# IDA: Informed Domain Adaptive Semantic Segmentation

Zheng Chen[1], Zhengming Ding[2], Jason M. Gregory[3], and Lantao Liu[1]

*Abstract*—Mixup-based data augmentation has been validated to be a critical stage in the self-training framework for unsupervised domain adaptive semantic segmentation (UDA-SS), which aims to transfer knowledge from a well-annotated (source) domain to an unlabeled (target) domain. Existing self-training methods usually adopt the popular region-based mixup techniques with a random sampling strategy, which unfortunately ignores the dynamic evolution of different semantics across various domains as training proceeds. To improve the UDA-SS performance, we propose an Informed Domain Adaptation (IDA) model, a self-training framework that mixes the data based on class-level segmentation performance, which aims to emphasize small-region semantics during mixup. In our IDA model, the class-level performance is tracked by an expected confidence score (ECS). We then use a dynamic schedule to determine the mixing ratio for data in different domains. Extensive experimental results reveal that our proposed method is able to outperform the state-of-the-art UDA-SS method by a margin of 1.1 mIoU in the adaptation of GTA-V to Cityscapes and of 0.9 mIoU in the adaptation of SYNTHIA to Cityscapes. Code link: **https://github.com/ArlenCHEN/IDA.git**

## I. INTRODUCTION

Semantic segmentation (SS) aims to learn a pixel-wise classification for a given image and plays a critical role in various applications such as infrastructure/industrial inspections [1], biomedical diagnoses [2], and vehicle autonomy [3]. Current mainstream segmentation models [4] [5] [6] heavily rely on deep neural networks (DNNs) which usually require a huge amount of manual annotations in order to achieve desirable performance, e.g., labeling for a single image might require more than 1.5 hours on average [7]. There exist some public datasets that provide dense annotations, e.g., Cityscapes [7], ACDC [8]. However, these existing datasets are far from providing sufficient coverage for other miscellaneous novel environments, leading to deep models that fail to generalize. In this case, how to transfer the knowledge in the easy-to-access data (e.g., data in simulators; existing public datasets) to boost the model generalization for unseen data of other domains is critical as there is usually a domain shift between the data used to train the model and the data encountered during deployment (see Fig. 1).

To achieve the model transfer, a broadly studied task is *domain adaptation* (DA), where we define the data with labels as a *source* domain, and the data to be processed and used for prediction as a *target* domain. In many real-world
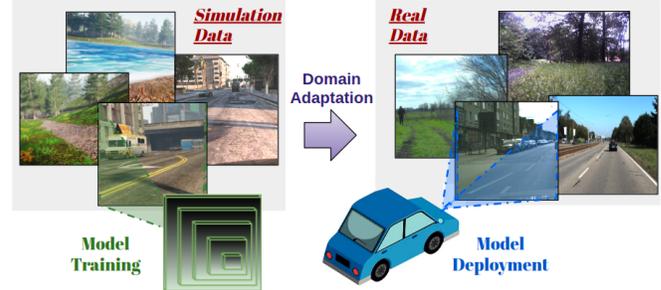


Fig. 1. Domain adaptation aims at diminishing the model performance drop due to the shift of different data domains, e.g., simulated data vs. real data.

scenarios, it can be challenging or costly to quickly create labels for target data. In this case, *unsupervised domain adaptation* (UDA) provides a means to solve the special DA setting where the target domain has no labels.

A critical step for UDA-SS is the mixup-based data augmentation which oftentimes utilizes random sampling where certain rectangular regions (e.g., in CutMix [9]) or class regions (e.g., in ClassMix [10]) are randomly selected with a predefined size and mixed. This strategy can generally mix the visual contents and mitigate the understanding shift of the DNNs. Unfortunately, this strategy ignores the dynamic evolution of the two domain data along the training progress, and can fail in leveraging important "informative" data that however are obscured amid the process. The "informativeness" of data can be defined as its significance to the final performance. In greater detail, in UDA-SS the segmentation of multiple classes are predicted such that the final performance of the model depends on the performance of each individual class. We conjecture that the model performance is mainly driven by the performance of "bottleneck" classes that usually have low-quality segmentations. Note also, the bottleneck classes may vary during the training process. The model performance can be significantly improved if we can identify what bottleneck classes are so as to inform us to specially improve their performance. Regions associated with those bottleneck classes are defined as informative data, and our proposed method is thus termed as Informed Domain Adaptation (IDA).

The key to the proposed IDA framework is a novel mixup technique — *Informed Mix* (IMix) built upon ClassMix [10]. Different from the ClassMix, our IMix bridges data from two domains according to the confidence values of training progress indicators. Thus the IMix is informed by and prioritizes, the regions indicated by a low indicator value which means the classes are balanced in both raw image space and label space, and the training will be unbiased, leading to an

increase of final performance. To comply with the training dynamics, we also propose a novel adaptation schedule for IDA. The proposed schedule adaptively determines the ratio of image regions from two domains for mixing because the dynamic changes in training progress are related not only to the type of bottleneck classes but also to the number of bottleneck classes. Setting a fixed number will either miss effective data or introduce possible noise. Our proposed schedule has three phases. In the first phase, IMix mainly selects easy classes from source images for adaptation. In the second phase, hard/bottleneck classes from source images are selected and the number of the selected classes decreases from a high level to a low level while the selections from target images increase from a low level to a high level. In the third phase, the numbers of selections from the source and target domains will be maintained at low and high levels, respectively. To summarize the contribution of this work,

- We propose a principled model, IDA, for the UDA-SS. The IDA model is a self-training framework that exploits the obscured informativeness of data, which has not been previously studied in DA.
- We propose a new mixup technique, IMix, that bridges the source and target domains according to the training progress defined by an expected confidence assessment.
- We propose a novel dynamic adaptation schedule which can adaptively adjust the mixing ratio for different domains to optimize the adaptation efficiency. We will make the code of this work public.

Finally, our extensive evaluations on popular datasets show that our IDA outperforms the current SOTA model HRDA [11] with a remarkable margin under the same settings.

## II. RELATED WORK

**Domain Adaptive Semantic Segmentation:** Mainstream methods for tackling UDA can be categorized into two classes — feature alignment (FA) [12] [13] [14] [15] [16] [17] [18] [19] [20] and self-training (ST) [21] [22] [23] [24] [11]. FA adapts the model by aligning the features from the source domain to the target domain using adversarial training, i.e., features from two domains are expected to be indistinguishable through a domain discriminator. However, FA suffers from two issues. First, FA aligns features from two domains in a *global* way by evaluating the domain discrimination using features of the whole image. This can be problematic for semantic segmentation as each image contains multiple classes. Aligning features globally cannot guarantee the class-level shift is eliminated, and even worse, it is possible that features are aligned at a global level but severely misaligned at a class level, causing the so-called negative transfer [25]. Recently, some class-level FA methods [16] [18] [20] [19] are proposed to consider a finer level of feature structure, but they still suffer from the lack of target labels and show a weak performance. Second, FA adopts an adversarial training paradigm which is known to be unstable to train [26]. On the other hand, ST tackles the UDA by a teacher-student framework[11] [24] [27], where the teacher is trained on the source domain and predicts *pseudo*

*labels* for target images. Then the student is supervised by those predicted pseudo labels. Recently ST [11] [24] has been prevalent since it constantly breaks the state-of-the-art record of UDA-SS due to the highly efficient feedback for adaptation from pseudo-labels.

**Mixup Data Augmentation:** mixup-based data augmentation has been demonstrated to be a vital step for UDA-SS as it can achieve adaptation directly in the raw input space and label space, by forcibly mixing different domains for each data sample. Region-based mixups, e.g., CutMix [9] and ClassMix [10] are two representative mixups used in UDA-SS. CutMix and ClassMix typically adopt a random sampling strategy when mixing data from two domains, i.e., the region in one domain is randomly selected with a predefined size while the rest regions are from the other domain. DACS [28] is recently proposed to apply the idea of ClassMix to domain adaptation. DACS randomly selects regions of half of the classes in source images and pastes the target data to the rest regions. DACS has been validated to be effective in recent ST methods [24], [11].

## III. METHODOLOGY

To organize the presentation, in Sect. III-A, we provide preliminary knowledge about the domain adaptive semantic segmentation. In Sect. III-B, we describe the general structure of our IDA model. In Sect. III-C, we first describe how we perform the identification of bottleneck classes on the fly along the training process. Then we introduce the IMix data augmentation by carefully considering the spatiotemporal changes of domain data during training.

### A. Preliminaries of UDA-SS

We consider a source domain distribution $\mathcal{S}$ and a target domain distribution $\mathcal{T}$ over the joint space of $\mathcal{X} \times \mathcal{Y}$, where $\mathcal{X}$ is the input space and $\mathcal{Y}$ is the label space, respectively. In UDA-SS, we have access to $N_s$ labeled samples $\left(\mathcal{X}_s = \{x_i^s\}_{i=1}^{N_s}, \mathcal{Y}_s = \{y_i^s\}_{i=1}^{N_s}\right)$ for the source domain, and only access to $N_t$ raw images $\mathcal{X}_t = \left\{x_j^t\right\}_{j=1}^{N_t}$ for the target domain. A neural network $g$ comprising of a feature extractor $f_\theta$ parameterized by $\theta$ and a segmentation head $h_\phi$ parameterized by $\phi$, i.e., $g_{\theta,\phi} = h_\phi(f_\theta)$, is usually adopted as the adaptation model. The expected error on the source domain is denoted by

$$L_\mathcal{S}(\theta, \phi) = \mathbb{E}_{(x,y)\sim\mathcal{S}} \left[l(g_{\theta,\phi}(x), y)\right], \qquad (1)$$

where $l(\cdot, \cdot)$ represents the loss function. In UDA-SS, typically the standard cross-entropy with one-hot ground-truth (gt) label is used to compute the training loss: $l(g_{\theta,\phi}(x), y) = -\sum_{c=1}^{C} \left[y^c \cdot \log g_{\theta,\phi}(x)^c\right]$, where $C$ is the class size.

Similarly, the expected error on the target domain is denoted by $L_\mathcal{T}(\theta, \phi)$, but we cannot obtain an expression for $L_\mathcal{T}(\theta, \phi)$ as we have no labels for target data. However, we have an indirect way to approximate $L_\mathcal{T}(\theta, \phi)$ by

$$L_\mathcal{T}(\theta, \phi) = \mathbb{E}_{x\sim\mathcal{T}} \left[l(g_{\theta,\phi}(x), \hat{y})\right], \qquad (2)$$

where $\hat{y}$ is a pseudo label generated by the model trained on the source domain, $\hat{y} = \texttt{one-hot}(\text{argmax}_c(g_{\theta_s,\phi_s}(x)))$,
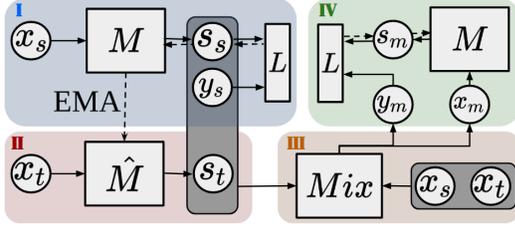
Fig. 2. Overall structure of our IDA model. Four training stages are involved. $x$ represents the input image; $y$ denotes the label; $s$ is the model prediction; $L$ means the standard cross-entropy loss; Model $M$ is the student model while $\hat{M}$ is the teacher model.

where $\theta_s$ and $\phi_s$ represent the neural network parameters trained on the source domain. Based on Eq. (1) and Eq. (2), we can obtain the adaptation objective as

$$\min_{\theta,\phi} L_{\mathcal{S}}(\theta,\phi) + L_{\mathcal{T}}(\theta,\phi). \qquad (3)$$

### B. Model Framework

We build the proposed IDA model based on the teacher-student model in self-training (see Fig. 2, $M$ represents the student model and $\hat{M}$ represents the teacher model). The teacher model shares the same network structure as the student's. Four stages are involved in one training iteration. In the beginning of each iteration, we use the *exponential moving average* (EMA) to update the teacher model's parameters by the ones of the student model such that the teacher model can be synchronized with the latest weights of the student model. Then the student model is trained with source labels. In the second stage, we use the teacher model to generate pseudo labels for target images without back-propagation. In the third stage, our proposed IMix module (described in Sect. III-C later) takes as input the source image, source prediction, source label, target image, and target prediction (pseudo label) to generate a new pair of data that mixes the data from the two domains. In the fourth stage, the student model is further trained using the newly generated mixed data.

In the proposed IDA model, our approximation to the expected error on the target domain differs from the Eq. (2), where the generated pseudo label $\hat{y}$ is directly used for training the model. Instead, we use the newly generated data pair $(x_m, y_m)$ to compute $L_{\mathcal{T}}$, as illustrated in Fig. 2. We denote the distribution of the mixed data as $\mathcal{M}$. The adaptation objective of our IDA model is

$$\min_{\theta,\phi} L_{\mathcal{S}}(\theta,\phi) + L_{\mathcal{M}}(\theta,\phi). \qquad (4)$$

Using $L_{\mathcal{M}}$ to approximate $L_{\mathcal{T}}$ has been validated in existing work [28] [24] [11]. The reason for this effectiveness lies in that Eq. (3) separates the supervision from the source domain and the target domain while Eq. (4) mixes supervision signals from the two domains in the term of $L_{\mathcal{M}}$. This mixture can efficiently guide the model to understand the target data with the accompanying source data at the sample level.

### C. Informed Mix

We propose Informed Mix (IMix) which is an important module in our IDA. The IMix considers dynamic temporal-spatial changes of data during the training process and is able
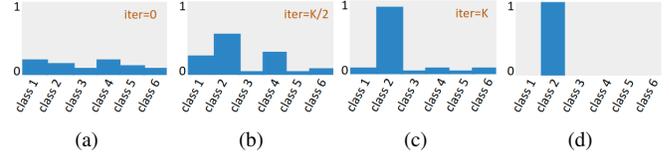


Fig. 3. (a) - (c) The changes of the categorical probability under the supervision of a (d) one-hot vector in different training iterations. $K$ represents the total number of training iterations.
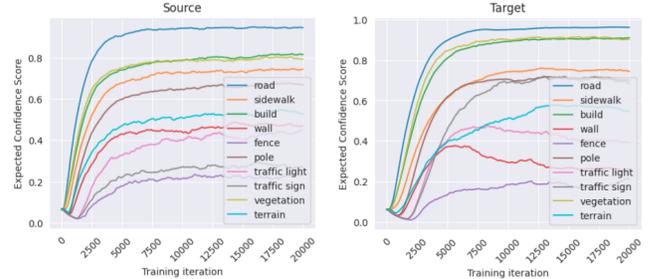


Fig. 4. Smoothed ECS values for source and target classes.

to adjust the mixture strategy accordingly. Different from previous mixup techniques [9] [10] [28] that use random sampling to select mix regions from images, our IMix informatively selects the regions based on the class-level performance during training.

*1) Class-level Performance Indicator:* To monitor the training progress, we usually observe how the loss changes over multiple epochs/iterations. However, we have no labels for target data in UDA-SS and thus are unable to use a loss-like indicator. In this work, we propose to use the *confidence score* (CS) of the predicted probability as the indicator. CS has been used by previous works as the uncertainty approximation, e.g., [29]. The metric is defined as $\mathrm{CS}(x) = \max(g_{\theta,\phi}(x))$, where the model $g$ usually has the `softmax` function as the last layer. The reason for using the confidence score is related to the standard cross-entropy (CE) loss function which assumes that the label is a one-hot vector and minimizing the CE loss is equivalent to maximizing the CS. In this case, a higher CS value can indicate a lower loss and thus better performance.

A simple illustration can be seen in Fig. 3, where Fig. 3(d) shows a one-hot vector label, the change of the corresponding probability over multiple training iterations can be seen from Fig. 3(a) to Fig. 3(c). The increased CE value can be a dual form of training error for indicating the training progress. By tracking the CS value during training, we are able to monitor the class-level performance on the fly, thus we can adaptively identify the data of the maximal informativeness.

In this work, we use the *expected confidence score (ECS)* as the class-level performance indicator. The ECS for class $c$ can be computed by

$$\mathrm{ECS}_c(x) = \mathbb{E}_{x \sim c}\left[\mathrm{CS}(x)\right], \qquad (5)$$

where with a slight abuse of notation, we use the first $c$ to conceptually represent the $c^{th}$ class, while the second $c$ to represent the *distribution* of the $c^{th}$ class. The ECS for the
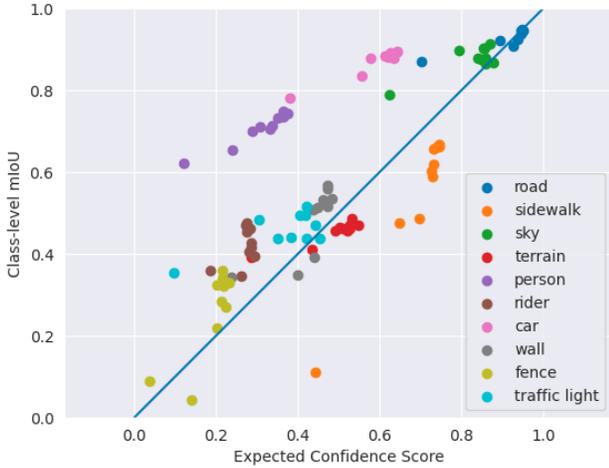
Fig. 5. Reliability diagram for class-level ECS.

source domain and the target domain can be expressed by

$$\mathrm{ECS}_c^s(x) = \mathbb{E}_{x \sim \mathcal{S}, x \sim c}[\mathrm{CS}(x)],$$
$$\mathrm{ECS}_c^t(x) = \mathbb{E}_{x \sim \mathcal{T}, x \sim c}[\mathrm{CS}(x)]. \quad (6)$$

To further validate the property of $\mathrm{ECS}_c^s$ and $\mathrm{ECS}_c^t$, we show the changes of the two ECSs during training in Fig. 4. Here we smooth the value of raw ECS by the EMA.

$$^j\mathrm{ECS}_c^s \leftarrow \tau \cdot {}^{j-1}\mathrm{ECS}_c^s + (1-\tau) \cdot {}^j\mathrm{ECS}_c^s,$$
$$^j\mathrm{ECS}_c^t \leftarrow \tau \cdot {}^{j-1}\mathrm{ECS}_c^t + (1-\tau) \cdot {}^j\mathrm{ECS}_c^t, \quad (7)$$

where $\tau$ is the smoothness weight; $j$ represents the $j^{th}$ iteration during training. As we can see from Fig. 4, ECS values of almost-all classes are monotonously increasing during training for both the source domain and target domain.

To further validate that the class-level ECS is well calibrated, we also show the reliability diagram for some of the representative classes from Cityscapes [7] in Fig. 5 where we show the relation between the ECS values and the widely-adopted segmentation performance metric *mean-Intersection-over-Union (mIoU)*. We can see that in general the mIoU is positively correlated with the value of ECS.

*2) Class Selection Strategy:* We propose our *Informed Mix* (IMix) based on the DACS [28] which is a variant of ClassMix [10], and build upon the idea of bridging images across domains. DACS is the first to apply the idea of ClassMix [10] to domain adaptation. In DACS, $x_a$ is selected from the source domain while $x_b$ is selected from the target domain. Our IMix follows the idea of DACS and applies a ClassMix-based data augmentation for domain adaptation. The differences between the proposed IMix and DACS lie in two aspects. First, we are not using random sampling to select classes instead, we select classes based on the ECS values (Eq. (7)). Second, we do not use a fixed ratio, e.g., 0.5, but a dynamic schedule to determine the value of the selection ratio. Details about ClassMix and IMix can be found in Algorithm 1.

Our IMix is based on the spatiotemporal change of data during training. We want to find a finer structure for this change to make the mixup more focused on bottleneck classes. The intuition is two-fold. First, both domains have

---

**Algorithm 1:** Functions of ClassMix and IMix

```
1  def ClassSample(i_a):
2      M=zero_like(i_a.shape)
3      u_a=unique(i_a)
4      û_a=randint(0, C, size=int(C/2))
5      M[i,j] = 1 if i_a[i,j] is in û_a
6      return M
7
8  def ISample(i_a, e, η):
9      M=zero_like(i_a.shape)
10     u_a=unique(i_a)
11     k = η·C
12     t_a = topk(e, k)
13     û_a = t_a.indices
14     M[i,j] = 1 if i_a[i,j] is in û_a
15     return M
16
17 # x_a, x_b: images from two domains
18 # F is the segmentation model
19 # C is the number of classes
20 # e ∈ℝ^C: ECS values for all classes
21 # η: Allocation ratio
22 def Mix(x_a, x_b, method, e, η):
23     y_a = F(x_a)  # y_a shape: [C, H, W]
24     y_b = F(x_b)
25     i_a = argmax(y_a, axis=0)
26     if method == 'ClassMix':
27         M = ClassSample(i_a)  # ClassMix
28     elif method == 'IMix':
29         M = ISample(i_a, e, η)  # IMix
30     x_m = M⊙x_a + (1-M)⊙x_b
31     y_m = M⊙y_a + (1-M)⊙y_b
32     return x_m, y_m
```

---

increasing performance as the training proceeds, the allocation ratio for mixing regions should incline to the target domain in a gradual manner as our goal is to boost the inference capability on the target domain. However, in the early phase of training, the ratio for source classes should be high as we still want to extract the main knowledge from the source domain. Second, we can empirically find that the performance of some classes is inferior to others during training. The overall performance might be significantly improved if those inferior classes are ameliorated. As we can learn, the dominating domain and the ratio for mixing should be adaptively adjusted along the training process.

For the convenience of analysis, we propose two concepts, *Source-Select-Target-Follow* (SSTF) and *Target-Select-Source-Follow* (TSSF). The difference between the two is the order of the class selection — which domain (selecting domain) provides the guaranteed selection of certain classes while the other one (following domain) acts accordingly. The data from the selecting domain is guaranteed to be exposed more during training, thus dominating the knowledge transfer process.

| Method | Road | S.Walk | Build | Wall | Fence | Pole | T. Light | Sign | Veg | Terrian | Sky | Person | Rider | Car | Truck | Bus | Train | MC | Bike | mIoU |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| APODA [30] | 85.6 | 32.8 | 79.0 | 29.5 | 25.5 | 26.8 | 34.6 | 19.9 | 83.7 | 40.6 | 77.9 | 59.2 | 28.3 | 84.6 | 34.6 | 49.2 | 8.0 | 32.6 | 39.6 | 45.9 |
| PatchAlign [31] | 92.3 | 51.9 | 82.1 | 29.2 | 25.1 | 24.5 | 33.8 | 33.0 | 82.4 | 32.8 | 82.2 | 58.6 | 27.2 | 84.3 | 33.4 | 46.3 | 2.2 | 29.5 | 32.3 | 46.5 |
| AdvEnt [17] | 89.4 | 33.1 | 81.0 | 26.6 | 26.8 | 27.2 | 33.5 | 24.7 | 83.9 | 36.7 | 78.8 | 58.7 | 30.5 | 84.8 | 38.5 | 44.5 | 1.7 | 31.6 | 32.4 | 45.5 |
| CBST [21] | 91.8 | 53.5 | 80.5 | 32.7 | 21.0 | 34.0 | 28.9 | 20.4 | 83.9 | 34.2 | 80.9 | 53.1 | 24.0 | 82.7 | 30.3 | 35.9 | 16.0 | 25.9 | 42.8 | 45.9 |
| MRKLD-SP [29] | 90.8 | 46.0 | 79.9 | 27.4 | 23.3 | 42.3 | 46.2 | 40.9 | 83.5 | 19.2 | 59.1 | 63.5 | 30.8 | 83.5 | 36.8 | 52.0 | 28.0 | 36.8 | 46.4 | 49.2 |
| BDL [32] | 91.0 | 44.7 | 84.2 | 34.6 | 27.6 | 30.2 | 36.0 | 36.0 | 85.0 | 43.6 | 83.0 | 58.6 | 31.6 | 83.3 | 35.3 | 49.7 | 3.3 | 28.8 | 35.6 | 48.5 |
| CADASS [33] | 91.3 | 46.0 | 84.5 | 34.4 | 29.7 | 32.6 | 35.8 | 36.4 | 84.5 | 43.2 | 83.0 | 60.0 | 32.2 | 83.2 | 35.0 | 46.7 | 0.0 | 33.7 | 42.2 | 49.2 |
| MRNet [34] | 89.1 | 23.9 | 82.2 | 19.5 | 20.1 | 33.5 | 42.2 | 39.1 | 85.3 | 33.7 | 76.4 | 60.2 | 33.7 | 86.0 | 36.1 | 43.3 | 5.9 | 22.8 | 30.8 | 45.5 |
| R-MRNet [35] | 90.4 | 31.2 | 85.1 | 36.9 | 25.6 | 37.5 | 48.8 | 48.5 | 85.3 | 34.8 | 81.1 | 64.4 | 36.8 | 86.3 | 34.9 | 52.2 | 1.7 | 29.0 | 44.6 | 50.3 |
| PIT [36] | 87.5 | 43.4 | 78.8 | 31.2 | 30.2 | 36.3 | 39.9 | 42.0 | 79.2 | 37.1 | 79.3 | 65.4 | 37.5 | 83.2 | 46.0 | 45.6 | 25.7 | 23.5 | 49.9 | 50.6 |
| SIM [37] | 90.6 | 44.7 | 84.8 | 34.3 | 28.7 | 31.6 | 35.0 | 37.6 | 84.7 | 43.3 | 85.3 | 57.0 | 31.5 | 83.8 | 42.6 | 48.5 | 1.9 | 30.4 | 39.0 | 49.2 |
| FDA [38] | 92.5 | 53.3 | 82.4 | 26.5 | 27.6 | 36.4 | 40.6 | 38.9 | 82.3 | 39.8 | 78.0 | 62.6 | 34.4 | 84.9 | 34.1 | 53.1 | 16.9 | 27.7 | 46.4 | 50.45 |
| CAG-UDA [39] | 90.4 | 51.6 | 83.8 | 34.2 | 27.8 | 38.4 | 25.3 | 48.4 | 85.4 | 38.2 | 78.1 | 58.6 | 34.6 | 84.7 | 21.9 | 42.7 | 41.1 | 29.3 | 37.2 | 50.2 |
| IAST [23] | 93.8 | 57.8 | 85.1 | 39.5 | 26.7 | 26.2 | 43.1 | 34.7 | 84.9 | 32.9 | 88.0 | 62.6 | 29.0 | 87.3 | 39.2 | 49.6 | 23.2 | 34.7 | 39.6 | 51.5 |
| DACS [28] | 89.9 | 39.7 | 87.9 | 30.7 | 39.5 | 38.5 | 46.4 | 52.8 | 88.0 | 44.0 | 88.8 | 67.2 | 35.8 | 84.5 | 45.7 | 50.2 | 0.0 | 27.3 | 34.0 | 52.1 |
| CorDA [40] | 94.7 | 63.1 | 87.6 | 30.7 | 40.6 | 40.2 | 47.8 | 51.6 | 87.6 | 47.0 | 89.7 | 66.7 | 35.9 | 90.2 | 48.9 | 57.5 | 0.0 | 39.8 | 56.0 | 56.6 |
| ProDA [41] | 87.8 | 56.0 | 79.7 | 46.3 | **44.8** | **45.6** | **53.5** | **53.5** | 88.6 | 45.2 | 82.1 | 70.7 | 39.2 | 88.8 | 45.5 | 59.4 | 1.0 | 48.9 | 56.4 | 57.5 |
| DAFormer [24] | 94.0 | 59.0 | 87.0 | 38.8 | 30.8 | 42.9 | 49.5 | 51.0 | 88.1 | **48.6** | 89.0 | 69.3 | 39.8 | 91.3 | 72.0 | 69.4 | 48.8 | 52.2 | 61.2 | 62.2 |
| HRDA [11] | 94.8 | 64.0 | **88.1** | **52.7** | 28.2 | 45.5 | 48.4 | 49.2 | **89.3** | 48.4 | 91.4 | 73.9 | 38.3 | 92.2 | 74.9 | 76.8 | **62.2** | **61.5** | 64.1 | 65.4 |
| IDA (ours) | **95.4** | **72.0** | 87.8 | 49.9 | 36.6 | 40.6 | 46.8 | 50.4 | 88.3 | 45.2 | **92.1** | **74.2** | **50.4** | **92.8** | **79.2** | **81.8** | 53.8 | 61.4 | **64.5** | **66.5** |

When selecting classes from one domain, we also need to decide whether well-performing classes or under-performing classes should be selected. Different types of classes in differing domains have different values. For example, under-performing classes in the source domain might indicate a strong signal for selection as those classes are bottleneck classes and the ground-truth label of those classes might boost the performance. On the contrary, the under-performing classes in the target domain might be out of choice as they can contain too much noise. The quality of each class is represented by the corresponding ECS value (Eq. (7)).

*3) Adaptation Schedule:* To account for the spatial change of the data during training, we propose a dynamic schedule to determine the value of $\eta$ in the function of **ISample** in Algorithm 1. The reason for using a dynamic schedule rather than a fixed ratio value is that any single fixed ratio might fail to capture the change throughout the training process. An extreme ratio, e.g., 0.1 or 0.9 can lead to a highly imbalanced mixing. We use the Kumaraswamy Cumulative Density Function (KCDF) as our basic scheduling function. The KCDF has a formulation of $y = 1 - (1 - x^a)^b$. Different KCDFs with different parameters $a$ and $b$ are shown as the solid curve in Fig. 6(a). Another variant of KCDF is expressed by $y = (1 - x^a)^b$, we denote this variant as Reversed KCDF (RKCDF), and show different RKCDFs as the dash curves in Fig. 6(a). Based on this basic function, we propose to use a truncated version of the functions to avoid extreme ratio values (see Fig. 6(b)).

## IV. EXPERIMENTS

### A. Evaluation Setup

**Datasets:** We test on three datasets. (1) **GTA-V** is a synthetic dataset collected in a simulated city environment. This dataset contains 24,966 synthetic frames with a resolution
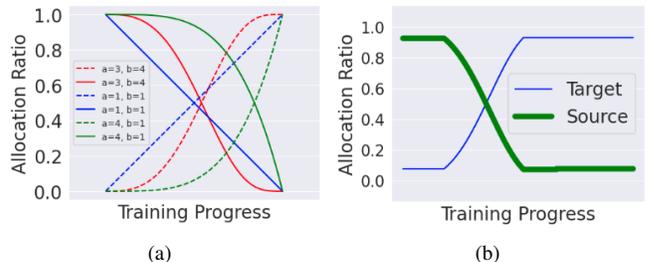


Fig. 6. (a) Kumaraswamy Cumulative Density function. (b) The schedule we use for assigning allocation ratio to the dominating domain (shown in a bold curve). In this work, the source domain dominates the allocation.

of $1914 \times 1052$. Images are provided with dense semantic annotations of 33 classes. (2) **SYNTHIA** is another city-like synthetic dataset that has 9,400 synthetic images with a resolution of $1280 \times 760$. Pixel-level semantic annotations for 13 classes are provided in SYNTHIA. (3) **Cityscapes** is a dataset containing 2,975 training images and 500 validation images with a resolution of $2048 \times 1024$. All images are collected in real European cities.

We perform two Sim2Real adaptations — one is the adaptation of GTA-V → Cityscapes and the other is the adaptation of SYNTHIA → Cityscapes. We evaluate segmentation performance with the standard *mean-Intersection-over-Union (mIoU)* metric. Evaluations for both adaptation scenarios are conducted on the 500 validation images in Cityscapes.

**Implementation Details:** We base our IDA framework on the self-training framework in HRDA [11]. We use a batch size of 1 and set the crop resolution as 952 due to the limited GPU memory. We compare our proposed IDA with recent SOTA methods [24] and [11]. To make the comparison fair, we also set the same batch size and image resolution for both

TABLE II

QUANTITATIVE COMPARISON FOR THE ADAPTATION OF SYNTHIA → CITYSCAPES.

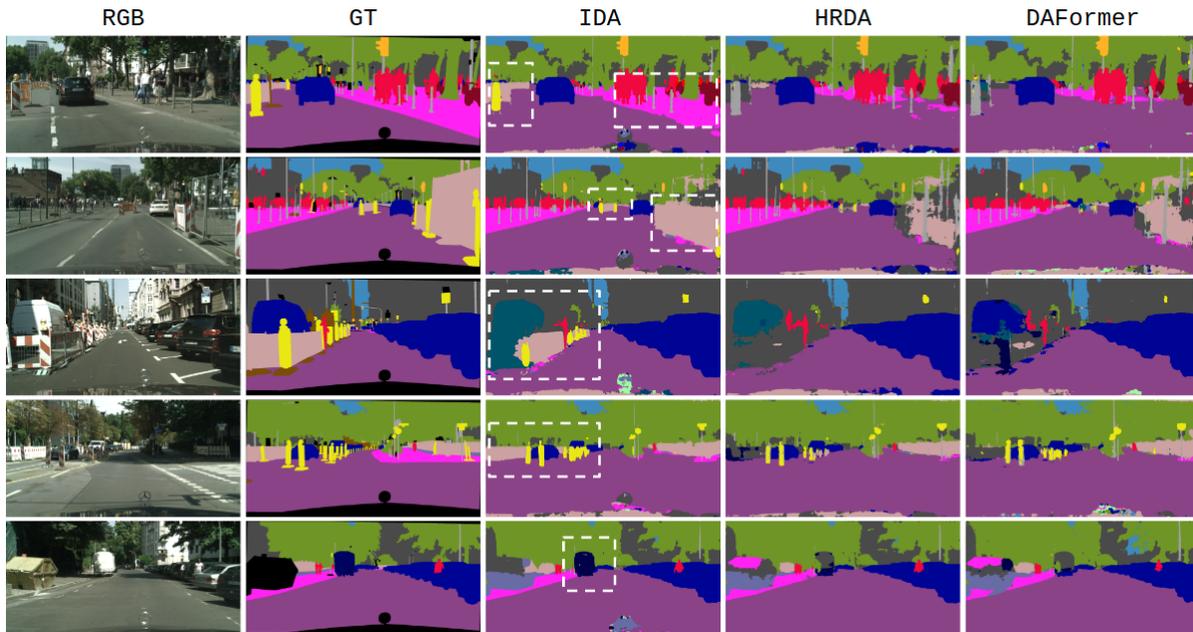| Method | Road | S.Walk | Build | Wall | Fence | Pole | T. Light | Sign | Veg | Sky | Person | Rider | Car | Bus | MC | Bike | mIoU |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PatchAlign [31] | 82.4 | 38.0 | 78.6 | 8.7 | 0.6 | 26.0 | 3.9 | 11.1 | 75.5 | 84.6 | 53.5 | 21.6 | 71.4 | 32.6 | 19.3 | 31.7 | 40.0 |
| AdvEnt [17] | 85.6 | 42.2 | 79.7 | 8.7 | 0.4 | 25.9 | 5.4 | 8.1 | 80.4 | 84.1 | 57.9 | 23.8 | 73.3 | 36.4 | 14.2 | 33.0 | 41.2 |
| CBST [21] | 68.0 | 29.9 | 76.3 | 10.8 | 1.4 | 33.9 | 22.8 | 29.5 | 77.6 | 78.3 | 60.6 | 28.3 | 81.6 | 23.5 | 18.8 | 39.8 | 42.6 |
| MRKLD [29] | 67.7 | 32.2 | 73.9 | 10.7 | 1.6 | 37.4 | 22.2 | 31.2 | 80.8 | 80.5 | 60.8 | 29.1 | 82.8 | 25.0 | 19.4 | 45.3 | 43.8 |
| MRNet [34] | 82.0 | 36.5 | 80.4 | 4.2 | 0.4 | 33.7 | 18.0 | 13.4 | 81.1 | 80.8 | 61.3 | 21.7 | 84.4 | 32.4 | 14.8 | 45.7 | 43.2 |
| R-MRNet [35] | 87.6 | 41.9 | 83.1 | 14.7 | 1.7 | 36.2 | 31.3 | 19.9 | 81.6 | 80.6 | 63.0 | 21.8 | 86.2 | 40.7 | 23.6 | 53.1 | 47.9 |
| PIT [36] | 83.1 | 27.6 | 81.5 | 8.9 | 0.3 | 21.8 | 26.4 | 33.8 | 76.4 | 78.8 | 64.2 | 27.6 | 79.6 | 31.2 | 31.0 | 31.3 | 44.0 |
| CAG-UDA [39] | 84.7 | 40.8 | 81.7 | 7.8 | 0.0 | 35.1 | 13.3 | 22.7 | 84.5 | 77.6 | 64.2 | 27.8 | 80.9 | 19.7 | 22.7 | 48.3 | 44.5 |
| IAST [23] | 81.9 | 41.5 | **83.3** | 17.7 | 4.6 | 32.3 | 30.9 | 28.8 | 83.4 | 85.0 | 65.5 | 30.8 | **86.5** | 38.2 | 33.1 | 52.7 | 49.8 |
| DACS [28] | 80.5 | 25.1 | 81.9 | 21.4 | 2.8 | 37.2 | 22.6 | 23.9 | 83.6 | **90.7** | 67.6 | 38.3 | 82.9 | 38.9 | 28.4 | 47.5 | 48.3 |
| DAFormer [24] | 82.3 | 36.9 | 76.1 | 41.8 | **6.0** | 44.5 | 45.1 | 46.5 | **85.9** | 82.9 | 68.0 | 44.4 | 84.9 | 47.3 | 49.1 | 57.5 | 56.2 |
| HRDA [11] | 83.0 | 43.9 | 76.5 | **49.8** | 4.3 | **51.7** | **55.8** | **52.8** | 85.2 | 80.0 | 68.2 | 43.0 | 80.7 | 56.7 | **59.1** | **59.7** | 59.4 |
| IDA (ours) | **88.9** | **44.2** | 78.2 | 49.1 | 4.9 | 48.6 | 52.3 | 49.3 | 84.9 | 88.2 | **70.1** | **47.0** | 85.3 | **58.2** | 58.7 | 56.9 | **60.3** |



Fig. 7. Comparison of different methods on the Cityscapes validation images that show challenging situations. IDA can still maintain a high segmentation quality even when objects are small, ambiguous, and highly irregular, e.g., regions marked by the white dotted boxes.

baselines as ours. More details about network structure and hyperparameters can be found in HRDA [11].

*B. Comparison*

We compare our proposed IDA model with the baseline UDA-SS methods both quantitatively and qualitatively. We first consider the adaptation of GTA-V → Cityscapes. The quantitative comparison can be seen in Table. I. Our IDA model exhibits the best overall mIoU performance among all the listed methods. IDA outperforms the current UDA-SS SOTA work HRDA [11] in the majority of classes (10 out of 19) and shows a considerable advantageous margin of 1.1% mIoU. The IDA model also shows superior performance for some challenging classes, e.g., Person, Rider, Car, Truck, Bus, Bike, etc. This superiority is consistent with our expectation for the IDA model as it especially aims to improve the performance of bottleneck classes during training. Note that the results of DAFormer and HRDA in Table I are different than the reported ones in the original works as we have

new hyperparameter settings for a fair comparison with our IDA model. The quantitative comparison for the adaptation of SYNTHIA→Cityscapes is shown in Table II. The IDA model can still outperform the strongest baseline HRDA by a margin of 0.9% mIoU and show advatanges on challenging classes such as Person, Rider and Bus. The effectiveness of our proposed IDA model can also be seen in visual examples in Fig. 7.

*C. Ablation Studies*

**Fixed selection ratios:** In IDA we adopt the SSTF-U mixing strategy where we first select the ground-truth regions of source under-performing classes to construct the mixing mask and then the target regions are selected according to the reverse source mask. Throughout the training process, we use a dynamic schedule to determine the value of the ratio for source class selection. Thus we conduct the ablation study showing issues with different fixed ratio values under different source class selection strategies, see Table. III.

TABLE III

COMPARISON OF USING DIFFERENT RATIO VALUES UNDER DIFFERENT
SELECTION STRATEGIES.

| | Selection | Class | Ratio | mIoU | $\delta$ (↑) | $\Delta$ (↑) |
|---|---|---|---|---|---|---|
| 0 | | | 0.1 | 37.2 | -29.3 | |
| 1 | | | 0.3 | 41.4 | -25.1 | |
| 2 | SSTF | W | 0.5 | 50.7 | -15.8 | -22.5 |
| 3 | | | 0.7 | 49.2 | -17.3 | |
| 4 | | | 0.9 | 41.3 | -25.2 | |
| 5 | | | 0.1 | 40.0 | -16.5 | |
| 6 | | | 0.3 | 45.4 | -11.1 | |
| 7 | SSTF | U | 0.5 | 58.2 | -3.3 | -8.32 |
| 8 | | | 0.7 | 55.2 | -2.3 | |
| 9 | | | 0.9 | 48.1 | -8.4 | |
| 10 | | | 0.1 | 50.6 | -15.9 | |
| 11 | | | 0.3 | 47.3 | -19.2 | |
| 12 | TSSF | W | 0.5 | 38.4 | -28.1 | -25.9 |
| 13 | | | 0.7 | 35.0 | -31.0 | |
| 14 | | | 0.9 | 31.4 | -35.1 | |
| 15 | | | 0.1 | 47.2 | -19.3 | |
| 16 | | | 0.3 | 49.3 | -17.2 | |
| 17 | TSSF | U | 0.5 | 45.6 | -20.9 | -25.4 |
| 18 | | | 0.7 | 33.4 | -33.1 | |
| 19 | | | 0.9 | 30.2 | -36.3 | |

TABLE IV

COMPARISON OF USING DIFFERENT SMOOTHNESS FOR ECS.

| Smoothness weight | 0 | 0.1 | 0.3 | 0.5 | 0.7 | 0.9 | 0.999 |
|---|---|---|---|---|---|---|---|
| mIoU | 62.1 | 61.5 | 60.8 | 63.2 | 64.9 | 66.2 | 66.5 |

Table. IV. It can be seen that the performance is generally increasing as the smoothness is lifted. The performance with the raw indicator values is the lowest.

The reason for this trend is the smoothed indicator is more stable than the raw values of the indicator which may change significantly among iterations, causing instability in the selection of classes and possibly a large distribution shift during training.
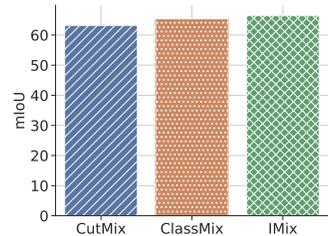


Fig. 8. Comparison of using different mixups.

**Different mixups:** In our work we propose a new mixup, IMix, for augmenting data in the IDA model. We compare the performance of the adaptation with different mixups, which is shown in Fig. 8. Compared with previous region-based mixups that use random sampling to generate new mixed data, our IMix considers dynamic changes in the data from the two domains. By capturing the fine structure of the adaptation, our IMix achieves the best performance among all listed mixups.

We show the testing performance on validation data of Cityscapes in the column of *mIoU* and the degradation from the best IDA model in the column of $\delta$. The column of $\Delta$ shows the mean of $\delta$ for each selection with a certain class type. We can see none of the settings in Table. III can achieve a positive $\Delta$. All methods with fixed ratios are degrading with a significant drop. The best result with a fixed ratio (SSTF-U-0.7) still has a gap of 2.3% mIoU compared with the best IDA.

From the value of $\Delta$ in Table. III, we can first validate that the SSTF is the most effective selection method. Based on the SSTF selection, we should give priority to selecting the region of *under-performing* classes in the source domain. Those underperforming classes can be treated as bottleneck classes for both domains, thus they can provide strong supervision to drive the improvement of the overall performance. An extreme value of the ratio might cause unacceptable damage to the adaptation. The performance with a mild value of the selection ratio, e.g., 0.5 or 0.7 can be better than other values, but still worse than the performance with dynamic scheduling. Our IDA model is even better than HRDA. The reason for this is we use the dynamic scheduling to balance the bias that we have introduced into the data, thus we are able to extract knowledge from under-performing classes while maintaining the data unbiased. One thing we need to note is that HRDA [11] uses a selection ratio of 0.5, but the performance of HRDA is better than IDA-SSTF-U-0.7. The reason for this drop is the bias we have injected by the selection strategy of SSTF-U to either well-performing classes or underperforming classes. On the contrary, HRDA uses random sampling such that the classes of the new data are not biased to any certain types, leading to better results. **Smoothness of the indicator:** We use the smoothed ECS as the class-level performance indicator. Here we show the necessity of using the smoothed values instead of the raw values. The quantitative comparison is shown in

## V. CONCLUSION

We present a principled model, Informed Domain Adaptation (IDA), for the un-supervised domain adaptive semantic segmentation. Our proposed IDA model is a self-training framework that exploits the obscured informativeness of data to improve the learning efficiency. To achieve this, we propose a new mixup technique, IMix, that bridges the source and target domains according to the training progress defined by an expected confidence assessment. We also propose a novel dynamic adaptation schedule which can adaptively adjust the mixing ratio for different domains. Extensive evaluations on popular datasets reveal that the IDA outperforms the SOTA model with a remarkable margin.

## REFERENCES

[1] Juan Jose Rubio, Takahiro Kashiwa, Teera Laiteerapong, Wenlong Deng, Kohei Nagai, Sergio Escalera, Kotaro Nakayama, Yutaka Matsuo, and Helmut Prendinger. Multi-class structural damage segmentation using fully convolutional networks. *Computers in Industry*, 112:103121, 2019.

[2] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015.

[3] Stephan R Richter, Vibhav Vineet, Stefan Roth, and Vladlen Koltun. Playing for data: Ground truth from computer games. In *European conference on computer vision*, pages 102–118. Springer, 2016.

[4] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. *Advances in Neural Information Processing Systems*, 34:12077–12090, 2021.

[5] Wenwei Zhang, Jiangmiao Pang, Kai Chen, and Chen Change Loy. K-net: Towards unified image segmentation. *Advances in Neural Information Processing Systems*, 34:10326–10338, 2021.

[6] Robin Strudel, Ricardo Garcia, Ivan Laptev, and Cordelia Schmid. Segmenter: Transformer for semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7262–7272, 2021.

[7] 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 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3213–3223, 2016.

[8] Christos Sakaridis, Dengxin Dai, and Luc Van Gool. ACDC: The adverse conditions dataset with correspondences for semantic driving scene understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021.

[9] Sangdoo Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6023–6032, 2019.

[10] Viktor Olsson, Wilhelm Tranheden, Juliano Pinto, and Lennart Svensson. Classmix: Segmentation-based data augmentation for semi-supervised learning. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1369–1378, 2021.

[11] Lukas Hoyer, Dengxin Dai, and Luc Van Gool. Hrda: Context-aware high-resolution domain-adaptive semantic segmentation. *arXiv preprint arXiv:2204.13132*, 2022.

[12] Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural networks. *The journal of machine learning research*, 17(1):2096–2030, 2016.

[13] Judy Hoffman, Dequan Wang, Fisher Yu, and Trevor Darrell. Fcns in the wild: Pixel-level adversarial and constraint-based adaptation. *arXiv preprint arXiv:1612.02649*, 2016.

[14] 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 *International conference on machine learning*, pages 1989–1998. Pmlr, 2018.

[15] Yi-Hsuan Tsai, Wei-Chih Hung, Samuel Schulter, Kihyuk Sohn, Ming-Hsuan Yang, and Manmohan Chandraker. Learning to adapt structured output space for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7472–7481, 2018.

[16] Kuniaki Saito, Kohei Watanabe, Yoshitaka Ushiku, and Tatsuya Harada. Maximum classifier discrepancy for unsupervised domain adaptation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3723–3732, 2018.

[17] Tuan-Hung Vu, Himalaya Jain, Maxime Bucher, Matthieu Cord, and Patrick Pérez. Advent: Adversarial entropy minimization for domain adaptation in semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2517–2526, 2019.

[18] Yawei Luo, Liang Zheng, Tao Guan, Junqing Yu, and Yi Yang. Taking a closer look at domain shift: Category-level adversaries for semantics consistent domain adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2507–2516, 2019.

[19] Zheng Chen, Durgakant Pushp, and Lantao Liu. Cali: Coarse-to-fine alignments based unsupervised domain adaptation of traversability prediction for deployable autonomous navigation. *arXiv preprint arXiv:2204.09617*, 2022.

[20] Haoran Wang, Tong Shen, Wei Zhang, Ling-Yu Duan, and Tao Mei. Classes matter: A fine-grained adversarial approach to cross-domain semantic segmentation. In *European conference on computer vision*, pages 642–659. Springer, 2020.

[21] Yang Zou, Zhiding Yu, BVK Kumar, and Jinsong Wang. Unsupervised domain adaptation for semantic segmentation via class-balanced self-training. In *Proceedings of the European conference on computer vision (ECCV)*, pages 289–305, 2018.

[22] Yang Zhang, Philip David, and Boqing Gong. Curriculum domain adaptation for semantic segmentation of urban scenes. In *Proceedings of the IEEE international conference on computer vision*, pages 2020–2030, 2017.

[23] Ke Mei, Chuang Zhu, Jiaqi Zou, and Shanghang Zhang. Instance adaptive self-training for unsupervised domain adaptation. In *European conference on computer vision*, pages 415–430. Springer, 2020.

[24] Lukas Hoyer, Dengxin Dai, and Luc Van Gool. Daformer: Improving network architectures and training strategies for domain-adaptive semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9924–9935, 2022.

[25] Wen Zhang, Lingfei Deng, Lei Zhang, and Dongrui Wu. A survey on negative transfer. *arXiv preprint arXiv:2009.00909*, 2020.

[26] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020.

[27] Nikita Araslanov and Stefan Roth. Self-supervised augmentation consistency for adapting semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15384–15394, 2021.

[28] Wilhelm Tranheden, Viktor Olsson, Juliano Pinto, and Lennart Svensson. Dacs: Domain adaptation via cross-domain mixed sampling. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1379–1389, 2021.

[29] Yang Zou, Zhiding Yu, Xiaofeng Liu, BVK Kumar, and Jinsong Wang. Confidence regularized self-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5982–5991, 2019.

[30] Jihan Yang, Ruijia Xu, Ruiyu Li, Xiaojuan Qi, Xiaoyong Shen, Guanbin Li, and Liang Lin. An adversarial perturbation oriented domain adaptation approach for semantic segmentation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 12613–12620, 2020.

[31] Yi-Hsuan Tsai, Kihyuk Sohn, Samuel Schulter, and Manmohan Chandraker. Domain adaptation for structured output via discriminative patch representations. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1456–1465, 2019.

[32] Yunsheng Li, Lu Yuan, and Nuno Vasconcelos. Bidirectional learning for domain adaptation of semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6936–6945, 2019.

[33] Jinyu Yang, Weizhi An, Chaochao Yan, Peilin Zhao, and Junzhou Huang. Context-aware domain adaptation in semantic segmentation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 514–524, 2021.

[34] Zhedong Zheng and Yi Yang. Unsupervised scene adaptation with memory regularization in vivo. *arXiv preprint arXiv:1912.11164*, 2019.

[35] Zhedong Zheng and Yi Yang. Rectifying pseudo label learning via uncertainty estimation for domain adaptive semantic segmentation. *International Journal of Computer Vision*, 129(4):1106–1120, 2021.

[36] Fengmao Lv, Tao Liang, Xiang Chen, and Guosheng Lin. Cross-domain semantic segmentation via domain-invariant interactive relation transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4334–4343, 2020.

[37] Zhonghao Wang, Mo Yu, Yunchao Wei, Rogerio Feris, Jinjun Xiong, Wen-mei Hwu, Thomas S Huang, and Honghui Shi. Differential treatment for stuff and things: A simple unsupervised domain adaptation method for semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12635–12644, 2020.

[38] Yanchao Yang and Stefano Soatto. Fda: Fourier domain adaptation for semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4085–4095, 2020.

[39] Qiming Zhang, Jing Zhang, Wei Liu, and Dacheng Tao. Category anchor-guided unsupervised domain adaptation for semantic segmentation. *Advances in Neural Information Processing Systems*, 32, 2019.

[40] Qin Wang, Dengxin Dai, Lukas Hoyer, Luc Van Gool, and Olga Fink. Domain adaptive semantic segmentation with self-supervised depth estimation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8515–8525, 2021.

[41] Pan Zhang, Bo Zhang, Ting Zhang, Dong Chen, Yong Wang, and Fang Wen. Prototypical pseudo label denoising and target structure learning for domain adaptive semantic segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12414–12424, 2021.