

Distribution Regularized Self-Supervised Learning for Domain Adaptation of Semantic Segmentation

Javed Iqbal^{a,*}, Hamza Rawal^a, Rehan Hafiz^a, Yu-Tseh Chi^b, Mohsen Ali^a

^aInformation Technology University, Lahore, 54000, Pakistan

^bFacebook, 1 Hacker Way, Menlo Park, CA 94025, USA.

Abstract

This paper proposes a novel pixel-level distribution regularization scheme (DRSL) for self-supervised domain adaptation of semantic segmentation. In a typical setting, the classification loss forces the semantic segmentation model to greedily learn the representations that capture inter-class variations in order to determine the decision (class) boundary. Due to the domain-shift, this decision boundary is unaligned in the target domain, resulting in noisy pseudo labels adversely affecting self-supervised domain adaptation. To overcome this limitation, along with capturing inter-class variation, we capture pixel-level intra-class variations through class-aware multi-modal distribution learning (MMDL). Thus, the information necessary for capturing the intra-class variations is explicitly disentangled from the information necessary for inter-class discrimination. Features captured thus are much more informative, resulting in pseudo-labels with low noise. This disentanglement allows us to perform separate alignments in discriminative space and multi-modal distribution space, using cross-entropy based self-learning for former. For later, we propose novel stochastic mode alignment method, by explicitly decreasing the distance between the target and source pixels that map to the same mode. The distance metric learning loss, computed over pseudo-labels and backpropagated from multi-modal modeling head, acts as the regularizer over the base network shared with the segmentation head. The results from comprehensive experiments on synthetic to real domain adaptation setups, i.e., GTA-V/SYNTHIA to Cityscapes, show that DRSL outperforms many existing approaches (a minimum margin of 2.3% and 2.5% in mIoU for SYNTHIA to Cityscapes).

Keywords: Semantic Segmentation, Self-supervised Learning, Domain Adaptation, Multi-modal distribution learning.

1. Introduction

In recent years, deep neural network based semantic segmentation models have achieved considerable success. This success is much reliant on the large pixel-level annotated dataset over which these models are trained. However, like many other deep neural network based models, semantic segmentation models suffer from considerable performance degradation when tested on images from the domain different than then one used in training. This problem, attributed to the domain shift, is exacerbated in semantic segmentation algorithms since many of them are trained on the synthetic dataset, due to lack of large real-world annotated datasets, and are tested over the real-world images. Retraining or fine-tuning for new domains is expensive, time consuming, and in many cases not possible due to the large number of ever-changing domains, especially in case of autonomous vehicles, and unavailability of annotated data.

To overcome domain shift, unsupervised domain adaptation (UDA), has been employed with reasonable success [1, 2, 3], but state-of-the-art is still lacking desired accuracy. Many un-

supervised domain adaptation algorithms for semantic segmentation [4, 5, 6, 7, 8, 9, 10] perform global marginal distribution alignment through adversarial learning to translate the input image or feature volume or output probability tensor from one domain to other. The adversarial loss looks at the whole tensor (image/feature or output probability) even when the objective is to improve the pixel-level label assignments [6], moreover aligning marginal distributions does not guarantee preserving the discriminative information across the domain [11]. Self-supervised learning methods [3, 12, 13, 14, 2, 7, 15] (either independently or along with adversarial learning) try to overcome this challenge by back-propagating the cross-entropy loss computed over pixel-level pseudo-labels generated by the source model. Quality of these pseudo-labels is dependent upon the generalization capacity of the classifier and effects overall adaptation process. The deep neural network based semantic segmentation model when trained by minimizing cross-entropy loss, greedily learns representations that capture inter-class variations. When optimally trained these inter-class variations should help map accurate decision boundary, projecting pixels from different classes to different sides of it (decision boundary). However, due to the domain shift, the decision boundary is not aligned in target domain, resulting in noisy pseudo-labels leading to poor self-supervised domain adaptation.

Previous works [16, 17] have shown discriminative clus-

*Corresponding author

Email addresses: javed.iqbal@itu.edu.pk (Javed Iqbal), mscs18004@itu.edu.pk (Hamza Rawal), rehan.hafiz@itu.edu.pk (Rehan Hafiz), jchi@fb.com (Yu-Tseh Chi), mohsen.ali@itu.edu.pk (Mohsen Ali)

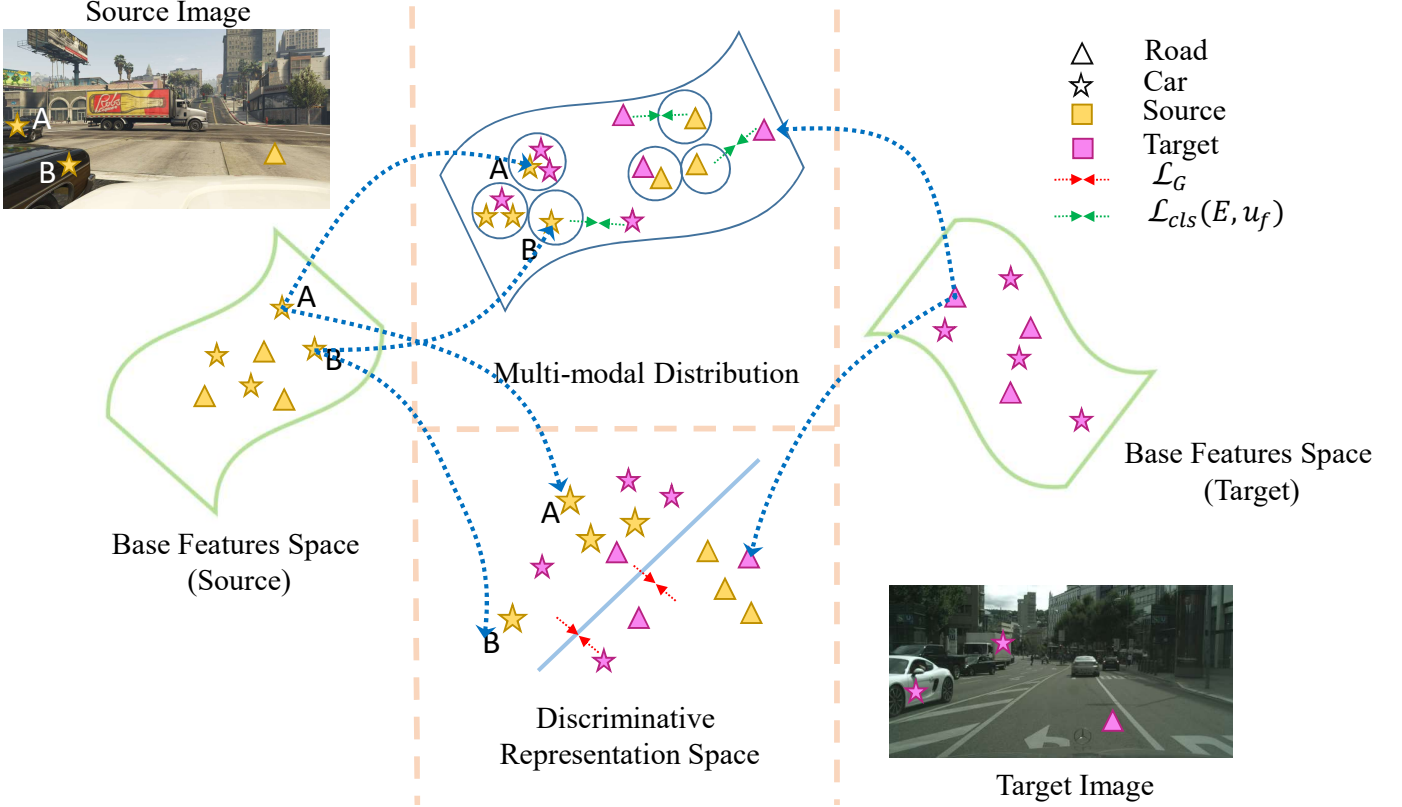


Figure 1: Separately capturing pixel-level intra-class variations and inter-class discriminative information, enables us to perform different alignment operations, i.e. *class aware mode alignment* & *cross-entropy based decision boundary alignment*. A & B are consistent in-term of being on same side of decision boundary, but variant enough to map to different modes.

tering on target data and moment matching across domains helps in adaptation. CAG-UDA [11] & [18] tried to align the class aware cluster centers across domains for better adaptation. However, visual semantic classes exhibit large set of variations, due to difference in texture, style, color, pose, illumination etc.. These variations are generally assumed to be across instance, e.g. two different types of cars, but they do manifest frequently in the same instance too, e.g. pixels belonging to different road locations or to different parts of car. Class aware single cluster based alignment might align centers of the source and target domain without aligning overall distribution, leaving classes with large variations vulnerable to misclassification in target domain. Learning to capture intra-class variations by representing each class with multiple modes and aligning the modes across domain might overcome these challenges.

Therefore, we propose a novel class aware multi-modal distribution alignment method for unsupervised domain adaption of semantic segmentation model. We combine together the ideas of distribution alignment and pseudo-label based adaptation, however, instead of just using discriminatively learned features during the adaptation, we explicitly learn representations separately. In addition to learning the inter-class variation through minimizing cross-entropy loss, i.e. the pixel-level intra-class features variations are captured by learning a multi-modal for each class (Fig. 1), resulting in a much more generalized representation. Both of these tasks have competing require-

ments, minimizing cross entropy loss results in learning inter-class discriminative representation along with intra-class consistency. Whereas multi-modal distribution learning intends to preserve information that can model intra-class variations. We disentangle these two information requirements by developing class-aware multi-modal distribution learning (MMDL) module, parallel to standard segmentation head. MMDL extracts the spatially low-resolution feature volume from the encoding block and maps to the spatially high-resolution embedding. Class aware multi-modal modeling is performed over these embedding using Distance metric learning [19]. Since both of these heads share the backbone, simultaneously decreasing loss on both act as a regularizer over the learned features, resulting in the less noisy pseudo-labels. During domain adaptation, the high quality pseudo-labels allow us to learn domain-invariant class discriminative feature representations in the discriminative space. At the same time, stochastic mode-alignment is performed across domains, by minimizing distance between representation of source pixels and target pixels mapping to same mode; thus preserving intra-class variations. Modes themselves are updated by increasing the posterior probability of target pixel belonging to the mode identified closest to target. During adaptation too, these losses computed parallelly act as regularizer over each other, hence dampening each others noise.

Our contributions are summarized as follows. First, we propose a multi-modal distribution alignment strategy for the self-

supervised domain adaptation. By designing a multi-modal distribution learning (MMDL) module parallel to standard segmentation head, with shared backbone, we disentangle inter-class discriminative and intra-class variation information; allowing them to be used during adaptation separately. We show that due to regularization of MMDL, the pseudo-labels generated over target domain are more accurate. Lastly, to perform stochastic mode alignment, we introduce the *cross domain consistency loss*. We present state-of-the-art performance for benchmark synthetic to real, e.g., GTA-V/SYNTHIA to Cityscapes adaptation.

2. Related Work

The domain shift between testing and training data deteriorates the model performance in most of the computer vision tasks like classification [20, 21, 22, 23, 24, 25], object detection [26, 27, 28] and semantic segmentation [5, 29, 6, 30, 31, 4, 32, 14, 33]. In this work, we focus on the domain shift problem for semantic segmentation with self-supervised learning. Our work is related to semantic segmentation, domain adaptation, and self-supervised learning.

Domain Adaptation for Semantic segmentation:

Recent works [34, 3, 14, 29, 5, 35, 13, 32, 1, 36, 9, 37] aiming to minimize the distribution gap between source and target domains are focused in two main directions. 1) adversarial learning and, 2) self-supervised learning for unsupervised domain adaptation (UDA) of semantic segmentation.

Adversarial Domain Adaptation: Adversarial learning is the most explored area for output space [35, 38, 34, 12, 9], latent/feature space [5, 39] and input space adaptation [4, 6, 40, 7]. We briefly describe the feature space/feature alignment, as our work is related to it. The authors in [41, 4, 42] used adversarial loss to minimize the distribution gap between the high level features representations of the source and target domain images. However, these methods do not align class-wise distribution shifts but instead match the global marginal distributions. To overcome this, [5, 6] combined category level adversarial loss (by defining class discriminators) with domain discriminator at feature space. [32] tried to regularize the segmentation network using weak labels along with latent space marginal distribution alignment for domain adaptation of semantic segmentation. Similarly, the authors in [33] investigated the robustness of the UDA of semantic segmentation and proposed a self-training augmented adversarial learning to improve the robustness to adversarial examples. Their approach resulted better performance in the presence of adversarial examples, however, reducing the performance over normal input images.

Self-supervised learning: Self-supervised learning for UDA is recently studied for major computer vision tasks like semantic segmentation and object detection [43, 3, 27, 13]. The authors in [1] proposed a self-paced self-training approach by generating class balanced pseudo-labels and class spatial priors extracted from the source dataset used to condition the pseudo-label generation. Zou et al. [2] extended the [1] with confidence regularization strategies and soft pseudo-labels for self-

training based UDA for semantic segmentation. LSE [14] further worked with self-generated scale-invariant examples and entropy based dynamic selection for self-supervised learning. The authors in [37] proposed a domain-aware meta-learning approach (MetaCorrection) to correct the segmentation loss and condition the pseudo-labels based on noise transition matrix. They report considerable mIoU gain especially when applied on pre-adapted model. In this work, we exploit a strategy similar to [1] to generate pseudo-labels for target domain images during adaptation.

Clustering Based Features Regularization: Some previous works also explored the effect of discriminative clustering on target data and moment matching across domains for target data adaptation [16, 17]. Recently, [11, 18] tried to define category anchors on the last feature volume of the segmentation model to align class aware centers across the source and target domains. Tsai et al. [44] tried to match the clustering distribution of discriminative patches from source and target domain images. Similarly, [3] and [13] exploited latent space and output space respectively by defining category based classification modules, forcing towards class-aware adaptation. However, these methods do not explore the intra-class variations present in source or target data but instead leverage the discriminative property to align the inter-class clusters. We specifically focus to capture the intra-class variations present in the source and target data by learning class-aware mixture models to help the adaptation.

3. Distribution Regularised Self-supervised Learning

In this section, we provide details of our distribution regularized self-supervised learning (DRSL) architecture. It employs DeepLab-v2 [45] as a baseline and embeds new components that enable the semantic segmentation model to be robust to domain shift.

3.1. Preliminaries

For supervised semantic segmentation, we have access to source domain images $\{x_s, y_s\}$ from $X_s \in \mathbb{R}^{H \times W \times 3}$ with corresponding ground truth labels $Y_s \in \mathbb{R}^{H \times W \times K}$. The $\{H, W\}$ shows the width and height of source domain images and K shows the number of classes. Let \mathcal{G} be a segmentation model with weights w_g that predicts the K channel softmax probability outputs. For a given source image x_s , the segmentation probability vector of class c at any pixel location (i, j) is obtained as $p(c|x_s, w_g)_{i,j} = \mathcal{G}(x_s)_{i,j}$. For fully labeled source data, the network parameters w_g are learned by minimizing the cross entropy loss (Eq. 1),

$$\mathcal{L}_{seg}^s(x_s, y_s) = - \sum_{i=1}^H \sum_{j=1}^W \sum_{c=1}^K y_s^{(c,i,j)} \log(p(c|x_s, w_g)_{c,i,j}) \quad (1)$$

where \mathcal{L}_{seg}^s is the source domain segmentation loss. For unsupervised domain adaptation of the target domain, we have access to the target domain images $\{x_t, -\}$ from $X_t \in \mathbb{R}^{H_t \times W_t \times 3}$ with no ground truths available. Thus, we adapt the iterative process used by [3, 1] to first generate pseudo-labels \hat{y}_t using

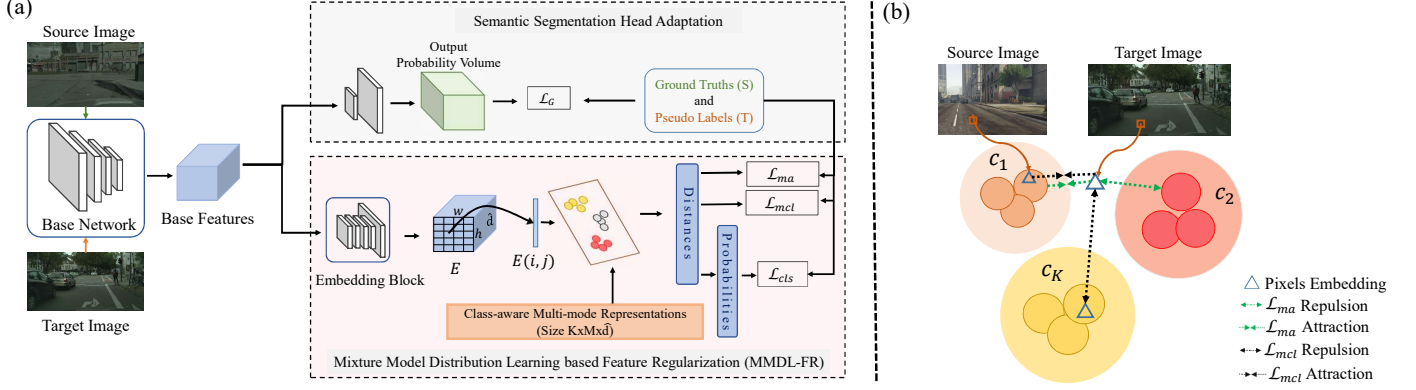


Figure 2: The proposed DRSL approach (a) Base features extracted from the base network are used for two separate tasks. The MMDL-FR module captures intra-class variations through multi-modal distribution learning. Semantic Segmentation head estimates the discriminative class boundaries necessary for the primary segmentation task. This disentanglement allows us simultaneous alignment in discriminative and multi-modal space, allowing MMDL-FR module to act as a regularizer over the Segmentation Head. (b) The proposed *Stochastic mode alignment*: Minimizing \mathcal{L}_{mcl} brings the source and target embeddings of the same mode of the same class closer than any source pixel’s embedding belonging to different class. \mathcal{L}_{ma} decreases the in-mode variance for the target samples by forcing them to come closer to the assigned mode and move away from other class’s modes.

the source trained model and then fine-tune the source trained model on target data using Eq. 2.

$$\mathcal{L}_{seg}^t(x_t, \hat{y}_t) = - \sum_{i=1}^{H_t} \sum_{j=1}^{W_t} b_t^{(i,j)} \sum_{c=1}^K \hat{y}_t^{(c,i,j)} \log(p(c|x_t, w_g)_{c,i,j}) \quad (2)$$

where \mathcal{L}_{seg}^t is the segmentation loss for target domain images with respect to generated pseudo-labels \hat{y}_t . b_t represents a binary mask with same resolution as \hat{y}_t to back-propagate loss for pixels which are assigned pseudo-labels. The total loss for the segmentation model is the combination of true labels based source domain loss and pseudo-labels based target domain loss and is given by Eq. 3,

$$\mathcal{L}_G(x_s, y_s, x_t, \hat{y}_t) = \mathcal{L}_{seg}^s(x_s, y_s) + \mathcal{L}_{seg}^t(x_t, \hat{y}_t) \quad (3)$$

3.2. Multi-Modal Distribution Learning

We propose to learn the complex intra-class variations through a multi-modal distribution learning (MMDL) framework where instead of a single cluster/anchor, each class is represented by multiple modes. This diverse representation of each class is used in the adaptation process to align the domains on fine-grained level. Furthermore, we disentangle the task of learning these intra-class variations (MMDL) from the main segmentation task by designing a separate module for it called multi-modal distribution learning based feature regularization (MMDL-FR). The proposed MMDL-FR module is model agnostic and can be appended at the encoder of any segmentation network.

The MMDL-FR module consists of mixture models based per-pixel classification augmented with distance metric learning (DML) based per-pixel embedding block. The input of the MMDL-FR module is the feature volume $F \in \mathbb{R}^{h \times w \times d}$, where $\{h, w, \text{and } d\}$ shows the spatial height, width and depth of the encoder output (base features) as shown in Fig. 2(a). The embedding block is comprised of 4 fully convolutional layers

with different dilation rates (similar to ones used in the last layer of the segmentation network) followed by an upsampling layer. The output of the embedding block \mathcal{E} is a feature volume $E = \mathcal{E}(F) \in \mathbb{R}^{h_o \times w_o \times \hat{d}}$, where $(h_o, w_o) = (H/2, W/2)$ (Sec.4.2.3) and $d \gg \hat{d}$ for any randomly selected source image.

To train the MMDL-FR module, we adapt a formulation similar to [19]. For each class c , a multi-modal distribution with M number of modes is learned. Let $e = E(i, j)$ be embedding for location (i, j) , a vector V_m^c represent the center of the mode m , ($m = 1, \dots, M$) of the class c , ($c = 1, \dots, K$) of the mixture models. In this work, these mode centers are formulated as the weights of a fully connected layer with size $K \cdot M \cdot \hat{d}$, and are reshaped into $(K \times M) \times \hat{d}$ producing $K \times M$ matrix for each input embedding vector e . This simple method makes it easy to flow back gradients to the fully connected layer and learn the segmentation backbone during training. To compute the classification probability for each embedding vector e , we compute the euclidean distance $D_m^c(e) = \|e - V_m^c\|_2^2$ between e and representative V_m^c and compute the posterior probabilities $q_m^c(e) \propto \exp(-(D_m^c(e))^2 / 2\sigma^2)$, where σ^2 is the variance of each mode and is set to 0.5. For class c posterior probability, we take the maximum over M modes of class c as, $Q(C = c|e) = \max_{m=1, \dots, M} q_m^c(e)$, where $C = c$ shows class c .

Loss Functions: To train the MMDL-FR module, two losses are used, i.e., triplet loss and the cross entropy loss. The triplet loss for *embedding block* is defined by Eq. 4,

$$\mathcal{L}_{emb}(E) = \sum_{e \in E} |\min_m D_m^{c^*}(e) - \min_{m, c^* \neq c} D_m^c(e) + \alpha|_+ \quad (4)$$

where $|\cdot|_+$ is the Relu function and α is the minimum margin between the distance of an embedding e to the closest mode representative $V_m^{c^*}$ of the true class c^* , and distance of embedding e to the closest mode representative of the incorrect class V_m^c . Similarly, the cross entropy loss for mixture models based

classification is given by Eq. 5,

$$\mathcal{L}_{cls}(E, u_f) = - \sum_{e \in E} \sum_{c=1}^K u_f^{(c)} \log(Q(C = c|e)) \quad (5)$$

where $u_f^{(c)}$ is the embedding classification label obtained from y_s^c or \hat{y}_t^c for class c . The triplet loss enforces the embedding block to learn representation that capture intra-class variation information, while cross entropy loss pushes them to not lose necessary class-specific information. Due to these two losses, the MMDL-FR module acts as a regularizer at latent space over the shared backbone, so that the shared features are much more informative if only segmentation head is used.

3.3. Stochastic Mode Alignment

One of the characteristic of domain generalization will be that the multi-modal distribution learning over one domain should result in the modes which are very close to the modes learned in the other domain. However, due to the domain shift, this is not generally true. That is in the target domain, the features of pixels assigned pseudo-label c might not be closer to the any of the modes belonging to the class c . In addition, features in target domain mapping to same mode might not be closer to each other, resulting in low posterior probability. We minimize two loss functions to perform *stochastic mode alignment*.

For first, we apply *domain invariant consistency loss*, ensuring that features of pixels mapped to same modes of same class should be near to each other regardless of the domain they are sampled from. Assume a batch consisting of arbitrary number of source and target images, $\{(x_s^i, y_s^i) | i = 0, 1 \dots, N_s, (x_t^i, \hat{y}_t^i) | i = 0, 1 \dots, N_t\}$, where \hat{y}_t^i are the pseudo-labels assigned to x_t^i . Embedding $E_t^i = \mathcal{E}(x_t^i)$ and $E_s^i = \mathcal{E}(x_s^i)$ are computed for all the target and source images in the batch. We randomly sample N_e number of embedding, $\{e_i^j | i = 0, 1 \dots, N_e\}$ from $\{E_i^j | i = 0, 1 \dots, N_t\}$, choosing only from the ones having valid pseudo-label.

For *domain invariant consistency*, we create a triplet $(e_t^i, e_s^i, \hat{e}_s^i)$ such that pseudo-label of e_t^i and ground-truth label of e_s^i is same class c , and both map to same mode m of class c . \hat{e}_s^i on the other hand is source pixel's embedding of any class $c^+ \neq c$. This loss when minimized brings e_t^i closer to e_s^i than any source pixel's embedding belonging to different class.

$$\mathcal{L}_{mcl} = \sum_i^{N_e} \|\|e_t^i - e_s^i\|_2^2 - \|e_t^i - \hat{e}_s^i\|_2^2 + \alpha 1\|_+ \quad (6)$$

Note: we could have chosen most closest source sample as negative, however, this would have been computationally prohibitive. Margin, $\alpha 1$, is set to 1, for all experiments.

The in-mode variance for the target samples is decreased by forcing them to come closer to the assigned mode and move away from the modes of the other classes. We sample T_e embeddings per image per class from both source and target images and create set E_s and E_t respectively. Eq. 7 minimizes the triplet loss for both the source and target embeddings simultaneously.

$$\mathcal{L}_{ma}(E_s, E_t) = \frac{1}{T_s} \mathcal{L}_{emb}(E_s) + \frac{1}{T_t} \mathcal{L}_{emb}(E_t) \quad (7)$$

where T_s^e and T_t^e represent cardinality of E_s and E_t , which might be different since samples from all classes might not be available.

3.4. Total Loss for Training and Adaptation

The DRSL model is trained using the combination of segmentation losses, mode consistency loss and MMDL-FR module losses. Let \mathcal{L}_{cls}^s and \mathcal{L}_{cls}^t represent call to Eq. 5 using source and target embeddings respectively. The source model with MMDL module is trained using Eq.8.

$$\mathcal{L}_{src} = \mathcal{L}_{seg}^s + \beta \mathcal{L}_{emb} + \eta \mathcal{L}_{cls}^s \quad (8)$$

During adaptation to target domain the loss functions in Eq.9 and Eq.10 are used.

$$\mathcal{L}_{DRSL} = \mathcal{L}_G + \beta \mathcal{L}_{ma} + \eta(\mathcal{L}_{cls}^s + \mathcal{L}_{cls}^t) \quad (9)$$

$$\mathcal{L}_{DRSL+} = \mathcal{L}_G + \beta \mathcal{L}_{ma} + \eta(\mathcal{L}_{cls}^s + \mathcal{L}_{cls}^t) + \gamma \mathcal{L}_{mcl} \quad (10)$$

where, β , η and γ are hyper-parameters to limit the effect of MMDL-FR module loss values.

4. Experiments and Results

We performed multiple experiments for domain adaptation of semantic segmentation and compare the obtained results with state-of-the-art methods.

4.1. Experimental Setup

Datasets: Following [13, 1, 3], we use the standard benchmark setting of *synthetic-to-real* setup for our experiments. Specifically we setup for, *GTA-V to Cityscapes* and *SYNTHIA to Cityscapes* dataset, where the prior is source domain dataset and the later is the target domain dataset.

Cityscapes [36] dataset is a known benchmark for the task of semantic segmentation and domain adaptation. The dataset have 5000 high resolution labeled images partitioned as, training (2975), validation (500) and testing (1125). However, the annotations are only available for training and validation sets. **GTA-V** dataset [48] is obtained from the video game and the images are densely labeled with similar classes to cityscapes. There are 24966 images with spatial resolution spatial resolution 1052×1914. The GTA-V dataset also covers the road scene imagery. **SYNTHIA** [49] is another synthetic labeled images collection having 16 classes similar to Cityscapes. The dataset have 9400 images each with a spatial size 760×1280. Contrary to GTA-V and Cityscapes, SYNTHIA dataset has more view-point variations where; the camera is not supposed to be on the top of a vehicle every time.

Network Architecture: Following [34, 35], we use ResNet-101 [50] backbone based DeepLab-v2 [45] as our baseline segmentation model. Parallel to the segmentation head is the multi-modal distribution learning based feature regularization (MMDL-FR) module consisting of a combination of DML based Embedding Block (EB) and multi-modal distribution learning. We call the DeepLab-v2 last block as the encoder

Table 1 Semantic segmentation performance for GTA-V to Cityscapes adaptation. The abbreviations “ A_I ”, “ A_F ” and “ A_O ” stand for adversarial training at input space, latent space, and output space. Similarly, “ S_T ” represents self-supervised learning.

GTA-V \rightarrow Cityscapes																								
Methods	Baseline	Appr.	Road	Sidewalk	Building	Wall	Fence	Pole	T. Light	T. Sign	Veg.	Terrain	Sky	Person	Rider	Car	Truck	Bus	Train	Motorcycle	Bicycle	mIoU	mIoU Gain	
Source [45]	DeepLab-v2	-	75.8	16.8	77.2	12.5	21.0	25.5	30.1	20.1	81.3	24.6	70.3	53.8	26.4	49.9	17.2	25.9	6.5	25.3	36.0	36.6	-	
MinEnt [34]		$A_O + S_T$	86.6	25.6	80.8	28.9	25.3	26.5	33.7	25.5	83.3	30.9	76.8	56.8	27.9	84.3	33.6	41.1	1.2	23.9	36.4	43.6	7.0	
FCAN [46]		$A_I + A_O$	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	46.6	10.0
IntraDA [12]		$A_O + S_T$	90.6	37.1	82.6	30.1	19.1	29.5	32.4	20.6	85.7	40.5	79.7	58.7	31.1	86.3	31.5	48.3	0.0	30.2	35.8	46.3	9.7	
PyCDA [13]		S_T	90.5	36.3	84.4	32.4	28.7	34.6	36.4	31.5	86.8	37.9	78.5	62.3	21.5	85.6	27.9	34.8	18.0	22.9	49.3	47.4	10.8	
LSE [14]		S_T	90.2	40.0	83.5	31.9	26.4	32.6	38.7	37.5	81.0	34.2	84.6	61.6	33.4	82.5	32.8	45.9	6.7	29.1	30.6	47.5	10.9	
Source [47]	ResNet-38	-	70.0	23.7	67.8	15.4	18.1	40.2	41.9	25.3	78.8	11.7	31.4	62.9	29.8	60.1	21.5	26.8	7.7	28.1	12.0	35.4	-	
CBST [1]		S_T	86.8	46.7	76.9	26.3	24.8	42.0	46.0	38.6	80.7	15.7	48.0	57.3	27.9	78.2	24.5	49.6	17.7	25.5	45.1	45.2	9.8	
CRST [2]		S_T	84.5	47.7	74.1	27.9	22.1	43.8	46.5	37.8	83.7	22.7	56.1	56.8	26.8	81.7	22.5	46.2	27.5	32.3	47.9	46.8	11.4	
Source [45]	DeepLab-v2	-	71.7	18.5	67.9	17.4	10.2	36.5	27.6	6.3	78.4	21.8	67.6	58.3	20.7	59.2	16.4	12.5	7.9	21.2	13.0	33.8	-	
MRENT [2]		S_T	91.8	53.4	80.6	32.6	20.8	34.3	29.7	21.0	84.0	34.1	80.6	53.9	24.6	82.8	30.8	34.9	16.6	26.4	42.6	46.1	12.3	
Ours (DRSL)		$A_I + S_T$	92.8	57.5	82.8	28.7	17.7	40.6	34.3	27.0	85.5	42.7	77.8	62.3	30.8	82.2	24.3	38.5	8.4	31.1	39.6	47.6	13.8	
Ours (DRSL+)		$A_I + S_T$	<u>92.6</u>	<u>55.9</u>	82.4	29.0	24.6	42.7	38.3	35.7	85.5	39.5	77.0	64.2	26.2	83.9	19.5	31.6	9.3	27.1	42.5	47.8	14.0	

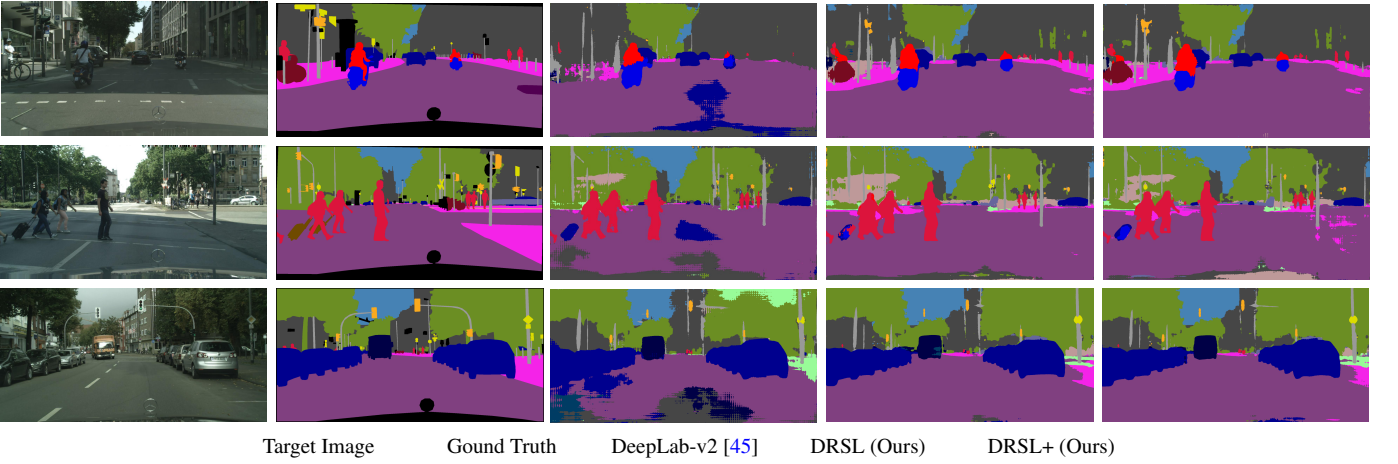


Figure 3: Semantic segmentation qualitative results for Cityscapes validation set when adapted from GTA-V dataset.

(base network) and the output feature-map as base features. For segmentation, these features are passed to segmentation layer while for MMDL-FR, these features are passed to the embedding block (Fig. 2). The embedding block consists of 4 fully convolutional layers with different dilation rates (similar to ones used in the segmentation layer of the segmentation network), producing an aggregated output. Unlike [19]’s fully-connected layers based DML for embedding generation, our strategy preserves the spatial structure necessary for segmentation and requires much less memory. The modes of the multi-modal are modeled with a fully connected layer as described in Sec. 3.2. and shown in Fig. 2. For each input the embedding block of the MMDL-FR module outputs an embedding volume E of size $(h \times w \times \hat{d})$. For an input image, we select a maximum of T_e embedding vectors per-class at random for further processing.

Implementation Details: To implement the proposed approach and conduct the experiments, we use PyTorch deep learning framework and a single GTX 1080ti GPU with a single Core-i5 machine with 32GB RAM. The ImageNet [51] trained weights for ResNet-101 [50] are used to train the DeepLab-v2 on source dataset. SGD optimizer with weight decay of 5×10^{-4} , momentum of 0.9, and initial learning rate of 2.5×10^{-4} for

source domain training and 5×10^{-5} during adaptation is used. In both source training and adaptation, we used a scale variance (0.5-1.5) and horizontal flipping randomly. For DML and mixture models based classification, the loss weights are set to $\beta = 0.25$ and $\eta = 0.1$ to limit the excessive gradient flow to segmentation model. Similarly for mixture models, the number of modes M is set to 3, and the number of embedding T_e per-class per-image is set to 300. For both source and target domain images, due to GPU memory limitations, small patches of size 512×512 cropped at random compared to original high-resolution images are processed.

The baseline segmentation model and the MMDL-FR module are initially trained with original source domain images, in-general called as source-only model. For self-supervised domain adaptation, selection of pixels as pseudo-labels is an important step as the adaptation process depends on the quality of pseudo-labels. We adapt an approach similar to [1], to generate pseudo-labels using the original source data trained model. For a given class c , we select δ confident pixels as pseudo-labels in the first round ($\delta = 20\%$) and increase this number of pixels ratio by 5% in each additional round. To further help the adaptation, we have obtained the translated version of the source domain datasets using CycleGan[4] and use these alongside orig-

inal source images during adaptation.

4.2. Experimental Results

In this section, we present experimental results of the proposed approach for semantic segmentation. We follow the standard synthetic to real adaptation setup.

4.2.1. Results on GTA-V to Cityscapes Adaptation

Table 1 presents domain adaptation performance for the task of semantic segmentation of the proposed DRSL approach compared to existing adversarial learning and self-supervised learning architectures. To have a fair comparison, the methods are divided into three groups where each comparing model is listed with its respective source model and backbone network. Fig. 3 shows example images to highlight the performance of the proposed DRSL qualitatively. The DRSL improves the performance for both objects and stuff classes, as shown in Fig. 3 (Column. 4). Small and far away objects like person, traffic light, and signboards are better adapted alongside near to camera objects and large area stuff classes like road, bus, and sidewalk. The cross domain mode alignment loss further penalizes the adaptation for small objects, further improving the performance for classes like bicycle, traffic sign, traffic light, pole, fence and person as shown in Table. 1 (DRSL+).

Overall, the proposed DRSL+ outperforms the latest self-supervised learning frameworks with clear gaps, surpassing the source dataset trained model with 14.0% gain in mIoU (last column of Table. 1). The DRSL+ performs well on both object classes as well as stuff classes compared to previous methods which may perform better on some classes but fail on other classes. Compared to CRST and MRENT [2] which regularizes the labels and models for high predictions, the proposed approach achieves a mIoU gain of 1.0 and 1.7% respectively. Similarly, the DRSL outperforms the PyCDA [13], which works on pyramid level labeling, and LSE [14] which incorporates scale invariances with class balancing strategies augmented with higher mIoU baseline models. Compared to composite adversarial learning-based methods like FCAN [46] and IntraDA [12], DRSL shows improvement with a minimum of 1% in mIoU and specifically with high margins in small objects. Similarly, compared to CAG-UDA[11] (mIoU=43.9% without warm-up training), the DRSL+ gains 3.9% in mIoU.

4.2.2. Results on SYNTHIA to Cityscapes Adaptation

Table 2 presents the proposed DRSL approach segmentation performance for SYNTHIA to Cityscapes adaptation. To have a fair comparison with existing methods, the comparing methods are divided into three groups and the respective source model results with different setups are shown. Moreover, for SYNTHIA to Cityscapes, we show the mIoU (16-classes) and mIoU* (13-classes) as shown by [3, 1]. Fig.4 shows qualitative results for DRSL and DRSL+ compared to baseline results. Row-1 and row-2 of Fig.4 focuses on objects like rider, bicycle, person, and the stuff classes, row-3 highlights the faraway objects and segmentation for road scene imagery.

The DRSL approach performs well on both stuff and object classes adaptation and shows an improvement of 11.7% in

mIoU and 12.9% in mIoU* compared to the baseline model (source). Compared to strong CBST[1] and MLST[3] self-supervised learning approaches, the DRSL shows a minimum improvement of 2.3% and 2.4% in mIoU and mIoU* respectively. Similarly, the DRSL shows significant improvement to existing regularization based models, like CRST [2] and entropy-based methods, e.g., LSE[14] and MinEnt [34]. Compared to CAG-UDA[11] (44.5% mIoU and 51.4% mIoU*), the DRSL+ gains 2.2% in mIoU and 1.9% in mIoU* respectively. The gaps can be more visible if compared with "without warm-up" training CAG-UDA.

4.2.3. Ablation Experiments

Ablation experiments are performed for GTA-V to Cityscapes.

Multi-Modal Distribution Learning based Regularization Module (MMDL-FR): During training and adaptation it's essential to understand the balance between the segmentation and different elements of MMDL-FR. We search over a range of values to identify (empirically) optimal values for the loss scaling factors, β and η (Table. 3). Based on the experiments, β and η are set to 0.25 and 0.1 respectively, for all the experiments including SYNTHIA to Cityscapes.

Effect of MMDL-FR Module on Adaptation Process: As described in Sec. 3.2 and Fig. 2, the MMDL-FR module regularizes the encoder (base-network) of the segmentation model with DML based embedding block and MMDL based classification. The MMDL-FR overall enhances the adaptation performance compared to the non-regularized version of the proposed method as shown in Table. 4.

Effect of Modes: As described in Sec. 3.2, it is very critical to select correct number of modes for multi-modal in MMDL. We have experimented with multiple number of modes (Table. 5) and selected $M=3$ for all the experiments.

Effect of Labels Reduction for MMDL-FR Module: The output of the embedding block in the MMDL-FR module is 8 times reduced compared to input image size. Embeddings needed to be upsampled 8 times if labels are not reduced requiring a lot of memory. Contrary to this, reducing labels 8 times introduces boxing effect. Based on these observations the scale factor 2 is used. A comparative performance of labels reduction is shown in Table. 6.

Table 6 Effect of label reduction ratio on mIoU.

GTA-V → Cityscapes				
Label Reduction Ratio	1	2	4	8
Embeddings Upsampling Ratio	8	4	2	1
Adaptation Performance (mIoU)	47.1	47.6	46.8	46.4

Pseudo-label Accuracy: To understand how the MMDL-FR results in more accurate pseudo-labels during the adaptation process, we compute mIoU of pseudo-labels for when MMDL-FR is not used (A) and when MMDL-FR is used (B). At the start of adaptation (round-0), we have same mIoU for both A & B (Table-7) since MMDL-FR will start to contribute when adaptation starts, i.e., during round-0. Due to MMDL-FR, the predictions by B after round-0 have much lower self-entropy

Table 2 Semantic segmentation performance of DRSL for SYNTHIA to Cityscapes adaptation. We present the mIoU (16-classes) and mIoU* (13-classes) comparison with existing state-of-the-art domain adaptation methods for the Cityscapes validation set.

SYNTHIA → Cityscapes																				
Methods	Baseline	Appr.	Road	Sidewalk	Building	Wall	Fence	Pole	T. Light	T. Sign	Veg.	Sky	Person	Rider	Car	Bus	Mcycle	Bicycle	mIoU	mIoU*
Source [45]	DeepLab-v2	-	64.3	21.3	73.1	2.4	1.1	31.4	7.0	27.7	63.1	67.6	42.2	19.9	73.1	15.3	10.5	38.9	34.9	40.3
CLAN [6]		A_O	81.3	37.0	80.1	-	-	-	16.1	13.7	78.2	81.5	53.4	21.2	73.0	32.9	22.6	30.7	-	47.8
Structure [8]		$A_F + A_O$	91.7	53.5	77.1	2.5	0.2	27.1	6.2	7.6	78.4	81.2	55.8	19.2	82.3	30.3	17.1	34.3	41.5	48.7
LSE [14]		S_T	<u>82.9</u>	<u>43.1</u>	78.1	9.3	0.6	28.2	9.1	14.4	77.0	83.5	58.1	25.9	71.9	38.0	29.4	31.2	42.6	49.4
CRST [2]		S_T	67.7	32.2	73.9	10.7	1.6	37.4	22.2	31.2	80.8	80.5	60.8	29.1	<u>82.8</u>	25.0	19.4	45.3	43.8	50.1
Source [47]	ResNet-38	-	32.6	21.5	46.5	4.81	0.03	26.5	14.8	13.1	70.8	60.3	56.6	3.5	74.1	20.4	8.9	13.1	29.2	33.6
CBST [1]		S_T	53.6	23.7	75.0	12.5	0.3	36.4	23.5	26.3	84.8	74.7	67.2	17.5	84.5	28.4	15.2	55.8	42.5	48.4
MLSL [3]		S_T	73.7	34.4	78.7	<u>13.7</u>	2.9	36.6	28.2	22.3	86.1	76.8	<u>65.3</u>	20.5	81.7	31.4	13.9	47.3	44.4	50.8
Source [45]	DeepLab-v2	-	69.2	26.6	66.5	6.5	0.1	33.2	4.1	18.0	80.5	80.0	55.3	15.1	67.5	20.1	6.8	14.0	35.2	40.3
DRSL		$A_I + S_T$	70.1	30.1	81.6	15.6	1.0	40.9	20.9	36.4	85.4	84.0	59.4	26.9	81.8	35.9	16.7	48.1	45.9	52.0
DRSL+		$A_I + S_T$	82.8	40.1	81.3	13.0	1.6	41.6	19.8	33.1	85.3	84.3	59.5	30.1	78.6	25.3	19.8	51.7	46.7	53.2

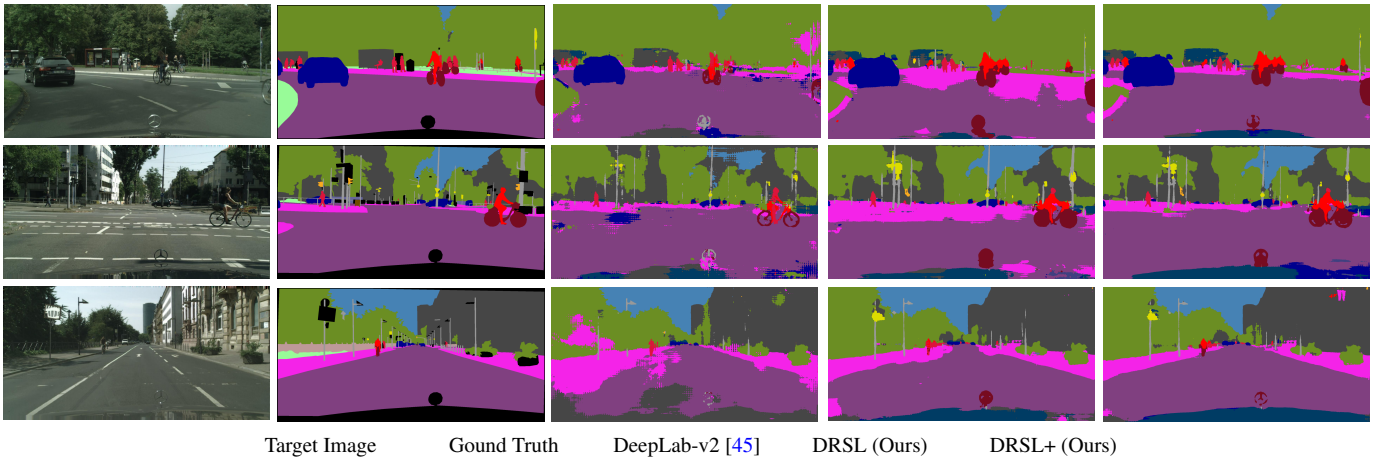


Figure 4: Semantic segmentation qualitative results for SYNTHIA to Cityscapes adaptation.

Table 3 Effect of (β, η) values of the MMDL-FR module.

β, η	(0.0, 0.0)	(0.1, 0.1)	(0.25, 0.1)	(0.5, 0.5)	(1.0, 1.0)
DRSL (mIoU)	44.9	46.1	47.6	45.9	46.0

Table 4 Effect of MMDL-FR module on adaptation.

Methods	Source [45]	Without MMDL-FR	With MMDL-FR
mIoU	33.6	44.9	47.6

Table 5 Effect to number of modes (M) in MMDL.

Number of Modes (M)	M=1	M=3	M=5
mIoU	44.7	47.6	46.2

and pseudo-labels have higher mIoU than the ones generated by model-A, thus improving self-supervised domain adaptation.

Table 7 Pseudo-labels with & without MMDL-FR module

Method	Start of Round-0		Start of Round-1	
	mIoU	Self-Entropy	mIoU	Self-Entropy
A: Without MMDL-FR {ST, ISA}	73.9	6.56×10^{-2}	76.4	1.57×10^{-2}
B: With MMDL-FR {ST, ISA, MMDL-FR}	73.9	6.56×10^{-2}	78.7	1.14×10^{-2}

Effect of Consistency Loss Weight: The cross domain mode consistency loss helps to make the embeddings of the source and target images belonging to the same mode of the same class closer, helping to better adapt the small object classes. However, its contribution in the whole loss needs to be limited to make the system stable. Our experiments suggests $\gamma = 0.1$ suits the DRSL+ as shown in Table. 8.

Table 8 Effect of cross domain mode consistency loss.

Loss weight γ	0.01	0.1	0.25
mIoU	46.0	47.8	45.3

Effect of Input Space Adaptation (ISA): Removing ISA module, mIoU decreases 1.6 points, from 47.6 (DRSL) to 46.0 (DRSL w/o ISA), indicating that ISA is needed but not vital for the effectiveness of the proposed model.

5. Conclusion

In this paper, we propose a distribution regularized self-supervised learning approach for domain adaptation of semantic segmentation. Parallel to the semantic segmentation decoding head, we employ a clustering based feature regulariza-

tion (MMDL-FR) module. Where segmentation head identifies what can differentiate a class, MMDL-FR explicitly models intra-class pixel-level feature variations, allowing the model to capture much richer representation of the class at pixel-level, thus improving model's generalization. Moreover, this disentanglement of information w.r.t tasks improves task dependent representation learning and allows performing separate domain alignments. Shared base-network enables MMDL-FR to act as regularizer over segmentation head, thus reducing the noisy pseudo-labels. Extensive experiments on the standard synthetic to real adaptation show that the proposed DRSL outperforms the state-of-the-art approaches.

References

- [1] 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 (ECCV), 2018, pp. 289–305. [1, 3, 5, 6, 7, 8](#)
- [2] Y. Zou, Z. Yu, X. Liu, B. Kumar, J. Wang, Confidence regularized self-training, in: Proceedings of the IEEE International Conference on Computer Vision, 2019, pp. 5982–5991. [1, 3, 6, 7, 8](#)
- [3] J. Iqbal, M. Ali, Msl: Multi-level self-supervised learning for domain adaptation with spatially independent and semantically consistent labeling, in: The IEEE Winter Conference on Applications of Computer Vision (WACV), 2020. [1, 3, 5, 7, 8](#)
- [4] J. Hoffman, E. Tzeng, T. Park, J. Zhu, P. Isola, K. Saenko, A. A. Efros, T. Darrell, Cycada: Cycle-consistent adversarial domain adaptation, in: Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholm, Sweden, July 10–15, 2018, 2018, pp. 1994–2003. [1, 3, 6](#)
- [5] Y.-H. Chen, W.-Y. Chen, Y.-T. Chen, B.-C. Tsai, Y.-C. F. Wang, M. Sun, No more discrimination: Cross city adaptation of road scene segmenters, in: Computer Vision (ICCV), 2017 IEEE International Conference on, IEEE, 2017, pp. 2011–2020. [1, 3](#)
- [6] Y. Luo, L. Zheng, T. Guan, J. Yu, Y. Yang, Taking a closer look at domain shift: Category-level adversaries for semantics consistent domain adaptation, in: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019. [1, 3, 8](#)
- [7] Y. Li, L. Yuan, N. Vasconcelos, Bidirectional learning for domain adaptation of semantic segmentation, in: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019. [1, 3](#)
- [8] W.-L. Chang, H.-P. Wang, W.-H. Peng, W.-C. Chiu, All about structure: Adapting structural information across domains for boosting semantic segmentation, in: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019. [1, 8](#)
- [9] T.-H. Vu, H. Jain, M. Bucher, M. Cord, P. Perez, Dada: Depth-aware domain adaptation in semantic segmentation, in: The IEEE International Conference on Computer Vision (ICCV), 2019. [1, 3](#)
- [10] J. Iqbal, R. Hafiz, M. Ali, Leveraging scale-invariance and uncertainty with self-supervised domain adaptation for semantic segmentation of foggy scenes, arXiv preprint arXiv:2201.02588 (2022). [1](#)
- [11] Q. Zhang, J. Zhang, W. Liu, D. Tao, Category anchor-guided unsupervised domain adaptation for semantic segmentation, in: Advances in Neural Information Processing Systems, 2019, pp. 433–443. [1, 2, 3, 7](#)
- [12] F. Pan, I. Shin, F. Rameau, S. Lee, I. S. Kweon, Unsupervised intra-domain adaptation for semantic segmentation through self-supervision, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020, pp. 3764–3773. [1, 3, 6, 7](#)
- [13] Q. Lian, F. Lv, L. Duan, B. Gong, Constructing self-motivated pyramid curriculums for cross-domain semantic segmentation: A non-adversarial approach, in: The IEEE International Conference on Computer Vision (ICCV), 2019. [1, 3, 5, 6, 7](#)
- [14] M. N. Subhani, M. Ali, Learning from scale-invariant examples for domain adaptation in semantic segmentation, in: Proceedings of the European Conference on Computer Vision (ECCV), 2020, pp. 1–17. [1, 3, 6, 7, 8](#)
- [15] M. A. Munir, M. H. Khan, M. Sarfraz, M. Ali, Ssal: Synergizing between self-training and adversarial learning for domain adaptive object detection, Advances in Neural Information Processing Systems 34 (2021). [1](#)
- [16] C. Chen, Z. Fu, Z. Chen, S. Jin, Z. Cheng, X. Jin, X.-S. Hua, Homm: Higher-order moment matching for unsupervised domain adaptation, in: Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 34, 2020, pp. 3422–3429. [1, 3](#)
- [17] A. Kumagai, T. Iwata, Unsupervised domain adaptation by matching distributions based on the maximum mean discrepancy via unilateral transformations, in: Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 33, 2019, pp. 4106–4113. [1, 3](#)
- [18] Z. Deng, Y. Luo, J. Zhu, Cluster alignment with a teacher for unsupervised domain adaptation, in: Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV), 2019. [2, 3](#)
- [19] L. Karlinsky, J. Shtok, S. Harary, E. Schwartz, A. Aides, R. Feris, R. Giryes, A. M. Bronstein, Repmet: Representative-based metric learning for classification and few-shot object detection, in: 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2019, pp. 5192–5201. [doi:10.1109/CVPR.2019.00534. 2, 4, 6](#)
- [20] E. Tzeng, J. Hoffman, K. Saenko, T. Darrell, Adversarial discriminative domain adaptation, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2017, pp. 7167–7176. [3](#)
- [21] P. O. Pinheiro, Unsupervised domain adaptation with similarity learning, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018, pp. 8004–8013. [3](#)
- [22] R. Xu, G. Li, J. Yang, L. Lin, Larger norm more transferable: An adaptive feature norm approach for unsupervised domain adaptation, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, 2019, pp. 1426–1435. [3](#)
- [23] Z. Deng, Y. Luo, J. Zhu, Cluster alignment with a teacher for unsupervised domain adaptation, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, 2019, pp. 9944–9953. [3](#)
- [24] A. Belal, M. Kiran, J. Dolz, L.-A. Blais-Morin, E. Granger, et al., Knowledge distillation methods for efficient unsupervised adaptation across multiple domains, Image and Vision Computing 108 (2021) 104096. [3](#)
- [25] S. Schrom, S. Hasler, J. Adamy, Improved multi-source domain adaptation by preservation of factors, Image and Vision Computing (2021) 104209. [3](#)
- [26] Y. Chen, W. Li, C. Sakaridis, D. Dai, L. Van Gool, Domain adaptive faster r-cnn for object detection in the wild, in: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2018. [3](#)
- [27] M. Khodabandeh, A. Vahdat, M. Ranjbar, W. G. Macready, A robust learning approach to domain adaptive object detection, in: Proceedings of the IEEE International Conference on Computer Vision, 2019, pp. 480–490. [3](#)
- [28] H.-K. Hsu, C.-H. Yao, Y.-H. Tsai, W.-C. Hung, H.-Y. Tseng, M. Singh, M.-H. Yang, Progressive domain adaptation for object detection, in: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, 2020, pp. 749–757. [3](#)
- [29] Y. Chen, W. Li, L. Van Gool, Road: Reality oriented adaptation for semantic segmentation of urban scenes, in: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2018. [3](#)
- [30] Y. Zhang, P. David, B. Gong, Curriculum domain adaptation for semantic segmentation of urban scenes, in: The IEEE International Conference on Computer Vision (ICCV), 2017. [3](#)
- [31] D. Dai, C. Sakaridis, S. Hecker, L. Van Gool, Curriculum model adaptation with synthetic and real data for semantic foggy scene understanding, International Journal of Computer Vision (2019) 1–23. [3](#)
- [32] J. Iqbal, M. Ali, Weakly-supervised domain adaptation for built-up region segmentation in aerial and satellite imagery, ISPRS Journal of Photogrammetry and Remote Sensing 167 (2020) 263–275. [3](#)
- [33] J. Yang, C. Li, W. An, H. Ma, Y. Guo, Y. Rong, P. Zhao, J. Huang, Exploring robustness of unsupervised domain adaptation in semantic segmentation, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, 2021, pp. 9194–9203. [3](#)
- [34] T.-H. Vu, H. Jain, M. Bucher, M. Cord, P. Pérez, Advent: Adversarial entropy minimization for domain adaptation in semantic segmentation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2019, pp. 2517–2526. [3, 5, 6, 7](#)
- [35] Y.-H. Tsai, W.-C. Hung, S. Schuster, K. Sohn, M.-H. Yang, M. Chandraker, Learning to adapt structured output space for semantic segmenta-

- tion, in: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2018. [3](#), [5](#)
- [36] 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: Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016. [3](#), [5](#)
- [37] X. Guo, C. Yang, B. Li, Y. Yuan, Metacorection: Domain-aware meta loss correction for unsupervised domain adaptation in semantic segmentation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2021, pp. 3927–3936. [3](#)
- [38] Z. Wang, M. Yu, Y. Wei, R. Feris, J. Xiong, W.-m. Hwu, T. S. Huang, H. 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, 2020, pp. 12635–12644. [3](#)
- [39] M. Mancini, L. Porzi, S. Rota Bulò, B. Caputo, E. Ricci, Boosting domain adaptation by discovering latent domains, in: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2018. [3](#)
- [40] Y. Zhang, Z. Qiu, T. Yao, D. Liu, T. Mei, Fully convolutional adaptation networks for semantic segmentation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018, pp. 6810–6818. [3](#)
- [41] M. Kim, H. Byun, Learning texture invariant representation for domain adaptation of semantic segmentation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2020, pp. 12975–12984. [3](#)
- [42] B. Zhang, S. Zhao, R. Zhang, Towards adaptive semantic segmentation by progressive feature refinement, in: 2020 IEEE International Conference on Image Processing (ICIP), IEEE, 2020, pp. 2221–2225. [3](#)
- [43] J. Zhang, C. Liang, C.-C. J. Kuo, A fully convolutional tri-branch network (fctn) for domain adaptation, in: 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, 2018, pp. 3001–3005. [3](#)
- [44] Y.-H. Tsai, K. Sohn, S. Schuler, M. Chandraker, Domain adaptation for structured output via discriminative patch representations, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, 2019, pp. 1456–1465. [3](#)
- [45] 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 (4) (2018) 834–848. [3](#), [5](#), [6](#), [8](#)
- [46] Y. Zhang, Z. Qiu, T. Yao, D. Liu, T. Mei, Fully convolutional adaptation networks for semantic segmentation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2018, pp. 6810–6818. [6](#), [7](#)
- [47] Z. Wu, C. Shen, A. Van Den Hengel, Wider or deeper: Revisiting the resnet model for visual recognition, Pattern Recognition 90 (2019) 119–133. [6](#), [8](#)
- [48] S. R. Richter, V. Vineet, S. Roth, V. Koltun, Playing for data: Ground truth from computer games, in: B. Leibe, J. Matas, N. Sebe, M. Welling (Eds.), European Conference on Computer Vision (ECCV), Vol. 9906 of LNCS, Springer International Publishing, 2016, pp. 102–118. [5](#)
- [49] G. Ros, L. Sellart, J. Materzynska, D. Vazquez, A. M. Lopez, The synthia dataset: A large collection of synthetic images for semantic segmentation of urban scenes, in: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016. [5](#)
- [50] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 770–778. [5](#), [6](#)
- [51] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al., Imagenet large scale visual recognition challenge, International journal of computer vision 115 (3) (2015) 211–252. [6](#)