

# Unseen Object Instance Segmentation with Fully Test-time RGB-D Embeddings Adaptation

Lu Zhang, Siqu Zhang, Xu Yang, Hong Qiao and Zhiyong Liu\*

**Abstract**—Segmenting unseen objects is a crucial ability for the robot since it may encounter new environments during the operation. Recently, a popular solution is leveraging RGB-D features of large-scale synthetic data and directly applying the model to unseen real-world scenarios. However, the domain shift caused by the sim2real gap is inevitable, posing a crucial challenge to the segmentation model. In this paper, we emphasize the adaptation process across sim2real domains and model it as a learning problem on the BatchNorm parameters of a simulation-trained model. Specifically, we propose a novel non-parametric entropy objective, which formulates the learning objective for the test-time adaptation in an open-world manner. Then, a cross-modality knowledge distillation objective is further designed to encourage the test-time knowledge transfer for feature enhancement. Our approach can be efficiently implemented with only test images, without requiring annotations or revisiting the large-scale synthetic training data. Besides significant time savings, the proposed method consistently improves segmentation results on the overlap and boundary metrics, achieving state-of-the-art performance on unseen object instance segmentation.

## I. INTRODUCTION

In recent years, a rising development trend in robotics is the transition from controlled labs to unstructured environments. When encountering new environments and objects, the robot system must have the ability to adjust itself and recognize unseen objects. Such capability is essential for robots to understand working environments better and perform various manipulation tasks. To achieve this goal, we approach the task of Unseen Object Instance Segmentation (UOIS) [1], [2], [3], [4], which aims to conduct instance-aware segmentation of unseen objects in tabletop scenes. In UOIS, the robot system needs to learn the concept of “object” and generalize it to unseen ones.

However, unlike ImageNet [5] and MS COCO [6], which have spurred significant development of the classification and object detection for natural images, a large-scale realistic dataset that contains sufficient objects for robotic manipulation scenes is currently unavailable [2]. As a result, existing UOIS methods [1], [2], [3], [4] generally resort to large-scale synthetic RGB-D data to train a perception model. For example, Xie *et al.* [1] propose using synthetic scenes that can be rendered into RGB-D images from various viewpoints with automatically generated labels. Previous works [1], [2], [3], [4] typically use synthetic depth or RGB-D images to train a

model that can separately segment each object instance, since the depth inputs have better generalization potential than the non-photorealistic RGB images. After training, the model is directly deployed on unseen realistic datasets. Though such a workaround achieves good performance, it has in turn led to the neglect of domain shift caused by the “sim2real gap”. For robot perception, the synthetic data fail to model many aspects of the real world, like the object texture, lighting conditions, depth noises, etc. This sim2real gap degrades the model’s performance on realistic data, especially on those in unseen environments.

Therefore, rather than betting the generalization ability to elaborate models as in previous works [1], [2], [3], [4], we place emphasis on the model’s adaptability across domains. Specifically, we propose a Fully Test-time RGB-D Embeddings Adaptation (FTEA) framework to mitigate the domain gap between synthetic and unseen realistic data. First, we propose a novel Non-parametric Entropy Objective (NEO) that can be calculated without the explicit classification layer to enable test-time adaptation for the open-world UOIS task. The NEO leverages non-parametric distributions to formulate Shannon entropy [7] as the learning objective, since the Shannon entropy is shown to be related to error and shift [8], *i.e.*, more confident predictions are generally more correct and have fewer shifts. Second, we design a Cross-modality Knowledge Distillation (CKD) module to encourage knowledge transfer during testing. Different from other KD methods [9], [10], CKD aims to distill the knowledge from multimodal (RGB-D) features to unimodal (RGB or depth) ones, thus enhancing the unimodal features for better fusion. Finally, we utilize the affine transformation provided by BatchNorm (BN) [11] layer as the modulation parameters to conduct fully test-time RGB-D embeddings adaptation.

Given the simulation-trained model, FTEA is independent of re-training on the large-scale synthetic data and does not introduce extra parameters, which establishes a key advantage of flexibility. Meanwhile, compared to the inference process (taking 1~10s per frame in recent state-of-the-arts [2], [3]), the computation overhead of the proposed adaptation is negligible (taking 0.04s per frame with a total of 500 iterations), making FTEA particularly efficient. The proposed method is evaluated on two real-world RGB-D image datasets for the unseen object instance segmentation, *i.e.*, OSD [12] and OCID [13]. Extensive experiments show that FTEA consistently improves segmentation performances and achieves state-of-the-art results on various evaluation metrics, demonstrating the effectiveness of our method.

All authors are with State Key Laboratory of Multimodal Artificial Intelligence Systems, Institute of Automation, Chinese Academy of Sciences, Beijing, China, and also with the School of Artificial Intelligence, University of Chinese Academy of Sciences, Beijing, China.

Zhiyong Liu is also with Nanjing Artificial Intelligence Research of IA, Jiangning District, Nanjing, 211100, Jiangsu, China.

{lu.zhang@ia.ac.cn, zhiyong.liu@ia.ac.cn}

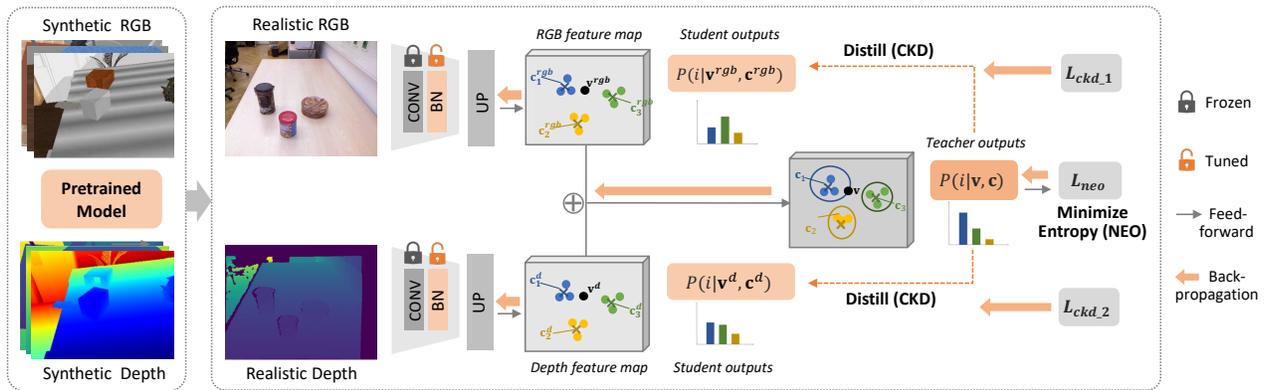


Fig. 1: Illustration of the proposed FTEA. We use the two-stream CNN architecture for RGB-D inputs, which is largely simplified for better visualization. During test time, all convolutional layers are frozen, only the affine parameters and normalization statistics of the BN layer are modulated based on two novel objectives, *i.e.*, NEO and CKD.

## II. RELATED WORK

**Unseen Object Instance Segmentation** UOIS aims to conduct instance-aware segmentation of unseen objects in tabletop environments, which is useful for robots to perform various manipulation tasks. As a pioneer, Xie *et al.* [1], [4] tackle this problem by proposing a two-stage framework. The framework first operates only on depth to produce rough initial segmentation masks and then refines those masks with RGB. Then, Xiang *et al.* [2] introduce a fully convolutional network based model called UCN that can be trained end-to-end. Different from previous approaches that mainly rely on depth for segmentation, UCN utilizes both depth images and non-photorealistic RGB images to produce feature embeddings for every pixel, which can be used to learn a distance metric for clustering to segment unseen objects. Recently, several methods are proposed to tackle specific challenges in UOIS. For instance, RICE [3] focuses on the occlusion problem in clutter scenes and utilizes a graph-based representation of instance masks to refine the outputs of previous methods. UOAIIS-Net [14] presents a new unseen object amodal instance segmentation (UOAIIS) task to emphasize the amodal perception for robotic manipulation in a cluttered scene, and introduce a large-scale photo-realistic synthetic dataset named UOAIIS-SIM to improve the sim2real transferability. Other works related to UOIS include clustering [15], [16], [17], novel category discovery [18], [19], [20], [21], and unsupervised semantic segmentation [22], [23], [24].

Though previous UOIS methods have been demonstrated to be effective, the domain shift problem is not explicitly concerned and tackled. It is worth noting that the photo-realistic data used by UOAIIS-Net are also generated with rendered scenes in the simulator, thus synthetic. Differently, we present a new perspective for the solution of the domain shift problem in UOIS and propose an efficient framework to conduct the adaptation during test time.

**Test-time Adaptation** When the model is deployed, it is inevitable to encounter unlabelled images that are not observed before. This is a key characteristic of robot per-

ception and the UOIS task. Therefore, developing strategies to adapt the model at test time is essential. Recently, test-time training (TTT, TTT+) [25], [26] uses unlabelled test instances to conduct self-supervised learning as the adaptation. But such methods heavily rely on the choice of proxy tasks and also need to visit training data that could be unavailable in practice. To address the above limitations, Wang *et al.* [8] propose Tent to reduce generalization error by test-time entropy minimization. Tent optimizes itself according to its own predictions, which is not relevant to proxy tasks. However, the aforementioned test-time adaptation method is mainly applied in traditional close-set tasks such as image classification, and exploits the output distributions of classifiers as the test-time objective optimization. For the UOIS task, the model generally can not be trained with an explicit discriminative layer with a certain number of classes, which poses a major obstacle for the test-time adaptation in UOIS.

## III. METHOD

### A. Overview

1) *Network Architecture*: To deal with RGB-D inputs, we adopt the two-stream CNN with late fusion as our basic network architecture. As shown in Figure 1, a pair of realistic RGB-D images are separately processed with CNN, and the RGB and depth feature maps are bilinearly upsampled to the full resolution as the input images. Then the late fusion is conducted for a joint representation. We use the UCN [2] as our simulation-trained model due to its conciseness and end-to-end fashion.

2) *The Pipeline*: First, we construct a non-parametric entropy objective (NEO) for test-time adaptation in an open-world setting, as described in Section III-B. Then, in Section III-C, a cross-modality knowledge distillation (CKD) module is further proposed to encourage test-time knowledge transfer for feature enhancement. Finally, we fix all convolutional layers and minimize the proposed NEO and CKD loss to modulate the affine parameters and normalization statistics of the BN layer, as detailed in Section III-D.

## B. Non-parametric Entropy Objective

To modulate features during test time, a learning objective based on the model’s predictions of test data is typically required. Recent works [8], [27] have demonstrated the effectiveness of using the entropy objective based on discriminative outputs (*e.g.*, classification probabilities) [8], [27]. However, unlike the standard recognition model, UOIS is conducted in an open-world setting, and thus does not explicitly train a discriminative layer that generate logits and probabilities for the direct calculation of entropy. To address this problem, we propose a novel non-parametric entropy objective (NEO). Specifically, NEO leverages the unsupervised clustering to obtain centroids that present pseudo instance labels, then calculates non-parametric classification probabilities to construct the entropy objective.

**Unsupervised Clustering** Given a bunch of RGB-D embeddings on a feature map, we aim to cluster all pixels into groups to segment unseen objects. But the number of unseen objects is uncertain, which prevents the usage of clustering algorithms with a known number of clusters such as  $k$ -means or spectral clustering. Thus, we follow UCN [2] to use the mean shift [28] clustering algorithm with the von Mises-Fisher (vMF) distribution [29], which automatically discovers the number of objects and generates a segmentation mask for each object. After the unsupervised clustering, we can calculate the centroid of each cluster. The  $i$ -th cluster’s centroid vector  $\mathbf{c}_i$  is obtained by averaging all feature map vectors  $\mathbf{v}_{x,y}$  which belong to the  $i$ -th cluster  $\mathbf{C}_i$  as

$$\mathbf{c}_i = \text{avg}(\mathbf{v}_{x,y}) \text{ for all } \mathbf{v}_{x,y} \in \mathbf{C}_i, i = 1, 2, \dots, n \quad (1)$$

where  $x$  and  $y$  denote locations on the feature map along the  $x$ -axis and  $y$ -axis. After performing the average operation for each cluster, we obtain the set of all centroids  $\mathbf{c}$  as

$$\mathbf{c} = \{\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_n\}, \quad (2)$$

where  $n$  is the number of objects, which is estimated by the unsupervised clustering algorithm.

**Non-parametric Classification Probability** Different from the classification with standard supervised learning, UOIS segments objects in the instance-level, thus having an uncertain number of “classes”, which prevents the usage of the parametric classifier. Inspired by recent self-supervised learning methods with instance discrimination [30], [31], we propose to calculate classification probabilities in a non-parametric way.

Without specifying the number of classes, we conduct non-parametric classification by using the metric between the candidate feature vector  $\mathbf{v}$  and the  $i$ -th cluster’s centroid vector  $\mathbf{c}_i$  as

$$P(i|\mathbf{v}, \mathbf{c}) = \frac{\exp(s(\mathbf{v}, \mathbf{c}_i))}{\sum_i \exp(s(\mathbf{v}, \mathbf{c}_i))}, s(\mathbf{v}_1, \mathbf{v}_2) = \frac{\mathbf{v}_1^\top \mathbf{v}_2}{\|\mathbf{v}_1\| \|\mathbf{v}_2\|}, i \in N_k, \quad (3)$$

where  $s(\cdot, \cdot)$  is the cosine similarity to measure how well  $\mathbf{v}$  matches the  $i$ -th cluster/object,  $N_k$  is the set of numbers to indicate candidate  $\mathbf{v}$ ’s corresponding cluster and  $k$  nearest clusters.

**Entropy Objective** As an unsupervised objective, the Shannon entropy [7] is widely used and demonstrated to be effective for assessing the output error and shift [8]. With the proposed non-parametric classification probability, we can calculate the Shannon entropy as

$$L_{neo} = - \sum_i P(i|\mathbf{v}, \mathbf{c}) \log_2 P(i|\mathbf{v}, \mathbf{c}), \quad (4)$$

where  $P(i|\mathbf{v}, \mathbf{c})$  is the probability that  $\mathbf{v}$  is recognized as the  $i$ -th cluster given a collection of centroids  $\mathbf{c}$ .

## C. Cross-modality Knowledge Distillation

Compared to the “perfectly” generated RGB-D data, real-world RGB-D images have inevitable noise, such as the black holes in the depth images. The fusion of RGB-D features generally makes the network more stable to such noise [2], [32]. Therefore, while testing on realistic RGB-D images, the *privileged* [33] multimodal (*i.e.*, RGB-D fused) features can be utilized to enhance the unimodal (*i.e.*, RGB or depth) networks. Based on this motivation, we propose to distill the knowledge from the multimodal features to the unimodal ones to encourage the test-time knowledge transfer for feature enhancement. This cross-modality knowledge distillation (CKD) constructs another objective for test-time adaptation. In CKD, the full multimodal network (two-stream CNN) is set as the teacher, while the student is the smaller partial unimodal network (one-stream CNN).

We use soft targets (*i.e.*, probabilities of the input belonging to the classes) as the distilled knowledge in CKD. For each candidate feature vector  $\mathbf{v}$  on the teacher feature map, the  $i$ -th soft target is

$$P(i|\mathbf{v}, \mathbf{c}, T) = \frac{\exp(s(\mathbf{v}, \mathbf{c}_i)/T)}{\sum_i \exp(s(\mathbf{v}, \mathbf{c}_i)/T)}, \quad (5)$$

where  $T$  is the temperature factor to control the importance of each soft target. Due to the absence of a parametric classifier, here we use the non-parametric probability as similar with Equation 3.

For the individual RGB and depth modality (*i.e.*, the student), we use the same spatial cluster assignment as the teacher, but they hold different representations  $\mathbf{v}^{rgb}$  and  $\mathbf{v}^d$ , thus producing different cluster centroids  $\mathbf{c}_i^{rgb}$  and  $\mathbf{c}_i^d$ . Then, the soft targets of RGB and depth modality, *i.e.*,  $P(i|\mathbf{v}^{rgb}, \mathbf{c}^{rgb}, T)$  and  $P(i|\mathbf{v}^d, \mathbf{c}^d, T)$ , can be obtained similarly to Equation 5. Finally, we formulate the cross-modality knowledge distillation objective as

$$L_{ckd} = \frac{1}{2} (KL(P(i|\mathbf{v}, \mathbf{c}, T), P(i|\mathbf{v}^{rgb}, \mathbf{c}^{rgb}, T)) + KL(P(i|\mathbf{v}, \mathbf{c}, T), P(i|\mathbf{v}^d, \mathbf{c}^d, T))), \quad (6)$$

where  $KL(\cdot, \cdot)$  denotes the Kullback-Leibler divergence loss. The gradients propagated through the multimodal distribution are stopped. By optimizing Equation 6, the response knowledge of the multimodal teacher network is distilled into the unimodal student network.

Method	OSD [12]							OCID [13]							
	Overlap			Boundary				F@.75	Overlap			Boundary			
	P	R	F	P	R	F	P		R	F	P	R	F	F@.75	
Mask RCNN [34]	74.4	72.7	73.4	53.1	48.1	49.8	-	80.8	73.9	76.1	68.2	58.4	61.8	-	
UOIS-Net-2D [1]	80.7	80.5	79.9	66.0	67.1	65.6	71.9	88.3	78.9	81.7	82.0	65.9	71.4	69.1	
UOIS-Net-3D [4]	85.7	82.5	83.3	75.7	68.9	71.2	73.8	86.5	86.6	86.4	80.0	73.4	76.2	77.2	
UCN [2]	84.3	88.3	86.2	67.5	67.5	67.1	79.3	86.0	92.3	88.5	80.4	78.3	78.8	82.2	
UCN+ [2]	87.4	87.4	87.4	69.1	70.8	69.4	83.2	91.6	92.5	91.6	<b>86.5</b>	87.1	86.1	89.3	
UOAIS-Net [14]	85.3	85.4	85.2	<b>72.7</b>	74.3	73.1	79.1	70.7	86.7	71.9	68.2	78.5	68.8	78.7	
<b>FTEA (Ours)</b>	85.8	<b>92.0</b>	88.6	69.2	75.7	71.7	87.3	86.2	<b>93.9</b>	89.5	79.5	79.5	79.1	85.1	
<b>FTEA+ (Ours)</b>	<b>89.9</b>	89.4	<b>89.5</b>	72.6	<b>76.0</b>	<b>73.8</b>	<b>88.3</b>	<b>92.0</b>	93.3	<b>92.3</b>	<b>86.5</b>	<b>88.0</b>	<b>86.7</b>	<b>91.1</b>	

TABLE I: The unseen object instance segmentation (UOIS) performances of the proposed FTEA and other state-of-the-art methods on OSD [12] and OCID [13] datasets. “+” denotes the zoom-in operation [2] to refine segmentation results.

#### D. Fully Test-time RGB-D Embeddings Adaptation

By minimizing the above entropy objective  $L_{neo}$  and distillation objective  $L_{ckd}$ , we can adapt our model fully in test time. However, tuning all parameters like that in the training phase is inefficient and could easily cause model instability [8]. As the channel of the feature map can be considered as a feature detector [35], [36], re-calibrating channel responses has been widely studied and utilized in network pruning [37], multimodal fusion [38], [39], representation learning [40], [41], etc. Besides, previous works [42], [43], [44], [45] on unsupervised domain adaptation (UDA) share a similar idea, *i.e.*, utilizing statistics of the BN layer as domain-related knowledge. However, the UDA methods require the use of source or target data to conduct the adaptation in advance, thus becoming difficult to be applied during test time or online. Differently, we use the channel-wise affine transformation provided by the BatchNorm (BN) [11] layer to stabilize the test-time adaptation and aim for an effective fusion of RGB-D embeddings.

In general, the BN layer is widely used in deep learning to eliminate internal covariate shift and improve generalization. It performs a linear transformation followed by the convolutional or fully-connected layers. We denote by  $\mathbf{x}_{m,l,k}$  the  $k$ -th channel for the  $l$ -th layer feature maps of  $m$ -th modality (RGB or depth in our setting), then the transformation of the BN layer can be written as

$$\mathbf{x}'_{m,l,k} = \gamma_{m,l,k} \frac{\mathbf{x}_{m,l,k} - \mu_{m,l,k}}{\sqrt{\sigma_{m,l,k}^2 + \epsilon}} + \beta_{m,l,k}, \quad (7)$$

where the scaling and shift factors  $\gamma_{m,l,k}$  and  $\beta_{m,l,k}$  are adjustable affine parameters, the normalization statistics  $\mu_{m,l,k}$  and  $\sigma_{m,l,k}$  are updated with momentum.

Thus, by adapting the scaling and shift factors  $\gamma_{m,l,k}$ ,  $\beta_{m,l,k}$  and updating the statistics  $\mu_{m,l,k}$ ,  $\sigma_{m,l,k}$  of the BN layer, the proposed method does not introduce new parameters. Given the simulation-trained model, our method is independent of re-training on the large-scale synthetic data. Meanwhile, the fully test-time adaptation re-calibrates channel-wise responses of RGB and depth feature maps, thus providing better weighted fused RGB-D embeddings for unseen object instance segmentation. Finally, the overall objective for test-time adaptation can be formulated as

$$L_{total} = \lambda_1 L_{neo} + \lambda_2 L_{ckd}, \quad (8)$$

where the two terms  $L_{neo}$  and  $L_{ckd}$  are weighted by the balancing parameters  $\lambda_1$  and  $\lambda_2$ .

## IV. EXPERIMENTS

### A. Datasets and Evaluation Metrics

**Datasets** The model is pre-trained on the synthetic Tabletop Object Dataset (TOD) [1], which consists of 40k synthetic scenes of cluttered objects in tabletop environments. The proposed method is evaluated and compared on two real-world datasets, named OSD [12] and OCID [13]. OSD consists of 111 images in tabletop environments with averaged 3.3 objects per image. OCID has 2,346 images on both tabletop and floor with averaged 7.5 objects per image. It is worth noting that OSD has manually labeled segmentation masks while OCID generates semi-automatically labeled annotations, which are easily influenced by the noise of depth images.

**Evaluation Metrics** We follow previous works [1], [2], [4] to use the object precision/recall/F-measure (Overlap P/R/F) metrics to evaluate the object segmentation performance, and the Boundary P/R/F to evaluate how sharp the predicted boundary matches against the ground truth boundary. Besides, F@.75 is used to denote the percentage of segmented objects with Overlap F-measure  $\geq 75\%$  [2]. Though the IoU is a standard metric in segmentation tasks, it is highly correlated to the overlap F-measure, thus is not reported in the UOIS task [3]. All P/R/F and F@.75 measures are reported in the range of [0, 100].

### B. Comparison with State-of-the-art Methods

In this section, we first compare the proposed method with several state-of-the-art methods. As shown in Table I, our method outperforms all competitors on both OSD and OCID datasets. Specifically, FTEA+<sup>1</sup> achieves 89.5 and 92.3 Overlap F-measure, 73.8 and 86.7 Boundary F-measure, 88.3 and 91.1 F@.75, on OSD and OCID respectively, which validates the effectiveness of the proposed approach. Additionally, when we use the end-to-end model without extra zoom-in refinement, *i.e.*, FTEA in the next-to-last line

<sup>1</sup>FTEA+ denotes adopting the zoom-in refinement strategy in [2].

$L_{neo}$	$L_{ckd}$	OSD [12]							OCID [13]						
		Overlap			Boundary			F@.75	Overlap			Boundary			F@.75
		P	R	F	P	R	F		P	R	F	P	R	F	
		84.3	88.3	86.2	67.5	67.5	67.1	79.3	86.0	92.3	88.5	80.4	78.3	78.8	82.2
✓		85.0	91.9	88.2	67.6	74.6	70.4	89.8	85.7	<b>94.0</b>	89.3	78.0	78.9	78.0	84.9
	✓	85.3	89.2	87.1	<b>70.1</b>	70.2	69.8	81.9	<b>86.2</b>	92.8	88.9	<b>80.5</b>	78.9	<b>79.2</b>	83.1
✓	✓	<b>85.8</b>	<b>92.0</b>	<b>88.6</b>	69.2	<b>75.7</b>	<b>71.7</b>	<b>87.3</b>	<b>86.2</b>	93.9	<b>89.5</b>	79.5	<b>79.5</b>	79.1	<b>85.1</b>

TABLE II: Ablation studies of the proposed FTEA. The ablation of test-time adaptation (TTA) is implicitly included in the first row since it should be in conjunction with at least one objective function, *i.e.*,  $L_{neo}$  or  $L_{ckd}$ .

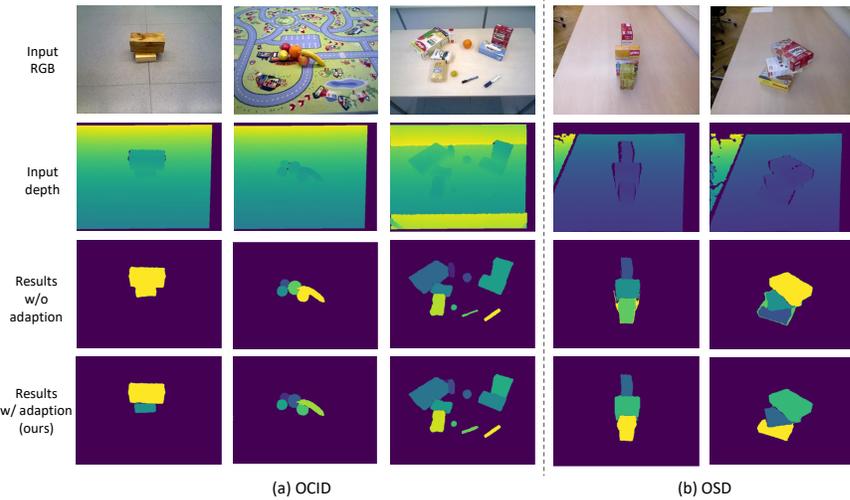
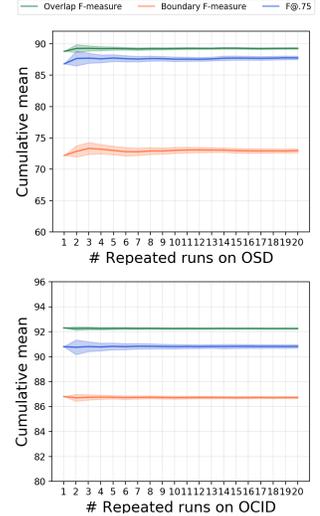


Fig. 2: Qualitative results of the proposed approach, (a) results on the OCID dataset, Fig. 3: 20 Repeated runs of different sampling orders.



in Table I, the improvement of our method is consistent. Compared to UCN [2], the proposed FTEA improves the Overlap F-measure by 2.4 and 1.0, the Boundary F-measure by 4.6 and 0.3, the  $F@.75$  by 5.1 and 1.8, on OSD and OCID respectively.

### C. Ablation Studies

**Non-parametric Entropy Objective** This objective is crucial for UOIS to conduct test-time adaptation. As shown in Table II, when equipped with the non-parametric entropy objective  $L_{neo}$ , the model’s performances on two real-world datasets are overall improved. Especially,  $L_{neo}$  significantly improves the overlap and boundary *recall*, which is desirable for discovering potential objects in unseen environments.

**Cross-modality Knowledge Distillation** The proposed cross-modality knowledge distillation provides another effective learning objective  $L_{ckd}$  for test-time feature enhancement. Table II shows that  $L_{ckd}$  works well on improving the overlap and boundary *precision* of segmentation results. Additionally, as shown in the last row in Table II, the cross-modality knowledge distillation  $L_{ckd}$  can be effectively combined with the former entropy objective  $L_{neo}$ , thus further improving the performance on most evaluation metrics for unseen object instance segmentation.

**Test-time Adaptation** Since the test-time adaptation (TTA) needs to be enabled with at least one proposed learning objective, *i.e.*, NEO or CKD, the ablation of TTA

is implicitly included in the first row in Table II. From Table II we can observe that whether TTA is in conjunction with the  $L_{neo}$  or  $L_{ckd}$ , the segmentation results on both OSD and OCID can be significantly improved.

### D. Discussions

**Adaptation Consumption** Besides the performance, we further investigate the computation consumption of the proposed method. The adaptation time only accounts for single backward pass of BN parameters, thus is way faster than the inference with multiple zoom-in operations [2]. Table III shows the averaged time consumption of the adaptation process and inference, *i.e.*, averaged 0.04s *v.s.* 1.24s per iteration. Based on the setting of only adapting 500 iterations as stated in Section IV-E, the total adaptation consumption for a new scenario is  $\sim 20s$  (0.04s per frame), which is efficient for practical application.

BatchSize	Inference time (avg)	Adaptation time (avg)
1	1.24±0.03s	0.04±0.00s

TABLE III: The time consumption of the inference and adaptation process, which is averaged over 500 iterations with 10 repeated runs.

**Qualitative Results** We visualize some qualitative results with and without the proposed FTEA in Figure 2. We can observe in Figure 2(a) that the proposed adaptation

Modalities		Overlap	Boundary	F@.75
RGB	Depth	F	F	
		87.4	69.4	83.2
✓		87.6	70.0	86.1
	✓	88.6	70.9	87.4
✓	✓	<b>89.5</b>	<b>73.8</b>	<b>88.3</b>

Parameters		Overlap	Boundary	F@.75
Conv	BN	F	F	
		87.4	69.4	83.2
✓		87.4	69.4	83.2
	✓	<b>89.5</b>	<b>73.8</b>	<b>88.3</b>
✓	✓	82.2	65.4	75.6

TABLE IV: Segmentation performances with different adaptation modalities and modulation parameters on OSD.

process mitigates the under-segmentation problem of two close objects. Besides, different from the smooth depth images in synthetic training data, realistic depth images are generally noisy, especially on the object’s boundary. This problem makes outputs of boundaries blurred and causes over-segmentation around the boundary, as illustrated in Figure 2(b). The last row in Figure 2(b) shows that, with the proposed adaptation process in FTEA, this problem can be largely alleviated.

**Robustness to the Sampling Order** Additionally, we train our models for multiple runs on different random sampling orders. Figure 3 illustrates the cumulative means of performance with 95% confidence intervals across 20 repeated runs on OCID and OSD. It can be observed that the overall performances of the model are stable, demonstrating the robustness of the proposed method. Note that the experiments on OCID also use random samples (not just the sampling order) since we conduct the adaptation with only 500 test-time images out of 2,346 images in OCID.

**Modalities for Adaptation** Table IV shows results with different adaptation modalities on OSD. First, segmentation performances can be improved whether RGB or depth modality is tuned. Second, by simultaneously tuning RGB and depth modalities, we can obtain better-fused RGB-D embeddings. Thus the performances become significantly better (see the last row in Table IV) due to the effective use of multimodal data.

**Modulation Parameters** We also conduct analysis on the modulation parameters, *i.e.*, the linear BN and the nonlinear convolution parameters. As shown in Table IV, using BN as modulation parameters achieves the best results. Since the nonlinear high dimensional convolution layers are difficult to optimize, tuning all model parameters could lead to a negative transfer.

**How Many BN Layers to Use** For brevity, we split all layers into four main building blocks as in ResNet [46]. Figure 4 illustrates the segmentation performances on OSD when we gradually increase the number of BN layers for adaptation. We can observe that the Overlap F-measure, Boundary F-measure, and F@.75 are consistently improved with more BN layers, demonstrating the effectiveness of the

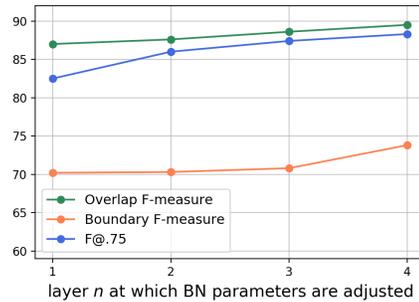


Fig. 4: Performances evolution on OSD when we gradually increase the number of BN layers for adaptation.

proposed method. On the other hand, the improvement is not plateaued, which implies that the adaptation process can be further enhanced by considering other effective parameters in addition to the BN layers of the model.

### E. Implementation Details

For fair comparisons, we follow previous work [2] to use a 34-layer, stride-8 ResNet (ResNet34-8s) as the backbone, and the full resolution  $640 \times 480$  feature map with embedding dimensions  $C = 64$  is obtained by bilinearly upsampling. To avoid noisy outliers, we set  $k = 1$  in Equation 3 (*i.e.*, select two clusters, the nearest cluster and the corresponding cluster) for the calculation of the non-parametric entropy objective  $L_{neo}$ . The temperature factor  $T$  in cross-modality knowledge distillation is 1. The weight factor for the overall loss is set as  $\lambda_1 = \lambda_2 = 1$ . During test time, our model is adapted with the SGD optimizer. We use batchsize=1 as in the typical inference phase. For the first 100 iterations (images), the learning rate is linearly warmed up to the base value  $lr = 0.005$ , then decayed with a cosine scheduler for another 400 iterations (images). We do not set the “epoch” number since there is no concept of “dataset” in the online test-time adaptation. We use the same learning schedule for the OSD and OCID. The data in test-time adaptation are not shuffled unless otherwise stated. All experiments are conducted on a single NVIDIA 2080Ti GPU with PyTorch.

## V. CONCLUSIONS

In this paper, we target the task of unseen object instance segmentation with an emphasis on the adaptation process for unseen realistic data. To mitigate the domain shift between the synthetic training and realistic testing data, a novel FTEA framework is proposed to conduct the fully test-time RGB-D embeddings adaptation. Specifically, during test time, we fix all convolutional layers and adjust the affine transformations provided by BN parameters via optimizing two novel unsupervised objectives, *i.e.*, the NEO and the CKD. NEO calculates the entropy of probability distributions of UOIS in a non-parametric way. CKD further encourages cross-modality knowledge transfer during test time. Extensive experiments on realistic RGB-D datasets OCID and OSD demonstrate the effectiveness of the proposed approach. We hope our work could draw attention to the test-time adaptation and reveal a promising direction for robot perception in unseen environments.

## VI. ACKNOWLEDGEMENT

We thank all anonymous reviewers for their constructive suggestions for improving our paper. This work was supported by the NSFC under Grant 62206288; in part by the National Key Research and Development Plan of China under Grant 2020AAA0108902; in part by the Strategic Priority Research Program of Chinese Academy of Science under Grant XDB32050100; in part by the Beijing Science and Technology Plan Project under Grant Z201100008320029; and in part by the Fujian Science and Technology Plan under Grant No. 2021T3003.

## REFERENCES

- [1] C. Xie, Y. Xiang, A. Mousavian, and D. Fox, "The best of both modes: Separately leveraging rgb and depth for unseen object instance segmentation," in *Conference on Robot Learning (CoRL)*, 2019.
- [2] Y. Xiang, C. Xie, A. Mousavian, and D. Fox, "Learning rgb-d feature embeddings for unseen object instance segmentation," in *Conference on Robot Learning (CoRL)*, 2020.
- [3] C. Xie, A. Mousavian, Y. Xiang, and D. Fox, "Rice: Refining instance masks in cluttered environments with graph neural networks," in *Conference on Robot Learning (CoRL)*, 2021.
- [4] C. Xie, Y. Xiang, A. Mousavian, and D. Fox, "Unseen object instance segmentation for robotic environments," *IEEE Transactions on Robotics*, 2021.
- [5] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009, pp. 248–255.
- [6] 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 *Proceedings of the European Conference on Computer Vision (ECCV)*, 2014, pp. 740–755.
- [7] C. E. Shannon, "A mathematical theory of communication," *The Bell System Technical Journal*, vol. 27, no. 3, pp. 379–423, 1948.
- [8] D. Wang, E. Shelhamer, S. Liu, B. Olshausen, and T. Darrell, "Tent: Fully test-time adaptation by entropy minimization," in *International Conference on Learning Representations (ICLR)*, 2020.
- [9] G. Hinton, O. Vinyals, J. Dean, *et al.*, "Distilling the knowledge in a neural network."
- [10] J. Gou, B. Yu, S. J. Maybank, and D. Tao, "Knowledge distillation: A survey," *International Journal of Computer Vision*, vol. 129, no. 6, pp. 1789–1819, 2021.
- [11] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *International Conference on Machine Learning (ICML)*, 2015, pp. 448–456.
- [12] A. Richtsfeld, T. Mörwald, J. Prankl, M. Zillich, and M. Vincze, "Segmentation of unknown objects in indoor environments," in *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2012, pp. 4791–4796.
- [13] M. Suchi, T. Patten, D. Fischinger, and M. Vincze, "Easylab: a semi-automatic pixel-wise object annotation tool for creating robotic rgb-d datasets," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 6678–6684.
- [14] S. Back, J. Lee, T. Kim, S. Noh, R. Kang, S. Bak, and K. Lee, "Unseen object amodal instance segmentation via hierarchical occlusion modeling," in *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022, pp. 5085–5092.
- [15] K. Tian, S. Zhou, and J. Guan, "Deepcluster: A general clustering framework based on deep learning," in *Joint European Conference on Machine Learning and Knowledge Discovery in Databases (ECML-KDD)*. Springer, 2017, pp. 809–825.
- [16] M. Caron, P. Bojanowski, A. Joulin, and M. Douze, "Deep clustering for unsupervised learning of visual features," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 132–149.
- [17] Y.-C. Hsu and Z. Kira, "Neural network-based clustering using pairwise constraints," *International Conference on Learning Representations Workshops (ICLRW)*, 2015.
- [18] Y.-C. Hsu, Z. Lv, and Z. Kira, "Learning to cluster in order to transfer across domains and tasks," in *International Conference on Learning Representations (ICLR)*, 2018.
- [19] Y.-C. Hsu, Z. Lv, J. Schlosser, P. Odom, and Z. Kira, "Multi-class classification without multi-class labels," in *International Conference on Learning Representations (ICLR)*, 2019.
- [20] K. Han, A. Vedaldi, and A. Zisserman, "Learning to discover novel visual categories via deep transfer clustering," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019, pp. 8401–8409.
- [21] K. Han, S.-A. Rebuffi, S. Ehrhardt, A. Vedaldi, and A. Zisserman, "Autonovel: Automatically discovering and learning novel visual categories," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [22] X. Ji, J. F. Henriques, and A. Vedaldi, "Invariant information clustering for unsupervised image classification and segmentation," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019, pp. 9865–9874.
- [23] Y. Ouali, C. Hudelot, and M. Tami, "Autoregressive unsupervised image segmentation," in *European Conference on Computer Vision (ECCV)*. Springer, 2020, pp. 142–158.
- [24] W. Van Gansbeke, S. Vandenhende, S. Georgoulis, and L. Van Gool, "Unsupervised semantic segmentation by contrasting object mask proposals," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2021, pp. 10052–10062.
- [25] Y. Sun, X. Wang, Z. Liu, J. Miller, A. Efros, and M. Hardt, "Test-time training with self-supervision for generalization under distribution shifts," in *International Conference on Machine Learning (ICML)*. PMLR, 2020, pp. 9229–9248.
- [26] Y. Liu, P. Kothari, B. van Delft, B. Bellot-Gurlet, T. Mordan, and A. Alahi, "Ttt++: When does self-supervised test-time training fail or thrive?" *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 34, 2021.
- [27] C. K. Mummadi, R. Huttmacher, K. Rambach, E. Levinkov, T. Brox, and J. H. Metzner, "Test-time adaptation to distribution shift by confidence maximization and input transformation," *arXiv preprint arXiv:2106.14999*, 2021.
- [28] D. Comaniciu and P. Meer, "Mean shift: A robust approach toward feature space analysis," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 24, no. 5, pp. 603–619, 2002.
- [29] T. Kobayashi and N. Otsu, "Von mises-fisher mean shift for clustering on a hypersphere," in *2010 20th International Conference on Pattern Recognition (ICPR)*. IEEE, 2010, pp. 2130–2133.
- [30] Z. Wu, Y. Xiong, S. X. Yu, and D. Lin, "Unsupervised feature learning via non-parametric instance discrimination," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 3733–3742.
- [31] C. Yang, Z. Wu, B. Zhou, and S. Lin, "Instance localization for self-supervised detection pretraining," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 3987–3996.
- [32] Y. Huang, C. Du, Z. Xue, X. Chen, H. Zhao, and L. Huang, "What makes multi-modal learning better than single (provably)," *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 34, 2021.
- [33] Z. Gao, J. Chung, M. Abdelrazek, S. Leung, W. K. Hau, Z. Xian, H. Zhang, and S. Li, "Privileged modality distillation for vessel border detection in intracoronary imaging," *IEEE transactions on medical imaging*, vol. 39, no. 5, pp. 1524–1534, 2019.
- [34] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 2961–2969.
- [35] M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," in *Proceedings of the European conference on Computer Vision (ECCV)*, 2014, pp. 818–833.
- [36] S. Woo, J. Park, J.-Y. Lee, and I. S. Kweon, "Cbam: Convolutional block attention module," in *European Conference on Computer Vision (ECCV)*, 2018, pp. 3–19.
- [37] Y. He, X. Zhang, and J. Sun, "Channel pruning for accelerating very deep neural networks," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 1389–1397.
- [38] Y. Wang, W. Huang, F. Sun, T. Xu, Y. Rong, and J. Huang, "Deep multimodal fusion by channel exchanging," *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 33, pp. 4835–4845, 2020.
- [39] L. Zhang, Z. Liu, S. Zhang, X. Yang, H. Qiao, K. Huang, and A. Hus-

- sain, “Cross-modality interactive attention network for multispectral pedestrian detection,” *Information Fusion*, vol. 50, pp. 20–29, 2019.
- [40] J. Hu, L. Shen, and G. Sun, “Squeeze-and-excitation networks,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 7132–7141.
- [41] W. Shao, S. Tang, X. Pan, P. Tan, X. Wang, and P. Luo, “Channel equilibrium networks for learning deep representation,” in *International Conference on Machine Learning (ICML)*. PMLR, 2020, pp. 8645–8654.
- [42] Y. Li, N. Wang, J. Shi, X. Hou, and J. Liu, “Adaptive batch normalization for practical domain adaptation,” *Pattern Recognition*, vol. 80, pp. 109–117, 2018.
- [43] W.-G. Chang, T. You, S. Seo, S. Kwak, and B. Han, “Domain-specific batch normalization for unsupervised domain adaptation,” in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 7354–7362.
- [44] Z. Zhuang, L. Wei, L. Xie, T. Zhang, H. Zhang, H. Wu, H. Ai, and Q. Tian, “Rethinking the distribution gap of person re-identification with camera-based batch normalization,” in *European Conference on Computer Vision (ECCV)*. Springer, 2020, pp. 140–157.
- [45] M. Klingner, J.-A. Termöhlen, J. Ritterbach, and T. Fingscheidt, “Unsupervised batchnorm adaptation (ubna): A domain adaptation method for semantic segmentation without using source domain representations,” in *Proceedings of the IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2022, pp. 210–220.
- [46] 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 (CVPR)*, 2016, pp. 770–778.