

Single-Stage Open-world Instance Segmentation with Cross-task Consistency Regularization

Xizhe Xue^{a,d}, Dongdong Yu^b, Lingqiao Liu^c, Yu Liu^b, Satoshi Tsutsui^d, Ying Li[‡], Zehuan Yuan^b, Ping Song^b, and Mike Zheng Shou^d

^aNorthwestern Polytechnical University

^bByteDance Inc.

^cUniversity of Adelaide

^dShow Lab, National University of Singapore

Abstract

Open-World Instance Segmentation (OWIS) is an emerging research topic that aims to segment class-agnostic object instances from images. The mainstream approaches use a two-stage segmentation framework, which first locates the candidate object bounding boxes and then performs instance segmentation. In this work, we instead promote a single-stage framework for OWIS. We argue that the end-to-end training process in the single-stage framework can be more convenient for directly regularizing the localization of class-agnostic object pixels. Based on the single-stage instance segmentation framework, we propose a regularization model to predict foreground pixels and use its relation to instance segmentation to construct a cross-task consistency loss. We show that such a consistency loss could alleviate the problem of incomplete instance annotation – a common problem in the existing OWIS datasets. We also show that the proposed loss lends itself to an effective solution to semi-supervised OWIS that could be considered an extreme case that all object annotations are absent for some images. Our extensive experiments demonstrate that the proposed method achieves impressive results in both fully-supervised and semi-supervised settings. Compared to SOTA methods, the proposed method significantly improves the AP_{100} score by 4.75% in UVO \rightarrow UVO setting and 4.05% in COCO \rightarrow UVO setting. In the case of semi-supervised learning, our model learned with only 30% labeled data, even outperforms its fully-supervised counterpart with 50% labeled data. The code will be released soon at: <https://github.com/showlab/SOIS>.

1 Introduction

Traditional instance segmentation [1, 2] methods often assume that objects in images can be categorized into a finite set of predefined classes (i.e., *closed-world*). Such an assumption, however, can be easily violated in many real-world applications, where models will encounter many new object classes that never appeared in the training data. Therefore, researchers recently attempted to tackle the problem of **Open-World Instance Segmentation (OWIS)** [3], which targets class-agnostic segmentation of all objects in the image.

Prior to this paper, most existing methods for OWIS are of two-stage [4, 5], which detects bounding boxes of objects and then segments them. Despite their promising performances, such a paradigm cannot handle and recover if object bounding boxes are not detected. In contrast, a single-stage approach called Mask2Former [6] has recently been introduced, yet only for closed-world instance

*This work was done when Xizhe Xue visited National University of Singapore. Email: xuexizhe@mail.nwpu.edu.cn

[†]Corresponding author

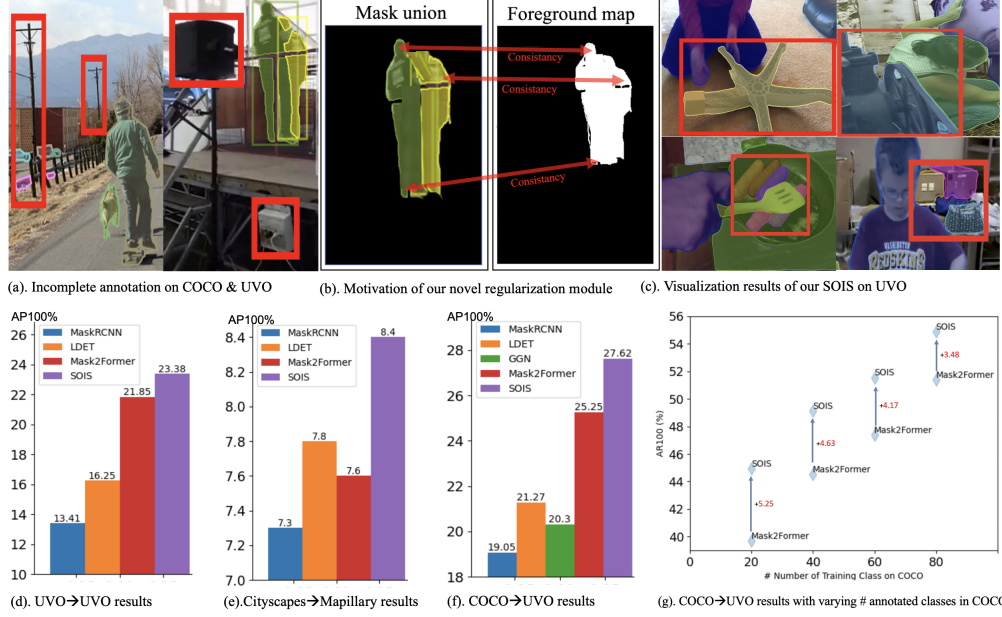


Figure 1: **(a).** Instances missing annotations in COCO and UVO datasets. The regions in **red boxes** are mistakenly annotated as background. **(b).** Motivation of our novel reg module (The consistency relationship between instance mask and foreground map). **(c).** Visualization results of our SOIS on UVO dataset. Here, the proposed SOIS is trained on COCO dataset and tested on UVO dataset. Our methods correctly segments many objects that are not labeled in COCO. **(d - f).** The $AP_{100}\%$ of our SOIS vs. SOTA methods on COCO→UVO, Cityscapes→Mapillary, COCO→UVO. **(g).** The $AR_{100}\%$ of our SOIS vs. baseline Mask2Former on COCO. From right to left, with the total number of classes decreases (i.e. more instance annotations missed), the gain of our SOIS over baseline becomes larger, thanks to the capability of our model to handle incomplete annotations.

segmentation. By extending it to open-world, we are **the first to develop a novel detection-free Single-stage Open-world Instance Segmentation method, dubbed as SOIS.**

Note that our work is not just a straightforward adaptation of Mask2Former from close-world to open-world. This is because unlike closed-world segmentation, where the object categories can be clearly defined before annotation, the open-world scenario makes it challenging for annotators to label all instances completely or ensure annotation consistency across different images because they cannot have a well-defined finite set of object categories. As shown in Figure 1(a), annotators miss some instances. **It still remains challenging that how to handle such incomplete annotations (i.e. some instances missed).**

Recent work LDET [5] addresses this problem by generating synthetic data with a plain background, but it based on a decoupled training strategy that can only be used in the two-stage method *while our method is of single-stage*. Another work called GGN [4] handles such incomplete instance-level annotation issue by training a pairwise affinity predictor for generating pseudo labels. But training such an additional predictor is complicated and time-consuming.

In contrast, *our proposed SOIS method is end-to-end and simpler*. We address this incomplete annotation issue via a **novel regularization module**, which is simple yet effective. Specifically, it is convenient to concurrently predict not only (1) *instance masks* but also a (2) *foreground map*. Ideally, as shown in Figure 1(b), the foreground region should be consistent with the union of all instance masks. To penalize their inconsistency, we devise a *cross-task consistency loss*, which can down-weight the adverse effects caused by incomplete annotation. This is because when an instance is missed in annotation, as long as it is captured by both our predictions of instance masks and foreground map, the consistency loss would be low and hence encourage such prediction. Experiments in Figure 1(g) show that such consistency loss is effective even when annotations miss many instances.

So far, like most existing methods, we focus on the fully-supervised OWIS. In this paper, we further extend OWIS to the semi-supervised setting, where some training images do not have any annotations at all. This is of great interest because annotating segmentation map is very costly. Notably, **our proposed regularization module can also benefit semi-supervised OWIS** – consider an unlabeled image as an extreme case of incomplete annotation where all of the instance annotations are missed. Specifically, we perform semi-supervised OWIS by first warming up the network on the labeled set and then continuing training it with the cross-task consistency loss on the mixture of labeled and unlabeled images.

Contributions. In a nutshell, our main contributions could be summarized as:

1. We propose a Single-stage Open-world Instance Segmentation (SOIS) for the first time while most OWIS methods are of two-stage.
2. We propose a novel cross-task consistency loss that mitigate the issue of incomplete mask annotations.
3. We further extend the proposed method into a semi-supervised OWIS model, which effectively makes use of the unlabeled images to help the OWIS model training .
4. Our extensive experiments demonstrate that the proposed method reaches the leading OWIS performance in the fully-supervised learning. (Figure 1(d-f)), and that our semi-supervised extension can achieve remarkable performance with a much smaller amount of labeled data.

2 Related Work

Closed-world instance segmentation (CWIS) [7, 8, 9, 10, 11] requires the approaches to assign a class label and instance ID to every pixel. Two-stage CWIS approaches, such as MaskRCNN, always include a bounding box estimation branch and a FCN-based mask segmentation branch, working in a ‘detect-then-segment’ way. To improve efficiency, one-stage methods such as CenterMask [9], YOLACT [10] and BlendMask [8] have been proposed, which remove the proposal generation and feature grouping process. To further free the CWIS from the local box detection, Wang et.al [11] proposed SOLO and obtained on par results to the above methods. In recent years, the methods [12, 13], following DETR [14], consider the instance segmentation task as an ensemble prediction problem. In addition, Cheng et al. proposed an universal segmentation framework MaskFormer [15] and its upgrade version Mask2Former [6], which even outperforms the state-of-the-art architectures specifically designed for the CWIS task.

Notably, two-stage method CenterMask preserves pixel alignment and separates the object simultaneously by integrating the local and global branch. Although introducing the global information in this way helps improve the mask quality in CWIS, it can not handle the open-world task very well. Because CenterMask multiplies the local shape and the cropped saliency map to form the final mask for each instance. There is no separate loss for the local shape and global saliency. When such method faces the incomplete annotations in OWIS tasks, the generated mask predictions corresponding to the unlabeled instances would still be punished during training, making it difficult to discover novel object at inference. The efficient way to jointly take advantages of global and local information in OWIS tasks deserves to be explored.

Open-world instance segmentation OWIS task [3] here focuses on the following aspects: (1) All instances (without stuff) have to be segmented; (2) Class-agnostic pixel-level results should be predicted with only instance ID and incremental learning ability is unnecessary. Several OWIS works have recently been developed. Yu et al. [16] proposed a two-stage segmentation algorithm, which decoupled the segmentation and detection modules during training and testing. This algorithm achieves competitive results on the UVO dataset thanks to the abundant training data and the introduction of effective modules such as cascade RPN [17], SimOTA [18], etc. Another work named LDET [5] attempts to solve the instance-level incomplete annotation problem. Specifically, LDET first generates the background of the synthesized image by taking a small piece of background in the original image and enlarging it to the same size as the original image. The instance is then matted to the foreground of the synthesized image. The synthesized data is used only to train the mask prediction branch, and the rest of the branches are still trained with the original data. Meanwhile, Wang et al. proposed GGN [4], an algorithm that combines top-down and bottom-up segmentation ideas to improve prediction accuracy by generating high-quality pseudo-labels. Specifically, a Pairwise Affinity (PA) predictor is trained first and a grouping module is used to extract and rank segments from predicted PA to generate pseudo-labels, which would be fused with groundtruth to train the segmentation model.

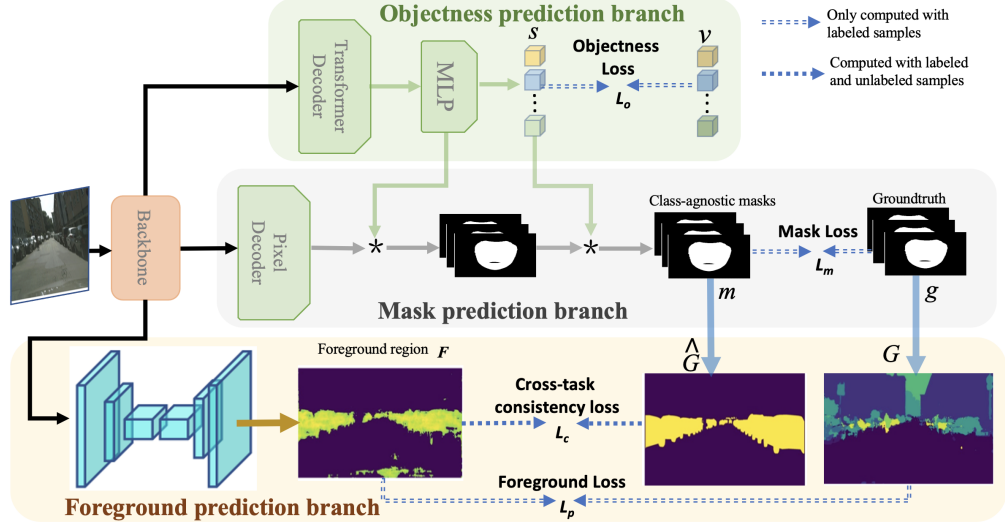


Figure 2: Overall framework of the proposed SOIS. The mask prediction branch generates the predicted masks, while the objectness prediction branch computes the objectness score for each mask. The foreground prediction branch segments a foreground region to guide the optimization of other two branches.

3 Methodology

In this section, we first define the OWIS problem with both fully and semi-supervised learning. Then the architecture of our SOIS and the proposed cross-task consistency loss are introduced in Section 3.3 and 3.4, respectively. Finally, Section 3.4 and 3.5 show how to optimize the SOIS in fully and semi-supervised way, respectively.

3.1 Problem Definition of OWIS

The open-world instance segmentation (OWIS) aims to segment all the object instances (things) of any class including those that did not appear in the training phase. Technically, OWIS is a task to produce a set of binary masks, where each mask corresponds to a class-agnostic instance. The pixel value of 1 in the mask indicates a part of an object instance while 0 indicates not.

3.2 Model Architecture

Our proposed SOIS framework consists of three branches to alleviate the incomplete annotating, as shown in Figure 6. Basically, we follow the design of one-stage Mask2Former [6]. The **objectness prediction branch** estimates the weighting score for each mask by applying a sequential Transformer decoder and MLP. The **mask prediction branch** predicts the binary mask for each instance. It first generates N binary masks with N ideally larger than the actual instance number K_i . Each mask is multiplied by a weighting score with a value between 0 and 1, indicating if a mask should be selected as an instance mask. This process generates the mask in an end to end way, which avoids to miss the instance because of poor detection bounding boxes and meanwhile reduces the redundant segmentation cost for each proposal. We refer to [6, 15] for more details including the training procedure.

The **foreground prediction branch** is a light-weight fully convolutional network to estimate the foreground regions that belong to any object instance. The more detailed design of the foreground prediction branch is in the Appendix. This guides the training of the mask branch through our cross-task consistency loss proposed in the following Sec. 3.3. Once training is done, we discard this branch and only use the objectness and mask prediction branch at inference time. Therefore, we would not introduce any additional parameter or computational redundancy, which benefits the running efficiency.

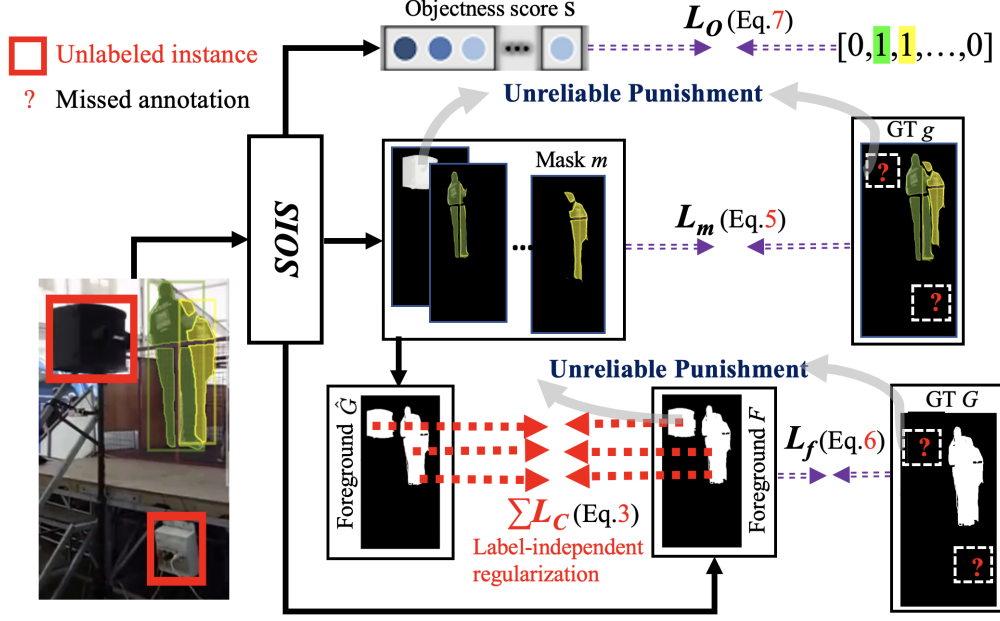


Figure 3: Working principle of consistency loss.

3.3 Learning with the cross-task consistency regularization

A critical limitation of the OWIS is the never-perfect annotations due to the difficulties in annotating class-agnostic object instances. Towards alleviating this issue, we propose a regularization to provide extra supervision to guide the OWIS model training under incomplete annotations.

We construct a branch to predict the foreground regions that belong to any of the object instance. Formally we create the the foreground annotation $G(x, y)$ calculated by

$$G(x, y) = \begin{cases} 0, & \text{if } \sum_{i=1}^K g^i(x, y) == 0 \\ 1, & \text{otherwise,} \end{cases} \quad (1)$$

where $g^i(x, y)$ is one of the K annotated object instances for the current image and the union of g^i defines the foreground object regions. Here (x, y) denotes a coordinate of a pixel in the an image. We use $G(x, y)$ as labels to train the foreground prediction branch.

Our consistency loss encourages the model outputs to have the relationship indicated in Eq 1, which states that the the foreground prediction should be the union of instance predictions. To do so, we use the following equation as an estimate of the foreground from the instance prediction:

$$\hat{G}(x, y) = \Phi \left(\sum_{j=1}^K m^j(x, y) \right), \quad (2)$$

where m^j means the confidence of pixels in j -th predicted mask, and Φ represents the Sigmoid function. Then, let the foreground prediction from the foreground prediction branch be F , our cross-task consistency loss is to make F and $\hat{G}(x, y)$ consistent, which finally leads to the following loss function.

$$L_c = \text{DICE}_{(\hat{G}, F)} + \text{BCE}_{(\hat{G}, F)}, \quad (3)$$

where DICE and BCE denote the dice-coefficient loss [19] and binary cross-entropy loss, respectively.

Consistency loss enjoys the following appealing properties. It is self-calibrated and independent with the incompleteness level of labels. As shown in Figure 3, for a instance mistakenly annotated as background, but the foreground prediction branch and mask prediction branch both correctly find

it, the model would be punished through mask loss and foreground loss. However, the consistency loss think this prediction is correct. In this way, consistency loss down-weights the adverse effects caused by other unreliable segmentation loss. The mitigation and the compensation factor synergize to relieve the overwhelming punishments on unlabeled instances.

3.4 Fully-supervised learning

The overall fully-supervised optimization of the proposed SOIS is carried out by minimizing the following joint loss formulation L_f ,

$$L_f = \alpha L_m + \beta L_p + \gamma L_c + \omega L_o, \quad (4)$$

$$\text{where } L_m = \text{BCE}_{(m,g)} + \text{DICE}_{(m,g)}, \quad (5)$$

$$L_p = \text{BCE}_{(F,G)} + \text{DICE}_{(F,G)}, \quad (6)$$

$$L_o = \text{BCE}_{(s,v)}, \quad (7)$$

where L_m , L_p and L_o denote the loss terms for mask prediction, foreground prediction, and objectness scoring, respectively. α , β , γ and ω are the weights of the corresponding losses. m and g represent the predicted masks and corresponding groundtruth, respectively. F and G is the foreground prediction result and the generated foreground groundtruth, while the estimated objectness score is denoted with s . v is a set of binary values that indicate whether each mask is an instance. Before computing the L_m , matching between the set of predicted masks and groundtruth has been done via the bipartite matching algorithm defined in [6].

3.5 Extension to semi-supervised learning

Due to the ambiguity of the instance definition in OWIS, it is much harder for the annotators to follow the annotation instruction, and this could make the annotations for OWIS expensive. It is desirable if we can use unlabeled data to help train OWIS models. In this regard, our proposed cross-task consistency loss only requires the outputs of both predictors to have a consistent relationship indicated in 1, and does not always need ground truth annotations. Thus, we apply this loss to unlabeled data, which becomes semi-supervised learning. Specifically, the easier-to-learn foreground prediction branch is able to learn well through a few labeled images in the warm-up stage. Then the resulted foreground map can serve as a constraint to optimize the open-world mask predictions with the help of our cross-task consistency loss, when the labels do not exist. In this way, our Semi-SOIS achieves a good trade-off between the annotation cost and model accuracy.

Semi-supervised learning process. Given a labeled set $D_l = \{(x_i, y_i)\}_{i=1}^{N_l}$ and an unlabeled set $D_u = \{x_i\}_{i=1}^{N_u}$, our goal is to train an OWIS model by leveraging both a large amount of unlabeled data and a smaller set of labeled data. Specifically, we initially use D_l to train the SOIS as a warm-up stage, giving a good initialization for the model. We then jointly train the OWIS model on the both labeled and unlabeled data. For the labeled data, we employ the loss function defined in Eq 4. For the unlabeled data, we apply only the cross-task consistency loss L_c .

4 Experiments

For demonstrating the effectiveness of our proposed SOIS, we compared it with other fully-supervised methods through intra-dataset and cross-dataset evaluations. We also performed ablation studies in these two settings to show the effect of each component. Moreover, we apply the proposed cross-task consistency loss for semi-supervised learning and test our method on the UVO validation set.

4.1 Implementation details and evaluation metrics

Implementation details Detectron2 [20] is used to implement the proposed SOIS framework, multi-scale feature maps are extracted from the ResNet-50 [21] or Swin Transformer [22] model pre-trained on ImageNet [23]. Our transformer encoder-decoder design follows the same architecture as in Mask2Former [6]. The number of object queries M is set to 100. Both the ResNet and Swin backbones use an initial learning rate of 0.0001 and a weight decay of 0.05. A simple data augmentation method, Cutout [24], is applied to the training data. All the experiments have been done on 8 NVIDIA V100 GPU cards with 32G memory.

Pseudo-labeling for COCO train set Pseudo-labeling is a common way to handle incomplete annotations. To explore the compatibility of our method and the pseudo-labeling operation, we employ a simple strategy to generate pseudo-labels for unannotated instances in the COCO train

Table 1: Results of **UVO-train** \rightarrow **UVO-val** intra-dataset evaluation.

Metric	Backbone	AP ₁₀₀ (%)	AP _s (%)	AP _m (%)	AP _l (%)	AR ₁₀₀ (%)	AR ₁₀ (%)
MaskRCNN	R-50	13.41	4.91	12.33	17.45	22.77	20.01
LDET	R-50	16.25	3.27	13.58	22.93	35.64	23.73
Mask2Former	R-50	21.85	6.16	16.82	31.65	41.18	28.26
SOIS (Ours)	R-50	23.38	6.59	17.35	34.23	41.94	29.24
Mask2Former	Swin-B	33.27	9.34	25.21	47.80	50.81	37.49
SOIS (Ours)	Swin-B	38.02	12.31	28.64	53.22	54.74	41.78

set [1] in our experiments. Specifically, we follow a typical self-training framework, introducing the teacher model and student model framework to generate pseudo-labels. These two models have the same architecture, as shown in Figure 6, but are different in model weights. The weights of the student model are optimized by the common back-propagation, while the weight of the teacher model is updated by computing the exponential moving averages (EMA) of the student model. During training, the image i is first fed into the teacher model to generate some mask predictions. The prediction whose confidence is higher than a certain value would be taken as a pseudo-proposal. The state S_{ij} of the pseudo-proposal p_{ij} is determined according to Equation (8).

$$S_{ij} = \begin{cases} \text{True,} & \text{if } \text{argmax}(\varphi(p_{ij}, g_i)) \leq \varepsilon, \\ \text{False,} & \text{otherwise,} \end{cases} \quad (8)$$

in which g_i means any ground truth instance in the image i . φ denotes the IOU calculating function, and ε is a threshold to further filter the unreliable pseudo-proposals. Finally, pseudo-proposals with states *True* would be considered as reliable pseudo-labels. Here, the confidence and IOU threshold ε for selecting pseudo-labels are set to 0.8 and 0.2, respectively. Then, we jointly use the ground truth and the pseudo-labels to form the training data annotations. If a region is identified as belonging to an instance in the pseudo-label, it will be considered as a positive sample during training.

Evaluation metrics The Mean Average Recall (AR) and Mean Average Precision (AP) [1] are utilized to measure the performance of approaches in a class-agnostic way.

4.2 Fully-supervised experimental setting

Intra-dataset evaluation UVO is the largest open-world instance segmentation dataset. Its training and test images are selected from the same domain, while they do not have any overlap. Here, we perform the leaning process of SOIS on the UVO-train subset and conduct the test experiments on the UVO-val subset. Besides, we split the COCO dataset into 20 seen (VOC) classes and 60 unseen (none-VOC) classes. We train a model only on the annotation of 20 VOC classes and test it on the 60 none-VOC class, evaluating its ability of discovering novel objects.

Table 2: Results of **COCO2017-train(VOC)** \rightarrow **COCO2017-val(none-VOC)** intra-set evaluation.

Metrics	AR ₁₀₀	AR _s	AR _m	AR _l
Mask2Former	9.21	4.56	8.79	19.30
SOIS	11.03	4.87	9.24	26.81

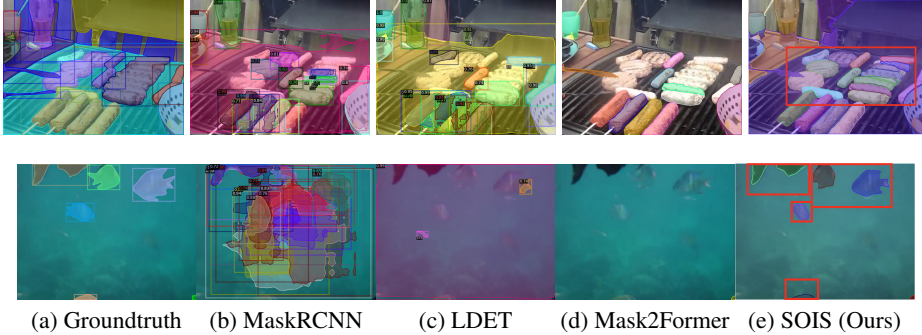
Cross-dataset evaluation Open-world setting assumes that the instance can be novel classes in the target domain. Therefore, it is essential for the OWIS method to handle the potential domain gap with excellent generalization ability. Cross-dataset evaluation, in which training and test data come from different domains, is necessary to be conducted. Here, we first train the proposed SOIS model and compared methods on the COCO-train subset, while testing them on the UVO-val dataset to evaluate their generalizability. Then we extend the experiments to an autonomous driving scenario, training the models on the Cityscapes [2] dataset and evaluating them on the Mapillary [25]. Cityscape have 8 foreground classes, while Mapillary contains 35 foreground classes including vehicles, animals, trash can, mailbox, etc.

4.3 Fully-supervised experimental results

Intra-dataset evaluation The results are illustrated in Table 1. The single-stage approaches based on the mask classification framework perform better than other two-stage methods. Among them, our proposed SOIS achieves a significant performance improvement over the Mask2Former baseline, which is 4.75% in AP_{100} and 3.93% in AR_{100} when using the Swin-B backbone. For VOC \rightarrow none-VOC setting, the experimental results are shown in Table 2, which verified that our proposed method can improve the performance for all instances, especially large ones.

Table 3: Results of COCO2017-train \rightarrow UVO-val cross-dataset evaluation.

Metric	Backbone	AR ₁₀₀	AP ₁₀₀	AP _s	AP _m	AP _l
MaskRCNN	R-50	38.17	19.05	6.27	13.15	28.05
LDET	R-50	42.63	21.27	5.66	17.52	18.38
GGN	R-50	43.30	20.30	8.70	18.20	27.30
Mask2Former	R50	48.71	25.24	6.46	16.09	40.37
SOIS (Ours)	R-50	51.28	27.62	7.80	18.61	43.42
Mask2Former	Swin-B	51.38	28.16	7.29	18.91	45.48
SOIS(Ours)	Swin-B	54.86	32.21	9.03	21.92	50.69

Figure 4: Visualization results of COCO \rightarrow UVO cross-dataset evaluation. The predicted boxes of two-stage methods MaskRCNN and LDET are also drawn. Proposed SOIS can discover both unlabeled object (first row) and unseen class of instances (second row) as shown in red boxes.

Cross-dataset evaluation For the COCO \rightarrow UVO task, according to Table 3, it is clear that the proposed SOIS outperforms all previous methods, achieving a new state-of-the-art AR_{100} at 54.86% which is 11.56% higher than previous state-of-the-art method GGN [4]. We also applied the proposed techniques to another classic one-stage method SOLO V2 [26]. The experimental results in Table 4 show that it improves AR_{100} and AP_{100} by 3.11% and 2.79% compared to SOLO V2. For the Cityscape \rightarrow Mapillary task, the overall AP and AR of SOIS still surpass the performance of other state-of-the-art methods (in Table 8), which demonstrates the effectiveness of our proposed techniques. We show some of the COCO \rightarrow UVO visualization results in Figure 7 to qualitatively demonstrate the superiority of our method. Please refer to the supplementary material for more qualitative examples.

Table 4: Results of SOIS with SOLOV2 structure (UVO-train \rightarrow UVO-val).

Metric	Backbone	AR ₁₀₀	AP ₁₀₀	AP _s	AP _m	AP _l
SOLO V2	R-50	39.41	22.25	5.56	14.18	34.12
SOLO V2SOIS	R-50	42.52	25.04	6.77	16.90	38.33

4.4 Ablation study

We perform cross-dataset and intra-dataset ablation studies to analyze the effectiveness of each component in the proposed SOIS, including the foreground prediction branch and the cross-task consistency loss. We also try combinations of the pseudo-label generation strategy and our cross-task consistency loss to investigate the individual and synergetic effects of them. Using the SwinB backbone, these models are trained on the COCO-train subset and the UVO-train subset, respectively. The metrics reported in Table 5 are tested on the UVO-val dataset.

Effectiveness foreground prediction branch Table 5 shows that although a separate foreground prediction branch can guide the method to optimize towards the direction of discovering foreground pixels, it only slightly boosts the performance.

Effectiveness of cross-task consistency loss Cross-task consistency loss has a positive effect on both sparse annotated (COCO) and dense annotated (UVO) training dataset. The values of AP_{100} and AR_{100} increase significantly (2.74% \uparrow and 2.49% \uparrow on COCO while 2.90% \uparrow and 3.35% \uparrow on UVO) after applying the cross-task consistency loss as well as the foreground prediction branches together. This result outperforms the SOIS counterpart with only pseudo-labeling, showing our effectiveness. In addition, jointly utilizing our cross-task consistency loss as well as the pseudo-labeling strategy

Table 5: Ablation results of the proposed components by cross-dataset and intra-dataset evaluations. Foreground prediction (FP), Cross-task consistency (CTC) loss, Pseudo label (PL).

Component			Train on COCO		Train on UVO	
FP	CTC loss	PL	AP ₁₀₀ (%)	AR ₁₀₀ (%)	AP ₁₀₀ (%)	AR ₁₀₀ (%)
✓		✓	28.65	51.54	35.12	51.39
			29.02	51.60	35.55	51.73
			30.09	52.97	32.94	51.64
✓	✓	✓	31.39	53.83	38.02	54.74
✓		✓	30.17	52.98	33.35	50.90
✓	✓	✓	32.21	54.86	37.71	52.27

Table 6: Results of our SOIS and classic semi-supervised method on UVO-val.

Training Data	UVO-train with 30% annotation			
	Fully-SOIS ₃₀	Mean Teacher	Pseudo Labeling	Semi-SOIS ₃₀
AP ₁₀₀ (%)	21.67	21.95	22.77	25.03
AR ₁₀₀ (%)	40.09	40.82	41.56	45.42

Table 7: Results of our SOIS and recent end to end method on UVO-val.

Training Data	UVO-train with 50% annotation			
	LDEF ₅₀	Mask2Former ₅₀	Fully-SOIS ₅₀	Semi-SOIS ₅₀
AP ₁₀₀ (%)	10.61	19.49	22.86	25.22
AR ₁₀₀ (%)	25.08	38.08	41.44	47.56

leads to performance improvements on two settings, which demonstrates the synergistic effect of both approaches.

Effectiveness of pseudo-labeling Pseudo-labeling is not always necessary and powerful for any types of datasets. As shown in Table 5, the AP_{100} and AR_{100} of the COCO trained model increase by 1.35% and 0.78%, respectively, after applying the pseudo-label generation. However, pseudo-labeling causes a performance degradation (e.g. 2.18%↓ in AP_{100}) to a model trained in the UVO dataset. Compared with COCO, the UVO dataset is annotated more densely. We conjecture that the background annotations of UVO are more reliable than those of COCO, where carefully selected pseudo-labels are more likely to represent unlabeled objects. The generated pseudo-labels of UVO contain higher noises than those of COCO. These additional noisy labels mislead the model training.

4.5 Semi-supervised learning experiment

Experimental setting We have divided the UVO-train dataset into the labeled subset D_l and the unlabeled subset D_u . Semi-supervised model Semi-SOIS is optimized as described in Section 3.6 on $D_l \cup D_u$, while the fully-supervised method Fully-SOIS is trained merely on the D_l . To ensure the comprehensiveness of the experiment, two different data division settings are included in our experiments: $\{D_L=30\%, D_u=70\%\}$ and $\{D_L=50\%, D_u=50\%\}$. The backbone applied here is Swin-B. We also implemented the classic Mean teacher model and a simple pseudo-label method based on the Mask2Former to perform comparison.

Results and analysis As presented in Figure 5, the Semi-SOIS₅₀ model trained on the UVO with 50% annotated data outperforms the Semi-SOIS₃₀ model leaning with 30% labeled training images. However, the performance increase between the Semi-SOIS₃₀ and Semi-SOIS₅₀ is slight. In addition, Semi-SOIS₃₀ improves Fully-SOIS₃₀ by 3.36% and 5.33% in AP_{100} and AR_{100} , respectively. Compared to Fully-SOIS₅₀, Semi-SOIS₅₀ still achieves significant advantages (2.36% in AP_{100} and 6.12% in AR_{100}). These results reflect that cross-task consistency loss has the ability to extract information from unlabeled data and facilitates model optimization in the semi-supervised setting. It is notable that the results of Semi-SOIS₃₀ are even better than those of Fully-SOIS₅₀. This illustrates that the information dug out by the cross-task consistency loss from the remaining 70% unlabeled data is more abundant than that included in 20% fully-labeled data. Therefore, our algorithm can achieve better

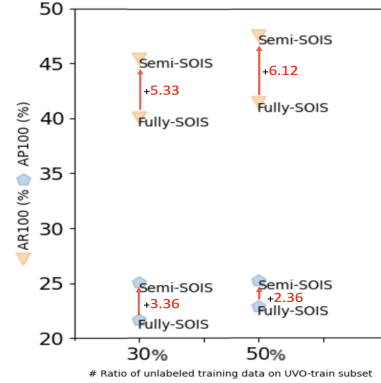


Figure 5: Comparison between fully-SOIS and semi-SOIS

performance with fewer annotations. This characteristic is promising in solving the OWIS problem. In addition, we also compared the semi-*SOIS* with classic semi-supervised method and recent end to end segmentation method. The results in Table 6 and 7 show our advantages over the compared methods.

5 Conclusion

This paper proposes the first single-stage framework (SOIS) for the open-world instance segmentation task. Apart from predicting the instance mask and objectness score, our framework introduces a foreground prediction branch to segment the regions belonging to any instance. Utilizing the outputs of this branch, we propose a novel cross-task consistency loss to enforce the foreground prediction to be consistent with the prediction of the instance masks. We experimentally demonstrate that this mechanism alleviates the problem of incomplete annotation, which is a critical issue for open-world segmentation. Our extensive experiments demonstrate that SOIS outperforms state-of-the-art methods by a large margin on typical datasets. We further demonstrate that our cross-task consistency loss can utilize unlabeled images to obtain some performance gains for a semi-supervised instance segmentation. This is an important step toward reducing laborious and expensive human annotation.

Table 8: Cross-set evaluation on autonomous driving scenes. Results of **Cityscapes** \rightarrow **Mapillary**.

Method	MaskRCNN	LDET	Mask2Former	OSIS(Ours)
AP(%)	7.3	7.8	7.6	8.4
AR ₁₀ (%)	6.1	5.5	7.0	7.5

Acknowledgments and Disclosure of Funding

Mike Zheng Shou is supported only by the National Research Foundation, Singapore under its NRFF award NRF-NRFF13-2021-0008.

A Appendix

In this appendix, we provide the architecture of the foreground prediction branch (in Figure 6) and detailed experimental settings first. Then some annotations in UVO dataset are visualized in Figure 7 to show the challenges of open world instance segmentation. Finally, additional visualization results of proposed SOIS are shown in Figure 8.

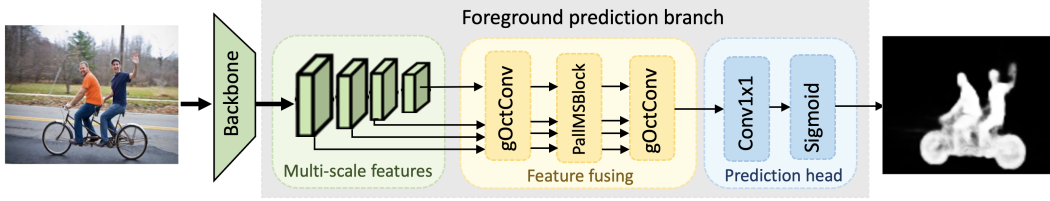


Figure 6: **Architecture of foreground prediction branch.** Multi-scale features extracted from backbone are fed into the feature fusing module to exchange and fuse the multi-scale information. Then a fused feature is sent to the prediction head to predict the final foreground map. Considering the efficiency, we follow [27] to introduce the gOctConv [27] and PalIMSBBlock [27] to perform feature fusing.

A.1 Detailed experimental settings

Implementation details For feature extracting, we obtain the multi-scale features through a sequential backbone network [21, 22], and FPN [28]. The multi-scale features contain D-dimensional feature maps with resolutions of 1/4, 1/8, 1/16, and 1/32. In the pixel decoder module, six MSDeformAttn layers are employed, while the transformer decoder have three layers with 100 queries by default.

In fully-supervised learning, the total loss L_f can be formulated as: $L_f = \alpha L_m + \beta L_p + \gamma L_c + \omega L_o$. We set the weight α of mask loss (L_m) to 5.0, the weight β of foreground loss (L_p) to 1.0, the weight γ of cross-task consistency loss (L_c) to 1.0 and the weight ω of objectness loss (L_o) to 2.0.

Training settings Specifically, AdamW [29] optimizer and the step learning rate schedule are applied to optimize our model. An initial learning rate of 0.0001 and a weight decay of 0.05 are utilized for all backbones. We set a learning rate multiplier of the backbone to 0.1 and we decay the learning rate at 0.9 and 0.95 fractions of the total number of training steps by a factor of 10. For data augmentation, we use the large-scale jittering (LSJ) augmentation with a random scale sampled from range 0.1 to 2.0 followed by a fixed size crop to 1024×1024 on COCO dataset and 640×640 on UVO dataset. Besides, a Cutout [30] strategy that randomly cuts out a region of size $[1/8 \cdot w, 1/8 \cdot h]$ to $[1/3 \cdot w, 1/3 \cdot h]$ is introduced during training. On COCO dataset, we train our models for 38×10^4 iterations with a batch size of 16, while on UVO dataset, we train our models for 12×10^4 iterations with the same batch size.

SOIS training process with pseudo-labeling on COCO dataset

Algorithm 1: SOIS training process with pseudo-labeling

Data: Image dataset

Result: Proposed SOIS Model M_u

```

1 initialization the student model  $M_u$ , and teacher model  $M_t = M_u.copy()$ ;
2 while Image  $i \notin \emptyset$  do
3   read image  $i$  and corresponding groundtruth  $gt_i$ ;
4   extract backbone feature  $X_i$ ;
5   pred_masks  $\leftarrow M_t.predictor(X_i)$ ;
6   pseudo_proposals  $\leftarrow filter\_masks\_with\_confidence(pred\_masks, confidence\_threshold)$ ;
7   pseudo_labels  $\leftarrow filter\_masks\_with\_IOU(pseudo\_proposals, IOU\_threshold)$ ;
8   training_labels  $\leftarrow merge(gt_i, pseudo\_labels)$ ;
9   aug_data  $\leftarrow Cutout(X_i, training\_labels)$ ;
10   $M_u \leftarrow M_u.training(aug\_data)$ ;
11   $M_t \leftarrow M_t.EMA\_update(M_t, M_u)$ 
12 end

```



Figure 7: Visualizations of UVO annotations. It is notable that the same class of object may be labeled as an instance or as background in different images. (as shown in the area highlighted by the ellipse). This inconsistency of annotations pose a great challenge to the algorithms.



Figure 8: Visualizations results of our proposed SOIS in UVO dataset. SOIS can discover many novel objects, as shown in regions in red boxes.

A.2 Visualization of annotations and our results on UVO dataset

Unlike in closed-world instance segmentation, where the object categories have been clearly defined, instance definition in OWIS is much more ambiguous and harder for annotators to follow. Inevitably, the instance annotation could become inconsistent across images, as shown in Figure ?? . Our method is motivated by this observation that the instance annotation in the existing datasets is very noisy. Our solution to this issue is to introduce a self-correcting mechanism to combat erroneous annotations, which provides additional guidance to both prediction tasks when the noisy annotations fail to provide correct supervision. The visualization results in Figure 7 demonstrate that our proposed SOIS can segment many novel objects that have not been unseen in the training set.

References

- [1] 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 *European conference on computer vision*, pp. 740–755, Springer, 2014.
- [2] 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. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [3] W. Wang, M. Feiszli, H. Wang, and D. Tran, “Unidentified video objects: A benchmark for dense, open-world segmentation,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 10776–10785, 2021.
- [4] W. Wang, M. Feiszli, H. Wang, J. Malik, and D. Tran, “Open-world instance segmentation: Exploiting pseudo ground truth from learned pairwise affinity,” *CVPR*, 2022.
- [5] K. Saito, P. Hu, T. Darrell, and K. Saenko, “Learning to detect every thing in an open world,” *arXiv preprint arXiv:2112.01698*, 2021.
- [6] B. Cheng, I. Misra, A. G. Schwing, A. Kirillov, and R. Girdhar, “Masked-attention mask transformer for universal image segmentation,” 2022.
- [7] K. He, G. Gkioxari, P. Dollár, and R. Girshick, “Mask r-cnn,” in *Proceedings of the IEEE international conference on computer vision*, pp. 2961–2969, 2017.
- [8] H. Chen, K. Sun, Z. Tian, C. Shen, Y. Huang, and Y. Yan, “Blendmask: Top-down meets bottom-up for instance segmentation,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 8573–8581, 2020.
- [9] J. Dai, K. He, and J. Sun, “Instance-aware semantic segmentation via multi-task network cascades,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3150–3158, 2016.
- [10] D. Bolya, C. Zhou, F. Xiao, and Y. J. Lee, “Yolact: Real-time instance segmentation,” in *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 9157–9166, 2019.
- [11] X. Wang, T. Kong, C. Shen, Y. Jiang, and L. Li, “Solo: Segmenting objects by locations,” in *European Conference on Computer Vision*, pp. 649–665, Springer, 2020.
- [12] Y. Fang, S. Yang, X. Wang, Y. Li, C. Fang, Y. Shan, B. Feng, and W. Liu, “Instances as queries,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 6910–6919, 2021.
- [13] B. Dong, F. Zeng, T. Wang, X. Zhang, and Y. Wei, “Solq: Segmenting objects by learning queries,” *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [14] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, “End-to-end object detection with transformers,” in *European conference on computer vision*, pp. 213–229, Springer, 2020.
- [15] B. Cheng, A. Schwing, and A. Kirillov, “Per-pixel classification is not all you need for semantic segmentation,” *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [16] Y. Du, W. Guo, Y. Xiao, and V. Lepetit, “1st place solution for the uvo challenge on video-based open-world segmentation 2021,” *arXiv preprint arXiv:2110.11661*, 2021.

- [17] T. Vu, H. Jang, T. X. Pham, and C. Yoo, “Cascade rpn: Delving into high-quality region proposal network with adaptive convolution,” *Advances in neural information processing systems*, vol. 32, 2019.
- [18] Z. Ge, S. Liu, F. Wang, Z. Li, and J. Sun, “Yolox: Exceeding yolo series in 2021,” *arXiv preprint arXiv:2107.08430*, 2021.
- [19] F. Milletari, N. Navab, and S.-A. Ahmadi, “V-net: Fully convolutional neural networks for volumetric medical image segmentation,” in *2016 fourth international conference on 3D vision (3DV)*, pp. 565–571, IEEE, 2016.
- [20] Y. Wu, A. Kirillov, F. Massa, W.-Y. Lo, and R. Girshick, “Detectron2,” <https://github.com/facebookresearch/detectron2>, 2019.
- [21] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- [22] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, “Swin transformer: Hierarchical vision transformer using shifted windows,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021.
- [23] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *2009 IEEE conference on computer vision and pattern recognition*, pp. 248–255, Ieee, 2009.
- [24] T. DeVries and G. W. Taylor, “Improved regularization of convolutional neural networks with cutout,” *arXiv preprint arXiv:1708.04552*, 2017.
- [25] G. Neuhold, T. Ollmann, S. Rota Bulò, and P. Kotschieder, “The mapillary vistas dataset for semantic understanding of street scenes,” in *Proceedings of the IEEE international conference on computer vision*, pp. 4990–4999, 2017.
- [26] X. Wang, R. Zhang, T. Kong, L. Li, and C. Shen, “Solov2: Dynamic and fast instance segmentation,” *Proc. Advances in Neural Information Processing Systems (NeurIPS)*, 2020.
- [27] S.-H. Gao, Y.-Q. Tan, M.-M. Cheng, C. Lu, Y. Chen, and S. Yan, “Highly efficient salient object detection with 100k parameters,” in *European Conference on Computer Vision*, pp. 702–721, Springer, 2020.
- [28] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, “Feature pyramid networks for object detection,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2117–2125, 2017.
- [29] I. Loshchilov and F. Hutter, “Decoupled weight decay regularization,” *arXiv preprint arXiv:1711.05101*, 2017.
- [30] T. DeVries and G. W. Taylor, “Improved regularization of convolutional neural networks with cutout,” *arXiv preprint arXiv:1708.04552*, 2017.