# Unsupervised Continual Semantic Adaptation through Neural Rendering

Zhizheng Liu[1*]  Francesco Milano[1*]  Jonas Frey[1,2]  Marco Hutter[1]  Roland Siegwart[1]
Hermann Blum[1†]  Cesar Cadena[1†]
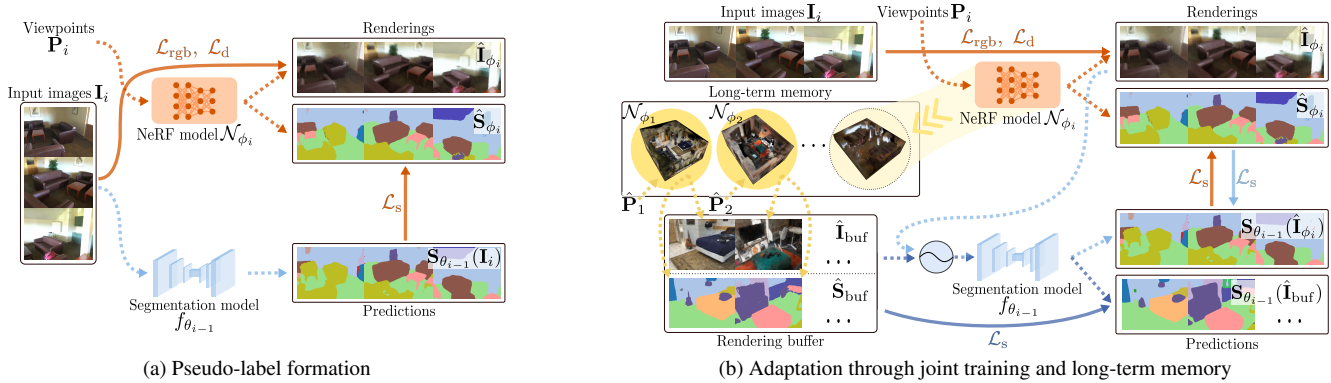[1]ETH Zurich  [2]Max Planck ETH Center for Learning Systems

Figure 1. We propose a method to continually adapt a semantic segmentation model $f$ in an unsupervised fashion across multiple scenes, using neural rendering. For each scene $\mathcal{S}_i$: a) RGB(-D) images $\mathbf{I}_i$ from multiple viewpoints $\mathbf{P}_i$ and their corresponding predictions $\mathbf{S}_{\theta_{i-1}}(\mathbf{I}_i)$ by the latest model $f_{\theta_{i-1}}$ are used to supervise a (Semantic-)NeRF model $\mathcal{N}_{\phi_i}$; b) Adaptation on $\mathcal{S}_i$ is performed through a *joint training*, in which the segmentation network is supervised using the 3D-aware, view-consistent pseudo-labels $\hat{\mathbf{S}}_{\phi_i}$ rendered from $\mathcal{N}_{\phi_i}$ and the NeRF model through the smooth predictions of $f_{\theta_{i-1}}$. For each scene, the NeRF model can be compactly stored in a long-term memory, from which images and pseudo-labels from arbitrary viewpoints $\hat{\mathbf{P}}$ can be rendered into a fixed-size rendering buffer and mixed with the renderings from the current scene to reduce forgetting. Bold and dotted lines denote supervision signals and inputs/outputs, respectively.

## Abstract

*An increasing amount of applications rely on data-driven models that are deployed for perception tasks across a sequence of scenes. Due to the mismatch between training and deployment data, adapting the model on the new scenes is often crucial to obtain good performance. In this work, we study continual multi-scene adaptation for the task of semantic segmentation, assuming that no ground-truth labels are available during deployment and that performance on the previous scenes should be maintained. We propose training a Semantic-NeRF network for each scene by fusing the predictions of a segmentation model and then using the view-consistent rendered semantic labels as pseudo-labels to adapt the model. Through joint training with the segmentation model, the Semantic-NeRF model effectively enables 2D-3D knowledge transfer. Furthermore, due to its compact size, it can be stored in a long-term memory and subsequently used to render data from arbitrary viewpoints to reduce forgetting. We evaluate our approach on Scan-Net, where we outperform both a voxel-based baseline and a state-of-the-art unsupervised domain adaptation method.*

## 1. Introduction

Data-driven models trained for perception tasks play an increasing role in applications that rely on scene understanding, including, *e.g.*, mixed reality and robotics. When deploying these models on real-world systems, however, mismatches between the data used for training and those encountered during deployment can lead to poor performance, prompting the need for an *adaptation* of the models to the new environment. Oftentimes, the supervision data required for this adaptation can only be obtained through a laborious labeling process. Furthermore, even when such data are available, a naïve adaptation to the new environment results in decreased performance on the original training data, a phenomenon known as *catastrophic forgetting* [21, 28].

In this work, we focus on the task of adapting a semantic segmentation network across multiple indoor scenes, under the assumption that no labeled data from the new environment are available. Similar settings are explored in the literature in the areas of *unsupervised domain adaptation (UDA)* [28, 46] and *continual learning (CL)* [21]. However, works in the UDA literature usually focus on a sin-

*Authors share first authorship.     † Authors share senior authorship.

gle source-to-target transfer where the underlying assumption is that the data from both the source and the target domain are available all at once in the respective training stage, and often study the setting in which the knowledge transfer happens between a synthetic and a real environment [6, 38, 39, 46]. On the other hand, the CL community, which generally explores the adaptation of networks across different *tasks*, has established the *class-incremental* setting as the standard for semantic segmentation, in which new classes are introduced across different scenes from the same domain and ground-truth supervision is provided [28]. In contrast, we propose to study network adaptation in a setting that more closely resembles the deployment of semantic networks on real-world systems. In particular, instead of assuming that data from a specific domain are available all at once, we focus on the scenario in which the network is sequentially deployed in multiple *scenes* from a real-world indoor environment (we use the ScanNet dataset [7]), and therefore has to perform multiple *stages* of adaptation from one scene to another. Our setting further includes the possibility that previously seen scenes may be revisited. Hence, we are interested in achieving high prediction accuracy on each new scene, while at the same time preserving performance on the previous ones. Note that unlike the better explored setting of class-incremental CL, in this setting, we assume a *closed set* of semantic categories, but tackle the covariate shift across scenes without the need for ground-truth labels. We refer to this setting as *continual semantic adaptation*.

In this work, we propose to address this adaptation problem by leveraging advances in neural rendering [32]. Specifically, in a similar spirit to [12], when deploying a pre-trained network in a new scene, we aggregate the semantic predictions from the multiple viewpoints traversed by the agent into a 3D representation, from which we then render pseudo-labels that we use to adapt the network on the current scene. However, instead of relying on a voxel-based representation, we propose to aggregate the predictions through a semantics-aware NeRF [32, 56]. This formulation has several advantages. First, we show that using NeRFs to aggregate the semantic predictions results in higher-quality pseudo-labels compared to the voxel-based method of [12]. Moreover, we demonstrate that using these pseudo-labels to adapt the segmentation network results in superior performance compared both to [12] and to the state-of-the-art UDA method CoTTA [51]. An even more interesting insight, however, is that due the differentiability of NeRF, we can jointly train the frame-level semantic network and the scene-level NeRF to enforce similarity between the predictions of the former and the renderings of the latter. Remarkably, this joint procedure induces better performance of both labels, showing the benefit of mutual 2D-3D knowledge transfer.

A further benefit of our method is that after adapting to a new scene, the NeRF encoding the appearance, geometry and semantic content for that scene can be compactly saved in long-term storage, which effectively forms a "memory bank" of the previous experiences and can be useful in reducing catastrophic forgetting. Specifically, by mixing pairs of semantic and color NeRF renderings from a small number of views in the previous scenes and from views in the current scene, we show that our method is able to outperform both the baseline of [12] and CoTTA [51] on the adaptation to the new scene and in terms of knowledge retention on the previous scenes. Crucially, the collective size of the NeRF models is lower than that of the explicit replay buffer required by [12] and of the teacher network used in CoTTA [51] up to several dozens of scenes. Additionally, each of the NeRF models stores a potentially infinite number of views that can be used for adaptation, not limited to the training set as in [12], and removes the need to explicitly keep color images and pseudo-labels in memory.

In summary, the main contributions of our work are the following: (i) We propose using NeRFs to adapt a semantic segmentation network to new scenes. We find that enforcing 2D-3D knowledge transfer by jointly adapting NeRF and the segmentation network on a given scene results in a consistent performance improvement; (ii) We address the problem of continually adapting the segmentation network across a sequence of scenes by compactly storing the NeRF models in a long-term memory and mixing rendered images and pseudo-labels from previous scenes with those from the current one. Our approach allows generating a potentially infinite number of views to use for adaptation at constant memory size for each scene; (iii) Through extensive experiments, we show that our method achieves better adaptation and performance on the previous scenes compared both to a recent voxel-based method that explored a similar setting [12] and to a state-of-the-art UDA method [51].

## 2. Related work

**Unsupervised domain adaptation for semantic segmentation.** Unsupervised domain adaptation (UDA) studies the problem of transferring knowledge between a source and a target domain under the assumption that no labeled data for the target domain are available. In the following, we provide an overview of the main techniques used in UDA for semantic segmentation and focus on those which are most closely related to our work; for a more extensive summary we refer the reader to the recent survey of [28].

The majority of the methods rely on auto-encoder CNN architectures, and perform network adaptation either at the level of the input data [3, 15, 23, 53, 54], of the intermediate network representations [4, 11, 15, 34, 54], or of the output predictions [3, 4, 11, 23, 27, 40, 41, 44, 49, 58, 59]. The main strategies adopted consist in: using adversarial

learning techniques to enforce that the network representations have similar statistical properties across the two domains [3, 4, 11, 15, 23, 34, 41], performing image-to-image translation to align the data from the two domains [3, 15, 23, 53, 54], learning to detect non-discriminative feature representations for the target domain [20, 40], and using self-supervised learning based either on minimizing the pixel-level entropy in the target domain [49] or on self-training techniques [5, 23, 27, 44, 55, 58, 59]. The latter category of methods is the most related to our setting. In particular, a number of works use the network trained on the source domain to generate semantic predictions on the unlabeled target data; the obtained *pseudo-labels* are then used as a self-supervisory learning signal to adapt the network to the target domain. While our work and the self-training UDA methods both use pseudo-labels, the latter approaches neither exploit the sequential structure of the data nor explicitly enforce multi-view consistency in the predictions on the target data. Furthermore, approaches in UDA mostly focus on single-stage, sim-to-real transfer settings, often for outdoor environments, and generally assume that the data from each domain are available all at once during the respective training stage. In contrast, we focus on a multi-step adaptation problem, in which data from multiple scenes from an indoor environment are available sequentially.

Within the category of self-training methods, a number of works come closer to our setting by presenting techniques to achieve *continuous*, multi-stage domain adaptation. In particular, the recently proposed CoTTA [51] uses a student-teacher framework, in which the student network is adapted to a target environment through pseudo-labels generated by the teacher network, and stochastic restoration of the weights from a pre-trained model is used to preserve source knowledge. ACE [52] proposes an input-level adaptation based on style transfer and counteracts forgetting by storing statistics that encode the style of previous domains. The method assumes access to the ground-truth source labels and due to its style-transfer approach it is only evaluated across different environmental conditions within the same (outdoor) scene, thus not being applicable to our setting. Finally, related to our method is also the recent work of Frey et al. [12], which addresses a similar problem as ours by aggregating predictions from different viewpoints in a target domain into a 3D voxel grid and rendering pseudo-labels, but does not perform multi-stage adaptation.

**Continual learning for semantic segmentation.** Continual learning for semantic segmentation (CSS) focuses on the problem of updating a segmentation network in a *class-incremental setting*, in which it is assumed that the domain is available in different *tasks* and that new classes are added over time in a sequential fashion [28]. The main objective consists in performing adaptation to the new task, mostly using only data from the current stage, while preventing

forgetting of the knowledge from the previous tasks. The methods proposed in the literature typically adopt a combination of different strategies, including distilling knowledge from a previous model [1, 10, 29, 31], selectively freezing the network parameters [29, 31], enforcing regularization of the latent representations [30], and generating or crawling data from the internet to replay [26, 35]. While similarly to CSS methods we explore a continual setting in which the network is sequentially presented with data from the same domain, we do not explore the class-incremental problem, and instead focus on a closed-set scenario with shifting distribution of classes and scene appearance. A further important difference is that while CSS methods assume each adaptation step to be supervised, in our setting no ground-truth labels from the current adaptation stage are available.

**NeRF-based semantic learning.** Since the introduction of NeRF [32], several works have proposed extensions to the framework to incorporate semantic information into the learned scene representation. Semantic-NeRF [56] first proposed jointly learning appearance, geometry, and semantics through an additional multi-layer perceptron (MLP) and by adapting the volume rendering equation to produce semantic logits. Subsequent works have further extended this framework along different directions, including combining NeRF with a feature grid and 3D convolutions to achieve generalization [48], interactively labeling scenes [57], performing panoptic segmentation [13, 19], and using pre-trained Transformer models to supervise few-shot NeRF training [16], edit scene properties [50], or distill knowledge for different image-level tasks [18, 47]. In our work, we rely on Semantic-NeRF, which we use to fuse predictions from a segmentation network and that we jointly train with the latter exploiting differentiability. We include the formed scene representation in a long-term memory and use it to render pseudo-labels to adapt the segmentation network.

## 3. Continual Semantic Adaptation

### 3.1. Problem definition

In our problem setting, which we refer to as *continual semantic adaptation*, we assume we are provided with a segmentation model $f_{\theta_0}$, with parameters $\theta_0$, that was pre-trained on a dataset $\mathcal{P} = (\mathbf{I}_{\mathrm{pre}}, \mathbf{S}^\star_{\mathrm{pre}})$. Here $\mathbf{I}_{\mathrm{pre}}$ is a set of input color images (potentially with associated depth information) and $\mathbf{S}^\star_{\mathrm{pre}}$ are the corresponding pixel-wise ground-truth semantic labels. We aim to adapt $f_{\theta_0}$ across a sequence of $N$ scenes $\mathcal{S}_i$, $i \in \{1, \ldots, N\}$ for each of which a set $\mathbf{I}_i$ of color (and depth) images, are collected from different viewpoints, but no ground-truth semantic labels are available. We assume that the input data $\{\mathbf{I}_{\mathrm{pre}}, \mathbf{I}_1, \ldots, \mathbf{I}_N\}$ originate from similar indoor environments (for instance, we do not consider simultaneously synthetic and real-world data) and that the classes to be predicted by the network belong to a

closed set and are all known from the pre-training. For each scene $\mathcal{S}_i$, $i \in \{1, \dots, N\}$, the objective is to find a set of weights $\theta_i$ of the network, starting from $\theta_{i-1}$, such that the performance of $f_{\theta_i}$ on $\mathcal{S}_i$ is higher than that of $f_{\theta_{i-1}}$. Additionally, it is desirable to preserve the performance of $f_{\theta_i}$ on the previous scenes $\{\mathcal{S}_1, \dots, \mathcal{S}_{i-1}\}$, in other words mitigate catastrophic forgetting.

The proposed setting aims to replicate the scenario of the deployment of a segmentation network on a real-world perception system (for instance a robot, or an augmented reality platform), where multiple sequential experiences are collected across similar scenes, and only limited data of the previous scenes can be stored on an on-board computing unit. During deployment, environments might be revisited over time, rendering the preservation of previously learned knowledge essential for a successful deployment.

### 3.2. Methodology

We present a method to address continual semantic adaptation in a self-supervised fashion (Fig. 1). In the following, $\boldsymbol{I}_i^k$ and $\boldsymbol{P}_i^k$ are the $k$-th RGB(-D) image collected in scene $\mathcal{S}_i$ and its corresponding camera pose, where $k \in \{1, \dots, |\mathbf{I}_i|\}$. We further denote with $\mathbf{S}_\theta(\boldsymbol{I}_i^k)$ the prediction produced by $f_\theta$ for $\boldsymbol{I}_i^{k}$[1]. With a slight abuse of notation, we use $\mathbf{S}_\theta(\mathbf{I}_i)$ in place of $\{\mathbf{S}_\theta(\boldsymbol{I}_i^k), \ \boldsymbol{I}_i^k \in \mathbf{I}_i\}$ and similarly for other quantities that are a function of elements in a set.

For each new scene $\mathcal{S}_i$, we train a Semantic-NeRF [56] model $\mathcal{N}_{\phi_i}$, with learnable parameters $\phi_i$, given for each viewpoint $\boldsymbol{P}_i^k$ the corresponding semantic label $\mathbf{S}_{\theta_j}(\boldsymbol{I}_i^k)$ predicted by a previous version $f_{\theta_j}$, $j < i$, of the segmentation model. From the trained Semantic-NeRF model $\mathcal{N}_{\phi_i}$ we render semantic pseudo-labels $\hat{\mathbf{S}}_{\phi_i}$ and images $\hat{\mathbf{I}}_{\phi_i}$. The key observation at the root of our self-supervised adaptation is that semantic labels should be *multi-view consistent*, since they are constrained by the scene geometry that defines them. While the predictions of $f$ often do not reflect this constraint because they are produced for each input frame independently, the NeRF-based pseudo-labels are by construction multi-view consistent. Inspired by [12], we hypothesize that this consistency constitutes an important prior that can be exploited to guide the adaptation of the network to the scene. Therefore, we use the renderings from $\mathcal{N}_i$ to adapt the segmentation network on scene $\mathcal{S}_i$, by minimizing a cross-entropy loss between the pseudo-labels and the network predictions. Crucially, we can use the NeRF and segmentation network predictions to supervise each other, allowing for joint optimization and adaptation of the two networks, which we find further improves the performance of both models.

To continually adapt the segmentation network $f$ to multiple scenes in a sequence $\mathcal{S}_1 \rightarrow \mathcal{S}_2 \rightarrow \cdots \rightarrow \mathcal{S}_N$ and prevent catastrophic forgetting, we leverage the compact

---

[1]Note that in our experiments $f_\theta$ does not use the depth channel of $\boldsymbol{I}_i^k$.

representation of NeRF by storing the corresponding model weights $\phi_i$ after adaptation in a long-term memory for each scene $\mathcal{S}_i$. Given that a trained NeRF can be queried from any viewpoint, this formulation allows generating for each scene a theoretically infinite number of views for adaptation, at the fixed storage cost given by the size of $\phi_i$. For each previous scene $\mathcal{S}_j$, images $\hat{\mathbf{I}}_{\phi_j}$ and pseudo-labels $\hat{\mathbf{S}}_{\phi_j}$ from both previously seen and novel viewpoints can be rendered and used in an experience replay strategy to mitigate catastrophic forgetting on the previous scenes. An overview of our method is shown in Fig. 1.

**NeRF-based pseudo-labels.** We train for each scene a NeRF [32] model, which implicitly learns the geometry and appearance of the environment from a sparse set of posed images and can be used to render photorealistic novel views. More specifically, we extend the NeRF formulation by adding a semantic head as in Semantic-NeRF [56], and we render semantic labels $\hat{\mathbf{S}}_\phi$ by aggregating through the learned density function the semantic-head predictions for $M$ sample points along each camera ray $\mathbf{r}$, as follows:

$$\hat{\mathbf{S}}_\phi(\mathbf{r}) = \sum_{i=1}^{M} T_i \alpha_i \mathbf{s}_i, \tag{1}$$

where $\alpha_i = 1 - e^{-\sigma_i \delta_i}$, $T_i = \prod_{j=1}^{i-1}(1 - \alpha_j)$, with $\delta_i$ being the distance between adjacent sample points along the ray, and $\sigma_i$ and $\mathbf{s}_i$ representing the predicted density and semantic logits at the $i$-th sample point along the ray, respectively.

We observe that if Semantic-NeRF is directly trained on the labels predicted by a pre-trained segmentation network on a new scene, the lack of view consistency of these labels can severely degrade the quality of the learned geometry, which in turn hurts the performance of the rendered semantic labels. To alleviate the influence of the inconsistent labels on the geometry, we propose to adopt several modifications. First, we stop the gradient flow from the semantic head into the density head. Second, we use depth supervision, as introduced in [8], to regularize the depth $\hat{d}(\mathbf{r}) = \sum_{i=1}^{N} T_i \alpha_i \delta_i$ rendered by NeRF via $\ell_1$ loss with respect to the ground-truth depth $d(\mathbf{r})$:

$$\mathcal{L}_{\mathrm{d}}(\mathbf{r}) = \left\| \hat{d}(\mathbf{r}) - d(\mathbf{r}) \right\|_1. \tag{2}$$

Through ablations in the Supplementary, we show that this choice is particularly effective at improving the quality of both the geometry and the rendered labels. Additionally, we note that since the semantic logits $\mathbf{s}_i$ of each sampled point are unbounded, the logits $\hat{\mathbf{S}}_\phi(\mathbf{r})$ of the ray $\mathbf{r}$ can be dominated by a sampled point with very large semantic logits instead of one that is near the surface of the scene. This could cause the semantic labels generated by the NeRF model to overfit the initial labels of the segmentation model and lose multi-view consistency even when the learned geometry is correct. To address this issue, we instead first apply softmax to the logits of each sampled point, so these are normal-

ized and contribute to the final aggregated logits through the weighting induced by volume rendering, as follows:

$$\hat{\mathbf{S}}'_\phi(\mathbf{r}) = \sum_{i=1}^N T_i \alpha_i \cdot \text{softmax}(\mathbf{s}_i), \ \hat{\mathbf{S}}_\phi(\mathbf{r}) = \hat{\mathbf{S}}'_\phi(\mathbf{r})/\|\hat{\mathbf{S}}'_\phi(\mathbf{r})\|_1. \tag{3}$$

The final normalized $\hat{\mathbf{S}}_\phi(\mathbf{r})$ is then a categorical distribution $(\hat{S}(\mathbf{r})_1, \cdots, \hat{S}(\mathbf{r})_C)$ over the $C$ semantic classes predicted by NeRF, and we use a negative log-likelihood loss to supervise the rendered semantic labels with the predictions of the semantic network:

$$\mathcal{L}_s(\mathbf{r}) = -\sum_{c=1}^C \log(\hat{S}(\mathbf{r})_c) \cdot \mathbb{1}_{c=c(\mathbf{r})}, \tag{4}$$

where $c(\mathbf{r})$ is the semantic label predicted by the segmentation network $f_\theta$. We train the NeRF model by randomly sampling rays from the training views and adding together the losses in (2) and (4), as well as the usual $\ell_2$ loss $\mathcal{L}_{\text{rgb}}(\mathbf{r})$ on the rendered color [32], as follows:

$$\mathcal{L} = \sum_{i=1}^R \mathcal{L}_{\text{rgb}}(\mathbf{r_i}) + w_d \mathcal{L}_d(\mathbf{r_i}) + w_s \mathcal{L}_s(\mathbf{r_i}), \tag{5}$$

where $R$ is the number of rays sampled for each batch and $w_d$, $w_s$ are the weights for the depth loss and the semantic loss, respectively. After training the NeRF model, we render from it both color images $\hat{\mathbf{I}}_\phi$ and semantic labels $\hat{\mathbf{S}}_\phi$, as *pseudo-labels* for adapting the segmentation network.

Being able to quickly fuse the semantic predictions and generate pseudo-labels might be of particular importance in applications that require fast, possibly online adaptation. To get closer to this objective, we adopt the multi-resolution hash encoding proposed in Instant-NGP [33], which significantly improves the training and rendering speed compared to the original NeRF formulation. In the Supplementary, we compare the quality of the Instant-NGP-based pseudo-labels and those obtained with the original implementation from [56], and show that our method is agnostic to the specific NeRF implementation chosen.

**Adaptation through joint 2D-3D training.** To adapt the segmentation network $f_{\theta_j}$ on a given scene $\mathcal{S}_i$ (where $i > j$), we use the rendered pseudo-labels $\hat{\mathbf{S}}_{\phi_i}$ as supervisory signal by optimizing a cross-entropy loss between the network predictions $\mathbf{S}_{\theta_j}$ and $\hat{\mathbf{S}}_{\phi_i}$, similarly to previous approaches in the literature [12, 51, 52]. However, we propose two important modifications enabled by our particular setup and by its end-to-end differentiability. First, rather than adapting via the segmentation predictions for the ground-truth input images $\mathbf{I}_i$, we use $\mathbf{S}_{\theta_j}(\hat{\mathbf{I}}_{\phi_i})$, that is, we feed the *rendered* images as input to $f$. This removes the need for explicitly storing images for later stages, allows the adaptation to use novel viewpoints for which no observations were made, and as we show in our experiments, results in improved performance over the use of ground-truth images.

Second, we propose to *jointly train* $\mathcal{N}_{\phi_i}$ and $f_{\theta_j}$ by itera-

tively generating labels from one and back-propagating the cross-entropy loss gradients through the other in each training step. In practice, to initialize the NeRF pseudo-labels we first pre-train $\mathcal{N}_{\phi_i}$ with supervision of the ground-truth input images $\mathbf{I}_i$ and of the associated segmentation predictions $\mathbf{S}_{\theta_j}(\mathbf{I}_i)$, and then jointly train $\mathcal{N}_{\phi_i}$ and $f_{\theta_j}$ as described above. We demonstrate the positive influence of this joint adaptation in the experiments, where we show in particular that this 2D-3D knowledge transfer effectively produces improvements in the visual content of both the network predictions and the pseudo-labels.

**Continual NeRF-based replay.** A simple but effective approach to alleviate catastrophic forgetting as the adaptation proceeds across scenes is to *replay* previous experiences, *i.e.*, storing the training data of each newly-encountered scene in a memory buffer, and for each subsequent scene, training the segmentation model using both the data from the new scene and those replayed from the buffer, as done for instance in [12]. In practice, the size of the replay buffer is often limited due to memory and storage constraints, thus one can only store a subset of the data for replay, resulting in a loss of potentially useful information. Unlike previous methods that save explicit data into a buffer, we propose storing the NeRF models in a long-term memory. The advantages of this choice are multifold. First, the memory footprint of multiple NeRF models is significantly smaller than that of explicit images and labels (required by [12]) or of the weights of the segmentation network, stored by [51]. Second, since the NeRF model stores both color and semantic information and attains photorealistic fidelity, it can be used to render a theoretically infinite amount of training views at a fixed storage cost (unlike [12], which fits semantics in the map, and could not produce photorealistic renderings even if texture was aggregated in 3D). Therefore, the segmentation network can be provided with images rendered from NeRF as input. As we show in the experiments, by rendering a small set of views from the NeRF models stored in the long-term memory, our method is able to effectively mitigate catastrophic forgetting.

## 4. Experiments

### 4.1. Experimental settings

**Dataset.** We evaluate our proposed method on the ScanNet [7] dataset. The dataset includes 707 unique indoor scenes, each containing RGB-D images with associated camera poses and manually-generated semantic annotations. In all the experiments we resize the images to a resolution of $320 \times 240$ pixels. Similarly to [12], we use scenes 11-707 in ScanNet to pre-train the semantic segmentation network, taking one image every 100 frames in each of these scenes, for a total of approximately $25\,000$ images. The pre-training dataset is randomly split into

a training set of 20k frames and a validation set of 5k frames. We use scene 1-10 to adapt the pre-trained model (cf. Sec. 4.3, 4.4, 4.5); if the dataset contains more than one video sequence for a given scene, we select only the first one. We select the first 80% of the frames (we refer to them as *training views*) from each sequence to generate predictions with the segmentation network and fuse these into a 3D representation, both by training our Semantic-NeRF model and with the baseline of [12]. The last 20% of the frames (*validation views*) are instead used to test the adaptation performance of the semantic segmentation model on the scene. We stress that this pre-training-training-testing setup is close to a real-world application scenario of the segmentation model, in which in an initial stage the network is trained offline on a large dataset, then some data collected during deployment may be used to adapt the model in an unsupervised fashion, and finally the model performance is tested during deployment on a different trajectory.

**Networks.** We use DeepLabv3 [2] with a ResNet-101 [14] backbone as our semantic segmentation network. To implement our Semantic-NeRF network, we rely on an open-source PyTorch implementation [45] of Instant-NGP [33]. Further details about the architectures of both networks can be found in the Supplementary. For brevity, in the following Sections we refer to Semantic-NeRF as "NeRF".

**Baselines.** As there are no previous works that explicitly tackle the continual semantic adaptation problem, we compare our proposed method to the two most-closely related approaches. The first one [12] uses per-frame camera pose and depth information to aggregate predictions from a segmentation network into a voxel map and then renders semantic pseudo-labels from the map to adapt the network. We implement the method using the framework of [42] and set a resolution of 5 cm for the voxels, so that the total size of the map is comparable to the memory footprint of the NeRF parameters (cf. Supplementary for further details). The second approach, CoTTA [51], focuses on continual test-time domain adaptation and proposes a student-teacher framework with label augmentation and stochastic weight restoration to gradually adapt the semantic segmentation model while keeping the knowledge on the source domain. We use the official open-source implementation, which we adapt to test its performance on the proposed setting.

**Metric.** For all the experiments, we report mean intersection over union (mIoU, in percentage values) as a metric.

### 4.2. Pre-training of the segmentation network

We pre-train DeepLab for 150 epochs to minimize the cross-entropy loss with respect to the ground-truth labels $\mathbf{S}^\star_{\mathrm{pre}}$. We apply common data augmentation techniques, including random flipping/orientation and color jitter. After pre-training, we select the model with best performance on the validation set for adaptation to the new scenes.

| | Pre-train | Mapping [12] | Ours | Ours Joint Training |
|---|---|---|---|---|
| Scene 1 | 41.1 | 48.9 | 48.8±0.7 | **54.8**±1.8 |
| Scene 2 | 35.5 | 33.9 | 36.2±0.8 | **38.3**±0.4 |
| Scene 3 | 23.5 | 25.1 | **27.1**±0.9 | 26.4±1.8 |
| Scene 4 | 62.8 | **65.3** | 62.9±0.5 | 65.0±1.1 |
| Scene 5 | 49.8 | 49.3 | **55.5**±1.3 | 46.6±0.2 |
| Scene 6 | 48.9 | **51.7** | 50.4±0.4 | 50.9±0.4 |
| Scene 7 | 39.7 | 41.2 | 40.4±0.5 | **41.7**±2.0 |
| Scene 8 | 31.6 | 34.8 | 34.0±0.4 | **39.0**±4.6 |
| Scene 9 | 31.7 | 33.8 | **35.6**±0.4 | 31.3±0.4 |
| Scene 10 | 52.5 | 55.8 | **56.4**±0.6 | 56.2±1.0 |
| Average | 41.7 | 44.0 | 44.7±0.7 | **45.0**±1.4 |

Table 1. Pseudo-label performance averaged over the training views and 10 different seeds for Ours pseudo-labels. "Pre-train" denotes the performance of the segmentation model $f_{\theta_0}$.

### 4.3. Pseudo-label formation

We train the NeRF network by minimizing (5) for 60 epochs using the training views. While with our method we can render pseudo-labels from any viewpoint, to allow a controlled comparison against [12] in Sec. 4.4 and 4.5, we generate the pseudo-labels from our NeRF model using the same training viewpoints. While the pseudo-labels of [12] are deterministic, to account for the stochasticity of NeRF, we run our method with 10 different random seeds and report the mean and variance over these. As shown in Tab. 1, the pseudo-labels produced by our method outperform on average those of [12]. A further improvement can be obtained by jointly training NeRF and the DeepLab model, which we discuss in the next Section.

### 4.4. One-step adaptation

As a first adaptation experiment, we evaluate the performance of the different methods when letting the segmentation network $f_{\theta_0}$ adapt in a single stage to each of the scenes 1-10. This setup is similar to that of one-stage UDA, and we thus compare to the state-of-the-art method CoTTA [51].

We evaluate our method in two different settings. In the first one, which we refer to as *fine-tuning*, we simply use the pseudo-labels rendered as in Sec. 4.3 to adapt the segmentation network through cross-entropy loss on its predictions. In the second one, we *jointly train* NeRF and DeepLab via iterative mutual supervision. For a fair comparison, in both settings we optimize the pre-trained NeRF for the same number of additional epochs, while maintaining supervision through color and depth images. In fine-tuning, we perform NeRF pre-training for 60 epochs, according to Sec. 4.3. In joint training, we instead first pre-train NeRF for 10 epochs, and then train NeRF concurrently with DeepLab for 50 epochs. We run each method 10 times and report mean and standard deviation across the runs. Given that the baselines do not support generating images from novel viewpoints, both in fine-tuning and in joint training we use images from the training viewpoints as input to DeepLab. Additionally, since our method allows *ren-*

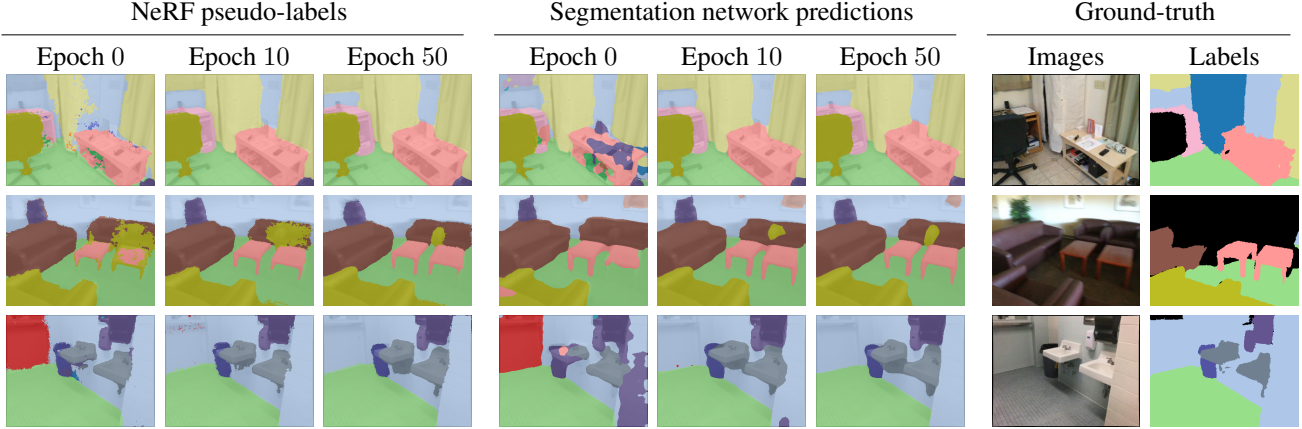|  | NeRF pseudo-labels | | | Segmentation network predictions | | | Ground-truth | |
|---|---|---|---|---|---|---|---|---|
|  | Epoch 0 | Epoch 10 | Epoch 50 | Epoch 0 | Epoch 10 | Epoch 50 | Images | Labels |

Figure 2. Effect of joint training over the pseudo-labels and the predictions of the segmentation network (DeepLab). Color-coded labels are overlaid on the corresponding color images. Black pixels in the ground-truth labels denote missing annotation. First scene: The noisy predictions of DeepLab are corrected and the segmentation results conform much better to the geometry of the scene. Second scene: The geometric details can be better recovered even for the legs of the table. Third scene: By enforcing multi-view consistency, the initial wrong predictions on the wall are corrected through the predictions from other views. Note that the obtained labels adhere accurately to the scene geometry, often even better than in the ground-truth annotations.

*dering* images, we evaluate the difference between feeding ground-truth images vs. NeRF renderings from the same viewpoints to DeepLab.

Table 2 presents the adaptation performance of the different methods on the validation views, which provides a measure of the knowledge transfer induced within the scene by the self-training. Since our method is unsupervised, in the Supplementary we additionally report the improvement in performance on the training views, which is indicative of the effectiveness of the self-supervised adaptation and is of practical utility for real-world deployment scenarios where a scene might be revisited from similar viewpoints.

As shown in Tab. 2, fine-tuning with our pseudo-labels results in improved performance compared to the pre-trained model, and outperforms both baselines for most of the scenes. Interestingly, using rendered images (NI + NL) consistently produces better results than fine-tuning with the ground-truth images (GI + NL). We hypothesize that this is due to the small image artifacts introduced by the NeRF rendering acting as an augmentation mechanism. We further observe that the mIoU can vary largely across the scenes. This can be explained with the variability in the room types and lighting conditions, which is also reflected in the scenes with more extreme illumination (and hence more challenging for NeRF to reconstruct the geometry) having a larger variance with our approach. However, the main observation is that jointly training NeRF and DeepLab (using rendered images as input) results in remarkably better adaptation on almost all the scenes. This improvement can be attributed to the positive knowledge transfer induced between the frame-level predictions of DeepLab and the 3D-aware NeRF pseudo-labels. As shown in Fig. 2, this strategy allows effectively resolving local artifacts in the

NeRF pseudo-labels through the smoothing effect of the DeepLab labels, while at the same time addressing inconsistencies in the per-frame outputs of the segmentation network due to its lack of view consistency.

## 4.5. Multi-step adaptation

To evaluate our method in the full scenario of continual semantic adaptation, we perform multi-step adaptation across scenes 1-10, where in the $i$-th step the segmentation network $f_{\theta_{i-1}}$ gets adapted on scene $\mathcal{S}_i$, resulting in $f_{\theta_i}$, and the NeRF model $\mathcal{N}_i$ is added to the long-term memory at the end of the stage. For steps $i \in \{2, \ldots, 10\}$, to counteract forgetting on the previous scenes we render images and pseudo-labels for each of the $\mathcal{N}_j$ models ($1 \leq j \leq i-1$) in the long-term memory. In practice, we construct a memory buffer of fixed size $100$, to which at stage $i$ each of the previous models $\mathcal{N}_j$ contribute equally with images $\hat{\mathbf{I}}_{\text{buf}}$ and pseudo-labels $\hat{\mathbf{S}}_{\text{buf}}$ rendered from $\lfloor 100/(i-1) \rfloor$ randomly chosen training views. Following [12], we additionally randomly select $10\%$ of the pre-training data and combine them to the data from the previous scenes, which acts as prior knowledge and prevents the model from overfitting to the new scenes and losing its generalization performance. This has a similar effect to the regularization scheme used by CoTTA [51] to preserve previous knowledge, namely storing the network parameters for the initial pre-trained model and the teacher network. Note that both the size of our memory buffer ($14\,\text{MB}$) and that of the replayed pre-training data ($65\,\text{MB}$) are much smaller than the size of two sets of DeepLab weights ($2 \times 225\,\text{MB}$), so our method actually requires less storage space than CoTTA [51]. A detailed analysis of the memory footprint of the different approaches is presented in the Supplementary; we show in particular

| | Pre-train | CoTTA [51] | Fine-tuning (GI + ML) | Ours Fine-tuning (GI + NL) | Ours Fine-tuning (NI + NL) | Ours Joint Training |
|---|---|---|---|---|---|---|
| Scene 1 | 43.9 | $44.0_{\pm0.0}$ | $46.3_{\pm0.3}$ | $46.2_{\pm1.0}$ | $47.1_{\pm1.0}$ | $\mathbf{50.0}_{\pm1.3}$ |
| Scene 2 | 41.3 | $41.2_{\pm0.0}$ | $39.4_{\pm0.3}$ | $39.5_{\pm1.0}$ | $44.2_{\pm1.0}$ | $\mathbf{47.1}_{\pm1.2}$ |
| Scene 3 | $\mathbf{23.0}$ | $22.8_{\pm0.0}$ | $21.6_{\pm0.1}$ | $21.9_{\pm0.7}$ | $21.5_{\pm1.0}$ | $19.9_{\pm2.3}$ |
| Scene 4 | 50.2 | $50.3_{\pm0.0}$ | $52.4_{\pm0.2}$ | $51.5_{\pm0.5}$ | $52.8_{\pm0.8}$ | $\mathbf{53.7}_{\pm2.4}$ |
| Scene 5 | 40.1 | $40.1_{\pm0.0}$ | $49.4_{\pm0.5}$ | $50.6_{\pm2.4}$ | $\mathbf{52.8}_{\pm2.9}$ | $42.7_{\pm1.0}$ |
| Scene 6 | 37.6 | $37.6_{\pm0.0}$ | $33.7_{\pm0.3}$ | $36.2_{\pm1.6}$ | $37.1_{\pm2.4}$ | $\mathbf{40.8}_{\pm1.2}$ |
| Scene 7 | 55.8 | $55.9_{\pm0.0}$ | $50.7_{\pm0.5}$ | $50.7_{\pm1.8}$ | $52.1_{\pm1.3}$ | $\mathbf{56.5}_{\pm4.8}$ |
| Scene 8 | $\mathbf{27.9}$ | $27.9_{\pm0.0}$ | $24.7_{\pm0.2}$ | $23.8_{\pm0.4}$ | $25.3_{\pm0.8}$ | $25.7_{\pm2.9}$ |
| Scene 9 | 54.9 | $54.9_{\pm0.0}$ | $62.2_{\pm1.3}$ | $57.6_{\pm5.3}$ | $52.1_{\pm2.7}$ | $\mathbf{63.7}_{\pm3.3}$ |
| Scene 10 | 73.5 | $73.5_{\pm0.0}$ | $\mathbf{73.8}_{\pm0.2}$ | $\mathbf{73.8}_{\pm0.2}$ | $73.5_{\pm0.4}$ | $73.7_{\pm0.5}$ |
| Average | 44.8 | $44.8_{\pm0.0}$ | $45.4_{\pm0.4}$ | $45.2_{\pm1.5}$ | $45.9_{\pm1.4}$ | $\mathbf{47.4}_{\pm2.1}$ |

Table 2. Performance of the segmentation network on the validation set of each scene after one-step adaptation. GI and NI denote respectively ground-truth color images and NeRF-rendered color images. ML and NL indicate adaptation using pseudo-labels formed respectively with the method of [12] and with our approach. In joint training, we use NeRF-based renderings and pseudo-labels.

| | | Step 1 | Step 2 | Step 3 | Step 4 | Step 5 | Step 6 | Step 7 | Step 8 | Step 9 | Step 10 | Average |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Pre-train | | 43.9 | 41.3 | 23.0 | 50.2 | 40.1 | 37.6 | 55.8 | $\mathbf{27.9}$ | 54.9 | 73.5 | 44.8 |
| Adapt | CoTTA [51] | $44.0_{\pm0.0}$ | $40.9_{\pm0.0}$ | $22.7_{\pm0.0}$ | $50.2_{\pm0.1}$ | $40.0_{\pm0.0}$ | $37.5_{\pm0.0}$ | $56.0_{\pm0.1}$ | $26.9_{\pm0.0}$ | $54.5_{\pm0.0}$ | $\mathbf{73.8}_{\pm0.0}$ | $44.7_{\pm0.0}$ |
| | Mapping [12] | $46.8_{\pm0.4}$ | $42.1_{\pm2.0}$ | $23.6_{\pm0.7}$ | $\mathbf{50.6}_{\pm2.6}$ | $\mathbf{44.0}_{\pm0.1}$ | $35.8_{\pm0.5}$ | $56.7_{\pm1.3}$ | $26.5_{\pm1.8}$ | $68.3_{\pm1.4}$ | $72.7_{\pm1.0}$ | $46.7_{\pm1.2}$ |
| | Ours ($\mathbf{I}_{pre}$ replay only) | $53.3_{\pm0.7}$ | $\mathbf{48.0}_{\pm2.4}$ | $20.5_{\pm0.1}$ | $49.0_{\pm1.5}$ | $43.4_{\pm0.0}$ | $39.0_{\pm1.4}$ | $\mathbf{62.1}_{\pm6.2}$ | $26.7_{\pm3.0}$ | $65.7_{\pm5.6}$ | $73.0_{\pm0.5}$ | $\mathbf{48.1}_{\pm2.1}$ |
| | Ours | $\mathbf{53.7}_{\pm1.3}$ | $46.3_{\pm0.7}$ | $\mathbf{24.3}_{\pm2.0}$ | $49.1_{\pm0.9}$ | $43.7_{\pm0.3}$ | $\mathbf{40.4}_{\pm1.5}$ | $55.8_{\pm0.8}$ | $26.2_{\pm0.9}$ | $\mathbf{68.9}_{\pm3.2}$ | $72.5_{\pm1.6}$ | $\mathbf{48.1}_{\pm1.3}$ |
| Previous | CoTTA [51] | – | $44.0_{\pm0.0}$ | $42.2_{\pm0.0}$ | $35.6_{\pm0.0}$ | $39.3_{\pm0.0}$ | $39.4_{\pm0.0}$ | $39.1_{\pm0.0}$ | $41.5_{\pm0.0}$ | $39.7_{\pm0.0}$ | $41.3_{\pm0.0}$ | $40.2_{\pm0.0}$ |
| | Mapping [12] | – | $46.5_{\pm0.1}$ | $42.8_{\pm1.0}$ | $37.3_{\pm0.9}$ | $40.4_{\pm0.6}$ | $40.9_{\pm0.7}$ | $39.9_{\pm1.1}$ | $42.2_{\pm0.5}$ | $40.0_{\pm0.4}$ | $42.8_{\pm0.7}$ | $41.4_{\pm0.7}$ |
| | Ours ($\mathbf{I}_{pre}$ replay only) | – | $52.3_{\pm0.3}$ | $47.5_{\pm1.1}$ | $38.6_{\pm1.1}$ | $40.8_{\pm0.7}$ | $42.4_{\pm0.3}$ | $41.5_{\pm0.7}$ | $\mathbf{44.3}_{\pm1.4}$ | $41.4_{\pm0.6}$ | $43.9_{\pm0.9}$ | $43.6_{\pm0.8}$ |
| | Ours | – | $\mathbf{53.2}_{\pm0.9}$ | $\mathbf{48.2}_{\pm0.8}$ | $\mathbf{41.5}_{\pm0.8}$ | $\mathbf{42.8}_{\pm0.8}$ | $\mathbf{43.2}_{\pm0.8}$ | $\mathbf{42.2}_{\pm0.8}$ | $44.1_{\pm0.2}$ | $\mathbf{41.7}_{\pm0.2}$ | $\mathbf{44.3}_{\pm0.3}$ | $\mathbf{44.6}_{\pm0.6}$ |

Table 3. Multi-step performance evaluated on the validation set of each scene. At Step $i$, Pre-train and Adapt denote respectively the performance of the pre-trained network $f_{\theta_0}$ and of the adapted network $f_{\theta_i}$ on the current scene $\mathcal{S}_i$, while Previous represents the average performance of $f_{\theta_i}$ on scenes $\mathcal{S}_1$ to $\mathcal{S}_{i-1}$. All Ours are with *joint training*.

that since our method is agnostic to the specific NeRF implementation, with the slower but lighter implementation of Semantic-NeRF [56] the storage comparison is in our favor up to 90 scenes. We deem this to be a realistic margin for real-world deployment scenarios (*e.g.*, it is hardly the case that an agent sequentially visits more than a few scenes during the same mission). For the baseline of [12] we use the same setup as our method, but with mapping-based pseudo-labels and ground-truth images in the memory buffer, due to its inability to generate images.

The multi-step adaptation results are shown in Tab. 3, where for each method the mean and standard deviation across 3 runs are reported. To better show the effect of NeRF-based replay, we also run our adaptation method with only replay from the pre-training dataset, without replaying from the old NeRF models (Ours ($\mathbf{I}_{pre}$ replay only)). Our method achieves the best average adaptation performance (Adapt) across the new scenes in the multi-step setting, improving by $\sim 3\%$ mIoU over the pre-trained model. Note that this improvement is consistent with the one observed in one-step adaptation (Tab. 2), which validates that our method can successfully adapt across multiple scenes, without the performance dropping after a specific number of steps. At the same time, while NeRF-based replay of the old scenes on average does not induce a positive forward transfer in the adaptation to the new scenes (Adapt), its usage can significantly alleviate forgetting compared to the case

with no replay. As a result, when using NeRF-based replay, our method is able to maintain in almost all the adaptation steps the best average performance over the previous scenes (Previous). Further in-detail results for each scene and after each adaptation step are reported in the Supplementary.

## 5. Conclusion

In this work, we present a novel approach for unsupervised continual adaptation of a semantic segmentation network to multiple novel scenes using neural rendering. We exploit the fact that the new scenes are observed from multiple viewpoints and jointly train in each scene a Semantic-NeRF model and the segmentation network. We show that the induced 2D-3D knowledge transfer results in improved unsupervised adaptation performance compared to state-of-the-art methods. We further propose a NeRF-based replay strategy which allows efficiently mitigating catastrophic forgetting and enables rendering a potentially infinite number of images for adaptation at constant storage cost. We believe this opens up interesting avenues for replay-based adaptation, particularly for use on real-world perception systems, which can compactly store collected experiences on board and generate past data as needed.

# References

[1] Fabio Cermelli, Massimiliano Mancini, Samuel Rota Buló, Elisa Ricci, and Barbara Caputo. Modeling the Background for Incremental Learning in Semantic Segmentation. In *CVPR*, 2020. 3

[2] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking Atrous Convolution for Semantic Image Segmentation. *CoRR:1706.05587*, 2017. 6

[3] Yun-Chun Chen, Yen-Yu Lin, Ming-Hsuan Yang, and Jia-Bin Huang. CrDoCo: Pixel-Level Domain Transfer With Cross-Domain Consistency. In *CVPR*, 2019. 2, 3

[4] Yi-Hsin Chen, Wei-Yu Chen, Yu-Ting Chen, Bo-Cheng Tsai, Yu-Chiang Frank Wang, and Min Sun. No More Discrimination: Cross City Adaptation of Road Scene Segmenters. In *ICCV*, 2017. 2, 3

[5] Jaehoon Choi, Taekyung Kim, and Changick Kim. Self-Ensembling With GAN-Based Data Augmentation for Domain Adaptation in Semantic Segmentation. In *ICCV*, 2019. 3

[6] 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 *CVPR*, 2016. 2

[7] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. ScanNet: Richly-Annotated 3D Reconstructions of Indoor Scenes. In *CVPR*, 2017. 2, 5

[8] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised NeRF: Fewer Views and Faster Training for Free. In *CVPR*, 2022. 4, 12

[9] Natalia Díaz-Rodríguez, Vincenzo Lomonaco, David Filliat, and Davide Maltoni. Don't forget, there is more than forgetting: new metrics for Continual Learning. In *NeurIPS Workshop*, 2018. 14

[10] Arthur Douillard, Yifu Chen, Arnaud Dapogny, and Matthieu Cord. PLOP: Learning without Forgetting for Continual Semantic Segmentation. In *CVPR*, 2021. 3

[11] Liang Du, Jingang Tan, Hongye Yang, Jianfeng Feng, Xiangyang Xue, Qibao Zheng, Xiaoqing Ye, and Xiaolin Zhang. SSF-DAN: Separated Semantic Feature Based Domain Adaptation Network for Semantic Segmentation. In *ICCV*, 2019. 2, 3

[12] Jonas Frey, Hermann Blum, Francesco Milano, Roland Siegwart, and Cesar Cadena. Continual Adaptation of Semantic Segmentation using Complementary 2D-3D Data Representations. *IEEE Robot. Autom. Lett.*, 2022. 2, 3, 4, 5, 6, 7, 8, 11, 12, 14, 15, 16, 17, 18

[13] Xiao Fu, Shangzhan Zhang, Tianrun Chen, Yichong Lu, Lanyun Zhu, Xiaowei Zhou, Andreas Geiger, and Yiyi Liao. Panoptic NeRF: 3D-to-2D Label Transfer for Panoptic Urban Scene Segmentation. In *3DV*, 2022. 3

[14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image Recognition. In *CVPR*, 2016. 6

[15] Judy Hoffman, Eric Tzeng, Taesung Park, Jun-Yan Zhu, Phillip Isola, Kate Saenko, Alexei Efros, and Trevor Darrell. CyCADA: Cycle-Consistent Adversarial Domain Adaptation. In *ICML*, 2018. 2, 3

[16] Ajay Jain, Matthew Tancik, and Pieter Abbeel. Putting NeRF on a Diet: Semantically Consistent Few-Shot View Synthesis. In *ICCV*, 2021. 3

[17] Diederik P Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. In *ICLR*, 2015. 11

[18] Sosuke Kobayashi, Eiichi Matsumoto, and Vincent Sitzmann. Decomposing NeRF for Editing via Feature Field Distillation. In *NeurIPS*, 2022. 3

[19] Abhijit Kundu, Kyle Genova, Xiaoqi Yin, Alireza Fathi, Caroline Pantofaru, Leonidas Guibas, Andrea Tagliasacchi, Frank Dellaert, and Thomas Funkhouser. Panoptic Neural Fields: A Semantic Object-Aware Neural Scene Representation. In *CVPR*, 2022. 3

[20] Seungmin Lee, Dongwan Kim, Namil Kim, and Seong-Gyun Jeong. Drop to Adapt: Learning Discriminative Features for Unsupervised Domain Adaptation. In *ICCV*, 2019. 3

[21] Timothée Lesort, Vincenzo Lomonaco, Andrei Stoian, Davide Maltoni, David Filliat, and Natalia Díaz-Rodríguez. Continual Learning for Robotics: Definition, Framework, Learning Strategies, Opportunities and Challenges. *Information Fusion*, 58, 2020. 1

[22] Yanghao Li, Naiyan Wang, Jianping Shi, Jiaying Liu, and Xiaodi Hou. Revisiting Batch Normalization For Practical Domain Adaptation. *CoRR:1603.04779*, 2016. 11

[23] Yunsheng Li, Lu Yuan, and Nuno Vasconcelos. Bidirectional Learning for Domain Adaptation of Semantic Segmentation. In *CVPR*, 2019. 2, 3

[24] 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 *ECCV*, 2014. 11

[25] David Lopez-Paz and Marc'Aurelio Ranzato. Gradient Episodic Memory for Continual Learning. In *NeurIPS*, 2017. 14

[26] Andrea Maracani, Umberto Michieli, Marco Toldo, and Pietro Zanuttigh. RECALL: Replay-based Continual Learning in Semantic Segmentation. In *ICCV*, 2021. 3

[27] Umberto Michieli, Matteo Biasetton, Gianluca Agresti, and Pietro Zanuttigh. Adversarial Learning and Self-Teaching Techniques for Domain Adaptation in Semantic Segmentation. *IEEE TIV*, 5, 2020. 2, 3

[28] Umberto Michieli, Marco Toldo, and Pietro Zanuttigh. Domain adaptation and continual learning in semantic segmentation. In *Advanced Methods and Deep Learning in Computer Vision*, chapter 8, pages 275–303. Elsevier, 2022. 1, 2, 3

[29] Umberto Michieli and Pietro Zanuttigh. Incremental Learning Techniques for Semantic Segmentation. In *ICCVW*, 2019. 3

[30] Umberto Michieli and Pietro Zanuttigh. Continual Semantic Segmentation via Repulsion-Attraction of Sparse and Disentangled Latent Representations. In *CVPR*, 2021. 3

[31] Umberto Michieli and Pietro Zanuttigh. Knowledge Distillation for Incremental Learning in Semantic Segmentation. *J. Comput. Vis. Image Understanding*, 205, 2021. 3

[32] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis. In *ECCV*, 2020. 2, 3, 4, 5

[33] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant Neural Graphics Primitives with a Multiresolution Hash Encoding. *ACM TOG*, 41(4):102:1–102:15, 2022. 5, 6, 11, 12, 16, 17

[34] Zak Murez, Soheil Kolouri, David Kriegman, Ravi Ramamoorthi, and Kyungnam Kim. Image to Image Translation for Domain Adaptation. In *CVPR*, 2018. 2, 3

[35] Mathieu Pagé Fortin and Brahim Chaib-draa. Continual Semantic Segmentation Leveraging Image-level Labels and Rehearsal. In *IJCAI*, 2022. 3

[36] Keunhong Park, Utkarsh Sinha, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Steven M. Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable Neural Radiance Fields . In *ICCV*, 2021. 17

[37] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-NeRF: Neural Radiance Fields for Dynamic Scenes. In *CVPR*, 2021. 17

[38] Stephan R. Richter, Vibhav Vineet, Stefan Roth, and Vladlen Koltun. Playing for Data: Ground Truth from Computer Games. In *ECCV*, 2016. 2

[39] German Ros, Laura Sellart, Joanna Materzynska, David Vazquez, and Antonio M. Lopez. The SYNTHIA Dataset: A Large Collection of Synthetic Images for Semantic Segmentation of Urban Scenes. In *CVPR*, 2016. 2

[40] Kuniaki Saito, Yoshitaka Ushiku, Tatsuya Harada, and Kate Saenko. Adversarial Dropout Regularization. In *ICLR*, 2018. 2, 3

[41] Swami Sankaranarayanan, Yogesh Balaji, Arpit Jain, Ser Nam Lim, and Rama Chellappa. Learning From Synthetic Data: Addressing Domain Shift for Semantic Segmentation. In *CVPR*, 2018. 2, 3

[42] Lukas Schmid, Jeffrey Delmerico, Johannes Schönberger, Juan Nieto, Marc Pollefeys, Roland Siegwart, and Cesar Cadena. Panoptic Multi-TSDFs: a Flexible Representation for Online Multi-resolution Volumetric Mapping and Long-term Dynamic Scene Consistency. In *ICRA*, 2022. 6

[43] Ken Shoemake. Animating Rotation with Quaternion Curves. In *SIGGRAPH*, 1985. 15

[44] Teo Spadotto, Marco Toldo, Umberto Michieli, and Pietro Zanuttigh. Unsupervised Domain Adaptation with Multiple Domain Discriminators and Adaptive Self-Training. In *ICPR*, 2021. 2, 3

[45] Jiaxiang Tang. Torch-ngp: a PyTorch implementation of instant-ngp. https://github.com/ashawkey/torch-ngp, 2022. 6, 11, 12, 16, 17

[46] Marco Toldo, Andrea Maracani, Umberto Michieli, and Pietro Zanuttigh. Unsupervised Domain Adaptation in Semantic Segmentation: a Review. *Technologies*, 8, 2020. 1, 2

[47] Vadim Tschernezki, Iro Laina, Diane Larlus, and Andrea Vedaldi. Neural Feature Fusion Fields: 3D Distillation of Self-Supervised 2D Image Representations. In *3DV*, 2022. 3

[48] Suhani Vora, Noha Radwan, Klaus Greff, Henning Meyer, Kyle Genova, Mehdi S. M. Sajjadi, Etienne Pot, Andrea Tagliasacchi, and Daniel Duckworth. NeSF: Neural Semantic Fields for Generalizable Semantic Segmentation of 3D Scenes. *TMLR*, 2022. 3

[49] Tuan-Hung Vu, Himalaya Jain, Maxime Bucher, Matthieu Cord, and Patrick Perez. ADVENT: Adversarial Entropy Minimization for Domain Adaptation in Semantic Segmentation. In *CVPR*, 2019. 2, 3

[50] Can Wang, Menglei Chai, Mingming He, Dongdong Chen, and Jing Liao. CLIP-NeRF: Text-and-Image Driven Manipulation of Neural Radiance Fields. In *CVPR*, 2022. 3

[51] Qin Wang, Olga Fink, Luc Van Gool, and Dengxin Dai. Continual Test-Time Domain Adaptation. In *CVPR*, 2022. 2, 3, 5, 6, 7, 8, 12, 14, 15, 16, 17, 18

[52] Zuxuan Wu, Xin Wang, Joseph E. Gonzalez, Tom Goldstein, and Larry S. Davis. ACE: Adapting to Changing Environments for Semantic Segmentation. In *ICCV*, 2019. 3, 5

[53] Yanchao Yang and Stefano Soatto. FDA: Fourier Domain Adaptation for Semantic Segmentation. In *CVPR*, 2020. 2, 3

[54] Yiheng Zhang, Zhaofan Qiu, Ting Yao, Dong Liu, and Tao Mei. Fully Convolutional Adaptation Networks for Semantic Segmentation. In *CVPR*, 2018. 2, 3

[55] Zhedong Zheng and Yi Yang. Rectifying Pseudo Label Learning via Uncertainty Estimation for Domain Adaptive Semantic Segmentation. *IJCV*, 2021. 3

[56] Shuaifeng Zhi, Tristan Laidlow, Stefan Leutenegger, and Andrew J. Davison. In-Place Scene Labelling and Understanding with Implicit Scene Representation. In *ICCV*, 2021. 2, 3, 4, 5, 8, 11, 12, 13, 16, 17

[57] Shuaifeng Zhi, Edgar Sucar, Andre Mouton, Iain Haughton, Tristan Laidlow, and Andrew J. Davison. iLabel: Interactive Neural Scene Labelling. *CoRR:2111.14637*, 2021. 3

[58] Yang Zou, Zhiding Yu, B.V.K. Vijaya Kumar, and Jinsong Wang. Unsupervised Domain Adaptation for Semantic Segmentation via Class-Balanced Self-Training. In *ECCV*, 2018. 2, 3

[59] Yang Zou, Zhiding Yu, Xiaofeng Liu, B.V.K. Vijaya Kumar, and Jinsong Wang. Confidence Regularized Self-Training. In *ICCV*, 2019. 2, 3

# Supplementary Material

The Supplementary Material is organized as follows. In Sec. A, we provide additional implementation details. In Sec. B, we present ablations on the NeRF-based pseudo-labels, showing the effect on their quality of different parameters and components of our method. In Sec. C, we report additional evaluations for the one-step adaptation experiments. In Sec. D we include in-detail results for the multi-step adaptation experiments and ablate on the replay-based strategy proposed by our method. In Sec. E we analyze in detail the memory footprint required by our method and by the different baselines that we compare against in the main paper. In Sec. F, we provide further visualizations, including examples of the pseudo-labels and network predictions produced by our method and the baselines. In Sec. G, we discuss limitations of our method and potential ways to address them. We will further release the code to reproduce our results.

Similarly to the main paper, in all the experiments we report mean intersection over union (mIoU, in percentage values) as a metric.

## A. Additional implementation details

**NeRF.** Following Instant-NGP [33, 45], to facilitate training of the hash encoding, we re-scale and re-center the poses used to train NeRF so that they fit in a fixed-size cube. For each ray that is cast from the training viewpoints, to render the aggregated colors and semantics labels we first sample 256 points at a fixed interval and then randomly select 256 additional points according to the density values of the initial points.

The base NeRF network uses a multi-resolution hash encoding with a 16-level hash table of size $2^{19}$ and a feature dimension of 2. Similarly to Semantic-NeRF [56], we implement the additional semantic head as a 2-layer MLP. In all the experiments, we train all the components of the Semantic-NeRF network concurrently, setting the hyperparameters in Eq. (5) from the main paper to $w_d = 0.1$ and $w_s = 0.04$ as suggested in [56], sampling 4096 rays for each viewpoint, and using the Adam [17] optimizer with a fixed learning rate of 1e−2.

In all the experiments in which the semantic segmentation model is trained using NeRF-rendered images, we use Adaptive Batch Normalization (AdaBN) [22] when performing inference on the ground-truth images, to improve the generalization ability of the model between NeRF-rendered images and ground-truth images.

**Dataset.** For convenience of notation, we re-map the scene indices in the dataset from $0000 - 0706$ to $1 - 707$ (so that we refer to scene 0000 as scene 1, to scene 0001 as scene 2, etc.). For sample efficiency, we downsample each sequence from the original 30 fps to 3 fps, resulting in a total of 100 to 500 frames for each video sequence.

**Pre-training.** To pre-train DeepLab on scenes $11 - 707$ from ScanNet, we initialize the model parameters with the weights pre-trained on the COCO semantic segmentation dataset [24]. We then run the pre-training on ScanNet using the Adam [17] optimizer with batch size of 4, and let the learning rate decay linearly from 1e−4 to 1e−6 over 150 epochs.

**One-step adaptation.** In all the one-step experiments with our method and with the baseline of [12], the semantic segmentation model is trained for 50 epochs with a fixed learning rate of 1e−5 and batch size of 4. Since CoTTA is an *online* adaptation method, in accordance with the settings introduced in the original paper, we adapt the segmentation network for a single epoch and with batch size 1, setting the learning rate to 2.5e−6. To prevent overfitting the semantic segmentation model to the training views of the new scene, we apply the same data augmentation procedure as in pre-training in each training step for our method and for [12]. Since CoTTA already implements a label augmentation mechanism for ensembling, we apply to the method only the augmentations used by its authors.

**Multi-step adaptation.** In the multi-step adaptation experiments, we use a batch size of 4 during training, where 2 samples come from the subset of the pre-training dataset used for replay (cf. main paper), and the other 2 data points are uniformly sampled from the training frames of the new scene and the replay buffer of the previous scenes.

**Hardware.** We train all our models using an AMD Ryzen 9 5900X with 32 GB RAM, and an NVIDIA RTX3090 GPU with 24 GB VRAM.

## B. NeRF-based pseudo-labels

In the following Section, we present ablations on the NeRF-based pseudo-labels, showing how the chosen NeRF implementation and the losses used in our method influence their segmentation accuracy.

### B.1. Comparison of NeRF frameworks

We compare the segmentation quality of the pseudo-labels obtained with our Instant-NGP [33, 45]-based implementation to that achieved with the original Semantic-NeRF [56] implementation, which we adapt to include the newly-introduced semantic loss (cf. Sec. 3.2 in the main paper and Sec. B.2). To this purpose, we train a semantics-aware NeRF model for scene 1 with both the methods, running the experiments 3 times for each method. In each run, we train the original implementation of Semantic-NeRF [56] for 200k steps and the one based on Instant-NGP [45] for 10 epochs (for a total of $10 \times 447 = 4470$ steps), which allows achieving a similar color reconstruction quality (measured as PSNR) for the two methods.

| Components | | Scene | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $\mathcal{L}_d$ | $\mathcal{L}_s$ | Scene 1 | Scene 2 | Scene 3 | Scene 4 | Scene 5 | Scene 6 | Scene 7 | Scene 8 | Scene 9 | Scene 10 | Average |
| ✗ | Semantic-NeRF [56] | 44.3±1.5 | 34.2±0.1 | 22.4±0.9 | **63.5**±1.2 | 52.3±1.2 | 47.3±0.5 | 38.9±0.6 | 33.8±0.4 | 32.4±0.5 | 53.3±0.6 | 42.2±0.7 |
| ✗ | Ours | 46.4±1.1 | 33.0±0.2 | 24.2±0.3 | 62.6±0.7 | 53.4±0.7 | 46.8±1.1 | 39.3±0.8 | **34.5**±0.6 | 33.8±0.6 | 55.8±0.2 | 43.0±0.6 |
| ✓ | Semantic-NeRF [56] | 44.0±0.6 | 34.8±0.5 | 22.8±0.9 | 63.1±0.7 | 55.8±2.0 | 49.1±1.2 | 39.0±0.8 | 33.9±0.5 | 33.0±1.5 | 55.1±0.6 | 43.1±0.9 |
| ✓ | Ours | **48.4**±0.9 | **36.0**±0.3 | **26.1**±0.4 | 61.6±0.5 | **57.0**±1.8 | **50.3**±0.6 | **39.8**±0.2 | 33.5±0.6 | **35.4**±0.7 | **57.4**±0.1 | **44.6**±0.7 |

Table 4. Effect of the $\ell_1$ depth loss $\mathcal{L}_d$ and of different types of semantic losses (either the original one proposed in [56] or ours) on the pseudo-label quality. The performance is evaluated on the training views of each scene and averaged over 3 runs.

| | Semantic-NeRF [56] | Instant-NGP [33] (impl. by [45]) |
|---|---|---|
| PSNR | 19.9±0.1 | 19.3±0.1 |
| mIoU | 50.0±0.5 | 48.4±0.9 |
| Model size (MB) | 4.9 | 50.0 |
| Training time / Step (s) | 0.19 | 0.06 |
| Total training time (min) | 633 | 5 |
| Inference time / Image (s) | 2.8 | 0.3 |

Table 5. Pseudo-label performance on the training views of scene 1, size of the associated model checkpoint, and the training and inference time using different NeRF frameworks. The implementation of [56] has been adapted to include the newly-introduced semantic loss (cf. Sec. 3.2 in the main paper). The results are averaged over 3 runs.

As shown in Tab. 5, the pseudo-labels produced by both implementations achieve a similar mIoU, with Semantic-NeRF slightly outperforming Instant-NGP. Furthermore, the size of the models produced by Semantic-NeRF is approximately 10 times smaller than the one required by Instant-NGP, at the cost however of longer training ($\sim 127\times$) and rendering ($\sim 9\times$) time.

Since in a real-world deployment scenario achieving fast adaptation might be of high priority, in the main paper we adopted the faster framework of Instant-NGP. However, the results above indicate that our method is agnostic to the specific NeRF framework chosen, and similar segmentation performance can be achieved by trading off between speed and model size depending on the main requirements. Further evaluations on the memory footprint in comparison also with the baselines of [12] and [51] are presented in Sec. E.

## B.2. Ablation on the NeRF losses

To investigate the effect of depth supervision [8] (through the $\ell_1$ depth loss $\mathcal{L}_d$) and of the proposed modifications to the semantic loss $\mathcal{L}_s$ (cf. Sec. 3.2 in the main paper), we evaluate on each scene the pseudo-labels produced by our method when ablating on these factors. For each scene, we train the NeRF model for 10 epochs without joint training, as we find training without semantic loss modifications is unstable for longer epochs. We run each experiment 3 times and report average and standard deviation across the runs. As shown in Tab. 4, both components induce a significant improvement of the pseudo-label quality. In particular, depth supervision and the use of our modified semantic loss instead of the one proposed in [56] produce an increase re-

spectively of $0.9\%$ mIoU and $0.8\%$ mIoU over the baseline with no modifications. The combined use of both ablated factors further increases the pseudo-label performance, resulting in a total improvement by $2.4\%$ mIoU.

The effect of the proposed modifications can also be observed in Fig. 3. In particular, as shown in Fig. 3a, the use of depth supervision is critical for properly reconstructing the scene geometry. The large number of artifacts in the reconstruction when the depth loss is not used are also reflected in the semantic pseudo-labels, which contain large levels of noise and often fail to assign a uniform class to each entity in the scene (Fig. 3b). Depth supervision applied together with the original semantic loss from [56] resolves some of the artifacts in the pseudo-labels, but still results in suboptimal quality. The combined use of depth supervision and of our modified semantic loss produces cleaner and smoother pseudo-labels, which also attain higher segmentation accuracy, as shown in Tab. 4.
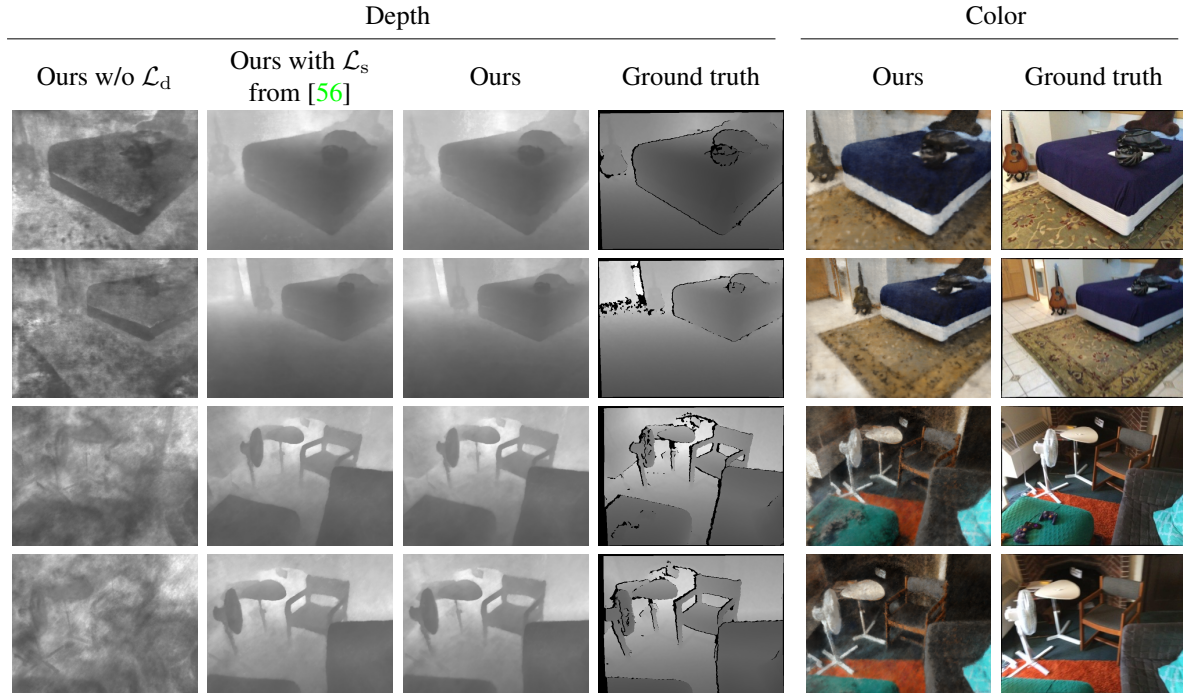
## C. One-step adaptation

In this Section, we report additional results on the one-step adaptation experiments.
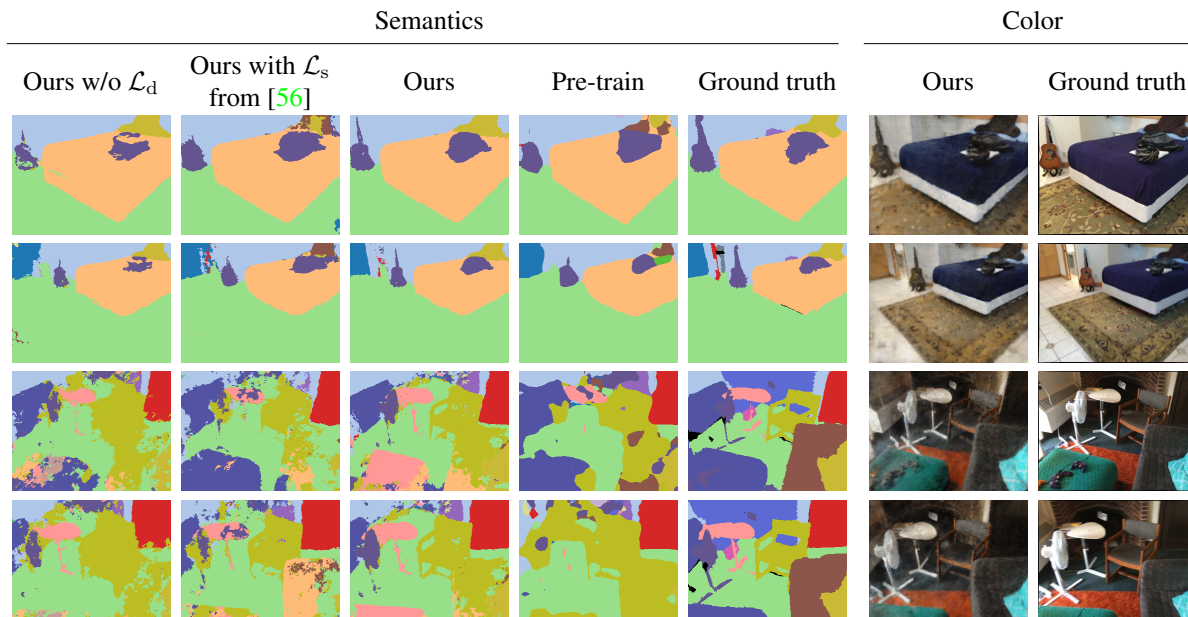
### C.1. One-step adaptation performance on the training set of each scene

Since in the scenario of a deployment of the semantic segmentation network on a real-world system a scene might be revisited from viewpoints similar to those used for training, in Tab. 6 we report the one-step adaptation performance evaluated on the training views. We compare our method to the baseline of CoTTA [51] and to fine-tuning, both with the pseudo-labels of [12] and with our NeRF-based pseudo-labels. For each method, we run the experiments 10 times and report average and standard deviation across the runs.

Similarly to the results obtained on the validation views (cf. main paper), our method with joint training obtains the best average performance across all scenes. Unlike what observed on the validation views, however, on the training views joint training does not result in an average performance improvement over fine-tuning with our NeRF-based pseudo-labels (NI + NL). We note however that these results are largely influenced by the outlier of Scene 5, where joint training achieves significantly lower segmentation accuracy. In Sec. G we analyze more in detail the failure cases

(a) Rendered depth



(b) Rendered semantics

Figure 3. Effect on the rendered depth and semantics of depth supervision and of the modification to the semantic loss. Black pixels in the ground-truth depth and ground-truth semantics denote respectively missing depth measurement and missing semantic annotation.

of our method and focus specifically also on Scene 5, which we find to contain several frames with extreme illumination conditions, which makes it particularly challenging to properly reconstruct the geometry of certain parts of the scene.

## D. Multi-step adaptation

In the following Section, we include in-detail results for the multi-step adaptation experiments, reporting addition-

| | Pre-train | CoTTA [51] | Fine-tuning (GI + ML) | Fine-tuning (GI + NL) | Fine-tuning (NI + NL) | Joint Training |
|---|---|---|---|---|---|---|
| Scene 1 | 41.1 | 41.9±0.0 | 50.6±0.1 | 50.1±0.6 | 50.7±0.5 | **55.5**±1.3 |
| Scene 2 | 35.5 | 35.6±0.0 | 33.5±0.1 | 35.7±0.8 | 36.6±0.3 | **39.5**±0.8 |
| Scene 3 | 23.5 | 23.7±0.0 | 24.4±0.1 | 26.9±1.0 | 27.1±1.2 | **27.5**±1.6 |
| Scene 4 | 62.8 | 63.0±0.0 | 66.1±0.3 | 63.2±0.6 | 66.1±0.8 | **67.7**±1.7 |
| Scene 5 | 49.8 | 49.8±0.0 | 51.2±0.1 | 57.1±1.2 | **59.9**±1.5 | 46.3±0.3 |
| Scene 6 | 48.9 | 48.9±0.0 | **53.1**±0.1 | 50.2±0.4 | 49.9±0.4 | 50.7±0.2 |
| Scene 7 | 39.7 | 39.8±0.0 | 41.4±0.1 | 40.8±0.6 | 42.1±0.8 | **43.8**±1.6 |
| Scene 8 | 31.6 | 31.7±0.0 | 36.2±0.2 | 34.4±0.5 | 33.9±0.4 | **38.1**±3.5 |
| Scene 9 | 31.7 | 31.7±0.0 | 32.7±0.1 | **35.5**±0.6 | 34.9±0.8 | 32.5±0.9 |
| Scene 10 | 52.5 | 52.7±0.0 | 57.8±0.1 | 57.1±0.6 | **58.4**±0.6 | 57.4±1.4 |
| Average | 41.7 | 41.9±0.0 | 44.7±0.1 | 45.1±0.7 | **45.9**±0.7 | **45.9**±1.3 |

Table 6. One-step adaptation performance on the training views of each scene. GI and NI denote respectively ground-truth color images and NeRF-rendered color images. ML and NL indicate adaptation using pseudo-labels formed respectively with the method of [12] and with our approach. In joint training, we use NeRF-based renderings and pseudo-labels. For each method, we run the experiments for 10 times and report average and standard deviation across the runs.

ally a set of standard metrics used in the continual learning literature. We further demonstrate the use, enabled by our method, of images and pseudo-labels rendered from *novel* viewpoints in previous scenes for multi-step adaptation. Remarkably, we find that this modification induces a further improvement in the retention of knowledge from the previous scenes.

## D.1. Detailed per-step evaluation

Table 8 reports the segmentation performance on the validation views of each scene after each step of adaptation, both for our method and for the baselines of [51] and [12]. For each method, we run the experiment 3 times and report main and standard deviation across the runs. The results complement Tab. 3 in the main paper, confirming in particular that in all the adaptation steps our method is the most effective at preserving knowledge on the previous scenes.

| | ACC Metric [25] | A Metric [9] | FWT [25] | BWT [25] |
|---|---|---|---|---|
| CoTTA [51] | 44.6±0.0 | 40.9±0.0 | **-0.2**±0.0 | **-0.1**±0.0 |
| Mapping [12] | 45.8±0.6 | 42.1±0.6 | -1.1±0.2 | -1.0±0.6 |
| Ours ($\mathbf{I}_{pre}$ replay only) | 46.8±0.8 | 43.7±0.6 | -1.4±0.7 | -1.4±0.7 |
| Ours | **47.2**±0.5 | **44.3**±0.2 | -1.1±0.2 | -0.9±0.4 |

Table 7. Continual learning metrics extracted from Tab. 8.

To facilitate the analysis of the results, in Tab. 7 we further report a set of metrics commonly used in the continual learning literature. Our method achieves the best performance both according to the ACC metric [25] and to the A metric [9], meaning that it obtains the best average mIoU across all previously visited scenes both at the final step and at any arbitrary adaptation step. The baseline of CoTTA [51] attains the best forward transfer (FWT) [25] and backward transfer (BWT) [25], which indicate respectively the influence that previous scenes have on the perfor-

mance on future scenes and the influence that adaptation on the current scenes has on the performance on the previous scenes (negative BWT corresponds to catastrophic forgetting). An important point to notice, however, is that the performance of CoTTA also does not vary significantly with respect to the pre-trained model, and in particular does not improve on average. Among the other methods, our method achieves the best FWT and BWT, which demonstrates the effectiveness of our NeRF-based replay buffer in alleviating forgetting and improving the generalization performance.

## D.2. "Replaying" from novel viewpoints

A key feature enabled by our method is the possibility of rendering both photorealistic color images and pseudo-labels from any arbitrary viewpoint inside a reconstructed scene. Crucially, this can include also *novel* viewpoints not seen during deployment and training, which can then be used for adaptation, at the fixed storage cost given by the size of the NeRF model parameters. In the following, we present an experiment demonstrating this idea in the multi-step adaptation scenario. Using the notation introduced in the paper, in each step $i \in \{1, \ldots, 10\}$, the semantic segmentation network $f_{\theta_{i-1}}$ is adapted on scene $\mathcal{S}_i$, and for each previous scene $\mathcal{S}_j$, $1 \leq j < i$ images and pseudo-labels rendered from viewpoints $\hat{\mathbf{P}}_j$ are inserted in a rendering buffer and mixed to the data from the current scene. However, unlike the experiments in the main paper, we do not enforce that for each scene $\mathcal{S}_j$ the viewpoints $\hat{\mathbf{P}}_j$ used for the rendering buffer coincide with those used in training $\mathbf{P}_j := \{\boldsymbol{P}_j^k\}_{k \in \{1, \cdots, |\mathbf{I}_j|\}}$, but instead allow novel viewpoints to be used, that is, $|\hat{\mathbf{P}}_j \backslash (\hat{\mathbf{P}}_j \cap \mathbf{P}_j)| > 0$.

Specifically, in the presented experiment we apply simple average interpolation of the training poses, and for each viewpoint $\hat{\boldsymbol{P}}_j^k \in \hat{\mathbf{P}}_j$ we compute its rotation component

| Method | Step | Scene 1 | Scene 2 | Scene 3 | Scene 4 | Scene 5 | Scene 6 | Scene 7 | Scene 8 | Scene 9 | Scene 10 | Average Prev. | Average |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Pre-training | – | 43.9 | 41.3 | 23.0 | 50.2 | 40.1 | 37.6 | 55.8 | 27.9 | 54.9 | 73.5 | – | 44.8 |
| CoTTA [51] | 1 | 44.0±0.0 | 40.9±0.0 | 22.9±0.0 | 50.3±0.0 | 40.1±0.1 | 37.5±0.0 | 55.9±0.0 | 27.6±0.0 | 54.7±0.0 | 73.6±0.0 | – | 44.7±0.0 |
| | 2 | 44.0±0.0 | 40.9±0.0 | 22.9±0.0 | 50.3±0.0 | 40.1±0.0 | 37.5±0.0 | 55.9±0.0 | 27.6±0.0 | 54.8±0.0 | 73.6±0.0 | 44.0±0.0 | 44.7±0.0 |
| | 3 | 43.6±0.1 | 40.7±0.1 | 22.7±0.0 | 50.1±0.1 | 39.9±0.0 | 37.5±0.0 | 56.1±0.0 | 27.3±0.0 | 54.6±0.0 | 73.7±0.0 | 42.2±0.0 | 44.6±0.0 |
| | 4 | 43.6±0.0 | 40.5±0.0 | 22.7±0.0 | 50.2±0.1 | 39.9±0.0 | 37.5±0.0 | 56.0±0.1 | 27.2±0.0 | 54.5±0.0 | 73.7±0.0 | 35.6±0.0 | 44.6±0.0 |
| | 5 | 43.7±0.1 | 40.5±0.0 | 22.7±0.0 | 50.2±0.1 | 40.0±0.0 | 37.5±0.0 | 55.9±0.1 | 27.1±0.0 | 54.6±0.0 | 73.7±0.0 | 39.3±0.0 | 44.6±0.0 |
| | 6 | 43.7±0.0 | 40.4±0.1 | 22.7±0.0 | 50.3±0.1 | 40.0±0.0 | 37.5±0.0 | 55.9±0.1 | 27.0±0.0 | 54.5±0.0 | 73.7±0.0 | 39.4±0.0 | 44.6±0.0 |
| | 7 | 43.7±0.1 | 40.4±0.1 | 22.7±0.1 | 50.3±0.1 | 39.9±0.1 | 37.6±0.1 | 56.0±0.1 | 26.9±0.0 | 54.5±0.0 | 73.7±0.0 | 39.1±0.0 | 44.6±0.0 |
| | 8 | 43.7±0.0 | 40.4±0.1 | 22.7±0.1 | 50.3±0.1 | 39.9±0.1 | 37.7±0.1 | 56.0±0.1 | **26.9**±0.0 | 54.5±0.0 | 73.7±0.0 | 41.5±0.0 | 44.6±0.0 |
| | 9 | 43.7±0.0 | 40.3±0.1 | 22.7±0.1 | 50.2±0.1 | 39.9±0.1 | 37.7±0.1 | 56.0±0.1 | 26.8±0.0 | 54.5±0.0 | 73.8±0.0 | 39.7±0.0 | 44.6±0.0 |
| | 10 | 43.7±0.1 | 40.2±0.1 | 22.7±0.1 | 50.3±0.1 | 39.9±0.1 | 37.6±0.0 | 56.1±0.1 | 26.8±0.0 | 54.4±0.1 | **73.8**±0.0 | 41.3±0.0 | 44.6±0.0 |
| Mapping [12] | 1 | 46.8±0.4 | 36.0±1.6 | 24.2±0.9 | 48.3±0.9 | 40.0±0.9 | 35.3±0.8 | 55.5±0.4 | 29.2±2.3 | 55.7±1.0 | 73.9±0.2 | – | 44.5±0.5 |
| | 2 | 46.5±0.1 | 42.1±2.0 | 23.6±0.9 | 48.4±1.3 | 41.3±1.0 | 35.5±0.7 | 54.8±1.1 | 28.3±0.8 | 56.5±0.9 | 73.7±0.2 | 46.5±0.1 | 45.1±0.2 |
| | 3 | 43.0±1.2 | 42.6±2.8 | 23.6±0.7 | 48.5±0.7 | 37.0±2.1 | 33.7±0.6 | 55.5±2.0 | 26.0±0.8 | 54.2±1.4 | 74.1±0.3 | 42.8±1.0 | 43.8±0.0 |
| | 4 | 45.5±0.3 | 42.9±2.2 | 23.5±0.8 | **50.6**±2.6 | 38.5±0.8 | 34.1±0.9 | 57.7±0.3 | 26.7±1.3 | 55.8±1.9 | 73.9±0.2 | 37.3±0.9 | 44.9±0.5 |
| | 5 | 44.9±0.6 | 42.9±1.2 | 23.5±0.7 | 50.2±2.5 | **44.0**±0.1 | 34.2±0.7 | 57.3±0.6 | 26.7±0.4 | 54.6±1.8 | 73.6±0.7 | 40.4±0.6 | 45.2±0.1 |
| | 6 | 44.8±1.1 | 43.5±0.6 | 22.8±0.9 | 49.6±2.4 | 43.9±0.4 | 35.8±0.5 | 57.9±1.3 | 25.7±0.0 | 56.1±1.7 | 73.3±0.6 | 40.9±0.7 | 45.3±0.3 |
| | 7 | 43.5±1.6 | 43.7±0.8 | 22.9±1.1 | 50.4±2.7 | 43.4±0.6 | 35.6±0.3 | 56.7±1.3 | 25.7±1.7 | 55.5±2.6 | 73.7±0.4 | 39.9±1.1 | 45.1±0.6 |
| | 8 | 42.0±0.7 | 43.5±1.2 | 23.0±0.7 | 50.3±2.5 | 43.8±0.1 | 35.9±1.5 | 57.1±0.1 | 26.5±1.8 | 56.1±2.9 | 73.9±0.5 | 42.2±0.5 | 45.2±0.2 |
| | 9 | 43.0±0.9 | 43.9±1.2 | 22.2±0.3 | 49.8±2.4 | 43.6±0.2 | 35.2±0.7 | 56.8±0.2 | 25.6±1.2 | 68.3±1.4 | 74.1±1.2 | 40.0±0.4 | 46.2±0.2 |
| | 10 | 42.5±0.7 | 43.5±1.3 | 22.5±0.3 | 49.7±2.5 | 43.6±0.2 | 35.6±1.1 | 55.6±1.0 | 26.2±1.4 | 65.8±4.0 | 72.7±1.0 | 42.8±0.7 | 45.8±0.6 |
| Ours (I_pre replay only) | 1 | 53.3±0.7 | 35.4±1.8 | 24.7±0.1 | 49.7±1.6 | 37.4±1.0 | 32.9±0.2 | 55.6±1.0 | 31.9±1.1 | 55.1±1.2 | 74.1±0.7 | – | 45.0±0.3 |
| | 2 | 52.3±0.3 | **48.0**±2.4 | 22.2±0.4 | 50.0±0.1 | 43.4±0.9 | 34.4±1.4 | 50.3±0.8 | 29.2±1.9 | 63.4±3.5 | 73.2±1.3 | 52.3±0.3 | 46.7±0.5 |
| | 3 | 51.8±1.9 | 43.2±1.6 | 20.5±0.1 | 48.6±0.9 | 40.0±2.1 | 33.1±1.9 | 55.3±0.6 | 27.7±1.5 | 57.8±4.7 | 73.7±0.6 | 47.5±1.1 | 45.2±0.6 |
| | 4 | 52.9±1.3 | 41.9±2.3 | 21.1±0.8 | 49.0±1.5 | 37.9±0.9 | 34.3±1.7 | 54.5±0.6 | 32.3±1.1 | 55.4±0.7 | 72.9±2.0 | 38.6±1.1 | 45.2±0.5 |
| | 5 | 51.5±0.8 | 41.7±1.0 | 21.2±0.9 | 48.8±1.2 | 43.4±0.0 | 35.2±0.5 | 56.4±1.0 | 29.3±0.1 | 53.2±2.4 | 72.2±1.0 | 40.8±0.7 | 45.3±0.5 |
| | 6 | 53.4±1.1 | 44.6±1.2 | 20.5±0.5 | 49.2±1.5 | 44.4±0.6 | 39.0±1.4 | 51.3±5.3 | 30.7±2.4 | 57.3±2.3 | 71.9±1.6 | 42.4±0.3 | 46.3±0.7 |
| | 7 | 52.1±0.5 | 45.5±2.5 | 21.1±0.3 | 49.7±1.1 | 44.0±0.4 | 36.6±1.8 | **62.1**±6.2 | 31.1±0.7 | 60.2±2.5 | 74.8±0.5 | 41.5±0.7 | 47.7±0.1 |
| | 8 | 50.7±2.1 | 47.1±2.5 | 21.0±0.7 | 49.3±1.6 | 44.3±1.7 | 38.2±1.5 | 59.6±7.2 | 26.7±3.0 | 57.0±0.9 | 74.2±0.4 | 44.3±1.4 | 46.8±1.1 |
| | 9 | 51.4±1.4 | 45.6±2.5 | 20.0±0.8 | 49.3±1.4 | 45.8±1.6 | 36.6±1.7 | 56.0±4.0 | 26.6±3.1 | 65.7±5.6 | 73.1±0.5 | 41.4±0.6 | 47.0±0.9 |
| | 10 | 48.7±1.5 | 44.5±3.9 | 21.1±0.3 | 50.1±1.5 | 44.2±1.0 | 35.5±1.9 | 56.8±3.5 | 28.3±3.2 | 65.8±5.4 | 73.0±0.5 | 43.9±0.9 | 46.8±0.8 |
| Ours | 1 | **53.7**±1.3 | 36.6±0.5 | 24.5±0.9 | 49.7±0.8 | 39.7±0.9 | 34.0±2.4 | 56.5±1.5 | 31.7±1.3 | 56.4±0.5 | 74.8±0.5 | – | 45.7±0.2 |
| | 2 | 53.2±0.9 | 46.3±0.7 | 23.2±0.5 | 48.5±1.1 | 41.9±0.9 | 33.7±1.5 | 56.4±1.7 | 30.4±1.2 | 59.1±0.5 | 74.1±0.5 | 53.2±0.9 | 46.7±0.2 |
| | 3 | 52.3±1.1 | 44.0±0.6 | 24.3±2.0 | 49.2±0.5 | 38.5±2.8 | 32.6±0.4 | 53.2±0.9 | 28.0±0.3 | 59.8±5.7 | 73.8±0.8 | 48.2±0.8 | 45.6±0.3 |
| | 4 | 53.5±0.6 | 46.3±1.4 | 24.7±2.9 | 49.1±0.9 | 37.3±3.4 | 34.8±2.5 | 54.8±2.0 | 29.8±1.2 | 59.3±4.0 | 72.9±0.4 | 41.5±0.8 | 46.3±0.6 |
| | 5 | 53.0±1.1 | 44.4±0.8 | 24.8±2.9 | 49.1±0.7 | 43.7±0.3 | 32.7±1.7 | 56.0±1.8 | 29.3±1.3 | 59.0±2.3 | 73.2±0.5 | 42.8±0.8 | 46.5±0.2 |
| | 6 | 53.0±0.9 | 45.0±0.9 | 24.8±2.5 | 49.0±0.2 | 44.1±0.5 | **40.4**±1.5 | 54.1±1.8 | 29.5±2.1 | 60.0±1.9 | 72.8±0.4 | 43.2±0.8 | 47.3±0.7 |
| | 7 | 51.6±0.4 | 44.7±0.5 | 23.8±2.6 | 49.6±0.5 | 44.1±0.3 | 39.2±2.0 | 55.8±0.8 | 28.6±1.8 | 62.1±6.4 | 73.7±0.3 | 42.2±0.8 | 47.3±0.8 |
| | 8 | 50.9±0.3 | 46.0±0.4 | 24.3±2.1 | 49.5±0.2 | 44.1±0.5 | 38.9±1.2 | 54.9±2.1 | 26.2±0.9 | 59.5±2.4 | 74.2±0.2 | 44.1±0.2 | 46.9±0.2 |
| | 9 | 51.6±0.3 | 46.4±1.5 | 23.6±2.1 | 49.0±0.3 | 44.1±0.3 | 37.4±1.4 | 55.4±2.8 | 25.9±0.4 | **68.9**±3.2 | 73.2±0.1 | 41.7±0.2 | 47.6±0.2 |
| | 10 | 50.8±0.4 | 44.6±1.1 | 23.7±2.1 | 49.4±0.1 | 43.8±0.5 | 37.0±1.9 | 54.8±1.8 | 26.1±0.7 | 69.6±1.0 | 72.5±1.6 | 44.3±0.3 | 47.2±0.5 |

Table 8. Detail of the multi-step performance evaluated on the validation set of each scene. At Step $i$, the performance of the adapted network $f_{\theta_i}$ on all the scenes is reported (for scenes $\mathcal{S}_j, j > i$ the values are greyed out). Pre-training denotes the performance of the pre-trained network $f_{\theta_0}$. For each Step $i$, we highlight: in **bold**, the performance of the method which achieves highest mIoU on the current scene $\mathcal{S}_i$, which is indicative of the adaptation performance; in underlined, for each scene $\mathcal{S}_j$, $1 \leq j \leq i-1$ the performance of the method which achieves highest mIoU on $\mathcal{S}_j$, which denotes the ability to preserve previous knowledge; in double-underlined, the performance of the method which achieves highest *average* mIoU on the previous scenes $\mathcal{S}_j, 1 \leq j < i$, which also provides an indication of the ability to counteract forgetting. For each method, the results are averaged over 3 runs. All Ours are with joint training.

through spherical linear interpolation [43] of the rotation components of $\boldsymbol{P}_j^k$ and $\boldsymbol{P}_j^{k+1}$, and its translation component as the average of the translation components of $\boldsymbol{P}_j^k$ and $\boldsymbol{P}_j^{k+1}$. The results of the experiment are shown in Tab. 9, which extends on Tab. 3 from the main paper. All the methods are run for 3 times and mean and standard deviation across the runs are reported.

As can be observed from the Adapt results, replaying from novel viewpoints achieves similar adaptation performance on the current scene as the other baselines of Ours, but with a slightly larger variance.

The crucial observation, however, is that this strategy

outperforms all the other baselines in terms of retention of previous knowledge (Previous) in almost all the steps, and improves on our method with replay of the training viewpoints on average by $0.7\%$ mIoU. This improvement can be attributed to the novel viewpoints effectively acting as a positive augmentation mechanism and inducing an increase of knowledge on the previous scenes. In other words, rather than simply counteracting forgetting, the model de facto keeps learning on the previous scenes, through the use of newly generated data points.

We believe this opens up interesting avenues for replay-based adaptation. In particular, more sophisticated strate-

| | | Step 1 | Step 2 | Step 3 | Step 4 | Step 5 | Step 6 | Step 7 | Step 8 | Step 9 | Step 10 | Average |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Pre-train | | 43.9 | 41.3 | 23.0 | 50.2 | 40.1 | 37.6 | 55.8 | **27.9** | 54.9 | 73.5 | 44.8 |
| Adapt | CoTTA [51] | $44.0_{\pm0.0}$ | $40.9_{\pm0.0}$ | $22.7_{\pm0.0}$ | $50.2_{\pm0.1}$ | $40.0_{\pm0.0}$ | $37.5_{\pm0.0}$ | $56.0_{\pm0.1}$ | $26.9_{\pm0.0}$ | $54.5_{\pm0.0}$ | $\mathbf{73.8}_{\pm0.0}$ | $44.7_{\pm0.0}$ |
| | Mapping [12] | $46.8_{\pm0.4}$ | $42.1_{\pm2.0}$ | $23.6_{\pm0.7}$ | $\mathbf{50.6}_{\pm2.6}$ | $\mathbf{44.0}_{\pm0.1}$ | $35.8_{\pm0.5}$ | $56.7_{\pm1.3}$ | $26.5_{\pm1.8}$ | $68.3_{\pm1.4}$ | $72.7_{\pm1.0}$ | $46.7_{\pm1.2}$ |
| | Ours ($\mathbf{I}_{pre}$ replay only) | $53.3_{\pm0.7}$ | $\mathbf{48.0}_{\pm2.4}$ | $20.5_{\pm0.1}$ | $49.0_{\pm1.5}$ | $43.4_{\pm0.0}$ | $39.0_{\pm1.4}$ | $\mathbf{62.1}_{\pm6.2}$ | $26.7_{\pm3.0}$ | $65.7_{\pm5.6}$ | $73.0_{\pm0.5}$ | $\mathbf{48.1}_{\pm2.1}$ |
| | Ours | $53.7_{\pm1.3}$ | $46.3_{\pm0.7}$ | $\mathbf{24.3}_{\pm2.0}$ | $49.1_{\pm0.9}$ | $43.7_{\pm0.3}$ | $\mathbf{40.4}_{\pm1.5}$ | $55.8_{\pm0.8}$ | $26.2_{\pm0.9}$ | $\mathbf{68.9}_{\pm3.2}$ | $72.5_{\pm1.6}$ | $\mathbf{48.1}_{\pm1.3}$ |
| | Ours (novel viewpoints) | $\mathbf{53.8}_{\pm0.4}$ | $46.7_{\pm2.1}$ | $23.2_{\pm3.3}$ | $49.0_{\pm1.0}$ | $42.9_{\pm0.4}$ | $40.1_{\pm0.7}$ | $58.0_{\pm8.5}$ | $23.2_{\pm2.0}$ | $66.7_{\pm7.1}$ | $71.5_{\pm2.2}$ | $47.5_{\pm2.8}$ |
| Previous | CoTTA [51] | – | $44.0_{\pm0.0}$ | $42.2_{\pm0.0}$ | $35.6_{\pm0.0}$ | $39.3_{\pm0.0}$ | $39.4_{\pm0.0}$ | $39.1_{\pm0.0}$ | $41.5_{\pm0.0}$ | $39.7_{\pm0.0}$ | $41.3_{\pm0.0}$ | $40.2_{\pm0.0}$ |
| | Mapping [12] | – | $46.5_{\pm0.1}$ | $42.8_{\pm1.0}$ | $37.3_{\pm0.9}$ | $40.4_{\pm0.6}$ | $40.9_{\pm0.7}$ | $39.9_{\pm1.1}$ | $42.2_{\pm0.5}$ | $40.0_{\pm0.4}$ | $42.8_{\pm0.7}$ | $41.4_{\pm0.7}$ |
| | Ours ($\mathbf{I}_{pre}$ replay only) | – | $52.3_{\pm0.3}$ | $47.5_{\pm1.1}$ | $38.6_{\pm1.1}$ | $40.8_{\pm0.7}$ | $42.4_{\pm0.3}$ | $41.5_{\pm0.7}$ | $44.3_{\pm1.4}$ | $41.4_{\pm0.6}$ | $43.9_{\pm0.9}$ | $43.6_{\pm0.8}$ |
| | Ours | – | $53.2_{\pm0.9}$ | $48.2_{\pm0.8}$ | $41.5_{\pm0.8}$ | $42.8_{\pm0.8}$ | $43.2_{\pm0.8}$ | $42.2_{\pm0.8}$ | $44.1_{\pm0.2}$ | $\mathbf{41.7}_{\pm0.2}$ | $44.3_{\pm0.3}$ | $44.6_{\pm0.6}$ |
| | Ours (novel viewpoints) | – | $\mathbf{54.8}_{\pm0.9}$ | $\mathbf{50.4}_{\pm2.1}$ | $\mathbf{41.8}_{\pm0.9}$ | $\mathbf{43.8}_{\pm0.8}$ | $\mathbf{43.4}_{\pm0.9}$ | $\mathbf{42.7}_{\pm1.0}$ | $\mathbf{44.8}_{\pm0.9}$ | $41.6_{\pm0.7}$ | $\mathbf{44.3}_{\pm0.2}$ | $\mathbf{45.3}_{\pm0.9}$ |

Table 9. Multi-step performance evaluated on the validation set of each scene. At Step $i$, Pre-train and Adapt denote respectively the performance of the pre-trained network $f_{\theta_0}$ and of the adapted network $f_{\theta_i}$ on the current scene $\mathcal{S}_i$, while Previous represents the average performance of $f_{\theta_i}$ on scenes $\mathcal{S}_1$ to $\mathcal{S}_{i-1}$. All Ours are with *joint training*. Our baseline with novel viewpoints used for replay (Ours (novel viewpoints)) is able to consistently retain knowledge better than the other methods.

gies to select the viewpoints from which to render could be designed, and further increase the knowledge retention on the previous scenes, without reducing the performance on the current scene.

## E. Memory footprint

In the following, we report the memory footprint of the different methods, denoting with $N$ the number of previous scenes at a given adaptation step.

For each previous scene, our method stores the corresponding NeRF model, which has a size of $50.0\,\mathrm{MB}$ with Instant-NGP [33, 45] and of $4.9\,\mathrm{MB}$ with Semantic-NeRF [56]. This results in either $(N \times 50.0)\mathrm{MB}$ or $(N \times 4.9)\mathrm{MB}$ of total data being stored in the *long-term* memory. Note however that during adaptation we only render data from a small subset of views to populate the replay buffer, hence the effective size of the data from the previous scenes that need to be stored in running memory during adaptation is $14.0\,\mathrm{MB}$. Additionally, we save one randomly selected data point every 10 samples in the pre-training dataset, taking up additional $64.6\,\mathrm{MB}$ of space.

Similarly to us, the method of [12] requires $14.0\,\mathrm{MB}$ for the replay buffer and $64.6\,\mathrm{MB}$ for the replay from the pre-training dataset, but stores voxel-based maps instead of NeRF models, taking up $71.8\,\mathrm{MB}$ for each scene. Importantly, since the voxel-based maps only include semantic information and cannot be used to render *color* images, the method of [12] additionally needs to save color images for the training viewpoints. In the 10 scenes that we used for our experiments, their size amounted on average to approximately $30.0\,\mathrm{MB}$ per scene, resulting in a total storage space of around $(N \times 101.8\,\mathrm{MB})$ required for the previous scenes.

In each step, in addition to the model that gets adapted on the current scene, CoTTA [51] requires storing the teacher model from which pseudo-labels for online adaptation are generated, and an additional version of the original, pre-
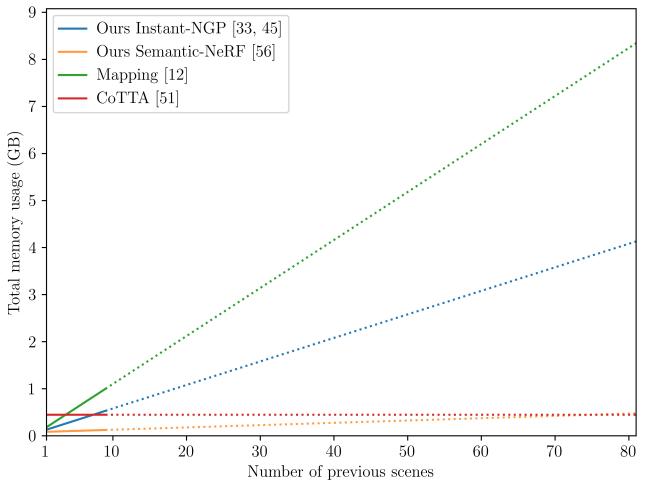


Figure 4. Memory footprint of the different methods as a function of the number of the previous scenes. Please refer to the text and to Tab. 10 for a detailed explanation. We use solid lines for the number of scenes used in our experiments.

trained model, to preserve source knowledge. The parameters of the DeepLab network used in our experiments have a size of $224.3\,\mathrm{MB}$, resulting in a total of $(2 \times 224.3)\mathrm{MB} = 448.6\,\mathrm{MB}$ of data that need to be stored.

A comparison of the memory footprint of the different methods as a function of the number of previous scenes can be found in tabular form in Tab. 10 and in graphical form in Fig. 4. For our method and for [12], we include in the total size both the data stored offline and the one inserted in the replay buffer.

Note that using the lighter implementation of Semantic-NeRF [56], the comparison is in our favour up to 75 scenes, and up to 91 scenes when only considering the size of the NeRF models.

| | | Previous scenes | | Source knowledge | Total |
|---|---|---|---|---|---|
| | | Offline | Online | | |
| Ours | Instant-NGP [33, 45] | $(N \times 49.9)$MB$^\star$ | 14.0 MB | 64.6 MB | $(78.6 + N \times 49.9)$MB |
| | Semantic-NeRF [56] | $(N \times 4.9)$MB$^\star$ | | | $(78.6 + N \times 4.9)$MB |
| CoTTA [51] | | – | 224.3 MB$^\dagger$ | 224.3 MB$^\dagger$ | 448.6 MB |
| Mapping [12] | | $\sim (N \times 101.8)$MB$^{\star\star}$ | 14.0 MB | 64.6 MB | $\sim (78.6 + N \times 101.8)$MB |

Table 10. Comparison of the memory footprint of different methods. $N$ denotes the number of previous scenes. $^\star$ The numbers refer to the storage cost required by the NeRF models. For actual adaptation (Online), only renderings from a subset of views are used, and inserted in a memory buffer of size 14.0 MB. $^{\star\star}$ The numbers refer to the storage cost required by each voxel-based map (71.8 MB), plus the explicit training views that need to be stored for each scene, which amount to an average of $\sim 30.0$ MB per scene. Similarly to Ours, for actual adaptation, a memory buffer of size 14.0 MB is used. $^\dagger$ CoTTA requires storing a teacher model for online adaptation, and an additional version of the original, pre-trained model, to preserve source knowledge.

## F. Further visualizations

In Fig. 5 we provide examples of the pseudo-labels produced on the training views by our method and by the different baselines. As previously observed by the authors of [12], the mapping-based pseudo-labels suffer from artifacts induced by the discrete voxel-based representation. Thanks to the continuous representation enabled by the coordinate-based multi-layer perceptrons, our NeRF-based pseudo-labels produce instead smoother and sharper segmentations. However, they occasionally fail to assign a uniform class label to each object in the scene (cf. last row in Fig. 5). This phenomenon, which we also observe in the mapping-based pseudo-labels, can be attributed to the inconsistent per-frame predictions of DeepLab, that cannot be fully filtered-out by the 3D fusion mechanism. By jointly training the per-frame segmentation network and the 3D-aware Semantic-NeRF, we are however able to effectively reduce the extent of this phenomenon, producing more uniform pseudo-labels.

Figure 6 further shows examples of the predictions returned by the segmentation network on the validation views after being adapted using the different methods. In accordance with what observed in the quantitative evaluations, while being able to preserve knowledge, CoTTA achieves limited improvements with respect to the initial performance. As a consequence, the predicted labels match very closely those of the pre-trained network. Fusing the predictions from multiple viewpoints into a 3D representation allows both the baseline of [12] and our method to reduce the amount of artifacts due to misclassifications in the per-frame predictions. The positive effect of this 3D fusion can be successfully transferred to the segmentation network through adaptation, as visible by comparing the predictions in the three rightmost columns of Fig. 6 to those of the pre-trained network (third column from the left in Fig. 6). We observe that fine-tuning with NeRF-based pseudo-labels instead of voxel-based pseudo-labels often results in a more consistent class assignment to different pixels of the same instance. This effect is amplified when using joint training, which often produces more accurate pseudo-labels compared to fine-tuning.

## G. Limitations

Since our approach relies on the assumption that a good reconstruction of the scene can be obtained, we find that our method achieves suboptimal performance when this assumption is not fulfilled. This is the case for instance for Scene 5 (cf. first and second row in Fig. 7), in which a large number of frames are overexposed and the ground-truth depth measurements are missing for a large part of the frame. Specular effects (cf. third row in Fig. 7) can further break the assumptions required by the volume rendering formulation of NeRF. Related to these problems is also the quality of the initial predictions of the segmentation network: Particularly when lighting conditions are poor, we observe that the predictions of the pre-trained segmentation network are very noisy (see, *e.g.*, first row in Fig. 7). The combination of these factors results in the pseudo-labels produced by our method assigning a uniform label to a large part of the scene and failing to correctly segment smaller details.

We observe that these degenerate cases can have a particularly large influence on the quality of the pseudo-labels and of the network predictions when jointly training the segmentation network and NeRF. We hypothesize that this might be due to the 2D-3D knowledge transfer enabled by our method inducing a negative feedback loop when poor segmentation predictions are combined with suboptimal reconstructed geometry. A possible way to tackle this problem in future work is by making use of regularization techniques, for instance by limiting large deviations of the predictions across adaptation steps, to avoid collapse, or by minimizing the entropy of the semantic predictions of both NeRF and the segmentation network.

A general limitation of our method is that it assumes scenes to be static. Extending the pipeline to handle dynamic scenes through the use of temporally-aware NeRFs [36, 37] is an interesting direction for future work.
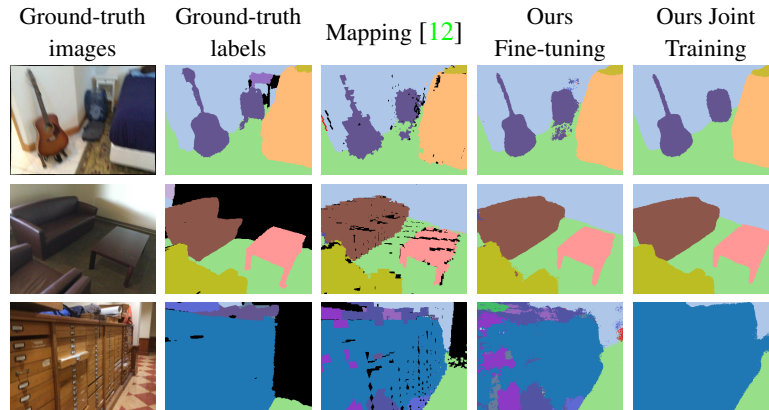
Figure 5. Comparison of example pseudo-labels obtained on the training views by the different methods.
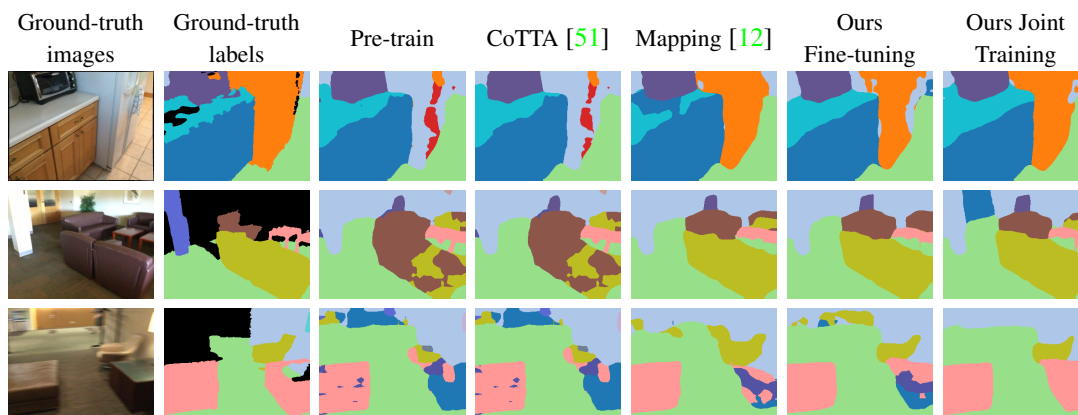


Figure 6. Comparison of the predictions of the semantic segmentation network on the validation views, when adapted using the different methods.
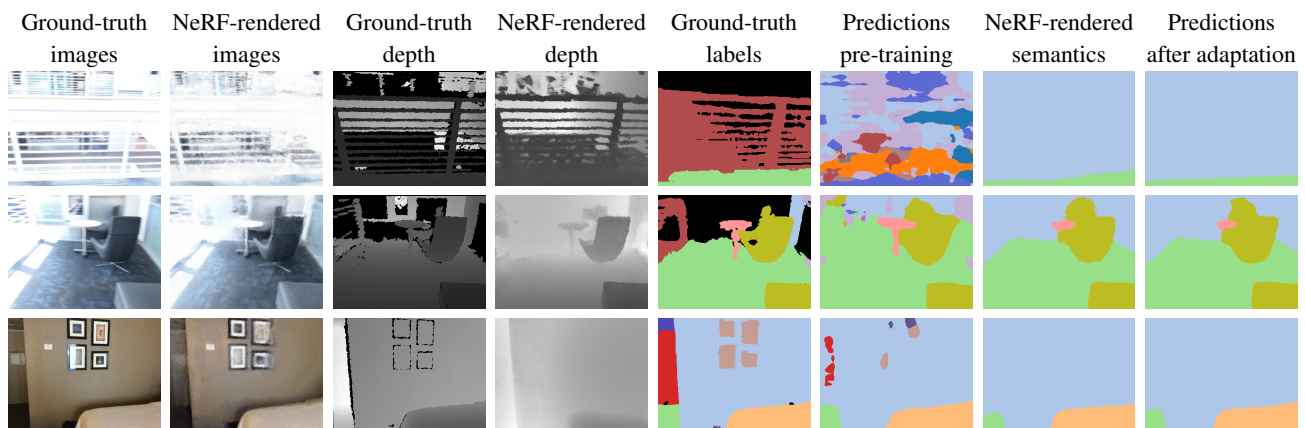


Figure 7. Examples of failure cases. First row: Poor lighting, incomplete ground-truth depth measurements, and noisy initial predictions result in both the rendered pseudo-labels and the predictions from the adapted network assigning uniform labels to large parts of the scene and failing to correctly segment fine details in the scene. Second row: Motion blur, diffusion lighting, shadows, and insufficient number of observations can also degrade the reconstruction quality and make the label propagate to the wrong objects. Third row: Specular effects can break the assumptions of the volume rendering formulation of NeRF; large flat areas with small variations in depth can be hard to reconstruct, resulting in smoothed-out, uniform labels for the background.