# *v*-CLR: View-Consistent Learning for Open-World Instance Segmentation

Chang-Bin Zhang[1]     Jinhong Ni[1]     Yujie Zhong[2]     Kai Han[1*]
[1]Visual AI Lab, The University of Hong Kong     [2]Meituan Inc.
{cbzhang, jhni}@connect.hku.hk     jaszhong@hotmail.com     kaihanx@hku.hk

## Abstract

*In this paper, we address the challenging problem of open-world instance segmentation. Existing works have shown that vanilla visual networks are biased toward learning appearance information, e.g., texture, to recognize objects. This implicit bias causes the model to fail in detecting novel objects with unseen textures in the open-world setting. To address this challenge, we propose a learning framework, called view-Consistent LeaRning (v-CLR), which aims to enforce the model to learn appearance-invariant representations for robust instance segmentation. In v-CLR, we first introduce additional views for each image, where the texture undergoes significant alterations while preserving the image's underlying structure. We then encourage the model to learn the appearance-invariant representation by enforcing the consistency between object features across different views, for which we obtain class-agnostic object proposals using off-the-shelf unsupervised models that possess strong object-awareness. These proposals enable cross-view object feature matching, greatly reducing the appearance dependency while enhancing the object-awareness. We thoroughly evaluate our method on public benchmarks under both cross-class and cross-dataset settings, achieving state-of-the-art performance. Project page:* https://visual-ai.github.io/vclr

## 1. Introduction

Modern object detectors [7, 22, 73, 78] and instance segmentors [14, 26, 44] have achieved many milestones. However, these detectors are based on the assumption of pre-defined taxonomy classes. Despite recent open-vocabulary detectors [23, 37] can be extended to larger taxonomy classes benefiting from the foundation model pre-trained on large-scale text-image pairs, these models are still limited by the finite taxonomy classes in the pre-trained data. In some realistic applications, models are required to identify out-of-taxonomy classes. Thus, recognizing objects in the open
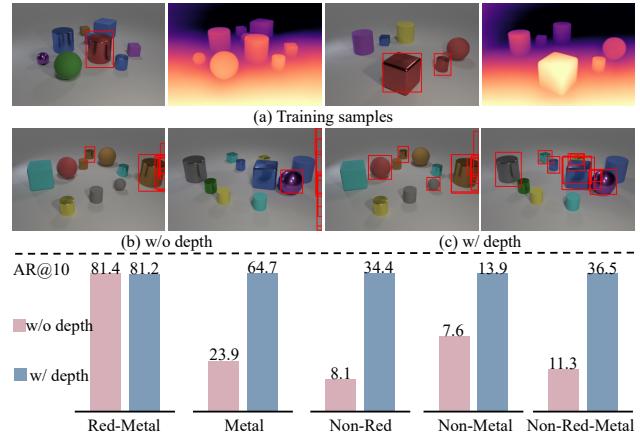


(a) Training samples

(b) w/o depth     (c) w/ depth

Figure 1. **Toy example on the CLEVR [33] dataset.** The model regards *red-metal* objects as the known class and is evaluated on different subsets in terms of AR@10. We train the model with and without incorporating depth image data, respectively. The prediction results are displayed in the middle row.

world has been increasingly interesting and challenging.

In open-world instance segmentation, models are trained on a set of predefined known classes and are evaluated to localize unknown objects during inference. Following [34, 36, 54, 68], we regard the open-world instance segmentor as a class-agnostic object discovery model. One straightforward solution is to train a class-agnostic detector on labeled instances of known classes, *i.e.,* performing binary object detection given ground truth labels from the known classes, and hope the models capture transferable features that generalize to unknown objects. However, various studies [1, 3, 18, 20] have demonstrated that neural networks exhibit a preference to capture texture information when recognizing objects. This hinders the model's ability to generalize in the open-world setting, especially to unknown objects with unseen textures.

To motivate the necessity of capitalizing appearance-invariant information, we showcase a toy open-world example run on the CLEVR [33] dataset in Fig. 1. In this example, we treat the *red metal* objects as the known class, and evaluate the model on detecting various other types of objects (with other colors or materials). We show an ex-

---

1

ample of training samples in Fig. 1 (a), where each sample consists of a natural image and a colorized depth map. We label the known class, *i.e.,* the red metal objects, with red bounding boxes. We then train a vanilla detector as the baseline model using only natural images as input, and a model incorporating colorized depth images. The evaluation results on various object subsets involving different colors and materials in Fig. 1 demonstrate that the model trained with depth images exhibits a much better generalization to novel objects. This toy example verifies the problem that the vanilla baseline models suffer from poor generalization due to the appearance bias, and emphasizes the importance of including appearance-invariant information to guide representation learning.

To overcome this challenge, we propose a *view-Consistent LeaRning* framework, dubbed *v*-CLR, to encourage the model to learn appearance-invariant representation for novel object discovery. To achieve this, we first transform images into multiple appearance-invariant views, from which we propose a feature-matching objective to enforce cross-view feature consistency. This objective alone would be insufficient as there is no guarantee that similar features correspond to objects, we thus adopt off-the-shelf general object proposals to ensure optimized representations are object-oriented. Specifically, we first exploit the appearance-invariant information by transforming the natural images into various other domains, *e.g.,* colorized depth images. Intuitively speaking, these transformations destroy or overwrite the appearance information from the natural image domain while preserving the original structures, thus encouraging the model to capitalize information other than appearance.

To facilitate appearance-invariant representation learning and effectively utilize training data containing multiple views, we build on top of DETR-like architectures [7, 74, 78], in which we enforce representation consistency across different views of the same image by matching similar queries. By doing this, we naturally circumvent the problem of implicit appearance bias by empowering the model to capture consistent cross-view information. However, naively enforcing such consistency may still fail in reality. The reason is that even if the model extracts similar representations across different views, it does not necessarily imply these representations are object-related. To sidestep this problem, we adopt pre-trained unsupervised instance detectors, *e.g.,* CutLER [64], to generate object proposals. These off-the-shelf instance detectors exhibit high instance awareness, for which we explicitly match the queries from different views with the object proposals to ensure these paired queries are object-oriented. To this end, we have devised a learning framework to allow models to capture object-related consistent appearance-invariant representations, enabling transferability to novel objects in open-world scenarios.

We conduct extensive experiments on various bench-

marks, including COCO 2017 [45], LVIS [24], UVO [60], and Objects365 [56], under cross-categories and cross-datasets settings. Our proposed learning framework consistently achieves state-of-the-art performance on several benchmarks in the open-world setting.

## 2. Related Work

**Object Detection and Instance Segmentation.** DETR [7] and its follow-up works [17, 43, 46, 50, 74, 75, 78] achieve an end-to-end detector with remarkable performance, improving transformer architecture [17, 46, 50, 78], training efficiency [10, 29, 32, 43, 73, 74] and label assignment [6, 47, 58]. MaskDINO [44] develop a unified model for object detection and instance segmentation. Benefiting from the powerful self-supervised learning [8, 28, 52, 62], unsupervised instance segmentation [63, 65] has received increasing interest by discovering pixel-level pseudo annotations automatically. Thanks to strong object-awareness from self-supervised pretrained models [52], CutLER [64] constructs a large-scale training set with pseudo masks, *e.g.,* ImageNet dataset, and train an instance segmentation model without any human annotation. In our work, we utilize the CutLER pre-trained on the ImageNet as a general objects proposal network.

**Open-world Instance Segmentation.** To promote the applications of modern object detectors in realistic scenarios, recent arts [36, 61] propose open-world instance segmentation. To avoid suppressing potential unknown objects in background regions, OLN [36] replaced the classification branch in Mask-RCNN [26] with a localization-aware score. LDET [54] proposed to synthesize training images by combining labeled objects and predefined background texture by copy-paste [21]. Segprompt [77] utilizes prompting designation to segment novel objects. Some other methods [30, 34, 61] design variant mechanisms to discover potential unknown objects in training images, including grouping pixels [61], leveraging prior mask [34] by MCG [53] and imposing geometry information [30]. SWORD [68] explores applying DETR-based model [78] on the open-world instance segmentation. SOS [66] propose to discover potential unlabeled objects by SAM [40] with DINOv2 [52] activation point as prompt. In our work, we conduct experiments based on the Deformable-DETR [78] and DINO-DETR [74], respectively.

**Texture-Invariant Representations.** Within the domain of generalization and adaptation, models are designed to utilize source domain training data to achieve effective performance on a different target domain, assuming that both domains share the same set of semantic categories. To successfully adapt to the target domain, which may exhibit different styles from the source, current approaches [15, 31, 38, 39, 42, 67, 71] incorporate style transfer techniques to modify training images to either the target

or an arbitrary style. Our approach emphasizes the transfer of knowledge across different semantic classes instead of across different domains. The challenges between domain shifts and semantic shifts are orthogonal [59] and the techniques for domain shifts are not suitable for semantic shifts. To this end, our method can leverage any transformation views with unified features among objects, extending beyond mere style transfer. For instance, depth images and edge maps can also be utilized to develop unified object representations, which are typically unsuitable for domain generalization due to their lack of class discriminability. In our approach, depth images serve as the primary transformation view, while stylized images and edge maps can be optionally used as auxiliary views.

## 3. Method

**Problem Statement.** Open-world instance segmentation aims to localize as many novel objects as possible during test time. Formally, the training labels are first divided into two sets of known classes ($\mathcal{C}_{base}$) and unknown classes ($\mathcal{C}_{novel}$), with no overlap between them (*i.e.,* $\mathcal{C}_{base} \cap \mathcal{C}_{novel} = \emptyset$). For each training sample image $I$ and its associated set of annotations $C$, we train the models only on the annotations of known classes, in a class-agnostic manner . During test time, we evaluate the model's capability of generalizing on the set of unknown classes ($\mathcal{C}_{novel}$).

### 3.1. Method Overview

**Architecture.** Inspired by instance segmentation models with transformer [14, 44], we decorate the Deformable-DETR [78] and DINO-DETR [74] into the instance segmentation model. Specifically, following [14, 44], each query predicts a prototype for a corresponding instance, and then the model will predict the instance segmentation map by computing the similarity between the output prototype and the pyramid features of the transformer encoder.

**Appearance-Invariant Transformation.** To enable such an appearance-invariant representation learning, we first leverage off-the-shelf image transformation to *overwrite the appearance from the natural images while leaving the overall structural contents intact*. The intuition is that we circumvent the texture bias [1, 3, 18, 20] by allowing the model to learn consistent and transferable representations from different image transformations. We adopt colorized depth maps [2] as the major transformation in this work, and with an additional auxiliary transformation *e.g.,* art-stylizing [69] and edge map [70], while we highlight that our method is not strictly bound by any transformation method so long as they suffice the aforementioned criteria. Complementing the two transformations with the natural images gives us three views, *i.e.,* natural images, colorized depth maps, and one additional auxiliary view, for each training sample, from which we randomly select one view per sample with equal

probability during training. To further destroy the appearance of objects, we apply random cropping and resizing to an image patch, subsequently integrating it with the original image. These various views play a crucial role in our method as described in the following section.

### 3.2. Appearance-Invariant Representation

Existing works have shown evidence that neural networks are biased toward learning appearance information, *e.g.,* texture, to differentiate different objects [1, 3, 18, 20]. This tendency of relying on appearance information inhibits the generalization ability to novel classes especially when unseen textures are presented during inference. To overcome this challenge, we devise a learning framework so that the model learns appearance-invariant representations complementing the appearance information and, thus are generalizable and unbiased during inference. Our proposed method is detailed below. Roughly speaking, the key to this learning framework is to enforce representation consistency by maximizing the query feature similarity between the transformed views and the natural image.

Our learning framework comprises two branches: the natural image branch, which always receives natural images as inputs; and the transformed image branch, which randomly processes any of the transformed images or the original natural image with equal probability. Both branches then utilize the adapted DETR transformer architectures [7, 74, 78] to make sets of predictions, where each prediction consists of a classification score, a predicted bounding box, and a predicted segmentation mask. We refer the readers to the *Model Architecture* paragraph in Sec. 3.1 for details regarding how we adopt detection transformers for instance segmentation. Following existing self-supervised learning frameworks [11–13, 27], to prevent feature collapsing, we update the transformer in the natural image branch as an exponential moving average (EMA) model of the transformed image branch.

**Object-centric Learning by Object Proposals.** At first glance, it seems to be feasible at this stage to ensure representation consistency on the query features outputted from the two branches. However, a high similarity between the matched queries does not necessarily imply the model learning informative representation. An example is when models capture shortcut solutions where the extracted representations are irrelevant to the objects. In the context of open-world learning, a lack of correlation with the objects can cause failure in generalization. Thanks to the high instance awareness of the large-scale pre-trained instance detectors [64], we sidestep the problem of the model falling into object-irrelevant solutions by leveraging these pre-trained detectors to provide object proposals. These object proposals serve as a medium to match object-related queries from both branches, thus ensuring our learning framework can learn meaningful object-oriented representation to
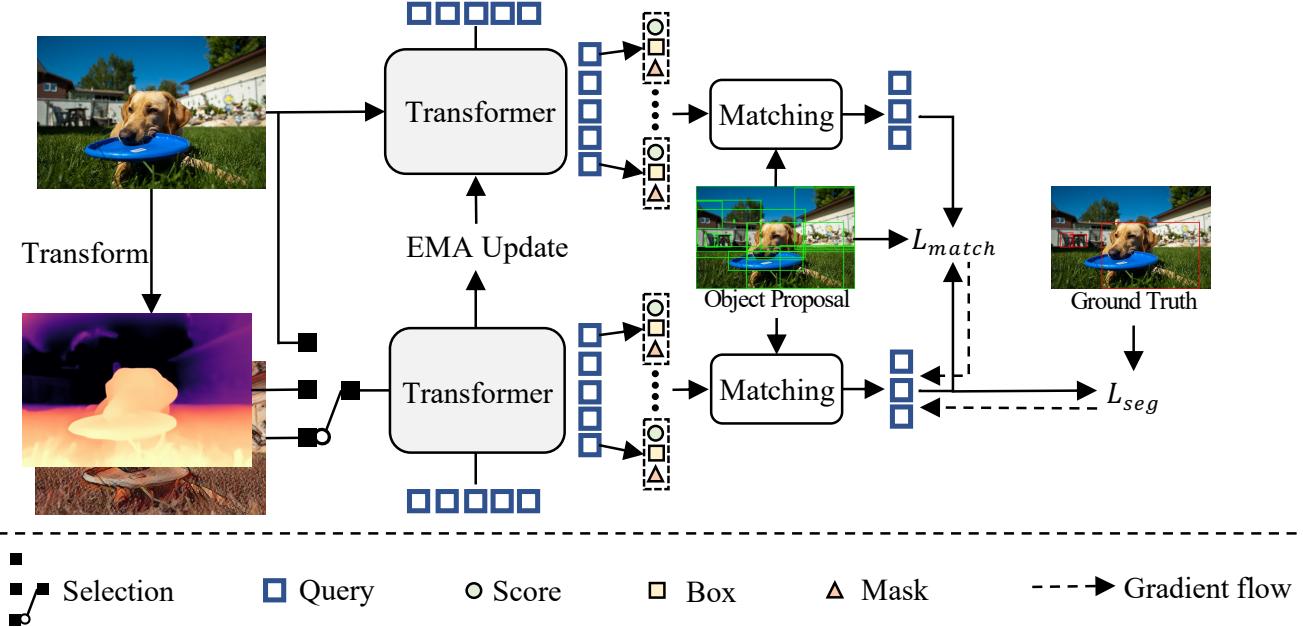
Figure 2. **Illustration of ν-CLR.** Our learning framework consists of two branches, the natural image branch (top) and the transformed image branch (bottom). Both branches adopt transformers to make predictions, which are then matched with the object proposals to obtain optimized object queries. We compute a matching loss $L_{match}$ which enforces the matched object-oriented query pairs from the two branches to be similar. We finally compute the ordinary segmentation loss $L_{gt}$ using the ground truth labels. The transformer in the natural image branch is updated as an EMA model of the transformed image branch.
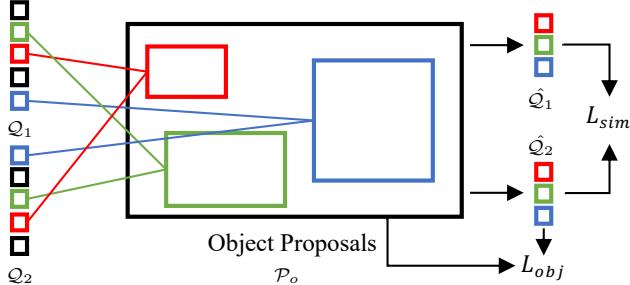


Figure 3. **Illustration of object feature matching in ν-CLR.** Let $\mathcal{Q}_1$ and $\mathcal{Q}_2$ represent the query outputs from the EMA teacher model and the student model, respectively. Predictions associated with object proposals demonstrating poor localization quality are removed, resulting in paired $\hat{\mathcal{Q}}_1$ and $\hat{\mathcal{Q}}_2$, and the objective $L_{sim}$ is utilized to maximize feature similarity between each pair. Concurrently, the student model is trained using these object proposals.

be successfully transferred to open-world settings.

**View-Consistent Learning.** Given the multiple transformed views of an image, we hope a model can learn to extract consistent characteristics shared across different views of the same image. To facilitate such training, we propose *view-Consistent LeaRning*. An overview of our method is illustrated in Fig. 2.

**Object Feature Matching.** We introduce the object feature matching in our view-consistent learning pipeline in detail. The overall illustration of the matching objective is shown in Fig. 3. Formally, denote the sets of predictions

from two branches as $\mathcal{P}_1$ and $\mathcal{P}_2$, and the set of extracted object proposals as $\mathcal{P}_o$, where each set $\mathcal{P} = \{(\hat{p}_i, \hat{b}_i, \hat{m}_i)\}$ consists of tuples of class score $\hat{p}_i$, bounding box $\hat{b}_i$, and segmentation mask $\hat{m}_i$, for $i = 1, \ldots, |\mathcal{P}|$. We also have the sets of queries $\mathcal{Q}_1$ and $\mathcal{Q}_2$ associated with the prediction sets, where we have $|\mathcal{Q}_i| = |\mathcal{P}_i|$ for $i = 1, 2$. Following the previous works [7, 41], for each proposal in $\mathcal{P}_o$, we find the optimal sets $\hat{\mathcal{P}}_1$ and $\hat{\mathcal{P}}_2$ for the two sets of predictions by minimizing the matching cost. The sets $\mathcal{P}_o$, $\hat{\mathcal{P}}_1$, and $\hat{\mathcal{P}}_2$ forms $\tilde{N}$ one-to-one triplets.

**Training Objectives.** We denote the optimal sets of queries as $\hat{\mathcal{Q}}_1$ and $\hat{\mathcal{Q}}_2$ corresponding to the sets of predictions $\hat{\mathcal{P}}_1$ and $\hat{\mathcal{P}}_2$, for which we compute the cosine similarity matching loss:

$$L_{sim} = \frac{1}{\tilde{N}} \sum_{q_i \in \hat{\mathcal{Q}}_i} \left(1 - \cos\left(q_1, q_2\right)\right),$$

where $\cos(q_1, q_2)$ denotes the cosine similarity between $q_1$ and $q_2$. Since we assume the object proposals to be reliably object-related, this may give us additional information for supervising the predicted boxes and segmentation maps. We thus compute the standard segmentation loss using the object proposals $L_{obj}$:

$$L_{obj} = \lambda_1 L_{dice} + \lambda_2 L_{mask} + \\ \lambda_3 L_{score} + \lambda_4 L_{box} + \lambda_5 L_{giou}, \quad (1)$$

where $\lambda_i$ from now on denotes the loss weight factor. The

4

total matching objective is computed as:

$$L_{match} = \lambda_{obj}L_{obj} + \lambda_{sim}L_{sim}.$$

The matching objective ensures the queries capture object-oriented appearance-invariant representations. We proceed to the regular segmentation loss using the ground truth labels. Formally, given the set of optimized transformed image queries $\hat{\mathcal{Q}}_2$ and the set of ground truth $\mathcal{G}$, we compute similar segmentation objective $L_{gt}$ as Eqn. (1) by replacing the object proposals $\mathcal{P}_o$ with $\mathcal{G}$. The total training objective is then:

$$L = \lambda_{match}L_{match} + \lambda_{gt}L_{gt}.$$

## 4. Experiments

### 4.1. Setup

**Datasets and Evaluations.** We conduct experiments in two popular open-world settings, cross-categories and cross-datasets, on the CLEVR [33], COCO 2017 [45], LVIS [24], UVO [60] and Objects365 [56] datasets. The prior setting divides the object classes into known and unknown classes, whereas the latter setting tests the generalization ability of the model on another dataset containing unseen object classes. Since the labels in validation images can not cover all objects, we apply the average recall (AR) over multiple IoU thresholds $[0.5, 0.95]$ to measure the model's performance, while ignoring the average precision (AP) as previous arts [36, 61, 68]. Following [34, 36, 68], *the most widely concerned metric in this task is AR@100*, which is denoted by $AR_{100}$ in our paper. As standard evaluation metrics on COCO, we use $AR^b$ and $AR^m$ to denote the results for predicted boxes and instance masks, respectively. We additionally report the performance for small, medium, and large objects, denoted by $AR_{s/m/l}$ respectively.

**Implementation Details.** We regard the model as a class-agnostic object detector in all experiments. We apply the DINO-DETR [74] with ResNet-50 [25] as the backbone to perform instance segmentation. We adopt the common settings in DETR-like models [7, 73, 74], *e.g.,* there are six layers in the transformer encoder and decoder, respectively. We set the number of denoising queries [43] as 300. Inspired by [14, 44, 68], we decorate the DINO-DETR with dynamic convolution for instance segmentation prediction. Following [68], we use 1500 and 1000 queries in the transformer decoder when training on VOC and COCO classes, respectively. We train the model for 8 epochs and the learning rate is decayed at the 7th epoch, while keeping other settings in the training schedule as fully-supervised object detectors. In our experiments, $\lambda_{sim}$, $\lambda_{obj}$ and $\lambda_{gt}$ is set to 1, and coefficients in Eqn. (1) are the same as DINO [74]. We use the pre-trained Cascade-Mask-RCNN [4] as the object proposal network without any fine-tuning, which is trained by CutLER [64] with ResNet-50 as the backbone.

| Method | $AR^b_{10}$ | $AR^b_{100}$ | $AR^m_{10}$ | $AR^m_{100}$ |
|---|---|---|---|---|
| Mask-RCNN [26] | 10.2 | 23.5 | 7.9 | 17.7 |
| CutLER [64] | 19.9 | 34.5 | - | - |
| OLN [36] | 18.0 | 33.5 | 16.9 | - |
| LDET [54] | 18.2 | 30.8 | 16.3 | 27.4 |
| GGN [61] | 17.3 | 31.6 | 16.1 | 28.7 |
| GGN + OLN [36] | 17.1 | 37.2 | 16.4 | 33.7 |
| UDOS [34] | - | 33.5 | - | 31.6 |
| GOOD$^\dagger$ [30] | - | 39.3 | - | - |
| Def-DETR [78] | 12.2 | 27.4 | 10.2 | 22.7 |
| SWORD [68] | 17.8 | 35.3 | 15.7 | 30.2 |
| *v*-CLR (Def-DETR) | 22.2 | 40.3 | 19.6 | 33.7 |
| DINO-DETR [74] | 13.2 | 31.1 | 9.7 | 22.0 |
| *v*-CLR (DINO) | **22.5** | **40.9** | **19.9** | **34.1** |

Table 1. **Evaluation results for novel classes in the VOC → Non-VOC setting.** The † denotes the model is trained with bounding boxes only.

| Method | $AR^b_{10}$ | $AR^b_{100}$ | $AR^m_{10}$ | $AR^m_{100}$ |
|---|---|---|---|---|
| Mask-RCNN [26] | 11.4 | 16.2 | 7.6 | 11.4 |
| LDET [54] | 16.0 | 31.9 | 12.3 | 25.2 |
| Def-DETR [78] | 13.5 | 33.5 | 9.5 | 25.3 |
| SWORD [68] | 16.8 | 43.1 | 13.3 | 34.9 |
| *v*-CLR (Def-DETR) | 20.3 | 45.8 | 16.1 | 34.6 |
| DINO-DETR [74] | 14.7 | 36.5 | 10.7 | 27.7 |
| *v*-CLR (DINO) | **21.0** | **47.2** | **16.8** | **35.9** |

Table 2. **Evaluation results for novel classes in the VOC→UVO setting.**

### 4.2. Main Results

To validate the effectiveness of our method, we conduct experiments in popular settings, including VOC → Non-VOC, COCO → LVIS, VOC → UVO, and COCO → Objects365, where $\mathcal{D}_A \rightarrow \mathcal{D}_B$ denotes training the model on dataset $\mathcal{D}_A$ and evaluating the transferability on the dataset $\mathcal{D}_B$.

**VOC → Non-VOC.** The VOC [16] dataset includes 20 common classes in natural images, for which we train the model on VOC classes to verify the generalization capability of our method. Specifically, the model is trained on the COCO 2017 training set with 20 VOC class labels, and tested on the other 60 Non-VOC classes on the COCO validation set. Following recent arts [36, 68], we also regard the prediction as a class-agnostic scheme, thus the most concerned evaluation metric is average recall (AR), especially AR@100. As shown in Tab. 1, we report the AR@10 and AR@100 on the Non-VOC classes, respectively. SOWRD [68] firstly explore adapting DETR-based detector to discover novel objects, and propose some techniques based on popular Deformable-DETR [78], including stop-gradient, IoU-based branch, and one-to-many assignment. However, we empirically find

| Method | $AR_{10}^b$ | $AR_{100}^b$ | $AR_{10}^m$ | $AR_{100}^m$ |
|---|---|---|---|---|
| Mask-RCNN [26] | 6.1 | 19.4 | 5.6 | 17.2 |
| GGN [61] | 7.6 | 22.4 | 7.2 | 20.4 |
| Def-DETR [78] | 6.3 | 19.4 | 5.5 | 16.4 |
| SWORD [68] | 8.8 | 23.5 | **8.0** | 20.4 |
| *v*-CLR (Def-DETR) | **9.4** | 27.2 | **8.0** | 22.3 |
| DINO-DETR [74] | 8.5 | 25.2 | 7.4 | 21.0 |
| *v*-CLR (DINO) | 9.3 | **28.4** | 7.9 | **23.6** |

Table 3. **Evaluation results for novel classes in the COCO →
LVIS setting.**

| Method | $AR_{10}^b$ | $AR_{100}^b$ | $AR_s^b$ | $AR_m^b$ | $AR_l^b$ |
|---|---|---|---|---|---|
| Mask-RCNN [26] | 19.3 | 32.8 | 18.2 | 36.4 | 43.5 |
| LDET [54] | 20.0 | 36.8 | 20.7 | 40.5 | 48.9 |
| Def-DETR [78] | 19.0 | 40.1 | 22.8 | 43.4 | 54.1 |
| SWORD [68] | **22.8** | 43.9 | 25.0 | 48.6 | 57.6 |
| *v*-CLR (Def-DETR) | 19.4 | 45.9 | 23.8 | 49.3 | 62.8 |
| DINO-DETR [74] | 19.0 | 46.4 | **28.8** | 50.0 | 58.6 |
| *v*-CLR (DINO) | 19.7 | **47.9** | 26.2 | **51.6** | **64.0** |

Table 4. **Evaluation results for novel classes in the COCO →
Objects365 setting.**

that vanilla DINO-DETR can achieve surprisingly strong performance with the help of denoising queries to accelerate training. Therefore, we conduct experiments based on Deformable-DETR [78] and DINO-DETR [74] for a fair comparison, respectively. Experimental results demonstrate that our method achieves state-of-the-art performance on all evaluation metrics in this setting.

**VOC → UVO.** The UVO dataset [60] is a large-scale dataset designed for open-world segmentation, covering many kinds of objects in the wild. To validate the cross-dataset generalization, we follow previous work to conduct experiments on the UVO dataset [54, 68]. Specifically, the model is trained on the 20 VOC classes of COCO 2017 training set, and is evaluated on the UVO dense v1.0 validation set. This split provides the category names of each instance, which allows us to split the novel classes and evaluate our model. We report the experimental results in Tab. 2. Compared with the previous state-of-the-art method, our method achieve a remarkable improvement of 2.7% in terms of $AR_{100}^b$ based on Deformable-DETR [78]. We argue that the baseline model tends to suffer from the bias on the limited appearance of known classes. Benefiting from learning appearance-invariant information, our method improves more than 10% both on $AR_{100}^b$ and $AR_{100}^m$.

**COCO → LVIS.** The LVIS dataset [24] enlarges the taxonomy of COCO, containing more than 1200 classes where a large number of classes are disjoint with COCO classes. In this setting, to verify the generalization ability on larger known taxonomy, the model is trained on 80 classes of COCO 2017 training set, and evaluated on other disjoint classes in LVIS validation set. As shown in Tab. 3, although vanilla DINO-DETR reaches a better performance than SWORD [68], our method can additionally improve the baseline by 3.2% in terms of $AR_{100}^b$. Our proposed method outperforms SWORD by about 3.7% and 1.9% in terms of $AR_{100}^b$ and $AR_{100}^m$. We argue this improvement arises from our proposed training framework, which encourages the model to learn appearance-invariant cues to discover potential objects.

**COCO → Objects365.** The Objects365 dataset [56] includes 365 common classes which is much larger than

COCO taxonomy. As shown in Tab. 4, the model is trained on COCO 80 classes and evaluated on the novel classes of Objects365. Since this dataset does not provide the instance mask annotation, we only evaluate the performance of bounding box prediction. The experimental results demonstrate that our method can outperform SWORD [68] by 2% in terms of $AR_{100}^b$. We also report the performance of different methods on the small, medium, and large objects, respectively. We observe that our method performs slightly worse on small objects than vanilla DINO-DETR and explore the potential reasons in Sec. 4.4.

### 4.3. Qualitative Results

We visualize prediction results of our method on the COCO 2017 validation set in Fig. 4. The model is trained on 20 VOC classes of the COCO 2017 training set. For each image, we show the *top*-10 predicted instances according to the prediction confidence.

### 4.4. Ablation Study

**Ablation study of components.** To validate the effectiveness of each component in our method, we conduct ablation studies in the VOC→Non-VOC setting, as shown in Tab. 5. Initially, incorporating general object proposals results in a 6% improvement over the vanilla DINO-DETR. Leveraging the colorized depth and auxiliary views introduced in our method, the detector achieves 40.0% in terms of $AR_{100}^b$, marking a 2.3% improvement over the strong baseline with $L_{obj}$ only. Based on this, our consistent constraint training objective yields an additional 0.2% improvement, raising $AR_{100}^b$ to 40.2%. To further enforce instance consistency, we filter paired object queries from the two branches before computing $L_{sim}$. This filtering results in a 0.7% improvement in $AR_{100}^b$, culminating in our final model with an $AR_{100}^b$ of 40.9%. Without CutLER [64] proposals, our method reaches 30.7% $AR_{100}^m$, achieving 8.6% improvement over the baseline model. Notably, general object proposals may be less effective when few or no unknown objects appear in the training images. In Tab. 5, following [30, 34], all experiments, except for the baseline model (first row), also utilize the unlabeled images in the training set in the VOC→Non-VOC

6

Figure 4. **Qualitative results of our method on COCO 2017 validation set.** The model is trained on 20 VOC classes. We show the top-10 predicted instances according to the prediction confidence.

| $L_{gt}$ | $L_{obj}$ | Transform. | $L_{sim}$ | filtering | $AR^b_{100}$ | $AR^m_{100}$ |
|---|---|---|---|---|---|---|
| ✔ | | | | | 31.1 | 22.0 |
| ✔ | ✔ | | | | 37.7 (+6.6) | 31.2 (+9.2) |
| ✔ | ✔ | ✔ | | | 40.0 (+8.9) | 33.2 (+11.2) |
| ✔ | ✔ | ✔ | ✔ | | 40.2 (+9.1) | 33.9 (+11.9) |
| ✔ | ✔ | ✔ | | ✔ | 35.9 (+4.8) | 30.7 (+8.6) |
| ✔ | ✔ | ✔ | ✔ | ✔ | 40.9 (+9.8) | **34.1** (+12.1) |

Table 5. **Ablation study of each component in our method.**

| Natural | Depth | Stylized | Edge | $AR^b_{100}$ | $AR^m_{100}$ |
|---|---|---|---|---|---|
| ✔ | | | | 38.5 | 32.0 |
| ✔ | ✔ | | | 40.5 | 33.3 |
| ✔ | | ✔ | | 40.2 | 33.5 |
| ✔ | ✔ | | ✔ | 40.5 | 33.7 |
| ✔ | ✔ | ✔ | | **40.9** | **34.1** |

Table 6. **Ablation study of different views used in our method.**

| $AR^b_{10}$ / $AR^b_{100}$ | Natural | Depth | Stylized |
|---|---|---|---|
| CutLER [64] | 19.9 / 34.5 | 10.3 / 17.5 | 11.6 / 22.4 |
| *v*-CLR (**ours**) | **22.5 / 40.9** | **18.8 / 35.7** | **21.0 / 35.2** |

Table 7. **Evaluation results on three different views in the VOC→Non-VOC setting.** Our method only uses natural images during inference, but it is also capable of processing multiple views.

both depth maps and stylized images perform similarly. By including an additional auxiliary view on top of the depth view, we observe a consistent improvement while adding stylized images perform slightly better than the edge map.

**Comparison with CutLER.** We leverage CutLER [64], which possesses a satisfactory object-identifying ability, in our work to generate object proposals. We compare the performance of CutLER as a detector versus our method on the novel Non-VOC classes in Tab. 7. While the performance margin is already 6.4% in $AR^b_{100}$ between CutLER and our method on natural images, it is noticeable that the performance of CutLER degrades rapidly on these transformed images, evidenced by around 15% performance gap on the two transformed views. These results demonstrate that CutLER may suffer from potential textual bias, thus emphasizing the strength of learning appearance-invariant representation.

**Application to vision transformers.** According to [51], vision transformers exhibit less texture bias compared to CNNs. We thus additionally investigate the applicability of our method to vision transformers. We present experimental results utilizing the Swin-Tiny backbone in Tab. 8.

setting. When CutLER object proposals are not applied, we use the trained baseline model to provide annotations for these unlabeled images to ensure a fair comparison.

**Image transformation.** We leverage colorized depth views with the help of additional auxiliary views to enforce the model to learn appearance-invariant representation. To study the impact of the transformed views, we apply the off-the-shelf model to generate different transformed views on the COCO 2017 validation set. We then study the impact of different views used and report the results in Tab. 6. The model is trained on VOC classes and evaluated on Non-VOC classes. When only one view is considered, we find that

|  | DINO-DETR [74] | Ours | w/o Transform. |
|---|---|---|---|
| $AR^b_{100}$ | 32.6 | **40.7** | 39.5 (-1.2) |
| $AR^m_{100}$ | 26.9 | **33.8** | 32.4 (-1.4) |

Table 8. **Experiments in the VOC→Non-VOC setting based on the Swin-Tiny [49] backbone.**

| Count | 1∼3 | 4∼6 | 7∼9 | ≥ 10 |
|---|---|---|---|---|
| N | 58.8 / 48.7 | 42.4 / 35.8 | 33.5 / 28.1 | 21.6 / 17.5 |
| N + D | **60.9** / 50.3 | 44.4 / 37.1 | 35.4 / 29.5 | 23.3 / 18.7 |
| N + S | 60.4 / 50.3 | 44.3 / 37.4 | 34.7 / 29.8 | 23.3 / 19.1 |
| N + D + S | 60.6 / **50.5** | **44.7 / 37.6** | **35.9 / 30.3** | **24.0 / 19.9** |

(a) Performance on scenarios with different number of instances

| Size | Small | Medium | Large | All |
|---|---|---|---|---|
| N | 16.6 / 12.3 | 45.3 / 38.1 | 73.9 / 63.5 | 38.5 / 32.0 |
| N + D | 17.4 / 12.5 | 49.1 / 41.1 | 75.3 / 64.8 | 40.5 / 33.3 |
| N + S | 17.1 / 12.7 | 48.8 / 41.1 | 75.2 / 65.0 | 40.2 / 33.5 |
| N + D + S | **17.6 / 13.1** | **49.6 / 42.2** | **75.5 / 65.5** | **40.9 / 34.1** |

(b) Performance on scenarios with different object sizes

Table 9. **Ablation study of view choices on different segmentation scenarios in the VOC→Non-VOC setting.** 'N', 'S', and 'D' denote natural images, stylized images, and depth maps, respectively. We report $AR^b_{100}$ / $AR^m_{100}$ in the table.

Our approach significantly surpasses DINO-DETR [74], emphasizing the necessity of transformed views for enhanced performance with transformer-based architectures. These experimental results indicate that our method is applicable and can enhance the performance of vision transformer backbones.

**View choices and segmentation scenarios.** We investigate the impact of incorporating various views on segmentation performance across different scenarios, including variations in object sizes and the number of instances. Specifically, as shown in Tab. 9a, we analyze the effect of different combinations of views and evaluate the model under scenarios with varying instance counts. The experimental results indicate that both depth maps and stylized images consistently improve performance across scenarios with diverse numbers of instances. Additionally, we assess the model's performance across different object sizes, as outlined in Tab. 9b, where objects are categorized into small, medium, and large based on the standard COCO dataset [45]. Our findings reveal that incorporating additional views significantly enhances performance on medium and large objects, while the improvements on small objects are relatively modest.

**Detailed performance on unknown and known classes.** To study the effect of our method on known and unknown classes, we train the model on a cross-dataset setting, VOC→UVO, and evaluate the model on known and unknown classes, respectively. As shown in Tab. 10, our method achieves performance comparable to the baseline model on known classes, while significantly improving recall

| Method | Known | | Unknown | | All | |
|---|---|---|---|---|---|---|
|  | $AR^b_{100}$ | $AR^m_{100}$ | $AR^b_{100}$ | $AR^m_{100}$ | $AR^b_{100}$ | $AR^m_{100}$ |
| DINO-DETR | 59.3 | **48.3** | 36.5 | 27.7 | 42.3 | 33.2 |
| v-CLR (**ours**) | **60.9** | 47.0 | **47.2** | **35.9** | **50.3** | **38.4** |

Table 10. **Evaluation results on known and unknown classes in the VOC→UVO setting.**

| Ratio | Small | Medium | Large |
|---|---|---|---|
| Ground-truth of Known Classes | 31.1% | 34.9% | 34.0% |
| + Proposals | 19.9% | 28.5% | 51.6% |

Table 11. **Ratio of small, medium and large objects in the supervision.** The ratio is measured under COCO→Objects365.

on unknown objects by 10.7% and across all classes by 8% in terms of $AR^b_{100}$. These results highlight the effectiveness of our method in discovering novel objects.

**Performance on small objects.** As shown in Tab. 4, our method exhibits unstable performance on small objects. Specifically, our method achieves an approximate 1% improvement on small objects when applied to Deformable-DETR [78], but leads to performance degradation when applied to DINO-DETR [74]. We attribute this instability to an imbalance in the ratio of objects with different sizes. In Tab. 11, we measure the size distribution of objects and observe that the ratio of small objects decreases significantly when incorporating proposals. This imbalance arises due to the CutLER [64] network's inherent preference for large objects, stemming from its pretraining process.

## 5. Conclusion

To encourage the model to utilize appearance-invariant cues to discover objects, we propose a learning framework, named view-Consistent LeaRning (v-CLR), for segmenting instances in an open world. Specifically, our method randomly picks one from natural images, depth images, and an auxiliary view as input during training. In this way, the model will tend to learn common features between the three views, which is beneficial for novel object discovery. Besides, to help the model learn appearance-invariant features, we design a consistent objective based on the general object proposals. The superiority of our approach is thoroughly validated with extensive experiments on cross-category and cross-dataset settings and consistently achieving state-of-the-art performance.

# References

[1] Pedro Ballester and Ricardo Araujo. On the performance of googlenet and alexnet applied to sketches. In *AAAI*, 2016. 1, 3

[2] Shariq Farooq Bhat, Reiner Birkl, Diana Wofk, Peter Wonka, and Matthias Müller. Zoedepth: Zero-shot transfer by combining relative and metric depth. *arXiv preprint arXiv:2302.12288*, 2023. 3, 12

[3] Wieland Brendel and Matthias Bethge. Approximating cnns with bag-of-local-features models works surprisingly well on imagenet. *arXiv preprint arXiv:1904.00760*, 2019. 1, 3

[4] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: Delving into high quality object detection. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2018. 5

[5] Zhaowei Cai and Nuno Vasconcelos. Cascade r-cnn: High quality object detection and instance segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 2019. 12

[6] Zhi Cai, Songtao Liu, Guodong Wang, Zheng Ge, Xiangyu Zhang, and Di Huang. Align-detr: Improving detr with simple iou-aware bce loss. *arXiv preprint arXiv:2304.07527*, 2023. 2

[7] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *Eur. Conf. Comput. Vis.*, 2020. 1, 2, 3, 4, 5

[8] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Int. Conf. Comput. Vis.*, 2021. 2

[9] Pratik Chaudhari, Anna Choromanska, Stefano Soatto, Yann LeCun, Carlo Baldassi, Christian Borgs, Jennifer Chayes, Levent Sagun, and Riccardo Zecchina. Entropy-sgd: Biasing gradient descent into wide valleys. *Journal of Statistical Mechanics: Theory and Experiment*, 2019. 12

[10] Qiang Chen, Xiaokang Chen, Gang Zeng, and Jingdong Wang. Group detr: Fast training convergence with decoupled one-to-many label assignment. *arXiv preprint arXiv:2207.13085*, 2022. 2

[11] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021. 3

[12] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020.

[13] Ze-Sen Chen, Gengshi Huang, Wei Li, Jianing Teng, Kun Wang, Jing Shao, Chen Change Loy, and Lu Sheng. Siamese detr. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2023. 3, 14, 15

[14] Bowen Cheng, Alex Schwing, and Alexander Kirillov. Per-pixel classification is not all you need for semantic segmentation. In *Adv. Neural Inform. Process. Syst.*, 2021. 1, 3, 5

[15] Sungha Choi, Sanghun Jung, Huiwon Yun, Joanne Taery Kim, Seungryong Kim, and Jaegul Choo. Robustnet: Improving domain generalization in urban-scene segmentation via instance selective whitening. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021. 2

[16] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *Int. J. Comput. Vis.*, 2010. 5

[17] Ziteng Gao, Limin Wang, Bing Han, and Sheng Guo. Adamixer: A fast-converging query-based object detector. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2022. 2

[18] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. Texture and art with deep neural networks. *Current opinion in neurobiology*, 2017. 1, 3

[19] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *The international journal of robotics research*, 2013. 12

[20] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A Wichmann, and Wieland Brendel. Imagenet-trained cnns are biased towards texture; increasing shape bias improves accuracy and robustness. *arXiv preprint arXiv:1811.12231*, 2018. 1, 3

[21] Golnaz Ghiasi, Yin Cui, Aravind Srinivas, Rui Qian, Tsung-Yi Lin, Ekin D Cubuk, Quoc V Le, and Barret Zoph. Simple copy-paste is a strong data augmentation method for instance segmentation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021. 2

[22] Ross Girshick. Fast r-cnn. In *Int. Conf. Comput. Vis.*, 2015. 1

[23] Xiuye Gu, Tsung-Yi Lin, Weicheng Kuo, and Yin Cui. Open-vocabulary object detection via vision and language knowledge distillation. *arXiv preprint arXiv:2104.13921*, 2021. 1

[24] Agrim Gupta, Piotr Dollar, and Ross Girshick. Lvis: A dataset for large vocabulary instance segmentation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019. 2, 5, 6

[25] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2016. 5

[26] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Int. Conf. Comput. Vis.*, 2017. 1, 2, 5, 6

[27] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020. 3

[28] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2022. 2

[29] Zhengdong Hu, Yifan Sun, Jingdong Wang, and Yi Yang. Dac-detr: Divide the attention layers and conquer. *Adv. Neural Inform. Process. Syst.*, 2024. 2

[30] Haiwen Huang, Andreas Geiger, and Dan Zhang. Good: Exploring geometric cues for detecting objects in an open world. In *Int. Conf. Learn. Represent.*, 2023. 2, 5, 6

[31] Wei Huang, Chang Wen Chen, Yong Li, Jiacheng Li, Cheng Li, Fenglong Song, Youliang Yan, and Zhiwei Xiong. Style projected clustering for domain generalized semantic segmentation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2023. 2

[32] Ding Jia, Yuhui Yuan, Haodi He, Xiaopei Wu, Haojun Yu, Weihong Lin, Lei Sun, Chao Zhang, and Han Hu. Detrs with

hybrid matching. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2023. 2

[33] Justin Johnson, Bharath Hariharan, Laurens Van Der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2017. 1, 5, 12

[34] Tarun Kalluri, Weiyao Wang, Heng Wang, Manmohan Chandraker, Lorenzo Torresani, and Du Tran. Open-world instance segmentation: Top-down learning with bottom-up supervision. *arXiv preprint arXiv:2303.05503*, 2023. 1, 2, 5, 6

[35] Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping Tak Peter Tang. On large-batch training for deep learning: Generalization gap and sharp minima. *arXiv preprint arXiv:1609.04836*, 2016. 12

[36] Dahun Kim, Tsung-Yi Lin, Anelia Angelova, In So Kweon, and Weicheng Kuo. Learning open-world object proposals without learning to classify. *IEEE Robotics and Automation Letters*, 2022. 1, 2, 5

[37] Dahun Kim, Anelia Angelova, and Weicheng Kuo. Region-aware pretraining for open-vocabulary object detection with vision transformers. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2023. 1

[38] Myeongjin Kim and Hyeran Byun. Learning texture invariant representation for domain adaptation of semantic segmentation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020. 2

[39] Sunghwan Kim, Dae-Hwan Kim, and Hoseong Kim. Texture learning domain randomization for domain generalized segmentation. In *Int. Conf. Comput. Vis.*, 2023. 2

[40] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross B. Girshick. Segment anything. In *Int. Conf. Comput. Vis.*, 2023. 2

[41] Harold W Kuhn. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 1955. 4

[42] Suhyeon Lee, Hongje Seong, Seongwon Lee, and Euntai Kim. Wildnet: Learning domain generalized semantic segmentation from the wild. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2022. 2

[43] Feng Li, Hao Zhang, Shilong Liu, Jian Guo, Lionel M Ni, and Lei Zhang. Dn-detr: Accelerate detr training by introducing query denoising. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2022. 2, 5

[44] Feng Li, Hao Zhang, Huaizhe Xu, Shilong Liu, Lei Zhang, Lionel M Ni, and Heung-Yeung Shum. Mask dino: Towards a unified transformer-based framework for object detection and segmentation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2023. 1, 2, 3, 5

[45] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Eur. Conf. Comput. Vis.*, 2014. 2, 5, 8, 14

[46] Shilong Liu, Feng Li, Hao Zhang, Xiao Yang, Xianbiao Qi, Hang Su, Jun Zhu, and Lei Zhang. Dab-detr: Dynamic anchor boxes are better queries for detr. *arXiv preprint arXiv:2201.12329*, 2022. 2

[47] Shilong Liu, Tianhe Ren, Jiayu Chen, Zhaoyang Zeng, Hao Zhang, Feng Li, Hongyang Li, Jun Huang, Hang Su, Jun Zhu, et al. Detection transformer with stable matching. *arXiv preprint arXiv:2304.04742*, 2023. 2

[48] Yun Liu, Ming-Ming Cheng, Xiaowei Hu, Kai Wang, and Xiang Bai. Richer convolutional features for edge detection. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2017. 12

[49] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Int. Conf. Comput. Vis.*, 2021. 8

[50] Depu Meng, Xiaokang Chen, Zejia Fan, Gang Zeng, Houqiang Li, Yuhui Yuan, Lei Sun, and Jingdong Wang. Conditional detr for fast training convergence. In *Int. Conf. Comput. Vis.*, 2021. 2

[51] Muzammal Naseer, Kanchana Ranasinghe, Salman Hameed Khan, Munawar Hayat, Fahad Shahbaz Khan, and Ming-Hsuan Yang. Intriguing properties of vision transformers. In *Adv. Neural Inform. Process. Syst.*, 2021. 7

[52] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 2, 12

[53] Jordi Pont-Tuset, Pablo Arbelaez, Jonathan T Barron, Ferran Marques, and Jitendra Malik. Multiscale combinatorial grouping for image segmentation and object proposal generation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 2016. 2

[54] Kuniaki Saito, Ping Hu, Trevor Darrell, and Kate Saenko. Learning to detect every thing in an open world. In *Eur. Conf. Comput. Vis.*, 2022. 1, 2, 5, 6

[55] Babak Saleh and Ahmed Elgammal. Large-scale classification of fine-art paintings: Learning the right metric on the right feature. *arXiv preprint arXiv:1505.00855*, 2015. 12

[56] Shuai Shao, Zeming Li, Tianyuan Zhang, Chao Peng, Gang Yu, Xiangyu Zhang, Jing Li, and Jian Sun. Objects365: A large-scale, high-quality dataset for object detection. In *Int. Conf. Comput. Vis.*, 2019. 2, 5, 6

[57] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and support inference from rgbd images. In *Eur. Conf. Comput. Vis.*, 2012. 12

[58] Yao Teng, Haisong Liu, Sheng Guo, and Limin Wang. Stageinteractor: Query-based object detector with cross-stage interaction. *arXiv preprint arXiv:2304.04978*, 2023. 2

[59] Hongjun Wang, Sagar Vaze, and Kai Han. Dissecting out-of-distribution detection and open-set recognition: A critical analysis of methods and benchmarks. *Int. J. Comput. Vis.*, 2024. 3

[60] Weiyao Wang, Matt Feiszli, Heng Wang, and Du Tran. Unidentified video objects: A benchmark for dense, open-world segmentation. In *Int. Conf. Comput. Vis.*, 2021. 2, 5, 6

[61] Weiyao Wang, Matt Feiszli, Heng Wang, Jitendra Malik, and Du Tran. Open-world instance segmentation: Exploiting pseudo ground truth from learned pairwise affinity. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2022. 2, 5, 6

[62] Xinlong Wang, Rufeng Zhang, Chunhua Shen, Tao Kong, and Lei Li. Dense contrastive learning for self-supervised visual pre-training. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2021. 2

[63] Xinlong Wang, Zhiding Yu, Shalini De Mello, Jan Kautz, Anima Anandkumar, Chunhua Shen, and Jose M Alvarez. Freesolo: Learning to segment objects without annotations. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2022. 2, 12

[64] Xudong Wang, Rohit Girdhar, Stella X Yu, and Ishan Misra. Cut and learn for unsupervised object detection and instance segmentation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2023. 2, 3, 5, 6, 7, 8, 12

[65] Yangtao Wang, Xi Shen, Yuan Yuan, Yuming Du, Maomao Li, Shell Xu Hu, James L Crowley, and Dominique Vaufreydaz. Tokencut: Segmenting objects in images and videos with self-supervised transformer and normalized cut. *IEEE Trans. Pattern Anal. Mach. Intell.*, 2023. 2, 12

[66] Christian Wilms, Tim Rolff, Maris Hillemann, Robert Johanson, and Simone Frintrop. Sos: Segment object system for open-world instance segmentation with object priors. In *Eur. Conf. Comput. Vis.*, 2024. 2

[67] Aming Wu and Cheng Deng. Single-domain generalized object detection in urban scene via cyclic-disentangled self-distillation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2022. 2

[68] Jiannan Wu, Yi Jiang, Bin Yan, Huchuan Lu, Zehuan Yuan, and Ping Luo. Exploring transformers for open-world instance segmentation. In *Int. Conf. Comput. Vis.*, 2023. 1, 2, 5, 6

[69] Xiaolei Wu, Zhihao Hu, Lu Sheng, and Dong Xu. Styleformer: Real-time arbitrary style transfer via parametric style composition. In *Int. Conf. Comput. Vis.*, 2021. 3, 12

[70] Saining Xie and Zhuowen Tu. Holistically-nested edge detection. *Int. J. Comput. Vis.*, 2015. 3

[71] Yanchao Yang and Stefano Soatto. Fda: Fourier domain adaptation for semantic segmentation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020. 2

[72] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 2021. 12

[73] Chang-Bin Zhang, Yujie Zhong, and Kai Han. Mr. detr: Instructive multi-route training for detection transformers. *arXiv preprint arXiv:2412.10028*, 2024. 1, 2, 5

[74] Hao Zhang, Feng Li, Shilong Liu, Lei Zhang, Hang Su, Jun Zhu, Lionel M Ni, and Heung-Yeung Shum. Dino: Detr with improved denoising anchor boxes for end-to-end object detection. *arXiv preprint arXiv:2203.03605*, 2022. 2, 3, 5, 6, 8, 12, 14

[75] Qiang Zhang, Zhang Zhang, Wei Cui, Jingkai Sun, Jiahang Cao, Yijie Guo, Gang Han, Wen Zhao, Jiaxu Wang, Chenghao Sun, et al. Humanoidpano: Hybrid spherical panoramic-lidar cross-modal perception for humanoid robots. *arXiv preprint arXiv:2503.09010*, 2025. 2

[76] Ying Zhang, Tao Xiang, Timothy M Hospedales, and Huchuan Lu. Deep mutual learning. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2018. 12

[77] Muzhi Zhu, Hengtao Li, Hao Chen, Chengxiang Fan, Weian Mao, Chenchen Jing, Yifan Liu, and Chunhua Shen. Segprompt: Boosting open-world segmentation via category-level prompt learning. In *Int. Conf. Comput. Vis.*, 2023. 2

[78] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020. 1, 2, 3, 5, 6, 8, 15

# Appendix

## A. Experimental Details

### A.1. Auxiliary Views

Our learning framework leverages multiple transformed views of the original natural image. Specifically, we apply off-the-shelf models to transform natural images into art-stylized and colorized depth images. For the art-stylized transformation, we utilize the pre-trained StyleFormer [69] model, which is trained on the WikiArt [55] dataset. For each natural image, we randomly select a target style from the WikiArt [55] dataset. For the colorized depth transformation, we employ the off-the-shelf ZoeDepth [2] model, pre-trained on the NYU Depth v2 [57] and KITTI [19] datasets. Additionally, for edge maps used in our ablation study, we apply the off-the-shelf RCF [48] model for edge detection. Examples of natural, art-stylized, and colorized depth images are shown in Fig. 5. Notably, no human annotations are used for generating depth maps or stylized images, ensuring that our method avoids any information leakage.

### A.2. Experiments on the CLEVR Dataset

The CLEVR dataset [33] is a synthetic dataset featuring objects characterized by four attributes:
- Size: large, small
- Shape: cube, sphere, cylinder
- Color: gray, red, blue, green, brown, purple, cyan, yellow
- Material: rubber, metal

In this work, we focus on two attributes—color and material—for illustrative simplicity. Specifically, we designate *red metal* objects as the known class, while objects with any other attribute combination are treated as unknown classes. The CLEVR dataset [33] comprises 70,000 training images and 15,000 validation images. We apply vanilla DINO-DETR [74] and train the model under two settings: with and without colorized depth images and stylized images. When using colorized depth images, the model randomly selects either a natural image, a depth map or a stylized image as input, each with equal probability. With 300 denoising queries in DINO-DETR [74], we train the model for 2,000 iterations with a batch size of 8, while retaining the remaining training configurations identical to those of vanilla DINO-DETR [74].

### A.3. Object Proposal Generation

Thanks to large-scale self-supervised learning, neural networks have shown remarkable capabilities in object recognition and localization [52, 65]. Leveraging this advancement, unsupervised instance segmentation [63, 64] has recently achieved significant progress. By benefiting from unsupervised training, these methods exhibit strong instance awareness, making them well-suited for generating object proposals in our work. Throughout this paper, we employ the ImageNet-pretrained Cascade R-CNN [5] from CutLER [64] to infer object proposals from the dataset. For each training image, we apply Non-Maximum Suppression (NMS) with a threshold of 0.7 and select the top-10 proposals based on prediction confidence.

---

**Algorithm 1** Pseudo-code of Parameter Perturbation in a PyTorch-like style.

```
# image: input image tensors
# model: the detector
# noise_std: the standard deviation of
    gaussian noise
def perturbation_forward(image, model,
    noise_std):
  # adding gaussian noise for each
      parameter
  for name, param in model.
      named_parameters():
    param += torch.randn_like(param) *
        noise_std
  output = model(image)
  return output
```

---

## B. Robustness against Parameter Perturbation

Numerous studies [9, 35, 72, 76] have demonstrated that neural networks trained with flatten minima exhibit superior generalization ability, *i.e.,* , the minima of the model should be in wide valleys rather than narrow crevices [9, 35, 72, 76]. In such cases, small perturbations to model parameters should not significantly degrade the performance of a model with strong generalization ability. Consequently, we can assess a model's generalization by introducing random perturbations to its parameters. Specifically, as detailed in Alg. 1, we inject Gaussian noise with varying standard deviations into all network parameters and evaluate the resulting performance. All models are trained on VOC classes and evaluated on Non-VOC classes. As illustrated in Fig. 6, we define the noise rate as the standard deviation of the Gaussian noise. With increasing noise rates, both our model and DINO-DETR experience performance degradation. However, at high noise rates, our method consistently outperforms the baseline by a substantial margin, demonstrating greater robustness to parameter perturbations.

## C. Robustness against Image Distortion

To verify the effectiveness of our method under different input perturbations, we evaluate our model under four popular distortions, Contrast, Gaussian Noise, Snow, and Frost. As shown in Fig. 7, we generate validation images with these distortions and evaluate the model's performance on them. We examine the robustness of our *v*-CLR approach against different types of image distortions. In Fig. 8, we plot the distribution of prediction scores for both the baseline DINO-DETR [74] and our method, with and without image
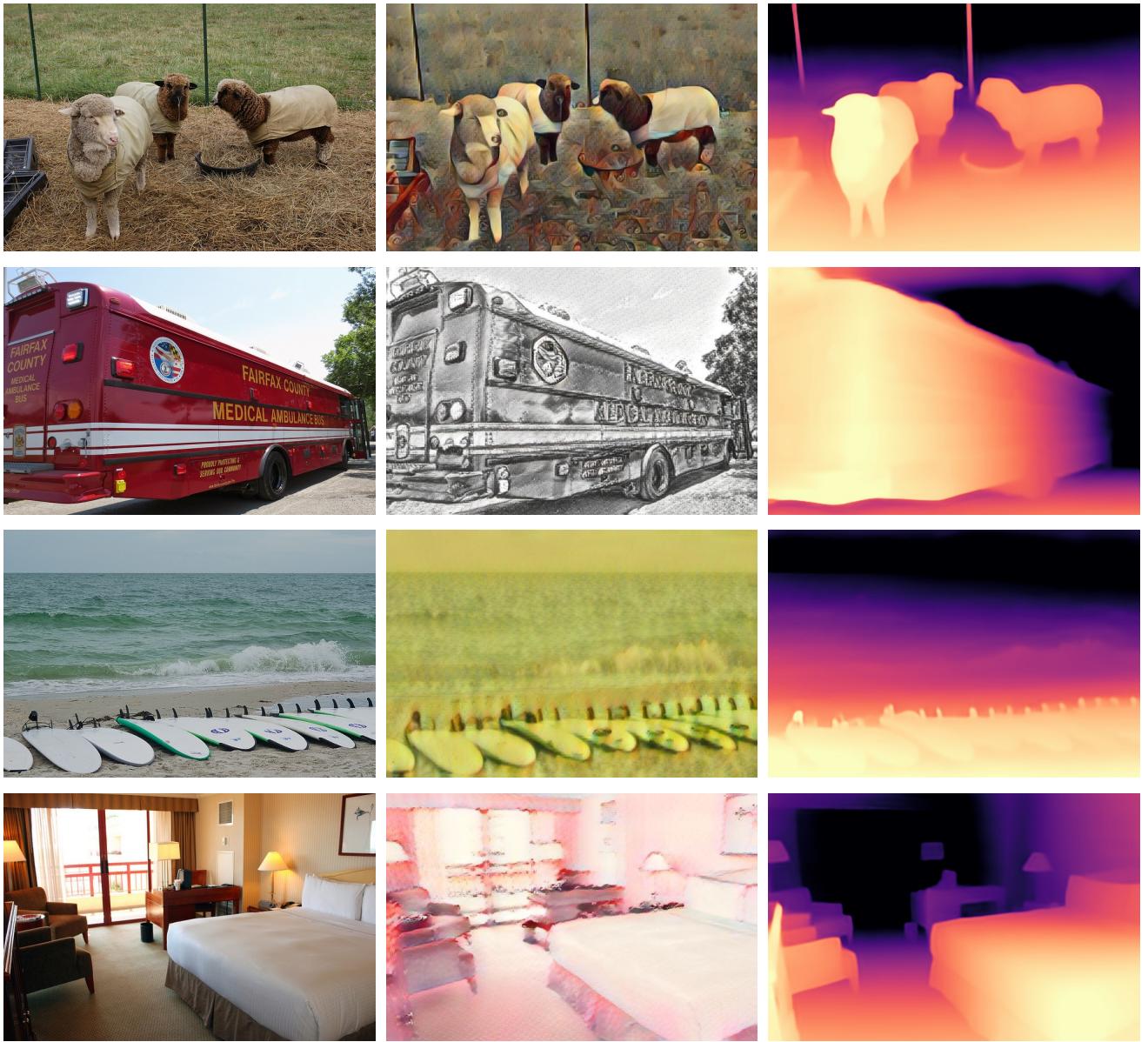
Figure 5. **Visualization of three views** used in our method, natural, art-stylized, and colorized depth images, respectively.
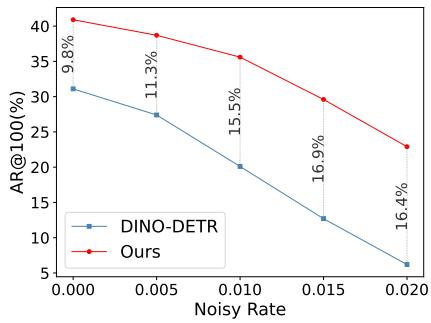


Figure 6. **AR$_{100}^{b}$ under different noisy rates.** All models are evaluated in the VOC→Non-VOC setting.

distortions. Our method (purple) consistently yields higher prediction scores than the baseline (green) on undistorted images. For distorted images, the distribution of the distorted baseline (red) exhibits a heavier right tail compared to the undistorted baseline (green), indicating that distortions reduce DINO-DETR's prediction confidence. In contrast, our method demonstrates greater robustness to image distortions, as the distributions of prediction scores for distorted and undistorted images show similar right-tail behavior. Surprisingly, the prediction score distribution for our method on distorted images exhibits even lower variance and a slightly higher mean than on undistorted images. This further suggests that image distortions have minimal impact on our
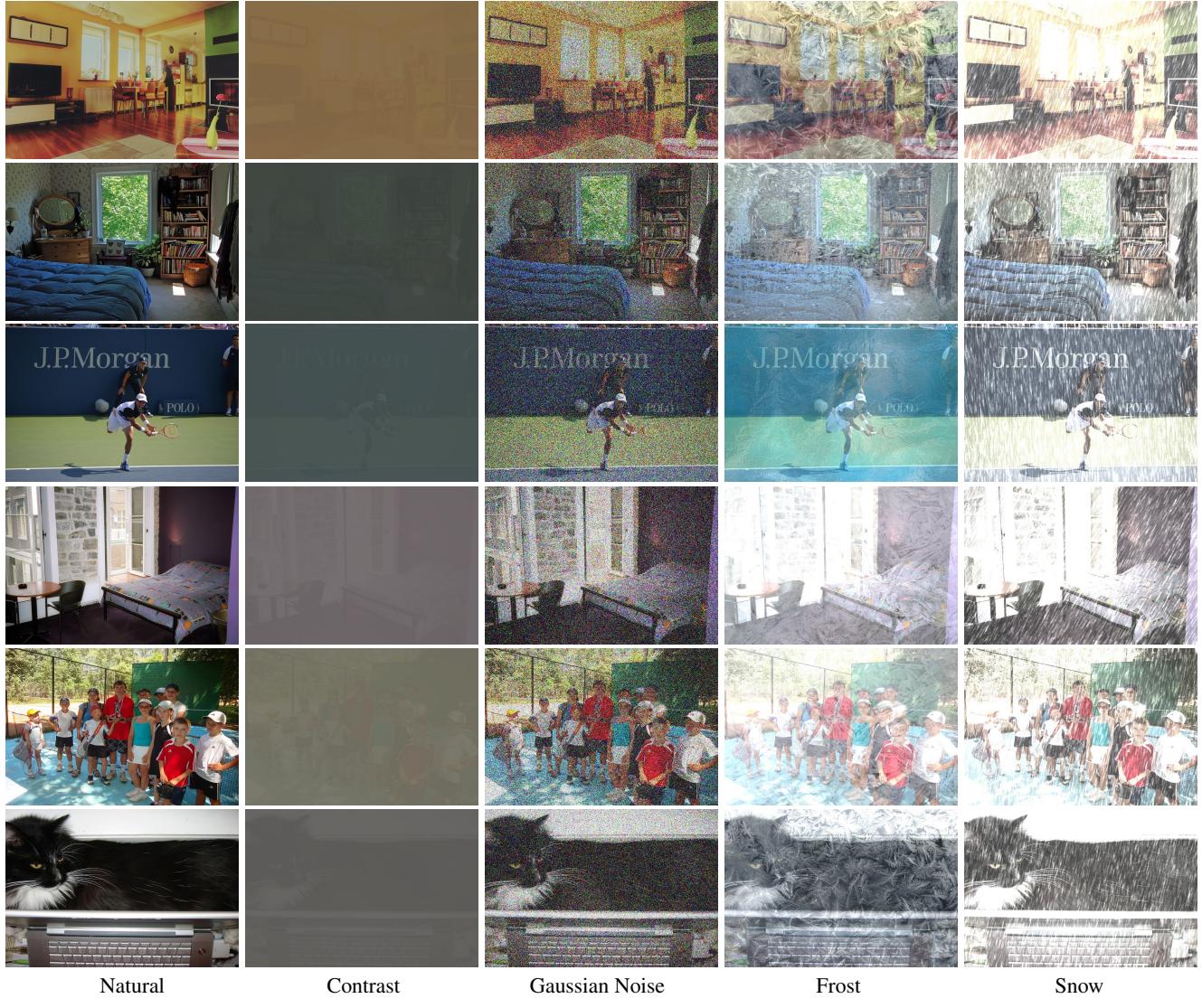
| Natural | Contrast | Gaussian Noise | Frost | Snow |

Figure 7. **Examples of distorted images on COCO 2017 [45] validation set.**



(a) Contrast      (b) Gaussian Noise      (c) Snow      (d) Frost
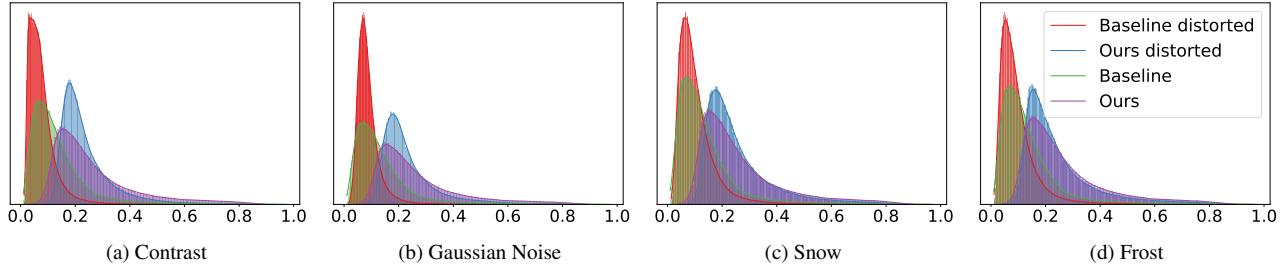
Figure 8. **The distribution of prediction scores from the baseline DINO-DETR [74] and our *v*-CLR under four types of image distortion.** For visualization clarity, we calculate the distribution of *top*-50 prediction scores.

model's prediction confidence.

## D. Comparison with SiameseDETR

We further compare our method with Siamese DETR [13], a recent self-supervised DETR-like object detector. Siamese DETR employs two augmented views to enforce instance-

14

| Method | $AR_1^b$ | $AR_{10}^b$ | $AR_{100}^b$ |
|---|---|---|---|
| SiameseDETR [13] | 12.4 | 23.0 | 30.7 |
| *v*-CLR (**ours**) | **16.8** | **42.7** | **60.2** |

Table 12. **Comparison with Siamese DETR [13].** For a fair comparison, all experiments are conducted on the Non-VOC→VOC setting with Deformable-DETR [78].

level consistency. Although it also utilizes transformations, its motivation differs substantially from ours, and the transformations in Siamese DETR do not specifically address texture bias. We evaluate both methods in the Non-VOC→VOC setting, as shown in Tab. 12, ensuring a fair comparison since VOC classes are unknown to both models. Experimental results reveal that our method surpasses Siamese DETR by a significant margin across all evaluation metrics, underscoring the effectiveness of our proposed framework.