

# Continual Coarse-to-Fine Domain Adaptation in Semantic Segmentation

Donald Shenaj, Francesco Barbato, Umberto Michieli, Pietro Zanuttigh

*Department of Information Engineering, University of Padova*

*{donald.shenaj, francesco.barbato, umberto.michieli, zanuttigh}@dei.unipd.it*

## Abstract

Deep neural networks are typically trained in a single shot for a specific task and data distribution, but in real world settings both the task and the domain of application can change. The problem becomes even more challenging in dense predictive tasks, such as semantic segmentation, and furthermore most approaches tackle the two problems separately. In this paper we introduce the novel task of coarse-to-fine learning of semantic segmentation architectures in presence of domain shift. We consider subsequent learning stages progressively refining the task at the semantic level; *i.e.*, the *finer* set of semantic labels at each learning step is hierarchically derived from the *coarser* set of the previous step. We propose a new approach (CCDA) to tackle this scenario. First, we employ the maximum squares loss to align source and target domains and, at the same time, to balance the gradients between well-classified and harder samples. Second, we introduce a novel coarse-to-fine knowledge distillation constraint to transfer network capabilities acquired on a coarser set of labels to a set of finer labels. Finally, we design a coarse-to-fine weight initialization rule to spread the importance from each coarse class to the respective finer classes. To evaluate our approach, we design two benchmarks where source knowledge is extracted from the GTA5 dataset and it is transferred to either the Cityscapes or the IDD datasets, and we show how it outperforms the main competitors.

*Keywords:* Coarse-to-Fine Learning, Unsupervised Domain Adaptation, Semantic Segmentation, Continual Learning, Deep Learning

## 1. Introduction

A fundamental aspect for any intelligent system is the ability to understand the surrounding environment. Semantic segmentation, *i.e.*, the dense (pixel-level) classification of an input image, is a key tool to achieve this target. Deep Convolutional Neural Networks, starting with the seminal work on FCNs by Long *et al.* [1] achieve impressive performance on this challenging task. Such architectures are very effective when trained on i.i.d. samples coming from a single domain, however, they suffer from an important shortcoming: they tend

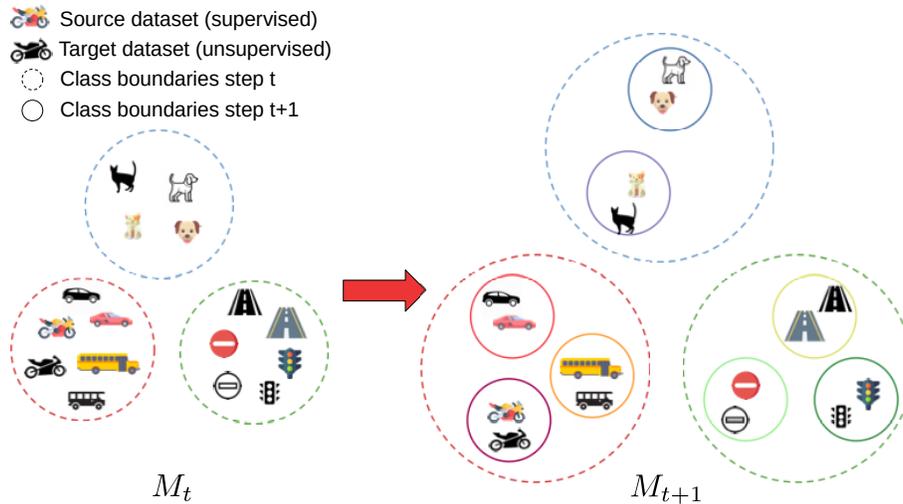


Figure 1: In semantic coarse-to-fine learning we train a model in subsequent steps to recognize incremental sets of finer classes over time, that are derived hierarchically from previous coarse ones. Similarly to UDA, labeled samples come from a source domain, which is different from the (testing) target domain, whose labels are unknown at training time.

to fail when even a small distribution shift occurs over the input domain distribution or over the output label space. Particularly, semantic segmentation solutions trained on a source domain show clear signs of *domain shift* when tested on a target test domain, limiting the performance [2]. A similar effect is found when training architectures in incremental steps, where *catastrophic forgetting* leads to a clear performance degradation [3]. Such effects impair significantly the employment of deep neural networks in real-world settings, where they are expected to adapt their knowledge to the specific encountered environment, while at the same time maintaining the ability of achieving subsequent levels of understanding of the scene incrementally. To address these undesirable effects, a multitude of Unsupervised Domain Adaptation (UDA) and Continual Learning (CL) approaches have been recently investigated to tackle domain shift and catastrophic forgetting, respectively.

In this work, we propose a joint continual and domain adaptive approach, placing an additional brick towards the development of deep learning architectures which can be deployed with intelligent agents. We consider coarse-to-fine learning, where the label sets are incrementally refined and classes are hierarchically derived from the previous ones, as it often occurs during the manual labeling by humans. A first, rough, segmentation comprised of few coarse classes is provided; then, it is further refined by splitting the coarse classes into finer categories. Moreover, similarly to continual learning, we assume that the set of classes to be recognized can change during training, as it is often unfeasible to retrain a segmentation architecture on the whole class set each time a new fine class is introduced. This incremental behavior is applied by masking classes

that have reached their maximum detail level in previous class set refinements, forcing the architecture to preserve previously acquired knowledge. Finally, to better mimic the deployment of an actual intelligent system, we evaluate our setup in the presence of domain shift. An overview of the considered task is shown in Figure 1, where we appreciate how the coarser classes are refined into finer categories by hierarchically splitting them, while adapting the network at both steps. The considered task, together with the mathematical notation needed to describe it, is reported in detail in Section 3. To tackle this task we introduce a novel approach named CCDA (Continual Coarse-to-fine Domain Adaptation), detailed in Section 4. In CCDA we employ a custom-designed coarse-to-fine knowledge distillation loss (see Section 4.2.1) and an unbiased coarse-to-fine weight initialization technique (see Section 4.2.2) together with an entropy minimization-based approach for unsupervised domain adaptation [4] which, can regularize the gradient evolution and better preserve less common classes during the adaptation procedure (see Section 4.1). Importantly, we apply the domain adaptation approach in all the incremental training steps, obtaining a network with improved generalization ability. This is due to the effect of self-supervision in earlier steps, which is very effective when a restricted number of classes is considered.

Given the novelty of the task, in order to evaluate our approach we design two new benchmarks, derived from common semantic segmentation synthetic-to-real unsupervised domain adaptation tasks. In particular, we reframed the GTA5 [5], Cityscapes [6] and IDD [7] datasets into class-hierarchical coarse-to-fine labeling datasets. These benchmarks ( $GTA5 \rightarrow Cityscapes$  and  $GTA5 \rightarrow IDD$ ) are then employed to evaluate our approach and compare it with other competing methods. The implementation details are discussed in Section 5, while the quantitative and qualitative results obtained are reported in Section 6.

## 2. Related Work

**Semantic Segmentation** has gained wide popularity over the last few years. A multitude of Convolutional Neural Network-based (CNN) segmentation architectures has been proposed for this task starting from FCN [1] to more recent approaches like SegNet [8], PSPNet [9] and DeepLab [10, 11, 12]. These approaches are supported by a fast-growing number of datasets such as Cityscapes [6], IDD [7], GTA5 [5], *etc.* Such architectures are very effective when trained on i.i.d. samples coming from a single domain distribution, however they tend to fail when distribution shifts occur over the input domain distribution or the output label space.

**Coarse-to-Fine Learning:** hierarchical composition of semantic representations has been explored in few previous work for part-based regularization [13, 14] and coarse-to-fine approaches where coarse-level classes are refined into finer categories [15, 16, 17, 18]. In particular, coarse-to-fine learning in semantic segmentation can be addressed by resorting either to a direct concatenation of features [15] or using a semantic embedding network to transfer object-level predictions during part-level training [18]. However, approaches for this task

all assume that the domain distribution remains unaltered at subsequent refinement stages and they are comprised of one single refinement stage.

**Continual Learning (CL)** is a more mature field where previous knowledge acquired on a subset of data is further expanded to learn new tasks (*e.g.*, class labels) [19, 20, 21]. Continual learning has recently been investigated also in semantic segmentation [3, 22, 23, 24, 25, 26]. The most popular strategy to retain knowledge from the previous learning step is the usage of knowledge distillation constraints at either the output [3, 23, 24, 22] or feature level [3, 25]. Another promising research direction involves feature-level regularization to increase separation among features from different classes [23]. Additionally, replay samples of past classes can also be either generated via a GAN or downloaded from the Web [26]. However, all these approaches do not consider classes hierarchically derived from macro-classes. The strongest similarity with most of these works lies in the special *background* class from which the novel classes are extracted in these approaches. Indeed, the *background* at a certain learning step could contain future classes, therefore at subsequent steps classes are derived from this unique macro-class. However, all these approaches deal only with the semantic shift of one class (*i.e.*, the background) and they do not tackle data domain shift (*i.e.*, source and target domains).

**Unsupervised Domain Adaptation** is a more challenging variant of domain adaptation [27, 2], where no labeled samples coming from the target domain are exploited for the adaptation procedure. These techniques deal with the reduction of the domain shift effect and, therefore, lead to a performance improvement in the target domain. While originally proposed to tackle such problems in other settings, like image classification, the rise of segmentation architectures made the transposition of adaptation techniques to such task a very desirable objective. Following the encoder-decoder structure of common semantic segmentation networks, a multitude of UDA techniques has been proposed acting on different network levels: *i.e.*, at the *input*, *feature* and *output* levels [2].

Works at input level usually resort to image-to-image translation [28, 29, 30, 31, 32, 33], or, more recently, to image-inpainting [34]. The former aims to reduce the domain shift by acting directly on the input images, either via style transfer or by generative architectures, while the latter exploits carefully designed cross-domain samples. Particularly, in [34] the network is adapted via a mixture of self-training and image inpainting.

Feature-level techniques reduce the domain shift by aligning the statistical distributions of the two domains exploiting the encoder-decoder latent representations. Recent works in this category [35, 36, 37] investigated various ways towards the goal, such as self-supervised and unsupervised orthogonal embeddings, source-target alignment objectives, prototype-driven clustering, vector sparsity and norm matching.

Finally, works at the output level fall into two major categories: entropy minimization and self-supervised training, both via direct approaches [38, 39] or by exploiting curriculum-learning techniques [40, 41]. In [38, 39] the network is forced to produce confident predictions in the target domain, similarly

to the behavior on the source domain, by acting on the logits produced by the segmentation architecture. In [40, 41] the network predictions on the target domain reinforce the training and prediction confidence. Other recent and high-performing techniques, such as [42], use a mixture of all the families of techniques. Here, the authors exploit image-inpainting, feature-level adaptation, self-supervised training and a modified segmentation architecture to surpass previous state-of-the-art approaches.

### 3. Task Formulation

In this section we provide a detailed overview of the proposed task along with the mathematical notation needed to describe it. We start by introducing the employed notation in Section 3.1, then in Section 3.2 we report an overview of the considered benchmarks describing their adaptation to our task; and in Section 3.3 we detail the hierarchical semantic class splits identified in such benchmarks.

#### 3.1. Problem Setup

In this work we consider semantic segmentation, which performs a pixel-level labeling of input images. Let us define the (input) image space as  $\mathcal{X} \subset \mathbb{R}^{H \times W \times 3}$ , which represents the set of RGB images of width  $W$  and height  $H$ . Together with the input we identify the corresponding pixel-level (output) labels, which belong to the space  $\mathcal{Y}_t \subset \mathcal{C}_t^{H \times W}$ , where  $t$  is a task index accounting for the changes in the label set  $\mathcal{C}$  over time. The semantic segmentation task can be expressed as the problem of finding a map  $M : \mathcal{X} \ni \mathbf{X} \mapsto \mathbf{P}_t = M_t(\mathbf{X}) \in [0, 1]^{H \times W \times |\mathcal{C}_t|}$  such that  $\hat{\mathbf{Y}} := \arg \max_{\mathbf{Y} \in \mathcal{Y}_t} \mathbf{P}_t = \mathbf{Y}$ . In recent semantic segmentation architectures the map  $M$  takes the form of an encoder-decoder network, *i.e.*,  $M_t = D_t \circ E_t$ . This allows to identify the latent features extracted by the encoder and used for the final predictions as  $\mathbf{F} = E_t(\mathbf{X})$  for any image  $\mathbf{X} \in \mathcal{X}$ . Indeed, as opposed to the standard supervised training where the whole class set is available from the beginning, in a continual setting the final class set  $\mathcal{C} = \bigcup_t \mathcal{C}_t$  is divided into subsets that are learned separately as tasks  $t = \{0, 1, \dots\}$ . In a coarse-to-fine setting, we can make a further separation: the class sets of each task can be divided as  $\mathcal{C}_t = \mathcal{C}_t^c \cup \mathcal{C}_t^f$ , where  $\mathcal{C}_t^c$  contains the classes that will be split into finer ones in the future and  $\mathcal{C}_t^f$  contains classes that have already reached the maximum level of refinement (*i.e.*, they will not be split anymore). The split between  $\mathcal{C}_t^c$  and  $\mathcal{C}_t^f$  is unknown at stage  $t - 1$ , but it becomes available when the subsequent training stage begins. At that time (task  $t$ ), the coarse set of classes of the previous task  $t - 1$  are mapped to a set of derived finer classes with the map  $S_t : \mathcal{C}_{t-1}^c \ni c \mapsto S_t(c) \subset \mathcal{C}_t$ . Therefore, at each step the last layer of the decoder  $D_t$  must change in order to accommodate the new classes, as we will detail in Section 4. Then, we can define the training set of task  $t$  as  $\mathcal{T}_t = \mathcal{T}_t^S \cup \mathcal{T}_t^T$ , where  $\mathcal{T}_t^S = \{(\mathbf{X}^S, Y^S) \stackrel{\mathcal{D}^S}{\sim} (\mathcal{X} \times \mathcal{Y}_t)\}$  is the (supervised) source

dataset, while  $\mathcal{T}_t^T = \{\mathbf{X}^T \overset{\mathcal{D}^T}{\sim} \mathcal{X}\}$  is the (unsupervised) target one. Note that, while the contents of  $\mathcal{T}_t^S$  and  $\mathcal{T}_t^T$  are relative to the same image space, the distributions according to which they are sampled varies (namely,  $\mathcal{D}^S$  for the source domain distribution and  $\mathcal{D}^T$  for the target one), leading to the presence of domain shift.

### 3.2. Datasets

For the evaluation of our strategy we use three different datasets, *i.e.*, *GTA5*, used as source dataset for supervised training and *Cityscapes* and *IDD*, used as target ones. This allows to build two synthetic-to-real semantic segmentation benchmarks: *GTA5*→*Cityscapes* and *GTA5*→*IDD*. In these benchmarks the segmentation architecture is trained in a supervised manner using samples exclusively from the *GTA5* [5] dataset and its performance is evaluated on *Cityscapes* [6] and *IDD* [7], respectively.

The **GTA5** [5] dataset provides 25000 synthetic images generated using the homonym video game, which is set in American cities. The images have varying resolutions close to  $1920 \times 1080$ , an aspect ratio of  $16 : 9$  and are paired with semantic segmentation labels with the same class set as *Cityscapes* (total of 35 classes, 19 of which are used for training).

The **Cityscapes** [6] dataset, instead, provides 2500 finely-labeled real images captured in multiple European cities and an additional set of 20000 coarsely-labeled samples from the same cities. The images have resolution  $2048 \times 1024$  (aspect ratio of  $2 : 1$ ) and are annotated using 35 classes.

Finally, the **IDD** [7] dataset provides 10004 finely-labeled samples captured in India. The images have a smaller resolution than the previous datasets ( $1280 \times 720$ , aspect ratio of  $16 : 9$ ) and the labels are provided in a different class set (36 classes, 27 of which are used in the training) leading to a different number of classes when performing closed-set domain adaptation: in particular the *train* and *terrain* classes are missing.

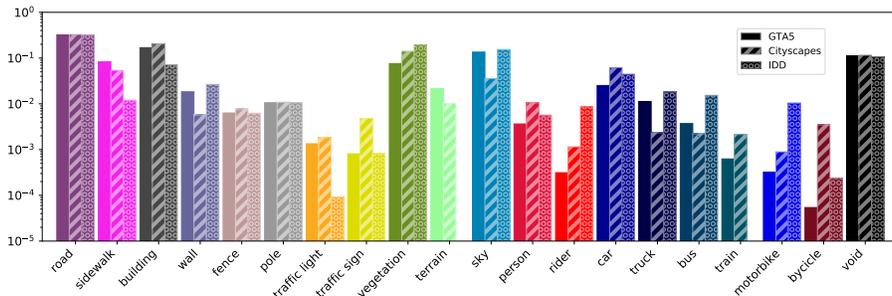


Figure 2: Relative pixel-wise class frequencies in each considered dataset (the plot uses a logarithmic scale). GTA5 in solid colors, Cityscapes in dashed colors, IDD in dotted colors (*train* and *terrain* are not present in IDD).



Figure 3: Samples from the three considered datasets highlight the large visual domain shift among them.

Autonomous driving datasets are inherently highly class-imbalanced due to the variable amount of pixels each class can have. We show the relative pixel class frequency of the three datasets in Figure 2, where the significant distribution shift and the missing IDD classes are clearly visible.

Even more, datasets acquired in different geographical regions and experimental scenarios have profound domain distribution shifts. In Figure 3 we report a couple of samples from each dataset, noticing how domain shift is evident from the visual aspect of the images. The images from the GTA5 dataset have a very crisp look, typical of synthetic datasets, even when the weather is cloudy (second row of first column); while the images coming from the real datasets are slightly blurry, side effect of the motion blur and of the real camera optics, which cannot focus perfectly on the whole scene. Even more evident is the difference in coloring between the datasets: GTA5 offers visually-pleasing well-saturated colors (which are inherited from the game engine used to generate the samples) while in IDD (second column) the colors are much dimmer. Even worse is the Cityscapes dataset (third column), where the white point is unbalanced and the resulting samples appear greenish and grayish.

### 3.3. Coarse-to-Fine Hierarchical Splits

To validate our setup, we devise a way of extending the benchmarks into class-hierarchical coarse-to-fine datasets. In particular, we focus on the 19 trainable classes of the Cityscapes dataset, developing a semantic grouping based both on the similarity of the classes and on common semantic concepts (*i.e.*, the *vehicle* class can be naturally split into finer classes like *two-wheels*, while simpler classes like *sky* reach the maximum level of detail earlier and cannot be further split). A graphical representation of our hierarchical split is reported in Figure 4, where different depths on the tree represent different incremental steps during the training procedure. The boxes highlighted with a blue border represent the original 19 classes from Cityscapes and, therefore, will be masked in subsequent steps, having reached their maximum refinement level. Classes

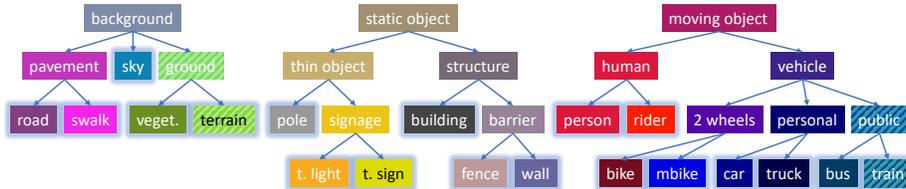


Figure 4: Hierarchical label splits: each level of depth in the tree corresponds to a different learning step. At each step, the provided ground truth segmentation maps contain only the labels of the given step, while all the other pixels are unlabeled. Boxes highlighted in blue contain the original Cityscapes classes. Classes in dashed boxes are not present in the IDD dataset (*terrain* and *train* are missing leading to the collapse of *ground* in *vegetation* and of *public transport* in *bus* at the second step).

in dashed boxes, instead, are not present in the IDD split (*i.e.*, *ground*, *terrain*, *public transport* and *train*). Such difference with respect to Cityscapes is due to the fact that *terrain* and *train* are missing, thus leading to the collapse of *ground* in *vegetation* and of *public transport* in *bus* at the second and third steps, respectively.

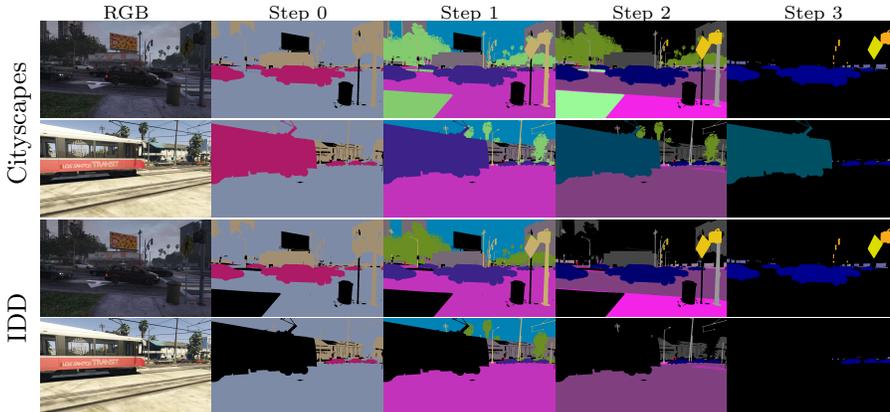


Figure 5: Continual coarse-to-fine annotation maps on the GTA5 dataset over different learning steps. These are the images used for supervised training. First two images reported in the Cityscapes class-split, second two in the IDD one.

In Figure 5 we show a qualitative example of the per-step labels provided to the architecture during training, notice how the fine classes are present only in one step and are masked in subsequent ones. Additionally, we can see how since the *terrain* in the first row and the *train* in the second row are missing in the IDD dataset, we masked them to *void* from the very first step when the IDD split is used, as can be verified by third and fourth row. On Cityscapes, instead, these classes are present (first two rows). This results in a different pixel-level class distribution between the two benchmarks.

Performing the training in an incremental manner dramatically changes the pixel-class distributions seen by the architecture during training and signifi-

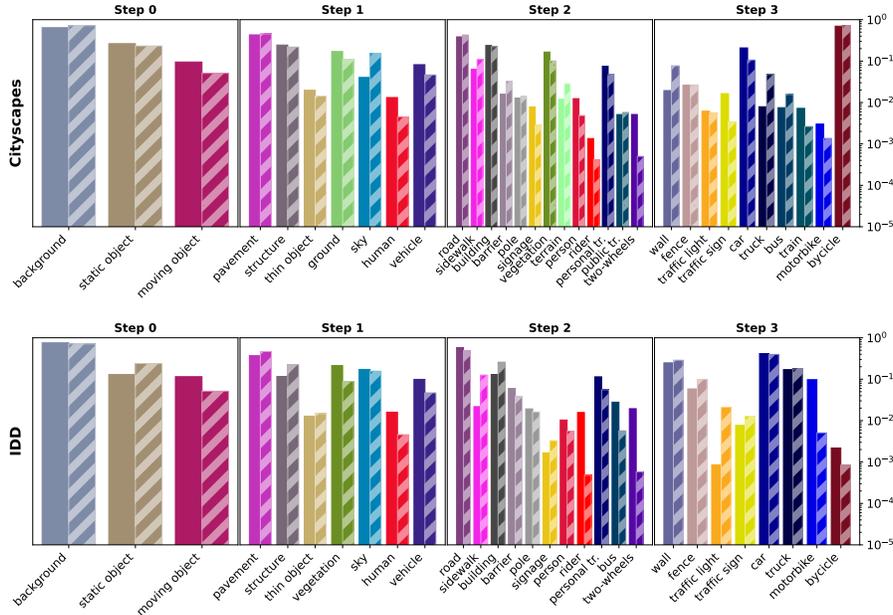


Figure 6: Relative frequency of pixels of given classes in each incremental step for both the considered target datasets. The classes at the incremental steps are refined according to the label splits of Figure 4

cantly bridges the domain gap. This is in part due to the effect of masking common old fine classes (*i.e.*, with many pixels), which allows to remove their contribution and leads to a much more uniform distribution, and in part due to the semantic grouping we propose, which is designed to reduce the distribution shift between source and target domain in each of the incremental steps. In Figure 6 we report the relative class frequency of each incremental step comparing the source and target distributions. From visual inspection it is clear that the distributions are much more aligned (both in the  $GTA5 \rightarrow Cityscapes$  and in the  $GTA5 \rightarrow IDD$  setups) than the complete distributions shown in Figure 2. This empirical finding is quantitatively confirmed when computing the KL-Divergence on the distributions. The original distribution scores 0.1361 and 0.2415 on  $GTA5 \rightarrow Cityscapes$  and  $GTA5 \rightarrow IDD$ , respectively. On the other hand, the incremental steps score  $\{0.0216, 0.0925, 0.0645, 0.1138\}$  (average 0.0731) and  $\{0.0604, 0.1374, 0.2556, 0.2286\}$  (average 0.1705), respectively. Such scores imply an average reduction in the distribution shift of 46.3% when using Cityscapes as target and of 29.4% when using IDD.

#### 4. Proposed Approach

In this section, we provide a detailed description of our proposed method: CCDA, Continual Coarse-to-fine Domain Adaptation. Our approach combines

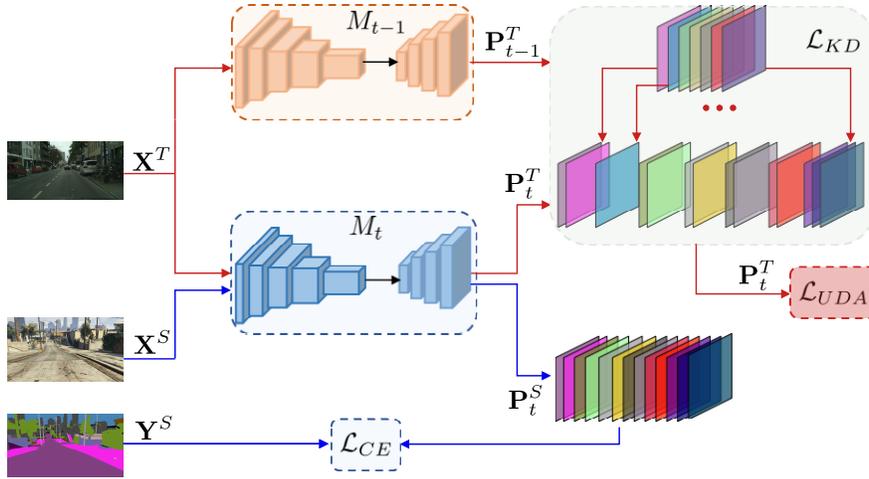


Figure 7: Overview of the proposed approach. The new refined classes are learned on the supervised source domain through  $\mathcal{L}_{CE}$ . At the same time, the predictions at current step are adapted to the unsupervised target domain with  $\mathcal{L}_{UDA}$ . Finally, the coarse-to-fine distillation module  $\mathcal{L}_{KD}$  is used to semantically adapt the knowledge from coarse to finer classes, while preserving the knowledge of the previously learned ones.

continual learning techniques with domain adaptation ones allowing incremental learning of semantically refined (finer) classes, while preserving the knowledge of older (coarse) classes and the semantic relation between any coarse class and its fine child classes. Moreover, to reduce the distribution gap between source and target domains, we include ad-hoc unsupervised domain adaptation techniques. Such strategies are employed jointly with the continual learning ones, allowing to incrementally adapt the architecture from simpler to harder settings (the number of classes increases throughout training) while at the same time improving the adaptation performance. A graphical representation of the proposed approach is shown in Figure 7, highlighting the three modules that make up our strategy. Here we see how the cross-entropy, knowledge distillation and domain adaptation objectives are used in conjunction with the previous step’s model and with source and target samples to optimize and adapt the current step’s architecture. The contribution of the various modules is summarized in the loss function:

$$\mathcal{L} = \mathcal{L}_{CE} + \lambda_{UDA} \mathcal{L}_{UDA} + \lambda_{KD}^c \mathcal{L}_{KD}^c + \lambda_{KD}^f \mathcal{L}_{KD}^f, \quad (1)$$

where  $\mathcal{L}_{CE}$  is the cross entropy loss, used to learn new classes in the supervised source domain. The other loss terms will be described in detail in the following subsections: the domain adaptation ( $\mathcal{L}_{UDA}$ ) loss based on maximum squares minimization is discussed in Section 4.1, while the coarse-to-fine knowledge distillation loss ( $\mathcal{L}_{KD}^c$  and  $\mathcal{L}_{KD}^f$ ) is introduced in Section 4.2.1. The hyper-

parameters  $\lambda_{UDA}$ ,  $\lambda_{KD}^c$  and  $\lambda_{KD}^f$  control the relative contribution of each component, and have been optimized on a validation set extracted from the training data of the target domains (as detailed in Section 5). Finally, the coarse-to-fine weight initialization rule is discussed in Section 4.2.2.

#### 4.1. Unsupervised Domain Adaptation

In this section we overview the UDA module of our method, presenting the problems that sparked its design. In a typical (closed-set) UDA setting we are interested in performing adaptation between two domains sharing a fixed set of classes in a single step. However, when the set of classes is large, when the target dataset does not provide enough distinctive cues to differentiate between classes or when such cues are too different from those learned on the source domain, the adaptation process has sub-optimal performances. A further cause of performance degradation is the inherent penalization of less common classes, which are learned through few very specific features (given the small number of samples and, therefore, the limited variability available) which can easily be missing or different in the target domain.

To tackle these issues, a good choice is to perform UDA in a curriculum learning fashion, starting from the adaptation of high level semantic concepts and transferring, in a multitude of steps, the finer and finer representations of the classes in a format compatible with the target dataset. The result of this approach is a more gradual and resilient adaptation, since the number of the new classes to be adapted is always smaller than in the joint case. This reduced number of classes can lead to significant improvements in the adaptation procedure, particularly when self-supervised techniques are employed. This is due to the reduced entropy of the prediction system, which having a restricted number of possible outputs has also a much smaller number of possible false predictions. A possible solution to these problems is the the maximum squares minimization loss [4], which we reframe for use in the proposed novel coarse-to-fine task. Formally, at each training step  $t$ , the objective to be optimized is:

$$\mathcal{L}_{UDA} = - \sum_{\mathbf{X} \in \mathcal{T}_t} \sum_{c \in \mathcal{C}_t} \frac{1}{2|\mathcal{C}_t|^\alpha (HW)^{1-\alpha}} (\mathbf{P}_t^T[c])^2, \quad (2)$$

where  $|\mathcal{C}_t|$  represents the number of classes at each training step  $t$ ,  $H$  and  $W$  are the image dimensions (*i.e.*,  $HW$  is the number of pixels), and  $\alpha$  is an hyper-parameter. Recall, also, that  $\mathbf{P}_t = M_t(\mathbf{X})$  for any given input image  $\mathbf{X}$ . At the end of each incremental step the model will be able to produce increasingly accurate predictions on the target domain, which are crucial for the successful domain-bridging effect of our knowledge distillation module (see Section 4.2.1).

#### 4.2. Coarse-to-Fine Learning

In this section we report in detail the two continual learning modules employed by our architecture: in Section 4.2.1 we present our novel cross-domain

coarse-to-fine knowledge distillation strategy, while in Section 4.2.2 we show our unbiased coarse-to-fine weight initialization strategy.

#### 4.2.1. Coarse-to-Fine Knowledge Distillation

Knowledge distillation is an established technique typically used in continual learning settings to preserve previous learned knowledge. Given its effectiveness, we extended the standard definition for use in our coarse-to-fine task. This technique assumes that at each training step the previous learned model is available, which will be used to guide the preservation of the learned knowledge on the previous classes. Storing the previous model, does not compromise the privacy of the training procedure, because we do not retain any previous image nor label used during training, and the model will be used just to perform inference on new training images (as it would do in any deployment setting).

Our key insight in improving the knowledge transfer across coarse and fine class sets, was to realize the importance of constraining the prediction of the fine classes into resembling their parent coarse class (*i.e.*, the one from which they have been derived). With this approach we benefit both in terms of accuracy in the recognition of new classes, and in the reduction of misclassification occurrences (particularly between semantically unrelated classes, *i.e.*, those that do not share a common ancestor in the hierarchical split, see Fig. 4).

Hence, at each training step  $t$ , we force the output of the previous model  $M_{t-1}$  at channel (coarse class)  $c$ , to match to the sum of the output channels  $f \in S_t(c)$  (*i.e.*, the set of  $c$ 's child classes) of the current model  $M_t$ . The matching is implemented with a cross entropy loss function, resorting to pseudo-labels computed on target domain samples, effectively enforcing an additional domain-bridging objective. This constraint results to be the natural coarse-to-fine extension of the simple approach proposed for standard continual learning in [24], where all fine classes are derived from a unique *background* macro class. More formally, the knowledge distillation  $\mathcal{L}_{KD}^c$  for the classes  $\mathcal{C}_{t-1}^c$  that are further split can be expressed as:

$$\mathcal{L}_{KD}^c = \frac{1}{|\mathcal{T}_t^T|} \sum_{\mathbf{x} \in \mathcal{T}_t^T} \sum_{c \in \mathcal{C}_{t-1}^c} \mathbf{P}_{t-1}^T[c] \sum_{f \in S_t(c)} \log \mathbf{P}_t^T[f]. \quad (3)$$

When considering classes that have reached the maximum level of semantic refinement (*i.e.*, that are not going to be further split), the expression collapses into the standard cross-entropy based knowledge distillation (since the sets  $S_t(c)$  contain only one class,  $c$  itself, the sum of logs disappears), which forces the prediction of each channel  $c$  in the previous step  $t-1$ , to approximate the same class  $c$  in the next step  $t$ . Mathematically:

$$\mathcal{L}_{KD}^f = \frac{1}{|\mathcal{T}_t^T|} \sum_{\mathbf{x} \in \mathcal{T}_t^T} \sum_{c \in \mathcal{C}_{t-1}^f} \mathbf{P}_{t-1}^T[c] \log \mathbf{P}_t^T[c]. \quad (4)$$

In our approach (CCDA), at each training step  $t$ , we can rely upon the knowl-

edge that the previous model  $M_{t-1}$  has been adapted to segment images coming from the target domain. For this reason, the knowledge distillation is performed exploiting target domain samples, leading to an additional domain bridging effect. As an additional study, we analyzed a variant of this approach SKDC (Source domain Knowledge-Distillation for Coarse-to-fine) which implements the same knowledge distillation module, but the pseudo-labels are computed from the source domain, removing the domain bridging component.

#### 4.2.2. Coarse-to-Fine Weight Initialization

When training on a new model, the weights of the classifier are usually initialized in a random fashion, since there is no *a priori* information to be exploited to identify a better value. In a continual learning setup, though, this is true only in the first incremental step ( $t = 0$ ). In the subsequent steps (*i.e.*,  $t > 0$ ) a better strategy for the initialization of the weights of the decoder  $D_t$  is to exploit the knowledge provided by the decoder of the previous step,  $D_{t-1}$  (recall that  $M_t = D_t \circ E_t$ ), which can be used to obtain weights relative to already-seen classes (reducing the number of random values to be used). Furthermore, failing to recognize this detail can lead to significant performance losses during the network optimization, as a total random initialization may compromise the effectiveness of the knowledge distillation objective. Our setup, differently from the incremental learning case (where classes are added without any semantic relationship between each other), provides additional *priors* which can be exploited to significantly improve the weights initialization: the hierarchical coarse-to-fine labels relationship. Such information allows to remove completely the need for random initialization, since all the finer classes predicted by  $D_t$  have a macro class in  $D_{t-1}$ . Our strategy, therefore, tracks and exploits this relationship to initialize weights for the new fine classes: the weights of each new child fine class are initialized with the value of its coarse parent class. Naturally, when a class reached the maximum level of refinement in the subsequent steps, its weights will be initialized with the same values of the previous step. Formally, the operations can be defined as:

$$\{\omega_t^f\} = \{\omega_{t-1}^c\} \quad \forall c \in C_{t-1}^c, \forall f \in S_t(c) \quad (5)$$

$$\{\omega_t^f\} = \{\omega_{t-1}^f\} \quad \forall f \in C_{t-1}^f. \quad (6)$$

For what concerns the biases initialization, we adopted the idea introduced in[24] for class incremental learning, again extending it to the more general case of coarse-to-fine learning. The idea consists in spreading the probability of each coarse class  $c \in C_{t-1}^c$  towards its respective set of fine classes  $S_t(c)$ . In formal terms, the bias update expression becomes:

$$\{\beta_t^f\} = \{\beta_t^c - \log(|S_t(c)|)\} \quad \forall c \in C_{t-1}^c, \forall f \in S_t(c) \quad (7)$$

$$\{\beta_t^f\} = \{\beta_t^f\} \quad \forall f \in C_{t-1}^f. \quad (8)$$

After applying a softmax layer to compute the output probabilities, this

updated rule leads all the fine classes to share an equal fraction of the original coarse class prediction probability.

For all the competing approaches, the results are reported initializing new weights according to Eq. (5) and (6) for fair comparison. Instead, bias initialization rule (Eq. (7) and (8)) is only employed in our approach.

In Table 6 of the ablation studies (Section 6.3) we analyzed the effects of this refined bias initialization (referred to as *Unbiased* weight initialization). Particularly, we observe that disabling it (*i.e.*, not subtracting the log in Eq. (7)) leads to a serious degradation of the performance. Our results confirm and extend the findings of [24] even in the more general case of coarse-to-fine learning, showing how an unbiased strategy can handle difficult class relationships.

## 5. Implementation Details

The proposed framework is implemented in PyTorch. The semantic segmentation network is a DeepLab-V3 with ResNet-101 as the backbone (pre-trained weights are available at <https://pytorch.org/vision/stable/models.html>). In order to perform hyper-parameters selection, we identified a validation set of respectively 496 and 50 samples within the original training set of Cityscapes and IDD. We recall that, at the beginning of each incremental step, the model is initialized with a new classifier, depending on the number of classes to be recognized. Following [43, 10] the model is trained using the Stochastic Gradient Descent (SGD) optimizer with momentum equal to 0.9 and weight decay of  $5 \times 10^{-4}$ . The learning rate is polynomially decreased with a decay power of 0.9. The initial value of learning rate for the step 0 is set to  $2.5 \times 10^{-4}$ , while in the incremental steps it is initialized to  $10^{-4}$ . In the  $GTA5 \rightarrow Cityscapes$  setup, the batch size per domain is set to 2 (total of 4 samples), images are resized to  $1280 \times 720$  and  $1280 \times 640$ , respectively for the source and target domains. Instead, in the  $GTA5 \rightarrow IDD$  setup, the batch size per domain is set to 1 (total of 2 samples) and images are resized to  $1280 \times 720$  for both domains. In both cases, we trained our models on a single NVIDIA RTX 3090, performing 50000 iterations in the first step and 25000 iterations for each incremental step. Following [4, 9], we perform data augmentation through random left-right flipping and Gaussian blur to improve generalization properties. The hyper-parameters, determined on the validation, are the same in both setups:  $\lambda_{UDA} = 0.1$ ,  $\lambda_{KD}^f = 10$  except for  $\lambda_{KD}^c$  which is equal to 10 for  $GTA5 \rightarrow Cityscapes$  and equal to 1 for  $GTA5 \rightarrow IDD$ . Finally, when a model makes use of maximum squares loss, we include a warm-up stage of 100 iterations, and we activate the maximum squares constraint thereafter.

## 6. Experimental Results

In this section we report the experimental results on two synthetic-to-real benchmarks:  $GTA5 \rightarrow Cityscapes$  and  $GTA5 \rightarrow IDD$ , comparing our approach with several methods proposed to solve the two tasks separately. Then, we

report some ablation studies to carefully validate the effects of the various components of our approach.

To evaluate our approach, we compare our strategy with several baselines. As a global upper bound, we report the results of the offline (single-shot) training on the target dataset, *i.e.*, the Joint Target Only, JTO. As an upper bound for all the steps of incremental learning approaches, we report the coarse-to-fine training on the target dataset, that we call Target Only Non-Continual, TNC; here, at each step the algorithm sees all the labels including the ones of the classes not refined in the current step. As a global lower bound for both incremental and UDA approaches we report the naïve fine-tuning optimization of the network on the source dataset (Source Only). Finally, we compare our approach with two competitors, namely: MaxSquareIW [4] (MSIW) as state-of-the-art UDA approach, which does not employ any incremental technique, and MiB [24] as state-of-the-art continual learning technique, without domain adaptation.

Our contributions are named as SKDC (Source Knowledge Distillation for Coarse-to-Fine) and CCDA (Continual Coarse-to-fine Domain Adaptation). SKDC is the continual counterpart of TNC trained on the source dataset, where our Coarse-to-fine Knowledge Distillation strategy is employed as a mean to preserve previously acquired knowledge. CCDA is our full approach and it combines our domain-adaptive approach and SKDC: here, both domain adaptation and knowledge distillation losses are employed in all steps. Differently from SKDC, though, the preserved knowledge comes from the target dataset, instead of the source one (*i.e.*, the pseudo-labels used in the knowledge distillation loss are generated via inference on the target dataset sample). This allows the knowledge distillation loss to perform additional entropy minimization and domain bridging tasks.

In Table 1 we report the mean quantitative results of all the approaches in both considered settings. In Tables 2 and 3 we report the detailed per-step and per-class IoU scores attained by our full approach (CCDA) in the benchmarks. All tables show the score on the target dataset.

### 6.1. *GTA5*→*Cityscapes* adaptation

On the *Cityscapes* dataset the upper bound on the mIoU score given by the target offline training (JTO) is 68.6%, while the coarse-to-fine training on the target dataset (TNC) leads to a reduction of 2.1% with a score of 66.5%.

When using only supervised source data (*GTA5*→*Cityscapes* setup) we can see that the Source Only method suffers greatly from both catastrophic forgetting and domain shift, obtaining very low score of 4.5%. The domain adaptation strategy of MSIW [4] is able to only marginally improve performances, reaching 6.9% with an improvement of 2.4%. Such results, shows how a continual learning approach is needed to properly tackle the task: indeed, MiB [24] reaches 26.5%, with an improvement of 22% points over the baseline. Nevertheless, both our strategies are able to surpass such performance: SKDC reaches a modest 30.4%, while our full approach (CCDA) further surpasses it with a mIoU score of 33.1%. This corresponds to an improvement of 6.6% with respect to the

	Model	mIoU <sub>0</sub>	mIoU <sub>1</sub>	mIoU <sub>2</sub>	mIoU <sub>3</sub>
GTA5 → Cityscapes	JTO	-	-	-	68.6
	TNC	92.1	83.7	72.6	66.5
	Source Only	65.0	56.7	25.1	4.5
	MiB [24]	65.0	56.1	27.8	26.5
	MSIW [4]	<b>85.4</b>	62.2	<b>38.5</b>	6.9
	SKDC (ours)	65.0	65.4	34.4	30.4
	CCDA (ours)	<b>85.4</b>	<b>67.9</b>	37.2	<b>33.1</b>
GTA5 → IDD	JTO	-	-	-	65.9
	TNC	87.0	78.8	73.2	64.2
	Source only	72.9	60.8	30.0	6.1
	MiB [24]	72.9	61.8	<b>44.7</b>	26.8
	MSIW [4]	<b>75.9</b>	62.7	30.2	8.9
	SKDC (ours)	72.9	60.9	38.7	32.4
	CCDA (ours)	<b>75.9</b>	<b>63.5</b>	40.3	<b>33.0</b>

Table 1: Quantitative results in terms of mIoU computed on the validation sets of Cityscapes and IDD, when employing source knowledge from the GTA5 dataset.

step	road	sidewalk	sky	vegetation	terrain	pole	traffic light	traffic sign	building	fence	wall	person	rider	bicycle	motorbike	car	truck	bus	train	mIoU
0			92.9						84.2							79.2				85.4
1	95.6	74.9	80.8			9.7		78.0				55.1				81.3				67.9
2	53.2	15.9	75.0	80.6	35.9	14.3	8.4	73.9	10.5			49.9	0.0	3.1		79.2		21.0		37.2
3	50.8	14.1	74.7	77.1	29.8	12.9	27.1	16.4	77.9	12.6	21.7	50.2	0.0	6.7	12.5	74.4	25.4	27.2	17.7	33.1

Table 2: *GTA5* → *Cityscapes*, CCDA per-class per-step IoU results on the Cityscapes test set.

best competitor (MiB). We further remark how the difference in score between CCDA and SKDC (2.7%) is higher than the one between MSIW and Source Only, meaning that our target self-supervised continual approach was able to improve the entropy minimization effect of MSIW, allowing to achieve better domain transfer capabilities.

In Table 1 we can also see that MSIW surpasses our strategy at step 2, this is due to the high number of new classes present in such step, which are generated by splitting of most existing classes in the previous step (there are 13 new classes, see Figure 6). This allows the improvement in domain shift to surpass the performance loss caused by the catastrophic forgetting. However, at the next step the capability of preserving old knowledge is required and the superior performance of our approach are clearly visible. In the detailed results reported in Table 2 we can see that an important component of the performance degradation derives from the *thin object* class, being it very difficult. Another cause of degradation is the *rider* class, which is confused for *person* leading to a drop in performance from 55.1% in class *human* at step 1 to 0% when the class is split in step 2.

From the qualitative results reported in Figure 8 we can appreciate how our

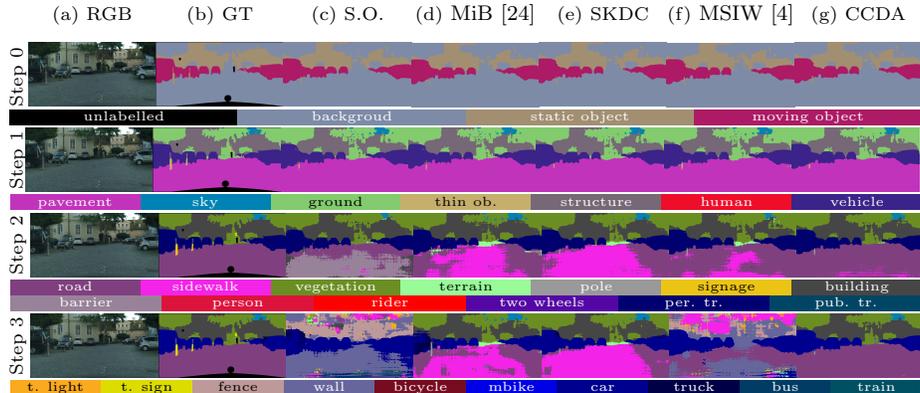


Figure 8: GTA5→Cityscapes qualitative results. GT refers to the complete ground truth for the given step.

step	road	sidewalk	sky	vegetation	pole	traffic light	traffic sign	building	fence	wall	person	rider	bicycle	motorbike	car	truck	bus	mIoU
0			93.1				56.8						77.9					75.9
1	92.2	85.7	75.6		20.6		51.8				50.6			68.3				63.5
2	84.3	21.8	82.4	71.7	9.8	29.2	44.0	21.2			21.8	18.2	31.7		62.8	24.9		40.3
3	78.2	16.9	79.1	70.8	7.1	2.4	35.1	25.0	5.4	19.1	19.1	2.8	14.0	39.4	74.4	43.3	28.9	33.0

Table 3: *GTA5*→*IDD*, CCDA per-class per-step IoU results on the Cityscapes test set.

CCDA strategy is consistently better than the competitors, in all the incremental steps. In particular, we can notice how in the third and last step CCDA is the only method that can correctly classify the *road* in the bottom half of the image, which the other methods confuse with *sidewalk*, *wall* or *car*. At the same time, we can notice how the shape of the parked cars on the sides of the image are much better defined in our method. Furthermore, also the visual results show a dramatic decrease in performance due to the catastrophic forgetting for Source Only and MSIW methods, and how such effect is strongly reduced in methods that exploit continual learning strategies (MiB, SKDC, CCDA).

## 6.2. *GTA5*→*IDD* adaptation

On the *IDD* dataset the upper bound given by JTO is 65.9%, *i.e.*, 2.7% lower than on Cityscapes (even with two classes less), indicating *IDD* as a more difficult dataset. The TNC method achieves a score of 64.2%, which is 1.7% lower than JTO, note how such difference is smaller than on the Cityscapes dataset (1.6% versus 1.3% relative loss, respectively). Moving to the domain adaptive setup (*GTA5*→*IDD*) the Source Only method reaches a mIoU of 5.7%, which is increased up to 7.4% by the UDA approach in the MSIW method. Also in this setup, continual learning strategies have much better performances: MiB reaches 26.8% of mIoU, surpassed by both SKDC (32.4%) and CCDA (33.0%).

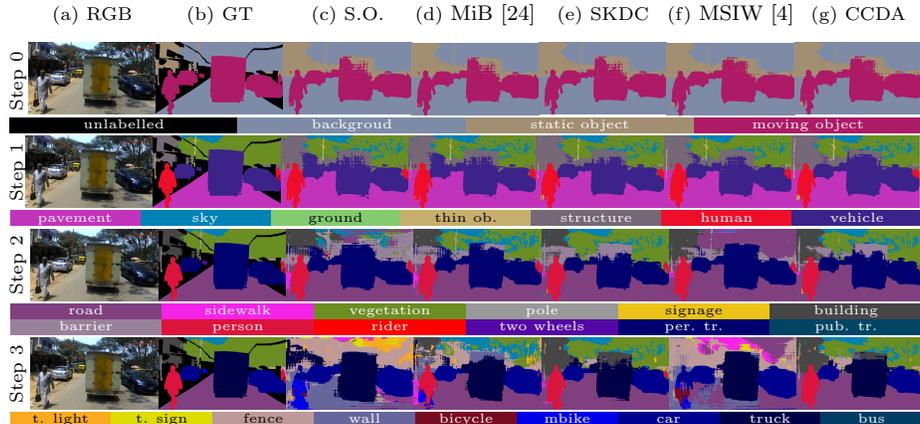


Figure 9: GTA5→IDD qualitative results. GT refers to the complete ground truth for the given step.

The improvement from MiB to CCDA is similar to the previous case, showing how the performance of the proposed approach are consistent across different datasets (6.6% for Cityscapes and 6.2% for IDD, respectively). In this setup the MiB method is the best-performing approach at step 2. We argue that this is due to the increased number of fine classes in step 1 due to the collapse of *ground* in *terrain*: this favors the naïve knowledge distillation approach of MiB, since there is no coarse-to-fine knowledge transfer. As before, though, the final results show how CCDA is able to preserve much better the knowledge during the third incremental step. Looking at Table 3 we can see that, even in this setup, one of the causes of decreased performance is the *pole* class, which after being split from *thin object* loses 10.8% mIoU (going from 20.6% to 9.8%). Other challenging classes are *fence* (losing 15.8%) and *rider* (with 32.4% loss).

In Figure 9 we report some qualitative results for the various methods in this benchmark. Here, we see the significant domain shift present between *GTA5* and *IDD* already from step 0 (first row): the *truck* (*moving object*) in the foreground is confused for *static object* in all strategies, even when only three classes are considered. From step 2 onward, though, the domain shift is overshadowed by the catastrophic forgetting effect: the Source Only and MSIW approaches overfit significantly on the new classes forgetting the old ones, up to the point that even the *sky* is confused as *road*. On the other hand, as for the *GTA5*→*Cityscapes* setting, our approach (CCDA) is consistently the best approach, maintaining a high precision and tight class borders in all steps. In particular, when we compare it to the MiB approach, we can see that the latter is overestimating the presence of the class set of step 3 (last row), for example it misleads part of the legs of the *person* on the left with a *motorbike*.

Model	mIoU <sub>0</sub>	mIoU <sub>1</sub>	mIoU <sub>2</sub>	mIoU <sub>3</sub>
Joint Source Only	-	-	-	30.8
Joint MaxSquareIW	-	-	-	38.9
SNC	65.0	56.2	29.0	24.8
TFT	92.1	83.7	67.9	17.5
TKDC (ours)	92.1	81.8	65.1	62.3
SKDC (ours)	65.0	65.4	34.4	30.4
CCDA (ours)	85.4	67.9	37.2	33.1

Table 4: Additional baseline experiments on variations of our full approach, conducted on the *GTA5*→*Cityscapes* setup.

Model	mIoU <sub>0</sub>	mIoU <sub>1</sub>	mIoU <sub>2</sub>	mIoU <sub>3</sub>
L1	<b>65.0</b>	51.5	27.2	18.7
L1 logits	<b>65.0</b>	61.3	31.2	23.4
L2	<b>65.0</b>	44.7	22.3	13.6
L2 logits	<b>65.0</b>	60.6	29.2	21.6
MiB [24]	<b>65.0</b>	56.1	27.8	26.5
SKDC (ours)	<b>65.0</b>	<b>65.4</b>	<b>34.4</b>	<b>30.4</b>

Table 5: Ablation experiments on knowledge distillation, performed on the *GTA5*→*Cityscapes* setup.

### 6.3. Ablation Studies

In this section we report some ablation results on the *GTA5*→*Cityscapes* benchmark, to better evaluate the components and design choices of our approach and compare its performance to other standard strategies. The section is divided into three parts: additional baseline results are reported in Table 4, ablation studies on the knowledge distillation loss (see Section 4.2.1) are shown in Table 5, and ablation studies on the weight initialization rule (see Section 4.2.2) are finally reported in Table 6.

As additional baseline results, in Table 4 we present three UDA-related methods (see Section 4.1) and two C2F-related methods (see Section 4.2). In the first group we placed methods tackling only the domain adaptation task with no incremental learning (*i.e.*, providing them the fully labeled data in a single shot): JSO (Joint Source Only), usually considered as baseline for standard UDA approaches; JMSIW (Joint MSIW), the upper bound for our CCDA approach given achieved by applying MSIW on the fully labeled dataset; SNC (Source Non Continual, the same as TNC but with source data instead of the target one), that is upper bound for approaches trained only on source data. In the second group we report: TFT (Target Fine Tuning), baseline for purely continual methods trained on the target domain, and TKDC, the counterpart of SKDC that is trained on target data. From the table we observe how our CCDA strategy achieves a score very close to JMSIW, with a loss of only 5.8% of mIoU, which is comparable to the performance reduction between JTO and TKDC (6.3%) and the one between JSO and SNC (6.0%).

In Table 5 we investigate the performance attainable by trying different

Model	Weights Init.	mIoU <sub>0</sub>	mIoU <sub>1</sub>	mIoU <sub>2</sub>	mIoU <sub>3</sub>
Source Only	Naïve	<b>65.0</b>	56.7	25.1	4.5
Source Only	Unbiased	<b>65.0</b>	60.4	28.2	5.5
SKDC (ours)	Naïve	<b>65.0</b>	60.1	32.6	27.9
SKDC (ours)	Unbiased	<b>65.0</b>	<b>65.4</b>	<b>34.4</b>	<b>30.4</b>

Table 6: Ablation experiments on weight initialization, performed on the *GTA5*→*Cityscapes* setup.

knowledge distillation strategies as described in detail in Section 4.2.1. We analyze four distance-based expressions and two cross-entropy based ones: as before, all training procedures were performed on the source dataset and tested on the target dataset. As distance-based strategies we considered the L1 and L2 norms, minimized before and after the *softmax* layer (identified as *logits* in the table). From the experimental results it is clear that matching the logits leads to the best results, gaining 4.7% and 8.0% mIoU in the L1 and L2 versions, respectively. A better approach, though, is to use cross-entropy based approaches. In this category we show MiB [24] and SKDC. The former is reported in the table since it could be considered a simplified case of SKDC where only one coarse class (*i.e.*, the *background*) is considered. This relation is made clear from the experimental results, where MiB loses 3.9% with respect to SKDC due to the missing coarse-to-fine class mapping.

Finally, in Table 6 we investigate the effect of our weight initialization strategy (described in Section 4.2.2), comparing the score attained by Source Only and SKDC when unbiased weight initialization is enabled or disabled. Recall that the difference lies in whether the bias is rescaled to preserve the original probabilities or not, in which case we fall back to the naïve coarse-to-fine weight initialization. In both cases, the unbiased weight initialization helps to improve mIoU since it reduces the class-confusion and helps the knowledge distillation and cross-entropy objectives. In particular, when the refined weight initialization rule is applied to the SKDC method the performance increases by 2.5%.

## 7. Conclusion

In this paper we introduced the novel task of continual coarse-to-fine unsupervised domain adaptation for semantic segmentation, which represents a step closer towards open-world learning. We started by defining in detail the task and we analyzed how it combines the continual learning and UDA tasks. Then, we proposed a new approach (CCDA) to address it. Our approach combines state-of-the-art UDA and continual learning strategies with specific provisions for the new challenging setting, including a novel knowledge distillation strategy and an unbiased weight initialization scheme for the incremental steps. We evaluated CCDA in two synthetic-to-real UDA benchmarks, extending them for the use in a continual coarse-to-fine setting. We compare our approach with several methods, including state-of-the-art UDA and CL techniques, outperforming

the competitors on both benchmarks. Further research will be devoted to the introduction of novel coarse-to-fine strategies and to the extension to different experimental scenarios with multiple source and target domains.

### Acknowledgement

This work was partially supported by the University of Padova Strategic Research Infrastructure Grant 2017: “*CAPRI: Calcolo ad Alte Prestazioni per la Ricerca e l’Innovazione*” and by the SID project “*Semantic Segmentation in the Wild*”.

### References

- [1] J. Long, E. Shelhamer, T. Darrell, Fully convolutional networks for semantic segmentation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2015, pp. 3431–3440.
- [2] M. Toldo, A. Maracani, U. Michieli, P. Zanuttigh, Unsupervised domain adaptation in semantic segmentation: a review, *Technologies* 8 (2) (2020).
- [3] U. Michieli, P. Zanuttigh, Incremental Learning Techniques for Semantic Segmentation, in: Proceedings of the International Conference on Computer Vision Workshops, 2019.
- [4] M. Chen, H. Xue, D. Cai, Domain adaptation for semantic segmentation with maximum squares loss, in: Proceedings of the International Conference on Computer Vision, 2019, pp. 2090–2099.
- [5] S. R. Richter, V. Vineet, S. Roth, V. Koltun, Playing for data: Ground truth from computer games, in: Proceedings of the European Conference on Computer Vision, 2016, pp. 102–118.
- [6] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, B. Schiele, The Cityscapes dataset for semantic urban scene understanding, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 3213–3223.
- [7] G. Varma, A. Subramanian, A. Namboodiri, M. Chandraker, C. Jawahar, Idd: A dataset for exploring problems of autonomous navigation in unconstrained environments, in: Proceedings of the Winter Conference on Applications of Computer Vision, 2019, pp. 1743–1751.
- [8] V. Badrinarayanan, A. Kendall, R. Cipolla, Segnet: A deep convolutional encoder-decoder architecture for image segmentation, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 39 (12) (2017) 2481–2495.
- [9] H. Zhao, J. Shi, X. Qi, X. Wang, J. Jia, Pyramid scene parsing network, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2017, pp. 2881–2890.

- [10] L. Chen, G. Papandreou, F. Schroff, H. Adam, Rethinking atrous convolution for semantic image segmentation, arXiv preprint arXiv:1706.05587 (2017).
- [11] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, A. L. Yuille, Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs, *IEEE Transactions on Pattern Analysis and Machine Intelligence* 40 (2018) 834–848.
- [12] L. Chen, Y. Zhu, G. Papandreou, F. Schroff, H. Adam, Encoder-decoder with atrous separable convolution for semantic image segmentation, in: *Proceedings of the European Conference on Computer Vision*, 2018, pp. 833–851.
- [13] J. Li, P. Zhou, C. Xiong, S. C. Hoi, Prototypical contrastive learning of unsupervised representations, arXiv preprint arXiv:2005.04966 (2020).
- [14] B. Yang, F. Wan, C. Liu, B. Li, X. Ji, Q. Ye, Part-based semantic transform for few-shot semantic segmentation, *IEEE Transactions on Neural Networks and Learning Systems* (2021).
- [15] M. Mel, U. Michieli, P. Zanuttigh, Incremental and multi-task learning strategies for coarse-to-fine semantic segmentation, *Technologies* 8 (1) (2020) 1.
- [16] O. Stretcu, E. A. Platanios, T. Mitchell, B. Póczos, Coarse-to-fine curriculum learning for classification, in: *Proceedings of the International Conference on Learning Representations Workshops*, 2020.
- [17] O. Stretcu, E. A. Platanios, T. M. Mitchell, B. Póczos, Coarse-to-fine curriculum learning, arXiv preprint arXiv:2106.04072 (2021).
- [18] U. Michieli, E. Borsato, L. Rossi, P. Zanuttigh, Gmnet: Graph matching network for large scale part semantic segmentation in the wild, in: *Proceedings of the European Conference on Computer Vision*, Springer, 2020, pp. 397–414.
- [19] G. I. Parisi, R. Kemker, J. L. Part, C. Kanan, S. Wermter, Continual lifelong learning with neural networks: A review, *Neural Networks* (2019).
- [20] T. Lesort, V. Lomonaco, A. Stoian, D. Maltoni, D. Filliat, N. Díaz-Rodríguez, Continual learning for robotics: Definition, framework, learning strategies, opportunities and challenges, *Information Fusion* 58 (2020) 52–68.
- [21] U. Michieli, M. Toldo, P. Zanuttigh, Unsupervised Domain Adaptation and Continual Learning in Semantic Segmentation, *Advanced Methods and Deep Learning in Computer Vision*, Elsevier (2021).

- [22] M. Klingner, A. Bär, P. Donn, T. Fingscheidt, Class-incremental learning for semantic segmentation re-using neither old data nor old labels, International Conference on Intelligent Transportation Systems (2020).
- [23] U. Michieli, P. Zanuttigh, Continual semantic segmentation via repulsion-attraction of sparse and disentangled latent representations, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2021, pp. 1114–1124.
- [24] F. Cermelli, M. Mancini, S. R. Bulò, E. Ricci, B. Caputo, Modeling the background for incremental learning in semantic segmentation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2020, pp. 9233–9242.
- [25] A. Douillard, Y. Chen, A. Dapogny, M. Cord, Plop: Learning without forgetting for continual semantic segmentation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2021, pp. 4040–4050.
- [26] A. Maracani, U. Michieli, M. Toldo, P. Zanuttigh, Recall: Replay-based continual learning in semantic segmentation, in: Proceedings of the International Conference on Computer Vision, 2021, pp. 7026–7035.
- [27] G. Csurka, Domain adaptation for visual applications: A comprehensive survey, arXiv preprint arXiv:1702.05374 (2017).
- [28] Y.-C. Chen, Y.-Y. Lin, M.-H. Yang, J.-B. Huang, Crdoco: Pixel-level domain transfer with cross-domain consistency, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2019, pp. 1791–1800.
- [29] J. Hoffman, E. Tzeng, T. Park, J.-Y. Zhu, P. Isola, K. Saenko, A. Efros, T. Darrell, Cycada: Cycle-consistent adversarial domain adaptation, in: Proceedings of the International Conference on Machine Learning, 2018, pp. 1994–2003.
- [30] J. Hoffman, D. Wang, F. Yu, T. Darrell, FCNs in the wild: Pixel-level adversarial and constraint-based adaptation, arXiv preprint arXiv:1612.02649 (2016).
- [31] Z. Murez, S. Kolouri, D. J. Kriegman, R. Ramamoorthi, K. Kim, Image to image translation for domain adaptation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018, pp. 4500–4509.
- [32] M. Toldo, U. Michieli, G. Agresti, P. Zanuttigh, Unsupervised domain adaptation for mobile semantic segmentation based on cycle consistency and feature alignment, Image and Vision Computing 95 (2020).

- [33] F. Pizzati, R. d. Charette, M. Zaccaria, P. Cerri, Domain bridge for unpaired image-to-image translation and unsupervised domain adaptation, in: Proceedings of the Winter Conference on Applications of Computer Vision, 2020, pp. 2990–2998.
- [34] W. Tranheden, V. Olsson, J. Pinto, L. Svensson, Dacs: Domain adaptation via cross-domain mixed sampling, in: Proceedings of the Winter Conference on Applications of Computer Vision, 2021, pp. 1379–1389.
- [35] M. Toldo, U. Michieli, P. Zanuttigh, Unsupervised domain adaptation in semantic segmentation via orthogonal and clustered embeddings, in: Proceedings of the Winter Conference on Applications of Computer Vision, 2021, pp. 1358–1368.
- [36] F. Barbato, M. Toldo, U. Michieli, P. Zanuttigh, Latent space regularization for unsupervised domain adaptation in semantic segmentation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops, 2021, pp. 2835–2845.
- [37] F. Barbato, U. Michieli, M. Toldo, P. Zanuttigh, Road scenes segmentation across different domains by disentangling latent representations, arXiv preprint arXiv:2108.03021 (2021).
- [38] Y. Zou, Z. Yu, B. Vijaya Kumar, J. Wang, Unsupervised domain adaptation for semantic segmentation via class-balanced self-training, in: Proceedings of the European Conference on Computer Vision, 2018, pp. 289–305.
- [39] Y. Zou, Z. Yu, X. Liu, B. V. Kumar, J. Wang, Confidence regularized self-training, in: Proceedings of the International Conference on Computer Vision, 2019, pp. 5982–5991.
- [40] Y. Zhang, P. David, B. Gong, Curriculum domain adaptation for semantic segmentation of urban scenes, in: Proceedings of the International Conference on Computer Vision, 2017, pp. 2020–2030.
- [41] Y. Zhang, P. David, H. Foroosh, B. Gong, A curriculum domain adaptation approach to the semantic segmentation of urban scenes, *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2020).
- [42] A. Cardace, P. Z. Ramirez, S. Salti, L. Di Stefano, Shallow features guide unsupervised domain adaptation for semantic segmentation at class boundaries, in: Proceedings of the Winter Conference on Applications of Computer Vision, 2022, pp. 1160–1170.
- [43] Y.-H. Tsai, W.-C. Hung, S. Schulter, K. Sohn, M.-H. Yang, M. Chandraker, Learning to adapt structured output space for semantic segmentation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018, pp. 7472–7481.