

# CoDEPS: Online Continual Learning for Depth Estimation and Panoptic Segmentation

Niclas Vödich<sup>1\*</sup>, Kürsat Petek<sup>1\*</sup>, Wolfram Burgard<sup>2</sup>, and Abhinav Valada<sup>1</sup>

<sup>1</sup>University of Freiburg    <sup>2</sup>University of Technology Nuremberg

**Abstract**—Operating a robot in the open world requires a high level of robustness with respect to previously unseen environments. Optimally, the robot is able to adapt by itself to new conditions without human supervision, e.g., automatically adjusting its perception system to changing lighting conditions. In this work, we address the task of continual learning for deep learning-based monocular depth estimation and panoptic segmentation in new environments in an online manner. We introduce CoDEPS to perform continual learning involving multiple real-world domains while mitigating catastrophic forgetting by leveraging experience replay. In particular, we propose a novel domain-mixing strategy to generate pseudo-labels to adapt panoptic segmentation. Furthermore, we explicitly address the limited storage capacity of robotic systems by leveraging sampling strategies for constructing a fixed-size replay buffer based on rare semantic class sampling and image diversity. We perform extensive evaluations of CoDEPS on various real-world datasets demonstrating that it successfully adapts to unseen environments without sacrificing performance on previous domains while achieving state-of-the-art results. The code of our work is publicly available at <http://codeps.cs.uni-freiburg.de>.

## I. INTRODUCTION

Deploying robots such as autonomous cars in urban scenarios requires a holistic understanding of the environment with a unified perception of semantics, instances, and depth. The joint solution of these tasks enables vision-based methods to generate a 3D semantic reconstruction of the scene, which can be leveraged for downstream applications such as localization or planning. While deep learning-based state-of-the-art approaches perform well when inference is done under similar conditions as used for training, their performance can drastically decrease when the new target domain differs from the source domain, e.g., due to environmental conditions [32], different sensor parameters [4, 7]. This domain gap poses a great challenge for robotic platforms that are deployed in the open world without prior knowledge about the target domain. Additionally, unlike the source domain where ground truth annotations are generally assumed to be known and can be used for the initial training, such supervision is not applicable to the target domain due to the absence of labels, rendering classical domain adaptation methods unsuitable. Unsupervised domain adaptation attempts to overcome these limitations. However, the vast majority of proposed approaches focuses on sim-to-real domain adaptation mostly in an offline manner [13, 24], i.e., a directed knowledge transfer without the need to avoid catastrophic forgetting and with access to

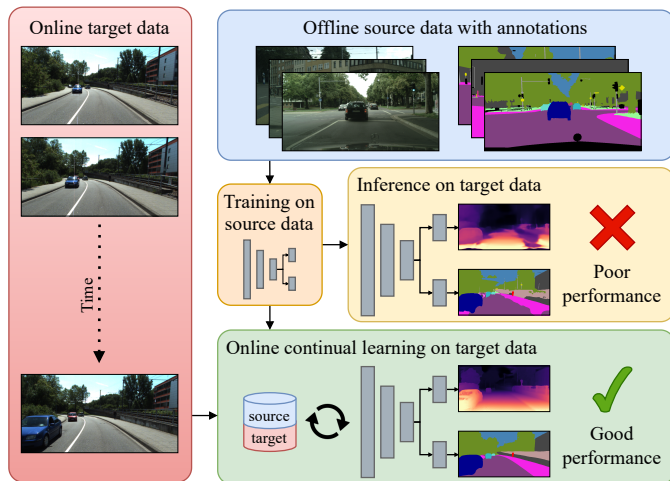


Fig. 1. Neural networks often perform poorly when deployed on a target domain that differs from the source domain used for training. To close this domain gap, we propose to continuously adapt the network by exploiting online target images. To mitigate catastrophic forgetting and enhance generalizability, we leverage a fixed-size replay buffer allowing the method to revisit data from both the source and target domains.

abundant target annotations. Additionally, such works might not consider limitations on a robotic platform, e.g., available compute hardware and limited storage capacity [18, 33].

In this work, we use online continual learning to address these challenges for depth estimation and panoptic segmentation in a multi-task setup. As shown in Fig. 1, we leverage images from an onboard camera to perform online continual learning enhancing performance during inference time. While a naive approach would result in overfitting to the current scene, our method CoDEPS mitigates forgetting by leveraging experience replay of both source data and previously seen target images. We combine a classical replay buffer with generative replay in the form of a novel cross-domain mixing strategy allowing us to exploit supervised ideas also for unlabeled target data. We explicitly address the aforementioned hardware limitations by using only a single GPU and restricting the replay buffer to a fixed size. We demonstrate that CoDEPS successfully improves on new target domains without sacrificing performance on previous domains.

The main contributions of this work are as follows:

- 1) We introduce CoDEPS, the first online continual learning approach for joint monocular depth estimation and panoptic segmentation.
- 2) We propose a novel cross-domain mixing strategy to adapt panoptic segmentation to unlabeled target data.

\*Equal contribution.

- 3) To address the storage restrictions of robotic platforms, we leverage a fixed-size replay buffer based on rare class sampling and image diversity.
- 4) We extensively evaluate CoDEPS and compare it to other methods in challenging real-to-real settings.
- 5) We release our code and the trained models at <http://codeps.cs.uni-freiburg.de>.

## II. RELATED WORK

In this section, we provide an overview of monocular depth estimation, panoptic segmentation, and unsupervised domain adaptation including continual learning.

*Monocular Depth Estimation:* Monocular depth estimation is the task of predicting a dense depth map from a single RGB image. While supervised approaches exploit measurements from range sensors to supervise the network predictions [30], unsupervised methods leverage geometric cues from temporal context [10, 38]. Most of the research on unsupervised learning tackles the limitations of the so-called photometric loss function that is usually employed for unsupervised depth learning, e.g., dynamic object handling [3, 5, 20], occlusion [11], and abrupt illumination changes [36]. In this work, we leverage Monodepth2 [11] for unsupervised depth learning and employ it similarly to Guizilini et al. [13] for the purpose of domain adaptation.

*Panoptic Segmentation:* Panoptic segmentation unifies the two tasks of semantic and instance segmentation by fusing the respective targets into a joint output. Furthermore, semantic classes are grouped into “stuff” classes, e.g., *road* or *building*, and “thing” classes, e.g., *car* or *pedestrian*. In particular, the goal of vision-based panoptic segmentation is to assign a semantic class to every pixel of an image and an additional instance label to each object belonging to the “thing” classes. Panoptic segmentation networks usually comprise a joint encoder and separate decoders for each subtask, whose outputs are subsequently merged by a panoptic fusion module. Existing works can be categorized into bottom-up [6, 28] and top-down [12, 27] approaches. Whereas bottom-up methods detect instances in a proposal-free manner from the semantic prediction, top-down methods include an additional proposal generation step. Contradictions to the semantic predictions are then resolved during post-processing. In this work, we build upon the bottom-up Panoptic-Deeplab [6] with changes to the semantic head according to Guizilini et al. [13].

*Unsupervised Domain Adaptation:* Domain adaptation aims to bridge the domain gap between a source domain  $\mathcal{S}$  used for training and a target domain  $\mathcal{T}$  used for inference to mitigate a loss in performance. An important aspect is whether the performance on the source domain must be maintained, linking domain adaptation to continual learning (CL) [23], where the objective of a task or the task itself can change over time. A CL system has to adapt to the new target objective while retaining the knowledge to solve the previous task(s), i.e., avoiding catastrophic forgetting. Ideally, the CL system can further achieve positive forward transfer, i.e.,

improve on future yet untrained tasks given the current task. In many real-world scenarios ground truth annotations for the target domain are not available, thus requiring unsupervised domain adaptation (UDA) methodology. Offline UDA assumes that abundant target data is accessible. However, in order to guarantee the continuous operation of a robot in new domains, UDA approaches have to work online without previous target data collection.

Offline UDA can leverage both annotated source data and abundant unlabeled target data, enabling learning a given task from  $\mathcal{S}$  while simultaneously adapting the network to  $\mathcal{T}$ . For depth estimation, DESC [24] adapts from a synthetic source domain containing RGB images and ground truth depth to a real-world target domain by performing source-to-target style transfer and using a consistency loss between depth predictions from RGB and semantic maps. GUDA [13] tackles UDA for semantic segmentation using depth estimation as a proxy task. A shared encoder with task-specific heads for depth estimation and semantic segmentation is trained via source supervision. Simultaneously, data from  $\mathcal{T}$  is used to update the encoder and depth head in an unsupervised manner. Due to the refined weights of the encoder, the semantic predictions on  $\mathcal{T}$  improve as well. Another common approach for adapting semantic segmentation is cross-domain sampling enabling partial supervision on  $\mathcal{T}$ . DACS [31] mixes images from  $\mathcal{S}$  and  $\mathcal{T}$  by copying the pixels of a source image to a target image based on the semantic labels [29]. The semantic prediction of the target image is updated with ground truth source labels for the same set of pixels. The network is then jointly trained on annotated source data and the pseudo-labeled mixing data. Recently, ConfMix [25] proposed a simple yet effective mixing strategy for object detection, where a target image is divided into rectangular image regions. The region with the most confident predictions is then copied onto a source image and the respective ground truth annotations. Finally, Huang et al. [16] propose a UDA method for panoptic segmentation by regularizing complementary features from semantic and instance segmentation. In this work, we extend the aforementioned mixing strategies to instance-based sampling and explicitly address differing camera parameters.

During online UDA, samples from  $\mathcal{T}$  can only be accessed in a consecutive manner resembling the image stream of a camera. Typically, a network is trained offline via supervision on  $\mathcal{S}$  and then adapted to  $\mathcal{T}$  during inference time. Such a setup rises two main challenges: first, incoming target samples originate from highly similar scenes and thus drastically reduce the diversity; second, this similarity of consecutive samples leads to a strong overfitting of the model to the scene [37]. Initial works for online UDA focused on depth learning [18, 37] and visual odometry [21, 33], for which unsupervised training schemes are already well established. Whereas Zhang et al. [37] propose novel network modules that are adapted via a meta-learning paradigm to mitigate forgetting, CoMoDA [18] employs a common CL strategy, i.e., experience replay to combine the online target sample with previously seen samples. Continual SLAM [33] also uses un-

supervised depth estimation as a proxy task to enhance visual odometry during inference time. Additionally, it demonstrates that incorporating samples from  $\mathcal{S}$  and previous target domains  $\mathcal{T}_i$  prevents catastrophic forgetting when revisiting domains. Similar settings involving multiple target domains, which are hence closely related to classical CL, are also addressed for semantic segmentation. CBNA [17] mixes statistics from  $\mathcal{S}$  and  $\mathcal{T}$  to update the batch normalization layers and showcases the efficacy of the approach on continually visited target domains. CoTTA [35] adapts the entire network without using source data but self-supervision. To tackle error accumulation, it uses an exponential moving average filter and student-teacher consistency when updating the network weights. Using depth estimation as a proxy task, Kuznietsov et al. [19] extend GUDA [13] to online UDA with experience replay and confidence regularization on the semantic predictions. To the best of our knowledge, we propose the first approach for online continual UDA for joint depth estimation and panoptic segmentation.

### III. TECHNICAL APPROACH

The setting investigated in this work consists of two steps. First, we train a neural network on the source domain  $\mathcal{S}$  partly using ground truth supervision. Second, to close the gap between domains, we continuously adapt the network during inference time on the target domain  $\mathcal{T}$  using a replay buffer and unsupervised training strategies.

#### A. Network Architecture and Source Domain Pretraining

In this section, we detail the network architecture and loss functions that we employ during the pretraining stage on the source domain.

*Architecture:* We build our network following a common multi-task design scheme, i.e., using a single backbone followed by task-specific heads. A high-level overview of the network architecture is shown in Fig. 2. In detail, we use a ResNet-101 [14] as the shared encoder for all three tasks including depth prediction, semantic segmentation, and instance segmentation. The depth head follows the design of Monodepth2 [11] comprising five consecutive convolutional layers with skip connections to the encoder. Additionally, we include a separate PoseNet consisting of a ResNet-18 encoder and a four-layer CNN to estimate the camera motion between two image frames. For panoptic segmentation, we follow the bottom-up method Panoptic-Deeplab [6], leveraging separate heads for semantic segmentation and instance segmentation, and slightly modify the semantic head [13]. Specifically, the instance head consists of two sub-heads to predict the center of an object and to associate each pixel of an image to the corresponding object or the background. Finally, a panoptic fusion module [6] assigns a semantic label to the class-agnostic instance predictions using majority voting over the semantic predictions of all pixels within an instance.

*Source Domain Pretraining:* During the initial training phase on the source domain, we assume to have access to image

sequences as well as ground truth panoptic segmentation annotations. In the following, we briefly describe the respective loss functions that we employ for training the three task-specific heads.

We train the depth estimation head using the common methodology of unsupervised training based on the photometric error [11]. In particular, we leverage an image triplet  $\{\mathbf{I}_{t_0}, \mathbf{I}_{t_1}, \mathbf{I}_{t_2}\}$  to predict depth  $\mathbf{D}_{t_1}$  and camera motion  $\mathbf{M}_{t_0 \rightarrow t_1}$  and  $\mathbf{M}_{t_1 \rightarrow t_2}$ . Afterwards, we compute the photometric error loss  $\mathcal{L}_{pe}^d$  as a weighted sum of the reprojection loss  $\mathcal{L}_{pr}^d$  and the image smoothness loss  $\mathcal{L}_{sm}^d$ :

$$\mathcal{L}_{pe}^d = \lambda_{pr} \mathcal{L}_{pr}^d + \lambda_{sm} \mathcal{L}_{sm}^d. \quad (1)$$

We train the semantic segmentation head in a supervised manner using the bootstrapped cross-entropy loss with hard pixel mining  $\mathcal{L}_{bce}^{sem}$  following Panoptic-Deeplab [6].

For training the instance segmentation head, we adopt the MSE loss  $\mathcal{L}_{center}^{ins}$  for the center head and the L1 loss  $\mathcal{L}_{offset}^{ins}$  for the offset head. The total loss to supervise instance segmentation is then computed as a weighted sum:

$$\mathcal{L}_{co}^{ins} = \lambda_{center} \mathcal{L}_{center}^{ins} + \lambda_{offset} \mathcal{L}_{offset}^{ins}. \quad (2)$$

#### B. Online Adaptation

After the described network has been trained on the source domain  $\mathcal{S}$  using the aforementioned losses, we aim to adapt it to the target domain  $\mathcal{T}$  in a continuous manner. That is, unlike other works, data from the target domain is revealed frame by frame resembling the online stream of an onboard camera. As depicted in Fig. 2, every adaptation iteration consists of the following steps:

- 1) Construct an update batch by combining online and replay data.
- 2) Generate pseudo-labels using the proposed cross-domain mixing strategy.
- 3) Perform backpropagation to update the network weights.
- 4) Update the replay buffer.

In this section, we first detail the structure of the utilized replay buffer and then propose adaptation schemes for both depth estimation and panoptic segmentation.

*Replay Buffer and Batch Generation:* Upon receiving a new image taken by the robot’s onboard camera, we construct a batch that is used to perform backpropagation on the network weights. In detail, a batch  $\mathbf{b}_t$  consists of the current online image  $\mathbf{I}_t \in \mathcal{T}$ , previously received target images  $\mathbf{I}_{\mathcal{T}_i} \in \mathbf{B}_{\mathcal{T}}$ , and fully annotated source samples  $\mathbf{I}_{\mathcal{S}_i} \in \mathbf{B}_{\mathcal{S}}$ . Here,  $\mathbf{B}_{\mathcal{T}} \subseteq \mathcal{T}$  and  $\mathbf{B}_{\mathcal{S}} \subseteq \mathcal{S}$  denote the respective replay buffers. Formally\*,

$$\mathbf{b}_t = \{\mathbf{I}_t, \mathbf{I}_{\mathcal{T}_0}, \mathbf{I}_{\mathcal{T}_1}, \dots, \mathbf{I}_{\mathcal{S}_0}, \mathbf{I}_{\mathcal{S}_1}, \dots\}. \quad (3)$$

By revisiting target images from the past, we increase the diversity in the loss signal on the target domain and hence mitigate overfitting to the current scene. This further accounts for situations in which the current online image suffers from visual

\*To improve readability, we omit in the notation that each image sample includes its two previous frames enabling unsupervised depth estimation.

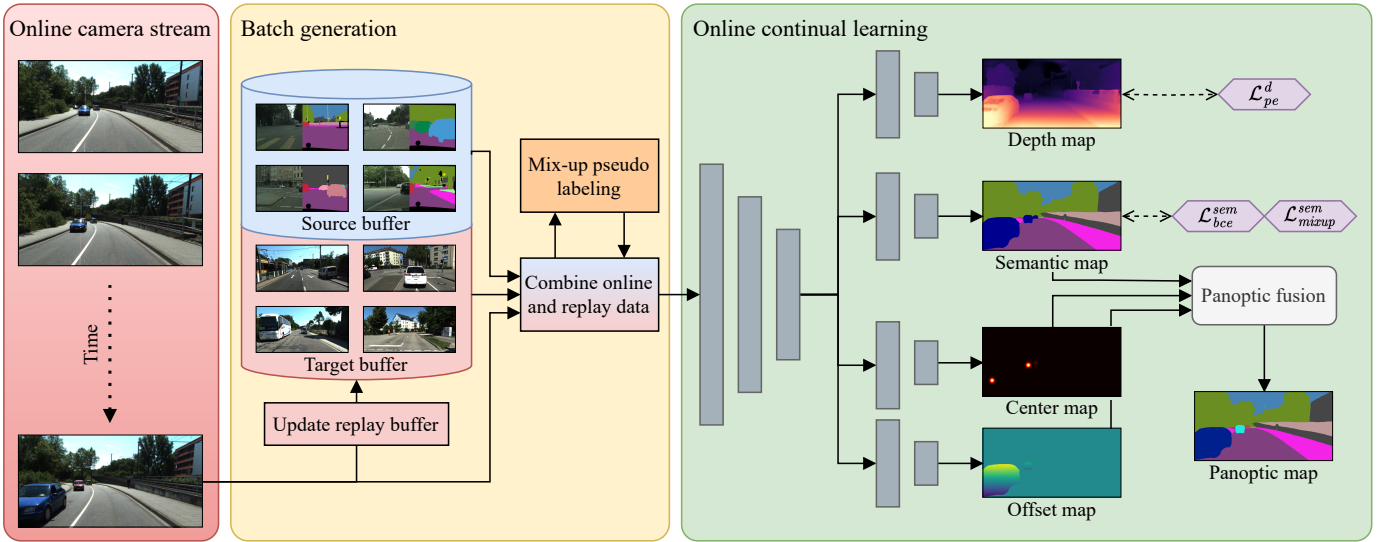


Fig. 2. Overview of our proposed CoDEPS. Unlabeled RGB images from an online camera stream are combined with samples from a replay buffer comprising both annotated source samples and previously seen target images. Cross-domain mixing enables pseudo-supervision on the target domain. The network weights are then updated via backpropagation using the constructed data batch. The additional PoseNet required for unsupervised monocular depth estimation is omitted in this visualization.

artifacts, e.g., overexposure. Similarly, revisiting samples from the source domain addresses the problem of catastrophic forgetting by ensuring that previously acquired knowledge can be preserved. Additionally, the annotations of the source samples enable pseudo-supervision on the target domain by exploiting cross-domain mixing strategies. For both the target and the source replay, we randomly draw multiple samples from the respective replay buffer and apply augmentation to stabilize the loss. In particular, we perform RGB histogram matching of the source images to the online target image, and all available source samples have to be selected once before repetition is allowed to ensure diverse source supervision.

Similar to previous works [1, 34], we explicitly consider limitations on the size of the replay buffer to closely resemble the deployment on a robotic platform, where disk storage is an important factor. This poses two questions: First, how to sample from  $\mathcal{S}$  to construct the fixed source buffer  $\mathbf{B}_{\mathcal{S}}$  that is prebuilt offline and, second, how to update the dynamic target buffer  $\mathbf{B}_{\mathcal{T}}$  during deployment? To construct  $\mathbf{B}_{\mathcal{S}}$ , we propose a refined version of rare class sampling (RCS) [15]. The frequency  $f_c$  of each class  $c \in \mathcal{C}$  is calculated based on the number of pixels with class  $c$ :

$$f_c = \frac{\sum_{\mathbf{I} \in \mathcal{S}} \sum_p \mathbf{1}_c(p_{c'})}{|\mathcal{S}| \cdot H \cdot W}, \quad (4)$$

where  $H$  and  $W$  denote the height and width of the images in  $\mathcal{S}$  and  $p_{c'} \in \mathbf{I}$  refers to a pixel with class  $c'$ . The indicator function is 1 if  $c'$  equals  $c$  and 0 otherwise. The probability of sampling a class is then given by

$$P(c) = \frac{e^{(1-f_c)/T}}{\sum_{c' \in \mathcal{C}} e^{(1-f_{c'})/T}}, \quad (5)$$

with temperature  $T$  controlling the smoothness of the distribution, i.e., a smaller  $T$  assigns a higher probability to rare

classes. In detail, we first sample a class  $c \sim P$  and then retrieve all candidate images containing pixels with class  $c$ . Instead of taking a random image from these candidates, we sample according to the number of pixels with class  $c$ . We repeat both steps  $|\mathbf{B}_{\mathcal{S}}|$  times without selecting the same image more than once. Using RCS ensures that  $\mathbf{B}_{\mathcal{S}}$  contains sufficiently many images with rare classes such that the performance on these classes will further improve during adaptation.

Since  $\mathcal{T}$  does not contain annotations and, particularly in the beginning, predictions are not reliable, we cannot use RCS for updating  $\mathbf{B}_{\mathcal{T}}$ . Instead, we invert the common methodology of loop closure detection for visual SLAM [33], i.e., the image  $\mathbf{I}_{\mathbf{t}}$  is only added to  $\mathbf{B}_{\mathcal{T}}$  if its cosine similarity with respect to all samples within the buffer is below a threshold.

$$\text{sim}_{\cos}(\mathbf{I}_{\mathbf{t}}) = \max_{\mathbf{I}_{\mathcal{T}_i} \in \mathbf{B}_{\mathcal{T}}} \cos(\text{feat}(\mathbf{I}_{\mathbf{t}}), \text{feat}(\mathbf{I}_{\mathcal{T}_i})), \quad (6)$$

where  $\text{feat}(\cdot)$  refers to the image features extracted from the final layer of the shared encoder, which is not adapted. If  $\mathbf{B}_{\mathcal{T}}$  is completely filled, we remove the following image to maximize image diversity:

$$\arg \max_{\mathbf{I}_{\mathcal{T}_i} \in \mathbf{B}_{\mathcal{T}}} \sum_{\mathbf{I}_{\mathcal{T}_j} \in \mathbf{B}_{\mathcal{T}}} \cos(\text{feat}(\mathbf{I}_{\mathcal{T}_i}), \text{feat}(\mathbf{I}_{\mathcal{T}_j})) \quad (7)$$

*Depth Adaptation:* To adapt the monocular depth estimation head along with the PoseNet, we exploit the fact that the photometric error loss (Eq. 1) does not require ground truth annotations. Hence, we can directly transfer it to the implemented continual adaptation. In particular, we compute  $\mathcal{L}_{pe}^d$  for the constructed batch  $\mathbf{b}_{\mathbf{t}}$  and average the loss such that each sample contributes by the same amount:

$$\mathcal{L}_{pe}^d(\mathbf{b}_{\mathbf{t}}) = \frac{\mathcal{L}_{pe}^d(\mathbf{I}_{\mathbf{t}}) + \sum_i \mathcal{L}_{pe}^d(\mathbf{I}_{\mathcal{T}_i}) + \sum_j \mathcal{L}_{pe}^d(\mathbf{I}_{\mathcal{S}_j})}{|\mathbf{b}_{\mathbf{t}}|}. \quad (8)$$

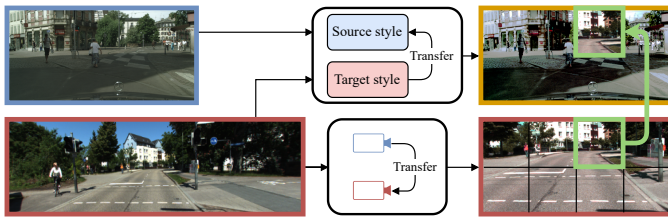


Fig. 3. Our proposed cross-domain mixing strategy first transfers the image style from the target to the source sample. Then it augments the target image to match the appearance of the source camera. Finally, a random image patch is copied from the target to the source image. The source annotations are retained and completed by the network’s estimate on the copied image patch. The result serves as pseudo-label, combining self-iterative learning with ground truth supervision.

Furthermore, if the predicted camera motion is below a threshold, i.e., the robot is presumably not moving, we do not compute the  $\mathcal{L}_{pe}^d(\mathbf{I}_t)$  and subtract 1 from the denominator to avoid adding a bias to the current scene.

*Panoptic Adaptation:* As described in Sec. III-A, panoptic segmentation is the fused output of a semantic head and an instance head. We observe that the decrease in performance on samples from unseen domains can mostly be attributed to the semantic head, while instance predictions remain stable. Cross-domain mixing strategies allow leveraging ideas from supervised training to an unsupervised setting, where ground truth annotations are unknown. In CoDEPS, we bootstrap annotated source samples and high-confident target predictions to artificially generate pseudo-labels for the target samples in an online fashion to supervise the semantic head. Similar to depth adaptation, we continue to compute  $\mathcal{L}_{bce}^{sem}$  on  $\{\mathbf{I}_{S_0}, \mathbf{I}_{S_1}, \dots\}$  to mitigate forgetting, and freeze the instance head.

We further design a mixing strategy combining pixels of images from both  $\mathcal{S}$  and  $\mathcal{T}$ , that considers multiple factors, which are unique to the online continual learning scenario: (1) the robust pretraining on a dedicated source dataset, which may result in significant performance degradation on the target dataset if the pre-trained weights are strongly adapted; (2) the existence of different cameras leading to significant changes in the field-of-view, geometric appearance of objects, resolution, and aspect ratio of the images; and (3) the continuously evolving visual appearance of street scenes during adaptation. To address these challenges, our cross-domain mixing approach employs a three-step method to generate the adaptation signal. First, we perform style transfer from the target image  $\mathbf{I}_{T_i}$  to the source image  $\mathbf{I}_{S_j}$  by aligning their pixel value histograms, as depicted in Fig. 3. This allows supervision with ground truth labels on images that are of similar visual appearance as the target image. Second, we apply a geometric transformation on  $\mathbf{I}_{T_i}$  based on the camera intrinsics of the source and target domains denoted by  $\mathbf{K}_S$  and  $\mathbf{K}_T$ , respectively. To this end, we assume a constant depth distribution over  $\mathbf{I}_{S_j}$ , lift the pixel values into Euclidean space via inverse camera projection, and project the lifted points back into the camera view of  $\mathbf{I}_{T_i}$  as follows:

$$\mathbf{I}'_T(\mathbf{p}_s) = \mathbf{I}_T \langle \mathbf{K}_T \mathbf{K}_S^{-1} \mathbf{p}_s \rangle, \quad (9)$$

where  $\langle \cdot \rangle$  denotes the bilinear sampling operator and  $\mathbf{p}_s$  is a

pixel coordinate in the source image. Equation 9 results in an adapted target image  $\mathbf{I}'_T$  with an adjusted field of view, resolution, and a geometric appearance of the scene similar to that of  $\mathbf{I}_S$ . The final step in the process involves separating  $\mathbf{I}'_T$  into multiple segments and randomly selecting one of them to be inserted into the style-transferred source image, see Fig. 3. To avoid providing a flawed supervision signal caused by geometrically unrealistic images, we only insert a single patch. Similarly, the ground truth labels of the pixels from  $\mathbf{I}_S$  are retained, and the semantic labels estimated by the network are used to label the inserted patch after intrinsics transformation. The generated image is then fed into the network and training is performed using the cross-entropy loss and the generated pseudo-labels of the mixed image. To mitigate the decline in performance commonly associated with self-iterative training on predicted pseudo-labels, often resulting in class collapse, we utilize an exponentially moving average (EMA) filter for updating the network weights. In detail, we create a duplicate with network weights  $w_{EMA}$  of the initial model with weights  $w$  and use this so-called EMA model to generate the semantic predictions. During continual learning, the weights  $w$  are updated via backpropagation on  $\mathbf{b}_t$ . Then, the EMA model is updated as follows:

$$w_{EMA} \leftarrow \alpha \cdot w_{EMA} + (1 - \alpha) \cdot w, \quad (10)$$

where  $\alpha$  denotes the contribution of the EMA model.

#### IV. EXPERIMENTAL EVALUATION

In the following sections, we provide further details on the pretraining step and the datasets that we evaluate on. We present extensive experimental results on the efficiency and efficacy of our proposed approach and include ablation studies on important design choices. Finally, we expand the experimental setup to multi-domain adaptation closely resembling classical continual learning settings.

We follow the evaluation protocol of Zhang et al. [37]. In detail, we compute the evaluation metrics on the frame of the current timestamp before using the same frame to perform backpropagation and update the model weights. Once 70% of a sequence is processed, we calculate the average of the accumulated metrics. Additionally, we report the scores on the remaining 30% of the same sequence without further weight updates to analyze the performance of the adapted model. In the tables, we refer to these types of evaluation by protocol 1 and protocol 2, respectively. We further denote the respective parts of a sequence by *adapt* and *eval*. Unlike Zhang et al. [37], we define our task in the context of continual learning. To measure knowledge retention and hence mitigate catastrophic forgetting, we introduce protocol 3 as evaluating the adapted model on the *val* split of the source dataset.

##### A. Datasets

To simulate data from a variety of domains, we employ our method on three datasets, namely Cityscapes [8], KITTI-360 [22], and SemKITTI-DVPS [2]. In particular, we utilize Cityscapes for pre-training and sequences of both KITTI-360

TABLE I  
EFFICACY OF THE NETWORK

Method	Dataset	mIoU $\uparrow$	RMSE $\downarrow$	Abs Rel $\downarrow$	$\delta_1$ $\uparrow$
GUDA	KITTI	—	4.42	0.11	0.88
CoDEPS		62.8	3.52	0.09	0.90
GUDA	Cityscapes	72.9	—	—	—
CoDEPS		72.9	10.16	0.19	0.78

Our utilized network is able to reproduce the performance of the baseline method GUDA [13] for both semantic segmentation (mIoU) and depth estimation (RMSE, Abs Rel,  $\delta_1$ ). The performance of GUDA is reported by the authors. To evaluate CoDEPS on KITTI, we use sequence 08 *eval* of SemKITTI-DVPS.

and SemKITTI-DVPS for adaptation. In the supplementary video we further provide qualitative results on in-house data recorded with our robotic platform.

*Cityscapes*: The Cityscapes Dataset [8] is a large-scale autonomous driving dataset that was recorded in 50 cities in Germany and bordering regions. It includes RGB images, panoptic annotations, and vehicle metadata. In this work, we utilize the fine panoptic labels to train the semantic and instance heads in a supervised manner. Additionally, we leverage the sequence image data of the left camera to train the depth prediction in an unsupervised fashion. Finally, we compute the depth error metrics using the provided disparity maps.

*KITTI-360*: The KITTI-360 Dataset [22] is a relatively recently released public dataset for the domain of autonomous driving, which was recorded in the city of Karlsruhe, Germany. It includes both 2D and 3D panoptic annotations for RGB images and LiDAR data. In this work, we predominantly utilize the RGB images to simulate an online image stream of an onboard camera. In particular, we use these images to adapt our network in a self-supervised manner. To compute evaluation metrics, we compare our predictions with the ground truth measurements and annotations of the dataset.

*SemKITTI-DVPS*: The SemKITTI-DVPS [30] is based on the odometry benchmark of the KITTI Dataset [9], which was recorded in Karlsruhe, Germany. We utilize the RGB images to simulate an onboard camera and to adapt our network to the new domain. Furthermore, we compute depth metrics based on the provided projected LiDAR points and the semantic/panoptic metrics using the extension SemanticKITTI [2].

*Semantic Labels*: As the aforementioned datasets use different labeling policies for the semantic annotations, we use the 19 classes of Cityscapes as the reference definition and remap classes of the other datasets. However, certain classes do not exist in the adaptation datasets (*wall*, *traffic light*, *bus*, *train*). For consistency across the datasets, we merge *wall* with *building* and remove the other three classes. Additionally, we merge *motorcycle* and *bicycle* into *two-wheeler* to increase the number of annotated pixels. Consequently, we consider nine “stuff” classes and five “thing” classes, listed in Table V. Note that *sky* is not included in SemKITTI-DVPS due to using LiDAR annotations and hence excluded in the evaluation on this dataset.

## B. Pretraining Protocol

The initial state of the network weights before adaptation is obtained by initializing the encoders using pretrained weights from the ImageNet dataset, followed by training the entire model on the Cityscapes dataset. In detail, we use the Adam optimizer with a constant learning rate  $lr = 0.0001$  and train the entire network for 250 epochs. In our experiments, we compare the performance of our approach to directly training on the target dataset, which can be considered as a theoretical upper limit having full target knowledge. Due to the unbalanced class distribution of KITTI-360, we train a copy of the network in two steps, using the Adam optimizer with  $lr = 0.0001$  on sequences 00-07. We train for 45 epochs while ignoring the most common classes *road*, *sidewalk*, *building*, and *vegetation*, followed by 55 epochs including all classes. Similarly, for SemKITTI-DVPS, we train another copy of the network on sequences 00-06, 09, and 10 for 30 epochs without the aforementioned classes plus *terrain* and *sky*, which is not included in the dataset, followed by 30 epochs including all classes. In Table I, we demonstrate that our implemented network is able to reproduce the performance of the baseline method GUDA [13].

## C. Online Adaptation

In this section, we extensively evaluate our proposed CoDEPS with respect to both adapting to a new domain and retaining knowledge to mitigate forgetting. In detail, for all presented experiments, we freeze the shared encoder following the study by McCraith et al. [26]. Based on the ablation study in Sec. IV-D, we use a buffer size of 300. For RCS, we follow Hoyer et al. [15] and set  $T = 0.01$ . Updating the EMA model is done with  $\alpha = 0.99$ .

In Table II, we assess the performance of CoDEPS on all sequences of the KITTI-360 dataset and compare it with the baseline method “only source”, which is also pretrained on Cityscapes but does not perform further adaptation to the target domain  $\mathcal{T}$ . This approach should be interpreted as a lower performance bound that must be improved. We demonstrate the key performance metrics of both protocols 1 and 2. As shown in Table II, CoDEPS achieves a performance boost across the board, as measured by the mIoU metric and all depth metrics of protocol 1. We attribute this improvement to the additional supervision signals incorporated into the segmentation head through our mixing strategy and the self-supervised reconstruction loss for depth adaptation. The improvement in semantic segmentation further enhances the panoptic segmentation metrics. With respect to protocol 2, CoDEPS reduces the depth errors on all sequences and improves the performance of semantic and panoptic segmentation on the vast majority of sequences. Note that on sequence 03 the panoptic metrics increase significantly despite the consistent mIoU, which we attribute to the more refined segmentation of objects due to our proposed cross-domain mixing strategy.

For the following experiments, we consider the case of using Cityscapes as the source domain and sequence 10 of KITTI-360 as the target domain. In Fig. 4, we illustrate

TABLE II  
ADAPTATION PERFORMANCE

Method	Sequence	Protocol 1						Protocol 2					
		mIoU $\uparrow$	PQ $\uparrow$	SQ $\uparrow$	RQ $\uparrow$	RMSE $\downarrow$	Abs Rel $\downarrow$	mIoU $\uparrow$	PQ $\uparrow$	SQ $\uparrow$	RQ $\uparrow$	RMSE $\downarrow$	Abs Rel $\downarrow$
Only source CoDEPS	00	51.61	39.10	72.72	50.48	6.54	0.36	49.94	35.29	72.14	45.50	6.08	0.34
		<b>53.76</b>	<b>40.72</b>	<b>72.90</b>	<b>52.51</b>	<b>5.09</b>	<b>0.19</b>	<b>52.08</b>	<b>36.08</b>	<b>72.58</b>	<b>46.08</b>	<b>4.34</b>	<b>0.15</b>
Only source CoDEPS	02	45.97	31.83	67.62	41.08	6.26	0.35	46.55	30.13	65.03	39.30	6.06	0.36
		<b>46.62</b>	<b>32.11</b>	<b>67.74</b>	<b>41.62</b>	<b>4.31</b>	<b>0.16</b>	<b>47.48</b>	<b>30.33</b>	<b>65.35</b>	<b>39.46</b>	<b>3.76</b>	<b>0.13</b>
Only source CoDEPS	03	46.63	28.15	57.41	35.23	<b>8.20</b>	0.34	<b>52.10</b>	28.20	56.67	35.77	7.34	0.29
		<b>47.94</b>	<b>29.05</b>	<b>58.07</b>	<b>36.10</b>	8.26	<b>0.33</b>	52.00	<b>31.13</b>	<b>61.51</b>	<b>39.65</b>	<b>6.98</b>	<b>0.18</b>
Only source CoDEPS	04	45.02	29.34	65.48	38.15	6.70	0.37	45.53	30.13	<b>70.85</b>	38.89	6.61	0.38
		<b>45.40</b>	<b>29.78</b>	<b>65.89</b>	<b>38.84</b>	<b>5.00</b>	<b>0.19</b>	<b>45.68</b>	<b>30.63</b>	66.18	<b>39.89</b>	<b>4.33</b>	<b>0.17</b>
Only source CoDEPS	05	48.94	32.19	66.80	41.37	6.76	0.37	<b>44.52</b>	<b>27.34</b>	<b>60.72</b>	<b>35.58</b>	5.93	0.43
		<b>49.26</b>	<b>32.96</b>	<b>66.98</b>	<b>42.40</b>	<b>5.25</b>	<b>0.21</b>	43.79	26.48	60.33	34.88	<b>4.68</b>	<b>0.25</b>
Only source CoDEPS	06	46.03	29.88	66.58	38.42	6.09	0.39	46.28	31.79	70.47	41.40	6.12	0.37
		<b>46.53</b>	<b>30.45</b>	<b>66.66</b>	<b>39.20</b>	<b>4.97</b>	<b>0.22</b>	<b>47.27</b>	<b>31.99</b>	<b>70.74</b>	<b>41.71</b>	<b>4.23</b>	<b>0.18</b>
Only source CoDEPS	07	40.54	28.48	66.52	34.42	7.83	0.34	59.07	27.62	45.88	35.41	9.64	0.38
		<b>41.46</b>	<b>29.30</b>	<b>67.64</b>	<b>35.58</b>	<b>6.50</b>	<b>0.22</b>	<b>60.57</b>	<b>30.91</b>	<b>50.25</b>	<b>39.79</b>	<b>6.48</b>	<b>0.20</b>
Only source CoDEPS	09	50.59	37.26	74.06	47.38	6.03	0.36	50.78	36.57	72.22	46.75	5.60	0.35
		<b>52.29</b>	<b>38.02</b>	<b>74.88</b>	<b>48.21</b>	<b>4.74</b>	<b>0.19</b>	<b>51.53</b>	<b>37.56</b>	<b>72.87</b>	<b>47.99</b>	<b>4.56</b>	<b>0.16</b>
Only source CoDEPS	10	51.94	32.60	71.27	32.60	8.06	0.35	45.74	30.62	69.56	39.49	7.90	0.33
		<b>53.02</b>	<b>33.50</b>	<b>71.53</b>	<b>33.50</b>	<b>7.19</b>	<b>0.22</b>	<b>49.91</b>	<b>31.91</b>	<b>70.68</b>	<b>40.95</b>	<b>5.57</b>	<b>0.15</b>

Comparison between our CoDEPS and the performance of the same architecture without performing online continual learning on the respective sequence of the KITTI-360 dataset. Thus, “only source” refers to the model weights after pretraining on Cityscapes. The listed metrics are mean intersection over union (mIoU) for semantic segmentation; panoptic quality (PQ), segmentation quality (SQ), and recognition quality (RQ) for panoptic segmentation; root mean squared error (RMSE) and absolute relative error (Abs Rel) for monocular depth estimation. Bold values denote the best result on each sequence.

TABLE III  
CONTINUAL LEARNING FOR MONOCULAR DEPTH ESTIMATION

Method	Batch current/target/source	Protocol 1					Protocol 2					Protocol 3				
		RMSE	Abs Rel	$\delta_1$	$\delta_2$	$\delta_3$	RMSE	Abs Rel	$\delta_1$	$\delta_2$	$\delta_3$	RMSE	Abs Rel	$\delta_1$	$\delta_2$	$\delta_3$
Only target	0 / 0 / 0	6.13	0.15	0.84	0.93	0.96	4.78	0.12	0.88	0.95	0.97	12.22	0.26	0.51	0.82	0.94
Only source	0 / 0 / 0	8.06	0.35	0.43	0.77	0.91	7.90	0.33	0.44	0.77	0.93	<b>10.16</b>	<b>0.19</b>	<b>0.78</b>	<b>0.93</b>	<b>0.97</b>
Online image	1 / 0 / 0	8.33	0.27	0.64	0.84	0.93	6.06	0.33	0.46	0.73	0.90	13.72	0.57	0.30	0.50	0.68
Target replay	1 / 2 / 0	<b>6.35</b>	<b>0.19</b>	<b>0.77</b>	<b>0.91</b>	<b>0.96</b>	<b>5.34</b>	<b>0.15</b>	<b>0.81</b>	<b>0.93</b>	<b>0.97</b>	12.48	0.44	0.34	0.68	0.88
CoDEPS	1 / 2 / 2	<u>7.19</u>	<u>0.22</u>	<u>0.73</u>	<u>0.89</u>	<u>0.94</u>	<u>5.57</u>	<b>0.15</b>	<b>0.81</b>	<b>0.93</b>	<b>0.97</b>	<u>11.38</u>	<u>0.21</u>	<u>0.75</u>	<u>0.91</u>	<u>0.96</u>

The root mean squared error (RMSE), absolute relative error (Abs Rel) as well as accuracies  $\delta_1 = \delta < 1.25$ ,  $\delta_2 = \delta < 1.25^2$ , and  $\delta_3 = \delta < 1.25^3$ , obtained by adapting Cityscapes to sequence 10 of the KITTI-360 dataset. Best results without access to ground truth target data (“only target”) in each category are in **bold**; second best are underlined.

the adaptation progress using unseen validation samples and compare the results to the ground truth. For depth, we visualize predictions generated by the network if it was only trained on  $\mathcal{S}$  and  $\mathcal{T}$ , respectively. For panoptic segmentation, the progressive adaptation on the target domain is particularly visible on the *sidewalk* and *terrain* image regions, which CoDEPS learns to differentiate from the similarly looking classes *road* and *vegetation*. Furthermore, instances become more pronounced, e.g., the cyclist in the right sample. Despite the enhancements on the target domain, CoDEPS successfully maintains its performance on the source domain with only minimum decreases in depth estimation.

*Depth Adaptation:* We present the results for monocular depth estimation in Table III. The first row “only target” shows the theoretical performance on  $\mathcal{T}$  (protocols 1 and 2) if the network would have been trained directly on this domain. Note that such a setup is infeasible in the real world when continuous operation must be guaranteed. The second row “only source” denotes the performance after pretraining on  $\mathcal{S}$  without performing online continual learning. Comparing the

absolute relative error as well as the accuracies  $\delta_1$ ,  $\delta_2$ , and  $\delta_3$  between these rows reveals the domain gap. Note that the opposite gap can be observed when evaluating on  $\mathcal{S}$  (protocol 3). While continual learning using the current online sample increases the accuracy of protocol 1, it also overfits to the current scene. That is, generalizability to the entire target domain is not achieved as shown by protocol 2. Introducing replay samples from the target buffer overcomes this issue and accounts for online samples of poor quality, improving protocols 1 and 2. However, both of the above result in catastrophic forgetting with respect to  $\mathcal{S}$  (protocol 3). The final CoDEPS adds additional source replay yielding low errors and high accuracy by compromising on both  $\mathcal{S}$  and  $\mathcal{T}$ .

*Panoptic Adaptation:* In Table IV, we also demonstrate the domain gap between  $\mathcal{S}$  and  $\mathcal{T}$  for semantic and panoptic segmentation. Similar to depth estimation, both “only target” and “only source” only perform well on their respective training domain without being able to generalize to the other. We further evaluate CoDEPS by comparing it with two competitive baselines that perform domain adaptation on segmentation

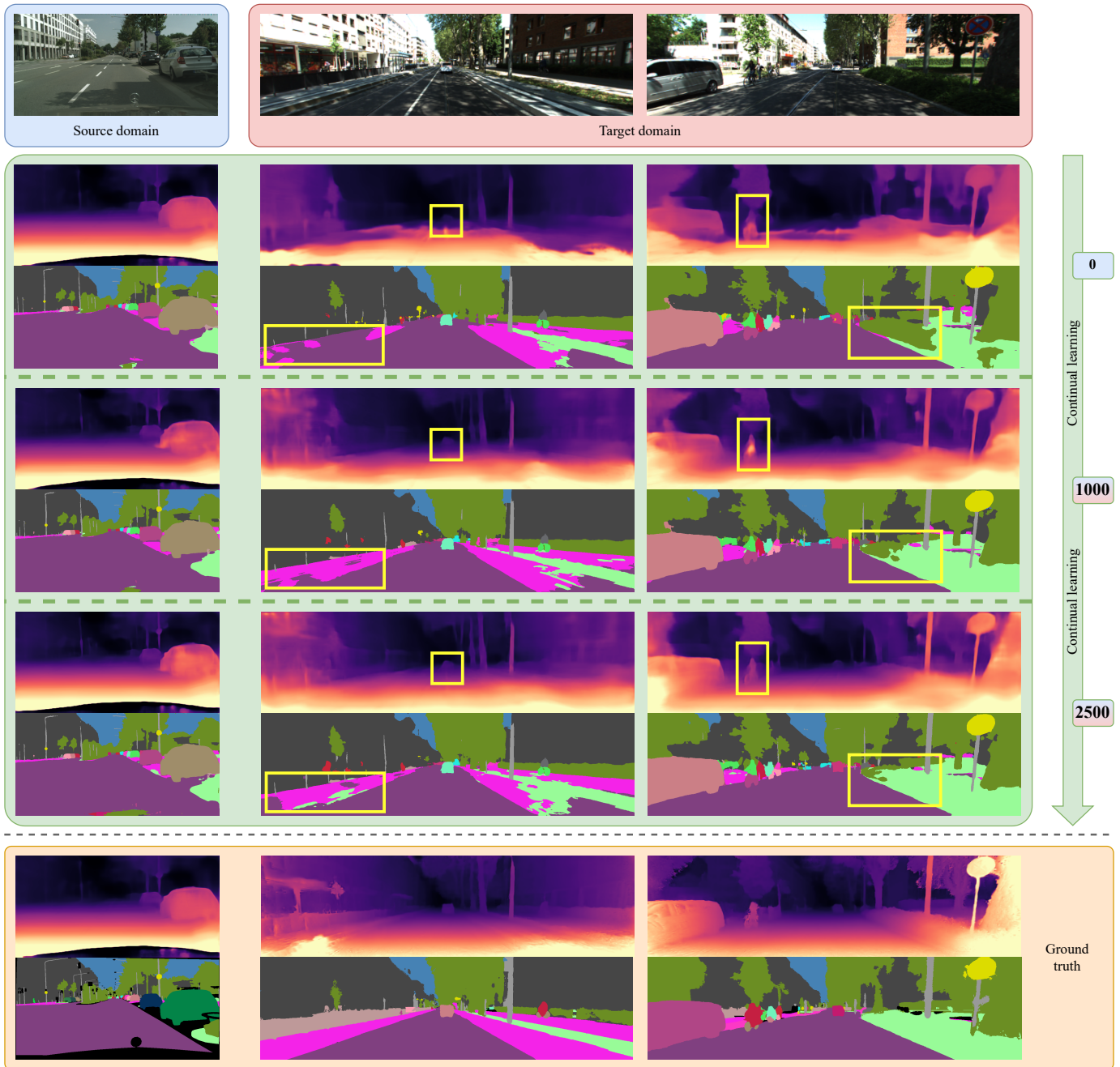


Fig. 4. Qualitative results for Cityscapes to KITTI-360 adaptation after pretraining on the source, i.e., 0 steps, and after having seen 1,000 and 2,500 frames. As shown in the left column, CoDEPS is able to avoid catastrophic forgetting on the source domain. The progressive adaptation on the target domain is particularly visible in the image areas highlighted by yellow boxes. “Stuff” classes of similar appearance like *sidewalk* vs. *road* (left image) and *terrain* vs. *vegetation* (right image) can be better distinguished by CoDEPS. Furthermore, instances become more pronounced as can be observed for the highlighted car (left image) and the cyclist (right image).



TABLE IV  
CONTINUAL LEARNING FOR PANOPTIC SEGMENTATION

Method	Protocol 1				Protocol 2				Protocol 3			
	mIoU	PQ	SQ	RQ	mIoU	PQ	SQ	RQ	mIoU	PQ	SQ	RQ
Only target	64.65	41.91	76.68	51.48	55.12	36.58	67.41	46.15	46.33	28.24	70.03	36.92
Only source	51.94	32.60	71.27	42.44	45.74	30.62	69.56	39.49	<u>72.87</u>	49.19	<u>77.45</u>	60.40
GUDA [13]	45.56	29.70	70.67	39.05	47.62	31.03	64.00	40.49	66.57	44.39	75.95	55.32
DACS [31]	51.14	32.09	71.12	42.23	45.24	29.05	69.47	38.11	72.66	49.27	77.33	60.60
CoDEPS (online image)	<b>53.22</b>	<u>33.46</u>	<b>71.63</b>	<u>43.46</u>	<u>49.51</u>	31.49	64.17	<u>40.71</u>	72.81	<b>49.83</b>	77.25	<b>61.49</b>
CoDEPS (random sampling)	52.36	33.24	<u>71.60</u>	43.25	48.78	<u>31.50</u>	<u>68.83</u>	40.56	72.05	49.11	77.18	60.52
CoDEPS	<u>53.02</u>	<b>33.50</b>	71.53	<b>43.62</b>	<b>49.91</b>	<b>31.91</b>	<b>70.68</b>	<b>40.95</b>	<b>72.90</b>	<u>49.76</u>	<b>77.49</b>	<u>61.22</u>

The mean intersection over union (mIoU), panoptic quality (PQ), semantic quality (SQ), and recognition quality (RQ) are obtained by adapting Cityscapes to sequence 10 of the KITTI-360 dataset. Best results without access to ground truth target data (“only target”) in each category are in **bold**; second best are underlined.

tasks: GUDA [13], which combines semantic segmentation and depth estimation, rendering their task comparable to ours, and DACS [31], which employs a class-mix strategy for offline domain adaptation of semantic segmentation. To ensure a fair comparison, both baselines are evaluated using the same settings as CoDEPS, including diversity sampling-based experience replay. The results in Table IV indicate that both approaches lead to a significant performance decrease across all three protocols. GUDA’s reliance on self-supervised feature alignment using depth training is not effective in the continual learning setting, as shown in the results. DACS also suffers from a decline in performance, likely due to the strong intervention of its mixing strategy into the pretrained network, which can already produce reasonable predictions on the target domain without adaptation.

These results imply that traditional approaches from offline sim-to-real adaptation may not perform well in the online continual learning scenario. To further assess the impact of target replay and our diversity-based buffer sampling, we selectively deactivate both components. Applying the proposed cross-domain mixing strategy results in an improvement in protocol 1. However, similar to depth adaptation, the results are not fully generalizable to the entire target domain, e.g., SQ of protocol 2. Instead of diversity-based sampling, we use random sampling when both creating the source buffer and when updating the target buffer. Compared to CoDEPS, the performance heavily degrades demonstrating the efficacy of the sampling method. Finally, we present the classwise evaluation of the segmentation performance in Table V, which demonstrates improvements of CoDEPS in the IoU metrics for most classes. In particular, we observe significant enhancements of the *two-wheeler* and *terrain* classes. The latter can also be observed in Fig. 4. In fact, CoDEPS outperforms even the model trained directly on the target domain using ground truth supervision for the latter class.

#### D. Ablation Study of the Replay Buffer

We extensively study different sizes of the replay buffer and the effect of diversity sampling as explained in Sec. III-B. We list our results in Table VI. Note that an infinite replay buffer contains 2,975 source and a maximum of 2,683 target samples in the employed setting, i.e., adapting from Cityscapes

TABLE V  
CLASSWISE EVALUATION

Class		Only target	Only source	CoDEPS
Stuff	Road	93	89	91
	Sidewalk	40	32	37
	Building	88	85	85
	Fence	43	14	22
	Pole	35	29	32
	Traffic sign	40	35	38
	Vegetation	78	73	75
	Terrain	54	21	39
Sky	82	79	81	
Thing	Person	47	38	38
	Rider	47	29	36
	Car	91	83	84
	Truck	1	4	2
	Two-wheeler	33	27	38
Mean	55.1	45.7	49.9	

The classwise mIoU is based on protocol 2 in Table IV. We compare CoDEPS against two baselines that were trained using source (“only source”) or target data (“only target”), respectively. CoDEPS provides a significant performance boost of 4.2% in terms of the mIoU metric.

*train* to KITTI-360 using sequence 10 *adapt*. Generally, a larger replay buffer yields higher performance with respect to both adaptation capability and avoiding catastrophic forgetting. Additionally, the proposed diversity sampling using semantic classes for the source and image features for the target samples increases the performance throughout the experiments. However, a greater buffer size increases the required storage posing a challenge for real-world deployment. Based on the presented results, we select a buffer size of 300 with active diversity sampling as for smaller sizes the performance of semantic segmentation on the target domain degrades.

#### E. Continual Adaptation

Finally, we evaluate the performance of CoDEPS in the context of multi-domain adaptation, i.e.,  $S \rightarrow \mathcal{T}_1 \rightarrow \mathcal{T}_2$ . In particular, we first adapt to sequence 10 of KITTI-360 followed by sequence 08 of SemKITTI-DVPS, then we invert the adaptation order. To analyze forward and backward transfer as defined for continual learning [23], we compute the metrics on the *val* split of the source and the *adapt* parts of the respective target domains. We report the results in Table VII. Note that we use  $\alpha_{S \rightarrow \mathcal{T}_1} = 0.9$  and  $\alpha_{\mathcal{T}_1 \rightarrow \mathcal{T}_2} = 0.7$  for updating the EMA model according to Eq. 10 since the network should adapt more strongly when deployed to  $\mathcal{T}_2$  due

TABLE VI  
ABLATION STUDY ON THE REPLAY BUFFER

Size	Div.	Protocol 2						Protocol 3					
		mIoU $\uparrow$	PQ $\uparrow$	SQ $\uparrow$	RQ $\uparrow$	RMSE $\downarrow$	Abs Rel $\downarrow$	mIoU $\uparrow$	PQ $\uparrow$	SQ $\uparrow$	RQ $\uparrow$	RMSE $\downarrow$	Abs Rel $\downarrow$
$\infty$		49.15	31.95	69.08	40.96	4.94	0.15	73.25	50.37	77.77	61.87	10.76	0.21
1000		49.11 $\pm$ 0.69	31.85 $\pm$ 0.25	66.82 $\pm$ 3.06	40.93 $\pm$ 0.06	5.04 $\pm$ 0.01	0.14 $\pm$ 0.00	72.84 $\pm$ 0.33	49.93 $\pm$ 0.20	77.51 $\pm$ 0.05	61.39 $\pm$ 0.28	11.35 $\pm$ 0.39	0.22 $\pm$ 0.01
1000	$\checkmark$	49.36	31.83	68.89	<b>41.01</b>	5.30	0.15	<b>73.50</b>	<b>50.05</b>	<b>77.67</b>	<b>61.48</b>	12.06	0.23
500		48.77 $\pm$ 0.39	31.54 $\pm$ 0.39	67.39 $\pm$ 2.16	40.66 $\pm$ 0.54	5.20 $\pm$ 0.20	0.15 $\pm$ 0.00	72.38 $\pm$ 0.26	49.48 $\pm$ 0.14	77.45 $\pm$ 0.16	60.90 $\pm$ 0.23	11.14 $\pm$ 0.54	0.22 $\pm$ 0.01
500	$\checkmark$	<u>49.56</u>	31.83	70.11	40.96	5.55	0.16	72.78	49.68	77.39	61.10	<u>11.30</u>	0.22
300		48.78 $\pm$ 0.05	31.50 $\pm$ 0.19	68.83 $\pm$ 2.31	40.56 $\pm$ 0.15	5.27 $\pm$ 0.16	0.15 $\pm$ 0.00	72.05 $\pm$ 0.25	49.11 $\pm$ 0.30	77.18 $\pm$ 0.07	60.52 $\pm$ 0.35	11.14 $\pm$ 0.22	0.22 $\pm$ 0.01
300	$\checkmark$	<b>49.91</b>	<b>31.91</b>	<b>70.68</b>	40.95	5.57	0.15	<u>72.90</u>	49.76	77.49	61.22	11.38	<b>0.21</b>
100		48.27 $\pm$ 0.84	30.71 $\pm$ 0.41	63.95 $\pm$ 0.41	39.79 $\pm$ 0.38	5.83 $\pm$ 0.15	0.16 $\pm$ 0.00	69.75 $\pm$ 1.77	47.94 $\pm$ 0.95	76.66 $\pm$ 0.25	59.39 $\pm$ 1.16	10.86 $\pm$ 0.56	0.22 $\pm$ 0.02
100	$\checkmark$	48.40	30.85	64.07	39.95	5.31	0.16	72.35	48.81	77.16	60.25	11.71	0.22
25		46.03 $\pm$ 1.03	29.62 $\pm$ 0.37	66.10 $\pm$ 2.26	38.48 $\pm$ 0.45	5.25 $\pm$ 0.26	0.14 $\pm$ 0.01	67.23 $\pm$ 0.85	45.90 $\pm$ 0.66	75.69 $\pm$ 0.38	57.21 $\pm$ 0.76	11.81 $\pm$ 0.22	0.22 $\pm$ 0.01
25	$\checkmark$	46.35	29.73	63.35	38.58	5.62	0.17	68.84	46.34	76.06	57.78	12.51	0.24

The numbers above are obtained by adapting Cityscapes to sequence 10 of the KITTI-360 dataset. Here, an infinite buffer size equals 2,975 source samples and a maximum of 2,683 target samples. Note that the effective size is two times the shown value as it refers to both source and target replay. The term ‘‘Div.’’ refers to diversity sampling. Where diversity sampling is not used, the same experiment is repeated three times with different random seeds to ensure a statistically reliable measure of performance. The results of these experiments are presented as the mean and standard deviation. Best results in each category are in **bold**; second best are underlined.

TABLE VII  
CONTINUAL LEARNING ON MULTIPLE DOMAINS

Domain	mIoU	PQ	SQ	RQ	RMSE	Abs Rel	mIoU	PQ	SQ	RQ	RMSE	Abs Rel	mIoU	PQ	SQ	RQ	RMSE	Abs Rel
	→ Pretraining on Cityscapes						→ Adaptation on KITTI-360						→ Adaptation on SemKITTI-DVPS					
Cityscapes	72.87	49.19	77.45	60.40	10.16	0.19	72.90	49.76	77.49	61.22	11.38	0.21	72.42	48.74	77.08	60.20	10.65	0.21
KITTI-360 seq. 10	45.74	30.62	69.56	39.49	7.90	0.33	49.91	31.91	70.68	40.95	5.57	0.15	49.26	32.32	64.08	40.95	5.23	0.15
SemKITTI-DVPS seq. 08	51.95	45.24	76.07	57.20	6.17	0.34	49.48	43.26	74.24	57.26	5.60	0.21	53.70	46.50	76.53	59.43	4.32	0.16
	→ Pretraining on Cityscapes						→ Adaptation on SemKITTI-DVPS						→ Adaptation on KITTI-360					
Cityscapes	72.87	49.19	77.45	60.40	10.16	0.19	72.75	49.01	77.36	60.35	10.82	0.22	72.51	48.87	76.98	60.28	11.41	0.21
KITTI-360 seq. 10	45.74	30.62	69.56	39.49	7.90	0.33	49.26	31.66	70.26	41.40	6.30	0.17	50.05	31.92	70.50	41.48	5.47	0.16
SemKITTI-DVPS seq. 08	51.95	45.24	76.07	57.20	6.17	0.34	52.31	44.29	75.58	56.87	4.56	0.16	53.83	47.29	76.55	60.01	4.25	0.16

CoDEPS is continually applied to three domains using Cityscapes as the initial source domain and then adapting to KITTI-360 and SemKITTI-DVPS. The listed numbers on the target domains are based on protocol 2.

to the larger amount of previously seen data. As shown in the first row of both adaptation orders, CoDEPS is able to mitigate catastrophic forgetting with respect to  $\mathcal{S}$  maintaining its performance. We make a similar observation when re-evaluating  $\mathcal{T}_1$  after the second adaptation step to  $\mathcal{T}_2$ . In particular, CoDEPS achieves positive backward transfer on SemKITTI-DVPS when adapting to KITTI-360. On the same adaptation order, we observe positive forward transfer for KITTI-360, i.e., the performance increases although CoDEPS was only adapted to SemKITTI-DVPS.

In Fig. 5, we illustrate the evolution of the performance metrics on SemKITTI-DVPS sequence 08 during adaptation (protocol 1). We compare the error without adaptation to directly adapting to SemKITTI-DVPS versus first adapting to KITTI-360. For both semantic segmentation and depth estimation, it can be clearly observed that the performance improves if more images have been seen. Additionally, adapting first to KITTI-360 results in a large performance increase for both semantic and panoptic segmentation. We account this to the fact that KITTI-360 sequence 10 leads to strongly improved performance, shown in Table VII, that can be transferred to the SemKITTI-DVPS domain.

## V. CONCLUSION

In this paper, we present CoDEPS as the first approach for online continual learning for joint monocular depth estimation and panoptic segmentation. CoDEPS enables the vision system of a robotic platform to continually enhance its performance

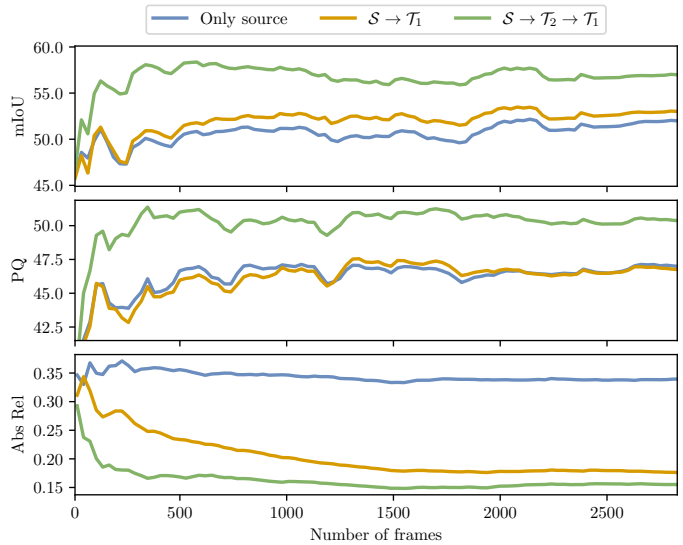


Fig. 5. Evolution of performance metrics on SemKITTI-DVPS sequence 08 during adaptation (protocol 1). The metrics are averaged until the given frame number. The target domains  $\mathcal{T}_1$  and  $\mathcal{T}_2$  refer to SemKITTI-DVPS and KITTI-360, respectively. It can be seen that there is positive forward transfer when first adapting on  $\mathcal{T}_2$ .

in an online fashion. In particular, we propose a new cross-domain mixing strategy to adapt panoptic segmentation combining annotated source data with unlabeled images from a target domain. To mitigate catastrophic forgetting, CoDEPS leverages experience replay using a buffer composed of source and target samples. We explicitly address the limited storage

capacity of robotic platforms by setting a fixed size for the replay buffer. To ensure distinct replay samples, we use rare class sampling on the source set and employ image-based diversity sampling when updating the target buffer. Using extensive evaluations, we demonstrate that CoDEPS outperforms competitive baselines while avoiding catastrophic forgetting in the online continual learning setting. Future work will explore cross-task synergies and the use of pretext tasks for domain adaptation.

#### ACKNOWLEDGMENT

This work was partly funded by the European Union’s Horizon 2020 research and innovation program under grant agreement No 871449-OpenDR and the Bundesministerium für Bildung und Forschung (BMBF) under grant agreement No FKZ 16ME0027.

#### REFERENCES

- [1] Rahaf Aljundi, Min Lin, Baptiste Goujaud, and Yoshua Bengio. Gradient based sample selection for online continual learning. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett, editors, *Conference on Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [2] Jens Behley, Martin Garbade, Andres Milioto, Jan Quenzel, Sven Behnke, Cyrill Stachniss, and Jurgen Gall. SemanticKITTI: A dataset for semantic scene understanding of LiDAR sequences. In *International Conference on Computer Vision*, 2019.
- [3] Borna Bešić and Abhinav Valada. Dynamic object removal and spatio-temporal RGB-D inpainting via geometry-aware adversarial learning. *IEEE Transactions on Intelligent Vehicles*, 7(2):170–185, 2022.
- [4] Borna Bešić, Nikhil Gosala, Daniele Cattaneo, and Abhinav Valada. Unsupervised domain adaptation for LiDAR panoptic segmentation. *IEEE Robotics and Automation Letters*, 7(2):3404–3411, 2022.
- [5] Vincent Casser, Soeren Pirk, Reza Mahjourian, and Anelia Angelova. Unsupervised monocular depth and ego-motion learning with structure and semantics. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2019.
- [6] Bowen Cheng, Maxwell D. Collins, Yukun Zhu, Ting Liu, Thomas S. Huang, Hartwig Adam, and Liang-Chieh Chen. Panoptic-DeepLab: A simple, strong, and fast baseline for bottom-up panoptic segmentation. In *IEEE/CVF Conference Computer Vision and Pattern Recognition*, pages 12472–12482, 2020.
- [7] Gong Cheng and James H. Elder. VCSeg: Virtual camera adaptation for road segmentation. In *IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1969–1978, 2022.
- [8] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The Cityscapes dataset for semantic urban scene understanding. In *IEEE/CVF Conference Computer Vision and Pattern Recognition*, pages 3213–3223, 2016.
- [9] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for Autonomous Driving? The KITTI Vision Benchmark Suite. In *IEEE/CVF Conference Computer Vision and Pattern Recognition*, pages 3354–3361, 2012.
- [10] Clément Godard, Oisín Mac Aodha, and Gabriel J Brostow. Unsupervised monocular depth estimation with left-right consistency. In *IEEE/CVF Conference Computer Vision and Pattern Recognition*, pages 270–279, 2017.
- [11] Clément Godard, Oisín Mac Aodha, Michael Firman, and Gabriel Brostow. Digging into self-supervised monocular depth estimation. In *International Conference on Computer Vision*, pages 3827–3837, 2019.
- [12] Nikhil Gosala and Abhinav Valada. Bird’s-eye-view panoptic segmentation using monocular frontal view images. *IEEE Robotics and Automation Letters*, 7(2):1968–1975, 2022.
- [13] Vitor Guizilini, Jie Li, Rareş Ambruş, and Adrien Gaidon. Geometric unsupervised domain adaptation for semantic segmentation. In *International Conference on Computer Vision*, pages 8537–8547, 2021.
- [14] Kaifeng He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE/CVF Conference Computer Vision and Pattern Recognition*, pages 770–778, 2016.
- [15] Lukas Hoyer, Dengxin Dai, and Luc Van Gool. DAFormer: Improving network architectures and training strategies for domain-adaptive semantic segmentation. In *IEEE/CVF Conference Computer Vision and Pattern Recognition*, pages 9924–9935, 2022.
- [16] Jiaying Huang, Dayan Guan, Aoran Xiao, and Shijian Lu. Cross-view regularization for domain adaptive panoptic segmentation. In *IEEE/CVF Conference Computer Vision and Pattern Recognition*, pages 10133–10144, 2021.
- [17] Marvin Klingner, Mouadh Ayache, and Tim Fingscheidt. Continual batchnorm adaptation (CBNA) for semantic segmentation. *IEEE Transactions on Intelligent Transportation Systems*, 23(11):20899–20911, 2022.
- [18] Yevhen Kuznietsov, Marc Proesmans, and Luc Van Gool. CoMoDA: Continuous monocular depth adaptation using past experiences. In *IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2907–2917, 2021.
- [19] Yevhen Kuznietsov, Marc Proesmans, and Luc Van Gool. Towards unsupervised online domain adaptation for semantic segmentation. In *European Conference on Computer Vision*, pages 261–271, 2022.
- [20] Hanhan Li, Ariel Gordon, Hang Zhao, Vincent Casser, and Anelia Angelova. Unsupervised monocular depth learning in dynamic scenes. In *Conference on Robot Learning*, pages 1908–1917. PMLR, 2021.
- [21] Shunkai Li, Xin Wang, Yingdian Cao, Fei Xue, Zike Yan, and Hongbin Zha. Self-supervised deep visual odometry with online adaptation. In *IEEE/CVF Conference Com-*

- puter Vision and Pattern Recognition*, pages 6339–6348, 2020.
- [22] Yiyi Liao, Jun Xie, and Andreas Geiger. KITTI-360: A novel dataset and benchmarks for urban scene understanding in 2d and 3d. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2022.
- [23] David Lopez-Paz and Marc' Aurelio Ranzato. Gradient episodic memory for continual learning. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Conference on Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- [24] Adrian Lopez-Rodriguez and Krystian Mikolajczyk. DESC: Domain adaptation for depth estimation via semantic consistency. *International Journal of Computer Vision*, 131(3):752–771, Mar 2023.
- [25] Giulio Mattolin, Luca Zanella, Elisa Ricci, and Yiming Wang. ConfMix: Unsupervised domain adaptation for object detection via confidence-based mixing. In *IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 423–433, January 2023.
- [26] Robert McCraith, Lukas Neumann, Andrew Zisserman, and Andrea Vedaldi. Monocular depth estimation with self-supervised instance adaptation. *arXiv preprint arXiv:2004.05821*, 2020.
- [27] Rohit Mohan and Abhinav Valada. Amodal panoptic segmentation. In *IEEE/CVF Conference Computer Vision and Pattern Recognition*, pages 20991–21000, 2022.
- [28] Rohit Mohan and Abhinav Valada. Perceiving the invisible: Proposal-free amodal panoptic segmentation. *IEEE Robotics and Automation Letters*, 7(4):9302–9309, 2022.
- [29] Viktor Olsson, Wilhelm Tranheden, Juliano Pinto, and Lennart Svensson. ClassMix: Segmentation-based data augmentation for semi-supervised learning. In *IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1368–1377, 2021.
- [30] Siyuan Qiao, Yukun Zhu, Hartwig Adam, Alan Yuille, and Liang-Chieh Chen. ViP-DeepLab: Learning visual perception with depth-aware video panoptic segmentation. In *IEEE/CVF Conference Computer Vision and Pattern Recognition*, pages 3996–4007, 2021.
- [31] Wilhelm Tranheden, Viktor Olsson, Juliano Pinto, and Lennart Svensson. DACS: Domain adaptation via cross-domain mixed sampling. In *IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1378–1388, 2021.
- [32] Abhinav Valada, Gabriel Oliveira, Thomas Brox, and Wolfram Burgard. Towards robust semantic segmentation using deep fusion. In *Robotics: Science and Systems Workshop, Are the Sceptics Right*, 2016.
- [33] Niclas Vödisch, Daniele Cattaneo, Wolfram Burgard, and Abhinav Valada. Continual SLAM: Beyond lifelong simultaneous localization and mapping through continual learning. In Aude Billard, Tamim Asfour, and Oussama Khatib, editors, *Robotics Research*, pages 19–35, Cham, 2023. Springer Nature Switzerland.
- [34] Niclas Vödisch, Daniele Cattaneo, Wolfram Burgard, and Abhinav Valada. CoVIO: Online continual learning for visual-inertial odometry. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2023.
- [35] Qin Wang, Olga Fink, Luc Van Gool, and Dengxin Dai. Continual test-time domain adaptation. In *IEEE/CVF Conference Computer Vision and Pattern Recognition*, pages 7201–7211, 2022.
- [36] Nan Yang, Lukas von Stumberg, Rui Wang, and Daniel Cremers. D3VO: Deep depth, deep pose and deep uncertainty for monocular visual odometry. In *IEEE/CVF Conference Computer Vision and Pattern Recognition*, pages 1281–1292, 2020.
- [37] Zhenyu Zhang, Stéphane Lathuilière, Elisa Ricci, Nicu Sebe, Yan Yan, and Jian Yang. Online depth learning against forgetting in monocular videos. In *IEEE/CVF Conference Computer Vision and Pattern Recognition*, pages 4493–4502, 2020.
- [38] Tinghui Zhou, Matthew Brown, Noah Snavely, and David G Lowe. Unsupervised learning of depth and ego-motion from video. In *IEEE/CVF Conference Computer Vision and Pattern Recognition*, pages 1851–1858, 2017.