panoptic segmentation," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2020.

[59] Y. Li, H. Zhao, X. Qi, L. Wang, Z. Li, J. Sun, and J. Jia, "Fully convolutional networks for panoptic segmentation," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, 2021.

[60] Y. Xiong, R. Liao, H. Zhao, R. Hu, M. Bai, E. Yumer, and R. Urtasun, "Upsnet: A unified panoptic segmentation network," in *Proc. IEEE Conf. Comp. Vis. Patt. Recogn.*, pp. 8818–8826, 2019.

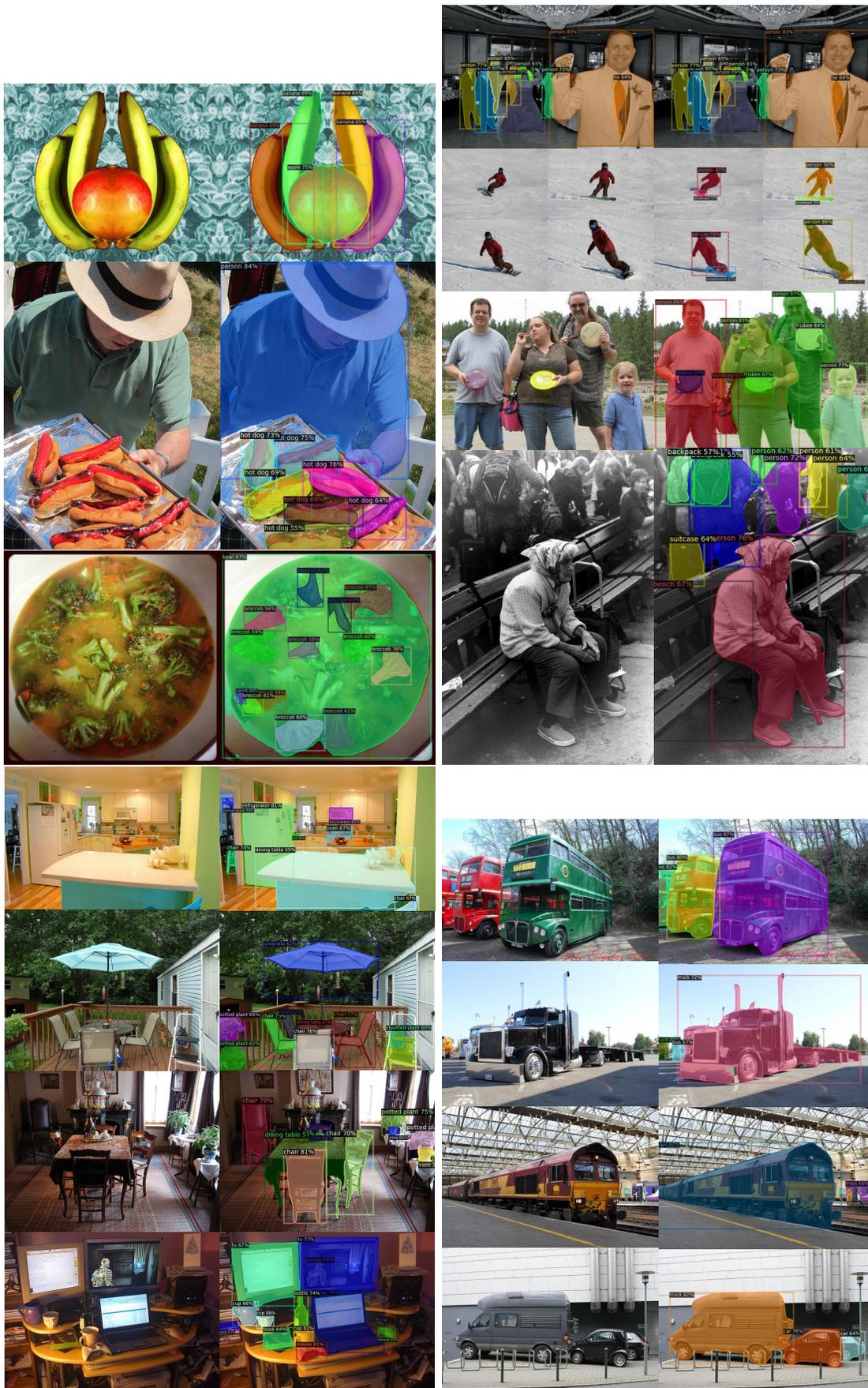**Authors' photograph and biography not available at the time of publication.**

# APPENDIX A

## VISUALIZATION OF RESULTS

Here we provide some visualization results of our model. Fig. A7 and Fig. A8 show some segmentation results of our model on COCO for instance segmentation and panoptic segmentation, respectively.

Fig. A9 show some results that our do not work very well for instance segmentation. In some cases, the COCO annotation is noisy, which may have caused confusion for our model. For example, for the third example in Fig. A9, the sailboat is incorrectly annotated. Occlusion in the last example also caused challenges.

Fig. A10 shows some panoptic results that our model does not perform well.

**Fig. A7** – More visualization of instance segmentation results on the COCO dataset (better viewed on screen). Color encodes categories and instances. Here the model is ResNet-101-DCN with BiFPN.
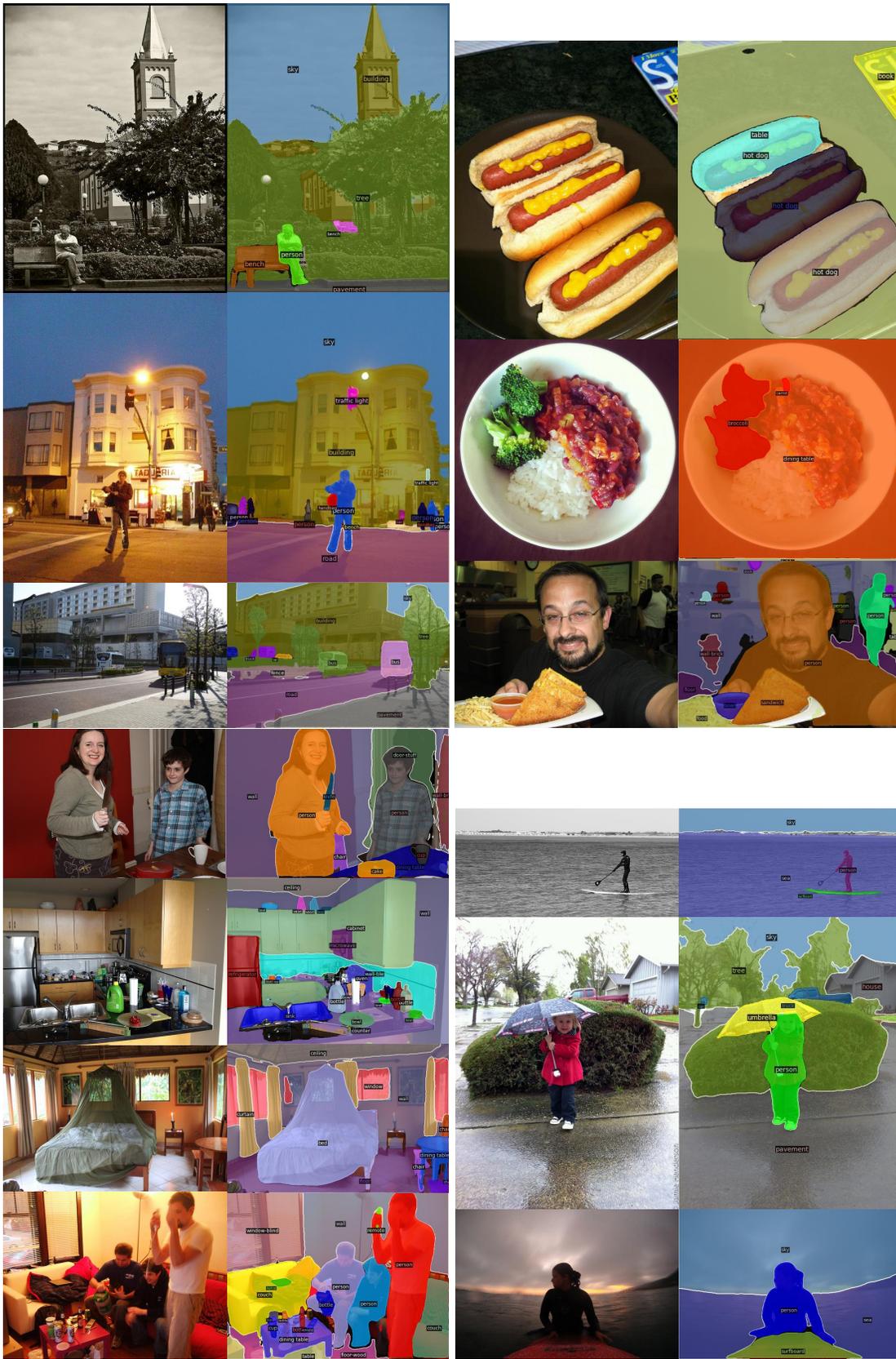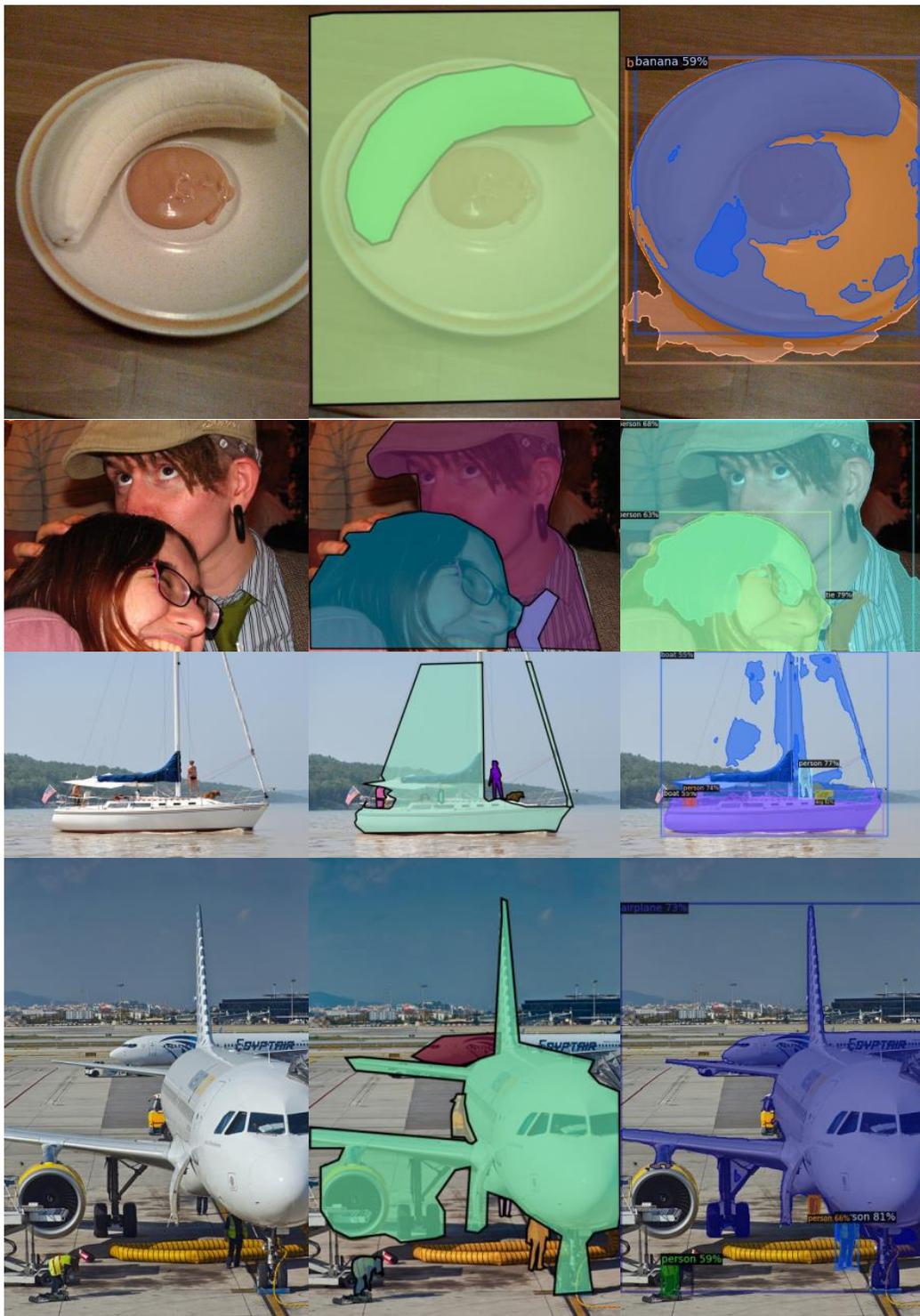
**Fig. A8** – More visualization of panoptic segmentation results on the COCO dataset (better viewed on screen). Here the model is ResNet-101-DCN with standard FPN.

**Fig. A9** – Some instance segmentation results that our model does not work very well, on the COCO dataset (better viewed on screen). Left to right: input image, ground-truth labels, model's predictions. In some cases (*e.g.*, the last two examples), the ground-truth annotation is incorrect or noisy.

**Fig. A10** – Some panoptic segmentation results that our model does not work very well, on the COCO dataset (better viewed on screen). Left to right: input image, ground-truth labels, model's predictions. On those challenging cases, our model makes plausible mistakes.