

000 001 Adaptive Feature Swapping for Domain 002 Invariant Learning 003

004 Anonymous ECCV submission
005

006 Paper ID 273
007

009 **Abstract.** The bottleneck of visual domain adaptation always lies in the
010 learning of domain invariant representations. In this paper, we present a
011 simple but effective technique named Adaptive Feature Swapping (AFS)
012 for learning domain invariant features in Unsupervised Domain Adap-
013 tation (UDA). AFS aims to select semantically irrelevant features from
014 labeled source data and unlabeled target data, and swap these features
015 with each other. Then the merged representations are also utilized for
016 training with prediction consistency constraint. In this way, the model
017 is encouraged to learn representations that are robust to domain-specific
018 information. We develop two swapping strategies including channel swap-
019 ping and spatial swapping. The former encourages the model to squeeze
020 redundancy out of features and pay more attention to semantic infor-
021 mation. The latter motivates the model to be robust to the background
022 and focus on objects. We conduct experiments on object recognition and
023 semantic segmentation in UDA setting and the results show that AFS
024 can promote various existing UDA methods.
025

026 **Keywords:** Domain Adaptation; Transfer Learning; Domain Adaptive
027 Semantic Segmentation
028

029 1 Introduction 030

031 Domain adaptation improves a target task with insufficient annotations by knowl-
032 edge transfer from a source domain with rich annotations. To achieve reliable
033 transfer, the domain discrepancy between the source and target domains should
034 be mitigated in domain adaptation. There are two main directions for mitigating
035 domain discrepancy including moment alignment [10,27,29,41,57,16] and adver-
036 sarial training [6,28,40,61]. The former aims to minimize domain discrepancy
037 estimated by maximum mean discrepancy (MMD) [27], correlation distance [41]
038 or other distance metric [16,29] calculated on task-specific features. The latter
039 direction of methods learn features that are indistinguishable to an additional
040 domain discriminator [6,28] or utilizes adversarial learning between two classi-
041 fiers for aligning conditional feature distributions [40,60,61,18].
042

043 Recently, domain-specific modeling and perturbation-tolerant learning in do-
044 main adaptation have drawn considerable attention. Modeling the domain-specific
045 features enhances the learning of domain invariant features [1,7,4]. For exam-
046 ple, Bousmalis et al. [1] proposed to learn domain-specific features via time-
047 consuming reconstruction at the pixel level and Cui et al. [4] adopted a single
048

045 fully connected layer whose responses are restricted to vanish. By introducing
046 auxiliary tasks like reconstruction or vanishing bridge responses, these meth-
047 ods learn the domain-specific features in implicit ways resulting in suboptimal
048 solutions with additional computation overhead.

049 Perturbation-tolerant learning imposes robustness over input perturbations.
050 Such robustness can be measured by the norm of the input-output Jacobian of
051 the network, and it correlates well with generalization [34]. To enhance the gen-
052 eralization in the target domain, several methods [5,39,52] imposed consistency
053 between the predictions of two perturbated samples by performing random data
054 augmentations over a target sample. With expansion assumption, Wei et al. [50]
055 proved that enforcing locally consistent prediction provides accuracy guarantees
056 on unlabeled target data for unsupervised domain adaptation. Typical consis-
057 tency constraints include L2 distance minimization [5], agreement maximiza-
058 tion [39] and mutual information maximization [52]. These methods usually
059 require two complete time-consuming forward computations of deep network.
060 Besides, the underlying information of samples from the two domains for gener-
061 ating perturbation is not well explored.

062 To handle these issues, we present Adaptive Feature Swapping to model the
063 domain-specific features and perform sample perturbation more efficiently and
064 effectively. The key insight is to swap semantically irrelevant features between
065 source and target data, obtaining perturbed representations that maintain the
066 semantic information for object recognition. By encouraging the consistency be-
067 tween predictions of the representations after swapping and the original ones,
068 the model attains robustness that benefits better generalization and more re-
069 liable knowledge transfer. It supports accuracy guarantees on unlabeled target
070 data as locally consistent prediction is ensured. Besides, enforcing the model to
071 predict the ground truth label over the source representations after swapping
072 encourages the model to pay more attention to the domain invariant features
073 against the domain-specific ones.

074 We develop two kinds of Adaptive Feature Swapping techniques including
075 Spatial Feature Swapping and Channel Feature Swapping. For Spatial Feature
076 Swapping, since the irrelevant information like the background surrounding an
077 object should not affect the prediction of the object, we attempt to swap the
078 irrelevant features of source samples with those of target samples. To enable
079 reliable Spatial Feature Swapping, we extract the spatial attention to locate the
080 irrelevant spatial features. We adopt activation-based attention [17] due to its
081 computation efficiency, and the features with less attention are regarded as ir-
082 relevant ones. Initially, the model may mistakenly put more attention to some
083 semantically irrelevant features. However, as the model is optimized during train-
084 ing, the attention is promoted to better locate the irrelevant features, which
085 makes Spatial Feature Swapping more effective.

086 Sharing the same spirit of Spatial Feature Swapping, Channel Feature Swap-
087 ping is conducted to squeeze the redundancy out of features, so the model tends
088 to pay more attention to the semantic information. Similarly, the accumulation
089 of feature responses of samples in a mini-batch is regarded as channel-wise at-

tention, serving to select irrelevant features for swapping. Since the two feature swapping strategies encourage the model to capture semantic information from different perspectives, they could cooperate well with each other which is verified in our experiments. Moreover, the representation with feature swapping randomly integrate feature from samples of other domain, exhibiting more diverse perturbation.

We evaluate our Adaptive Feature Swapping on several benchmarks for object classification and semantic segmentation in UDA setting. Experimental results show that the proposed Adaptive Feature Swapping improves source-only model by a considerable margin. Besides, as a light-weight technique, Adaptive Feature Swapping can be easily plugged into various domain adaptation methods and boost their performances.

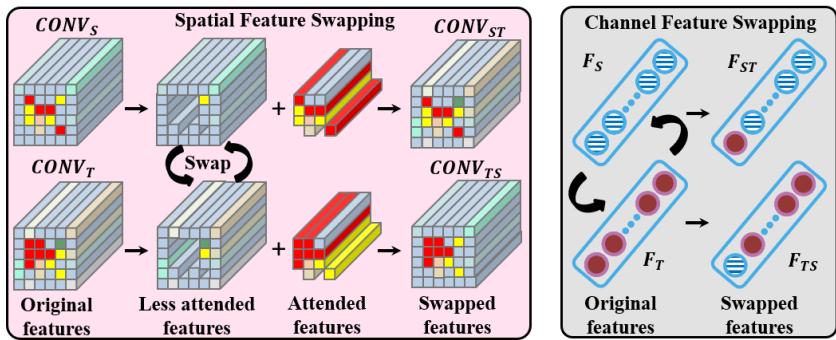
2 Related Work

The unsupervised domain adaptation methods can be mainly divided into moment alignment based methods [27,29,41,57,16] and adversarial training based methods [6,28,40,61]. Moment alignment based methods try to estimate the domain discrepancy related to distribution moments over deep representations. [45] firstly adopted the MMD [9] in the deep neural network and later DAN [27] utilized the multiple kernel variant of MMD over multiple layer representations for better moment matching. JAN [29] considered the joint distribution of features and logits, which made distribution alignment more effective. CAN [16] proposed to optimize a new metric which explicitly models the intra-class domain discrepancy and the inter-class domain discrepancy. Enhanced Transport Distance (ETD) [21] builds an attention-aware transport distance, which can be viewed as the prediction feedback of the iteratively learned classifier, to measure the domain discrepancy.

Early Adversarial training based methods introduced a domain discriminator [6,28] that is trained to distinguish source and target representations while the backbone network is trained to generate representations indistinguishable to the domain discriminator. DANN [6] performed domain adversarial learning on features while [28] performed domain adversarial learning on the multilinear mapping of feature and classifier response. Recent methods [40,60,61,18] utilized adversarial learning between two classifiers for aligning conditional feature distributions. MCD [40] alternatively maximized and minimized the L1 distance of output probabilities of two symmetric classifiers to better utilize the task-specific decision boundaries. In [61], the authors used two asymmetrical classifiers to estimate the conditional feature distributions with marginal loss. In [51], the authors propose to make the domain alignment proactively serve classification via feature decomposition and alignment with the prior knowledge induced from the classification task.

Other methods focus on some characteristics of specific layers in a deep neural network for domain adaptation such as adaptively increasing the feature norms in AFN [54], maintaining both discriminability and diversity in BNM [3] and

135 minimizing inter-class confusion of unlabeled target data in MCC [15]. Besides,
 136
 137



149 **Fig. 1.** Illustration of Spatial Feature Swapping and Channel Feature Swapping. For
 150 the convolutional activations of a source (target) sample, the semantic features (in red
 151 and yellow) are kept and combined with the rest irrelevant features from convolutional
 152 activations of a target (source) sample, obtaining new convolutional activations. The
 153 Channel Feature Swapping exchanges some irrelevant elements between source and
 154 target features.

155
 156 Bousmalis et al. [1] developed a split model that models the shared domain
 157 invariant features, and domain-specific features supporting effective reconstruc-
 158 tion. Gong et al. [7] proposed a domain flow generation model to generate a
 159 continuous sequence of intermediate domains with various domain-specific infor-
 160 mation. Both methods require time-consuming reconstruction at the pixel level.
 161 Recently, Cui et al. [4] proposed to use a single fully connected layer, named grad-
 162 ually vanishing bridge (GVB) to capture domain-specific features. Being simple
 163 and effective, additional parameters and careful design of gradually vanishing
 164 loss are still required in GVB. In contrast, our method learns domain-specific
 165 features via simple feature swapping with prediction consistency constraint.

166 For the perturbation-tolerant learning, two random data augmentations are
 167 performed over a target sample and the consistency between the predictions
 168 of the two augmented samples is imposed in [5]. In [39], the authors proposed
 169 to minimize Min-Entropy Consensus loss that univocally selects a pseudo-label
 170 that maximizes the agreement between two perturbed versions of the same tar-
 171 get sample. In [52], the authors explicitly maximized the mutual information
 172 between the rotated image and the label. Our method differs from these meth-
 173 ods by implementing the perturbation via feature swapping. The other related
 174 work is [49] that exchanged channel activations between the two convolutional
 175 layer activations extracted from two modalities of a sample for better fusing
 176 the two modalities. [49] selected feature according to the magnitude of Batch-
 177 Normalization scaling factor. On the contrary, our method swaps both **channel**
 178 and **spatial** features from two **independent samples** according to carefully
 179 designed **aggregated attention**.

180 3 Methodology

181
 182 We illustrate our method in object recognition task under unsupervised domain
 183 adaptation (UDA) setting for simplicity. Suppose that there are a source domain
 184 $D_s = \{(x_i^s, y_i^s)\}_{i=1}^{n_s}$ of n_s labeled samples and a target domain $D_t = \{(x_i^t)\}_{i=1}^{n_t}$
 185 of n_t unlabeled samples. The two domains distribute differently that the input
 186 marginal distributions $p(x)$ vary. In UDA, both the two domains share the same
 187 K categories. The goal is to learn a domain invariant model that generalizes well
 188 in target domain.

189 We propose Adaptive Feature Swapping for learning domain invariant rep-
 190 resentations. As shown in Fig. 1, Adaptive Feature Swapping involves spatial
 191 swapping and **channel** swapping. Spatial swapping aims to select irrelevant fea-
 192 tures like background for source and target samples and swap these features
 193 with each other. Similarly, channel swapping selects irrelevant features along the
 194 channel and swaps these features between source and target samples. The pre-
 195 diction over the original feature and the one after swapping is enforced to be
 196 consistent. Both feature swapping strategies encourage the model to pay more
 197 attention to the semantic information that is invariant across domains.

198 3.1 Spatial Feature Swapping

199 We perform Spatial Feature Swapping on a specific convolutional layer ac-
 200 tivations. Suppose that for a batch of source samples and a batch of target
 201 samples both with size B , their convolutional layer activations are $CONVs \in \mathcal{R}^{B \times C \times W \times H}$ and $CONv_T \in \mathcal{R}^{B \times C \times W \times H}$. To perform feature swapping, we
 202 need to determine which spatial feature is semantically irrelevant. Inspired by
 203 activation-based attention [17], we first map the batch of convolutional layer ac-
 204 tivations into an attention map and then select features with less attention for
 205 swapping. Concretely, the spatial attention $M_{sp}(A)$ for an activation tensor A is
 206 obtained via

$$207 \quad (M_{sp}(A))_{i,j} = \sum_{b=1}^B \sum_{c=1}^C |A_{b,c,i,j}|^p, \quad (1)$$

208 where $i \in \{1, 2, \dots, H\}$, $j \in \{1, 2, \dots, W\}$ are spatial indexes and p is set to 1. b
 209 and c are the instance index and the channel index, respectively.

210 To prevent a semantically relevant feature in $CONVs$ to be swapped to
 211 $CONv_T$ that affects the semantic of swapped $CONv_{TS}$, we select the feature
 212 with less aggregated attention response ($M_{sp}(CONVs) + M_{sp}(CONv_T)$). The
 213 features corresponding to the least r_1 percentage of aggregated attention are
 214 selected for swapping.

215 3.2 Channel Feature Swapping

216 In CNN-based classifier, the feature before the final classifier layer is usually
 217 of high dimension that involves redundancy. The network may tend to remem-
 218 ber irrelevant features for recognition, especially for a small scale labeled source
 219 domain. This issue motivates us to adopt Channel Feature Swapping for encour-
 220 aging the model to ignore semantically irrelevant information.

Suppose that $F_S \in \mathcal{R}^{B \times L}$ and $F_T \in \mathcal{R}^{B \times L}$ are the features before the final classifier for a batch of source samples and a batch of target samples both with size B. We construct channel attention to determine which channel feature is semantically irrelevant. Channel attention can be modeled in several ways like global average pooling (GAP), global max/min pooling, or learned via FC layers [14]. Among these solutions, GAP is simple yet effective and widely used in NIN [25], GoogLeNet [42], and ResNet [11]. Therefore, we use GAP to extract channel attention for selecting redundant features. Specifically, given an activation tensor $V \in \mathcal{R}^{B \times L}$, the channel attention can be computed as

$$(M_{ch}(V))_c = \sum_i^B |V_{i,c}|, \quad (2)$$

where $i \in \{1, 2, \dots, B\}$ and c is the channel index.

Similar to Spatial Feature Swapping, we need to select channels that contribute little to predicting objects. Specifically, we construct the aggregated attention as $(M_{ch}(F_S) - M_{ch}(F_T))$ assuming that the channel that the model pays nearly equal attention for both domains delivers less semantic information. Then the aggregated attention is sorted and the features corresponding to the least r_2 percentage of aggregated attention are selected for swapping.

3.3 Learning with Feature Swapping

In unsupervised domain adaptation, a model is trained by minimizing a classification loss on the labeled source domain and an additional transfer loss, to learn both discriminative and domain invariant representations. The classification loss in \mathcal{D}_S can be computed as:

$$\mathcal{L}_{cls} = \frac{1}{N_s} \sum_{i=1}^{N_s} L_{ce}(p_i^s, y_i^s), \quad (3)$$

where L_{ce} is the cross-entropy loss function and p_i^s, y_i^s are classification response and label for the source sample x_i^s .

The additional transfer loss can be norm increasing loss [54], MMD loss [27], MCC loss [15] or adversarial discriminator score [6].

For the model learning with Adaptive Feature Swapping, the transfer loss can be formulated as:

$$\mathcal{L}_{cst} = \frac{1}{N_s} \sum_{i=1}^{N_s} L_{ce}(p_i^{st}, y_i^s) + \lambda_{cst} \frac{1}{N_t} \sum_{j=1}^{N_t} \|p_j^t - p_j^{ts}\|_2^2, \quad (4)$$

where p_i^{st} is the classification response over representation of source sample x_i^s that is swapped with some features from target sample. And p_j^{ts} is the classification response over the representation of target sample x_j^t that is swapped with some features from a source sample. p_j^t is the classification response of the original target sample x_j^t .

The transfer loss for Adaptive Feature Swapping encourages consistency between model predictions of the original representation and the representation with feature swapping for a sample. For source samples, the consistency is achieved by cross-entropy loss with corresponding labels. As for unlabeled target samples, the consistency is simply achieved by mean square loss.

The total loss of model learning with Adaptive Feature Swapping is

$$\mathcal{L} = \mathcal{L}_{cls} + \mathcal{L}_{cst}. \quad (5)$$

It is worth noting that the proposed Adaptive Feature Swapping can cooperate with existing methods. Suppose that the loss of existing methods is $\mathcal{L}_{cls} + \lambda_{tran}\mathcal{L}_{tran}$, one can easily add \mathcal{L}_{cst} to the original loss obtaining

$$\mathcal{L} = \mathcal{L}_{cls} + \lambda_{tran}\mathcal{L}_{tran} + \mathcal{L}_{cst}. \quad (6)$$

3.4 Theoretical Insight

We provide the connections between our method and the theory proposed in [50]. Let F denote a learned scoring functions (a neural network), G be the label predictor that $G(x) = \arg \max_i F(x)_i$, and G^* be the ground-truth classifier. Define \mathcal{B} the set of perturbated versions x' for a sample x where the perturbation is the proposed feature swapping. We can see that $\|x' - x\| \leq r_* \|x\|$, where r_* is swapping ratio. The neighborhood of x can be defined as $\mathcal{N} = \{x' : \mathcal{B}(x) \cap \mathcal{B}(x') \neq \emptyset\}$. The consistency loss can be formulated as:

$$R_{\mathcal{B}}(G) = \mathbb{E}_P[1(\exists x' \in \mathcal{B}(x) \text{ such that } G(x') \neq G(x))] \quad (7)$$

which is strongly reminiscent of our consistency loss \mathcal{L}_{cst} .

In accordance with [50], we assume P , a distribution of unlabeled examples over input space \mathcal{X} , is \mathcal{B} -separated with probability $1 - \mu$ by G^* that $R_{\mathcal{B}}(G^*) \leq \mu$. Besides, let $\mathcal{M}(G_{pl}) = \{x : G_{pl}(x) \neq G^*(x)\}$ be the set of examples that mistakenly labeled by a pseudolabeler G_{pl} . The expansion assumption stated in [50] requires that $P_i(\mathcal{N}(V)) \leq \min(cP_i(V), 1)$ for all $V \subseteq X$ and all class i with $P_i(V) \leq a$, where P_i is the class-conditional distribution, $a = \max_i\{P_i(\mathcal{M}(G_{pl}))\}$ and a is required to be less than $1/3$ and $c > 3$.

Then it is proved in [50] that for any minimizer \widehat{G} of $\min_G \frac{c+1}{c-1} L_{0-1}(G, G_{pl}) + \frac{2c}{c-1} R_{\mathcal{B}}(G) - Err(G_{pl})$, we have

$$Err(\widehat{G}) \leq \frac{2}{c-1} Err(G_{pl}) + \frac{2c}{c-1} \mu \quad (8)$$

where $L_{0-1}(G, G') = \mathbb{E}[1(G(x) \neq G'(x))]$, and $Err(G) = L_{0-1}(G, G^*)$ is the error rate of G .

Therefore, our method could achieve the above bounded target error (Eqn.(8)) for unsupervised domain adaptation under suitable assumptions. Specifically, during the training process, the current model has relatively low $Err(G_{pl})$, which in turn lower the error rate $Err(\widehat{G})$ of the updated model. With stronger base model like MCC [15] and BNM [3], the $Err(G_{pl})$ is lower. Therefore, the trained model gets lower error rate, as verified experimentally in Section 4.

315 4 Experiments

316 4.1 Datasets and Setup

317 We use two standard object classification benchmarks and a standard semantic
 318 segmentation benchmark for unsupervised domain adaptation to evaluate the
 319 proposed methods. Office-Home [46] is a challenging benchmark with 65 cat-
 320 egories in four domains: *Art* (Ar), *Clipart* (Cl), *Product* (Pr) and *Real-world*
 321 (Rw), which contains 12 transfer tasks. VisDA-C [36] is a large-scale benchmark
 322 that contains the images from 12 categories of two very distinct domains, the
 323 synthetic domain, and the real-world domain. The synthetic domain contains
 324 152,397 images and the latter contains 55,388 images. Following standard pro-
 325 tocol, we focus on the synthetic-to-real transfer task.

326 For the segmentation task, we evaluate our method by adapting the segmen-
 327 tation from game scenes, GTA5 [38] dataset, to real scene, the Cityscapes [2]
 328 dataset. GTA5 contains 24,966 images with the resolution of 1914×1052 . The
 329 Cityscapes dataset contains 2,975 training images and 500 images for validation
 330 with the resolution of 2048×1024 .

331 Following standard protocol, we use the training set of GTA5 as the la-
 332 beled source domain and the unlabeled training images from Cityscapes as the
 333 target domain. We conduct evaluations on the validation set of Cityscapes and
 334 adopt the Intersection-over-Union (IoU) of each class and the mean-Intersection-
 335 over-Union (mIoU) as performance metrics. For the task GTA5 \rightarrow Cityscapes, we
 336 report the results on the common 19 classes.

337 4.2 Implementation Details

338 **Classification.** The methods are implemented based on PyTorch [35], and for
 339 fair comparisons, we use the ResNet-50/ResNet-101 [11] as the backbone of net-
 340 work, fixing the batch-size of 36 for all experiments. We use the same settings
 341 as [24]. We adopt mini-batch SGD optimizer with momentum 0.9 to train the
 342 model. For Office-Home, we use ResNet-50 as the backbone, and the hyper-
 343 parameters r_1, r_2, λ_{cst} are set to 1/10, 1/15 and 40, respectively. The performances
 344 with respect to various hyperparameters can be referred to the Supplementary.
 345 For VisDA-C we adopt ResNet-101 as backbone for fair comparisons. The hyper-
 346 parameters r_1, r_2 and λ_{cst} are set to 1/4, 1/3 and 40, respectively. The Channel
 347 Feature Swapping is performed on the feature after mean pooling over the last
 348 residual block of ResNet [11]. The spatial Feature swapping is performed on the
 349 activations of the third residual block. We evaluate our approach by applying
 350 Adaptive Feature Swapping (AFS) to the source-only model and several main-
 351 stream DA methods including CDAN [28], BNM [3], MCC [15], and CST [26]
 352 based on their open-source codes. λ_{tran} for BNM, CDAN and MCC are kept the
 353 same as their original implementation. The reimplemented MCC is marked with
 354 \dagger as MCC \dagger . Codes can be found in the Supplementary.

355 **Segmentation.** We evaluate our approach by applying Channel Feature Swap-
 356 ping to the source-only model and the state-of-the-art domain adaptive segmen-
 357 tation method ProDA [58] and DAFormer [13] based on their open-source codes.

Table 1. Accuracy (%) On VisDA-C for UDA using the ResNet-101 backbone. The best accuracy is indicated in bold and the second best one is underlined.

Method	plane	bicycl	bus	car	horse	knife	mcycle	persn	plant	sktb	train	trunk	Avg
DANN [6]	81.9	77.7	82.8	44.3	81.2	29.5	65.1	28.6	51.9	54.6	82.8	7.8	57.4
DAN [27]	87.1	63.0	76.5	42.0	90.3	42.9	85.9	53.1	49.7	36.3	85.8	20.7	61.1
MCD [40]	87.0	60.9	83.7	64.0	88.9	79.6	84.7	76.9	88.6	40.3	83.0	25.8	71.9
AFN [54]	93.6	61.3	84.1	70.6	94.1	79.0	91.8	79.6	89.9	55.6	89.0	24.4	76.1
DRMEA [32]	92.1	75.0	78.9	75.5	91.2	81.9	89.0	77.2	93.3	77.4	84.8	35.1	79.3
ATDOC [24]	93.7	83.0	76.9	58.7	89.7	95.1	84.4	71.4	89.4	80.0	86.7	55.1	80.3
DTA [19]	93.7	82.2	85.6	83.8	93.0	81.0	90.7	82.1	95.1	78.1	86.4	32.1	81.5
STAR [30]	95.0	84.0	<u>84.6</u>	73.0	91.6	91.8	85.9	78.4	94.4	84.7	87.0	42.2	82.7
SHOT [23]	94.3	<u>88.5</u>	80.1	57.3	93.1	94.9	80.7	80.3	91.5	89.1	86.3	<u>58.2</u>	82.9
RWOT [53]	95.1	80.3	83.7	90.0	92.4	68.0	92.5	82.2	87.9	78.4	90.4	68.2	84.0
ResNet [11]	67.7	27.4	50.0	61.7	69.5	13.7	85.9	11.5	64.4	34.4	84.2	19.2	49.1
+AFS	95.1	59.3	83.9	72.9	<u>95.3</u>	66.2	91.9	63.4	93.6	62.3	84.1	17.0	73.8
CDAN [28]	85.2	66.9	83.0	50.8	84.2	74.9	88.1	74.5	83.4	76.0	81.9	38.0	73.9
+AFS	93.0	73.5	80.1	69.5	91.8	<u>97.0</u>	82.9	79.0	83.3	72.4	84.5	39.7	78.1
BNM [3]	95.6	86.7	77.4	48.9	92.0	83.3	76.2	76.3	88.0	64.2	84.7	46.3	76.6
+AFS	<u>95.9</u>	87.0	78.1	53.6	93.7	90.5	78.8	78.8	87.4	84.9	83.4	47.0	79.9
MCC [15]†	91.5	82.6	75.9	61.5	91.4	92.9	81.6	79.6	85.8	88.4	69.5	51.8	80.2
+AFS	93.1	78.5	81.1	79.0	93.9	94.2	88.1	<u>85.1</u>	91.0	85.6	85.2	41.6	83.1
CST+SAM [26]†	95.2	89.0	75.8	89.6	95.7	98.5	87.8	<u>84.1</u>	94.4	<u>90.9</u>	82.4	54.5	<u>86.5</u>
+AFS	97.2	89.0	83.8	<u>87.5</u>	95.7	94.0	87.9	82.8	96.6	94.7	82.4	53.1	87.1

For the segmentation task, we also use PyTorch [35] to implement our methods. Our training is carried out on 2 TITAN RTX GPUs. To train our models using SGD, the initial learning rate is 1e-4, and adjusted according to the ‘poly’ learning rate scheduler with a power of 0.9. The weight decay is set to 0.0005. These settings are the same as ProDA [58] for fair comparisons. In accordance with [13], we train our model with AdamW, a base learning rate of 6e-5 for the encoder and 6e-4 for the decoder, a weight decay of 0.01, linear learning rate warmup with 1.5k iteration, and linear decay afterwards. The data augmentation and other hyperparameters are kept the same as [13]. The hyperparameters r_2 , λ_{cst} are set to 1/15 and 15, respectively.

Note that ProDA [58] adopts pseudo label for training target samples, so the consistency loss on target domain is directly implemented as the cross entropy loss of prediction over representations with feature swapping and its corresponding pseudo labels. Note that only Channel Feature Swapping is utilized since swapping spatial features changes their semantic in segmentation task. Besides, due to the limited computation resources, we only implemented ProDA for its first stage with a batch of 2 images (marked with †).

4.3 Classification Results

Results on VisDA-C are reported in Tab.1. We can observe that our method consistently improves the generalization ability of all baseline methods. Specifically, our model achieves the highest accuracy 83.1% on average, 2.9% higher

Table 2. Accuracy (%) On Office-Home for UDA using the ResNet-50 backbone. The best accuracy is indicated in bold and the second best one is underlined.

Method	Ar	→Cl	Ar	→Pr	Ar	→Rw	Cl	→Ar	Cl	→Pr	Cl	→Rw	Pr	Rw	→Ar	Rw	→Cl	Rw	→Pr	Avg
DAN [27]	43.6	57.0	67.9	45.8	56.5	60.4	44.0	43.6	67.7	63.1	51.5	74.3	56.3							
DANN [6]	45.6	59.3	70.1	47.0	58.5	60.9	46.1	43.7	68.5	63.2	51.8	76.8	57.6							
JAN [29]	45.9	61.2	68.9	50.4	59.7	61.0	45.8	43.4	70.3	63.9	52.4	76.8	58.3							
EntMin [8]	51.0	71.9	77.1	61.2	69.1	70.1	59.3	48.7	77.0	70.4	53.0	81.0	65.8							
AFN [54]	52.0	71.7	76.3	64.2	69.9	71.9	63.7	51.4	77.1	70.9	57.1	81.5	67.3							
SymNet[60]	47.7	72.9	78.5	64.2	71.3	74.2	64.2	48.8	79.5	74.5	52.6	82.7	67.6							
MDD [61]	54.9	73.7	77.8	60.0	71.4	71.8	61.2	53.6	78.1	72.5	60.2	82.3	68.1							
ECT [21]	51.3	71.9	85.7	57.6	69.2	73.7	57.8	51.2	79.3	70.2	57.5	82.1	67.3							
GVB [4]	57.0	74.7	79.8	64.6	74.1	74.6	65.2	55.1	81.0	74.6	59.7	84.3	70.4							
TCM [56]	58.6	74.4	79.6	64.5	74.0	75.1	64.6	56.2	80.9	74.6	60.7	84.7	70.7							
SHOT [23]	57.1	78.1	81.5	68.0	78.2	78.1	67.4	54.9	82.2	73.3	58.8	84.3	71.8							
AAA [20]	56.7	78.3	82.1	66.4	78.5	79.4	67.6	53.5	81.6	74.5	58.4	84.1	71.8							
SENTRY [37]	61.8	77.4	80.1	66.3	71.6	74.7	66.8	63.0	80.9	74.0	66.3	84.1	72.2							
ResNet [11]	34.9	50.0	58.0	37.4	41.9	46.2	38.5	31.2	60.4	53.9	41.2	59.9	46.1							
+AFS	52.5	71.4	76.7	58.0	68.6	69.6	57.5	48.1	76.1	70.2	53.9	80.7	65.3							
CDAN [28]	50.7	70.6	76.0	57.6	70.0	70.0	57.4	50.9	77.3	70.9	56.7	81.6	65.8							
+AFS	56.3	73.8	78.4	62.7	72.4	71.4	64.3	54.4	80.0	72.9	60.2	83.0	69.2							
BNM [3]	52.3	73.9	80.0	63.3	72.9	74.9	61.7	49.5	79.7	70.5	53.6	82.2	67.9							
+AFS	56.8	75.5	79.6	63.4	76.0	74.4	63.0	53.3	80.3	73.0	59.8	83.4	69.9							
MCC [15]†	55.8	77.4	80.8	66.3	75.7	75.8	65.4	53.4	80.7	74.0	59.0	84.4	70.7							
+AFS	59.0	77.5	81.4	67.2	76.8	76.0	65.8	57.3	81.9	74.0	61.8	85.0	72.0							
CST [26]†	59.0	79.6	83.4	68.4	77.1	76.7	68.9	56.4	83.0	75.3	62.2	85.1	73.0							
+AFS	58.7	80.2	<u>83.3</u>	67.6	79.0	76.7	68.9	57.1	<u>82.6</u>	<u>75.1</u>	65.5	85.7	73.4							

than MCC. Compared with source only model, integrating our AFS brings 24.7% accuracy improvement on average. The results show that our methods can directly improve the performance without any other tricks. As for CST+SCM [26], our simple model surpasses it by 0.6%. Note that input-consistency via random augmentation is involved in CST+SCM, leaving little room for our method for improvement. However, we can observe that our method could still boost CST+SCM, showing its promising potential to boost existing unsupervised domain adaptation methods. Based on these experiments, we can infer that AFS can stably enhance the transfer ability of classifier.

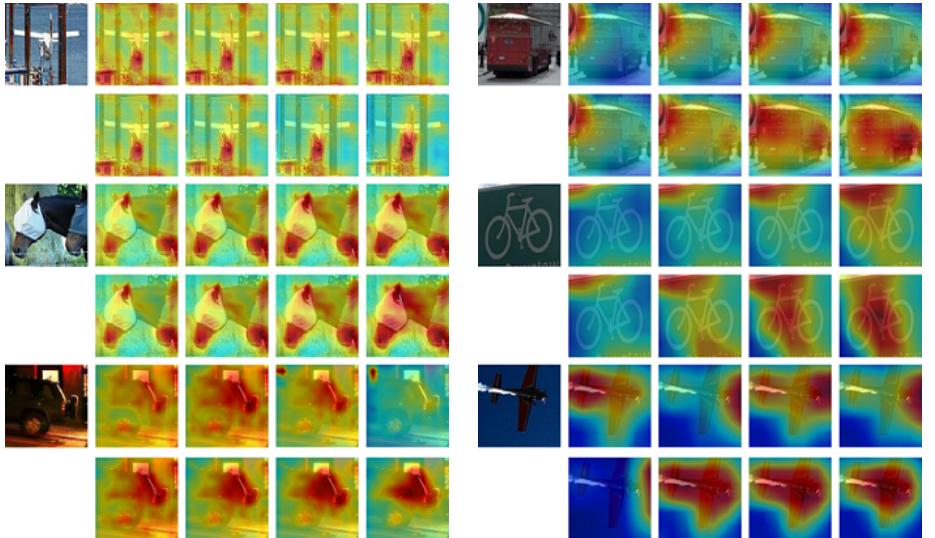
Table 3. Ablation Study On VisDA-C for UDA using the ResNet-101 backbone. The best accuracy is indicated in bold and the second best one is underlined.

Method	plane	bicyl	bus	car	horse	knife	mcycle	persn	plant	sktb	train	trunk	Avg
MCC [15]†	91.5	82.6	75.9	61.5	91.4	92.9	81.6	79.6	85.8	88.4	69.5	51.8	80.2
MCC+SS	92.7	<u>83.0</u>	78.2	69.2	91.7	94.3	85.1	80.7	88.8	86.7	<u>81.7</u>	46.3	81.5
MCC+CS	93.5	85.4	<u>78.9</u>	<u>70.6</u>	93.4	93.0	<u>86.0</u>	<u>81.6</u>	92.5	85.8	80.7	46.0	<u>82.3</u>
MCC+SS+CS	<u>93.1</u>	78.5	81.1	79.0	93.9	<u>94.2</u>	88.1	85.1	<u>91.0</u>	85.6	85.2	41.6	83.1

Results on Office-Home are summarized in Tab. 2. Office-Home is a challenging dataset due to its large domain discrepancy. Utilizing the proposed AFS into source only model brings 19.2% accuracy improvement on average. Compared with CST [26], our simple model surpasses it by 0.4% on average. Our AFS also boosts the CDAN, BNM and MCC baselines by considerable margins. Based on

450 these promising results, we can infer that AFS can stably explore truly useful
 451 semantic information to better adapt the classifier.

452 **Ablation Study.** In Tab. 3, we validate the influence of the proposed spatial
 453 swapping (SS) and channel swapping (CS) in VisDA-C dataset. All the results
 454 reported in Tab. 3 show that our methods can stably improve the transferability
 455 of the model.



474 **Fig. 2.** The 1st column are images from
 475 VisDA-C dataset. The 2nd to fifth columns
 476 show the attention maps extracted from
 477 the activations of **third residual block**
 478 of model trained with 1K, 3K, 5K, and
 479 10K iterations, respectively. The model
 480 that produces attention maps of the 1st
 481 (2nd), 3rd (4th), and 5th (6th) rows is
 482 MCC† baseline (MCC+AFS). The model
 483 pays more attention to the regions in red
 484 than regions in blue.

474 **Fig. 3.** The 1st column are images from
 475 VisDA-C dataset. The 2nd to fifth columns
 476 show the attention maps extracted from
 477 the activations of the **last residual block**
 478 of model trained with 1K, 3K, 5K, and
 479 10K iterations, respectively. The model
 480 that produces attention maps of the 1st
 481 (2nd), 3rd (4th), and 5th (6th) rows is
 482 MCC† baseline (MCC+AFS). The model
 483 pays more attention to the regions in red
 484 than regions in blue.

485 Firstly, we validate the effectiveness of using spatial swapping. Compared
 486 with baseline MCC, integrating spatial swapping (MCC+SS) shows a gain of
 487 1.3% accuracy on average. Based on it, we can infer that spatial swapping helps
 488 the model to pay more attention to objects instead of the background. Secondly,
 489 the large improvement is achieved by adopting channel swapping into baseline
 490 MCC (MCC+CS) that encourages the model to focus on semantic information.
 491 Finally, we also utilize both two swapping strategies into baseline MCC and
 492 it performs better than MCC+SS and MCC+CS. Therefore, the two feature
 493 swapping strategies can benefit each other and help the model achieve better
 494 performance.

Attention visualization. We visualize the attention maps of three images from VisDA-C dataset for baseline MCC and MCC+AFS with Adaptive Feature Swapping in Fig. 2 and Fig. 3.

Specifically, the image is input to the trained network, and the activations of a convolutional layer of ResNet-101 is extracted, named $A \in \mathcal{R}^{1024 \times W \times H}$. We sum over the first dimension of A and obtain $M \in \mathcal{R}^{W \times H}$. Then we interpolate M into 224×224 to get the final attention map. For Fig. 2, attention map is extracted on the activations of the third residual block, where the spatial feature swapping is performed. As for Fig. 3, attention map is extracted from the activations of the last residual block. These activations are global average pooled and the pooled results are performed with channel feature swapping. Therefore, the attention map of this block could reflect where the model attends when trained with channel feature swapping.

We can observe that, as the model is optimized with increasing iterations, MCC+AFS pays more attention to the semantic region of Plane or Horse. And the background regions attain less attention in MCC+AFS. While in baseline MCC, some background regions are still highlighted with considerable attention. The superiority of MCC+AFS over MCC verifies that learning with Adaptive Feature Swapping encourages the model to capture semantic information. More results can be referred to the Supplementary.

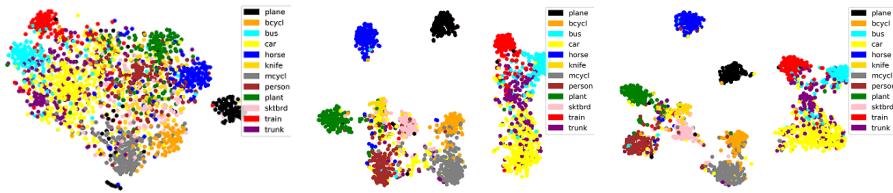


Fig. 4. The 1st, 2nd and 3rd figures correspond to the t-SNE embedding visualization of Source-only model, MCC, and our models MCC+AFS on VisDA-C. Different colors indicate different categories.

Representation visualization. To better illustrate that our method can improve the discriminability of target representations, we visualize the representations generated by source-only model, MCC†, and MCC+AFS on VisDA-C. We employ the t-SNE method [33] and map the representations into 2D points. We randomly select 2000 samples across 12 categories from real-world domain in VisDA-C. As shown in Fig. 4, compared with the MCC†, our method better separates the target samples and the representations within the classes are well clustered and more compact. And the mis-classification of features in each cluster is reduced, especially for “train”, “bus” and “car” categories.

Hyper-parameter Sensitivity Hyper-parameter r_1 and r_2 control the perturbation strength of the representations with feature swapping and λ_{cst} balances the consistency loss with other losses. To evaluate the parameter sensitivity of AFS, we conduct experiments with MCC+AFS on VisDA-C. To reduce the complexity of evaluation, we set $r = r_1 = r_2$ and vary $r \in \{1/2, 1/3, 1/4, 1/5, 1/10\}$

, $1/15, 1/20\}$ and $\lambda_{cst} \in \{10, 20, 30, 40, 50\}$. Fig. 5 shows that, except for combinations of large $r = \{1/2, 1/3\}$ and large $\lambda_{cst} = \{40, 50\}$, MCC+AFS is not that sensitive to r_1 , r_2 and λ_{cst} , and can achieve competitive results under a wide range of hyper-parameter values.

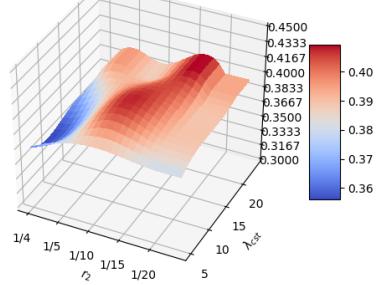
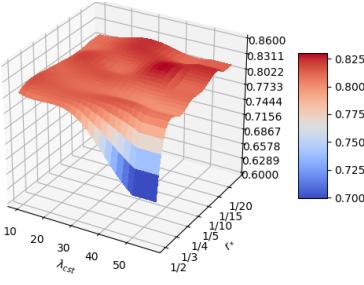


Fig. 5. Hyper-parameter sensitivity analysis of AFS on VisDA-C based on MCC [15]. **Fig. 6.** mIoUs of various Source only+AFS on GTA2Cityscape.

4.4 Segmentation Results

Table 4. Comparison results of GTA5 \rightarrow Cityscapes adaptation in terms of mIoU. The best accuracy is indicated in bold and the second best one is underlined.

Method	road	sidewalk	building	wall	fence	pole	light	sign	vegetation	terrain	sky	person	rider	car	truck	bus	train	motorcycle	bicycle	mIoU	
AdaStruct [43]	86.5	25.9	79.8	22.1	20.0	23.6	33.1	21.8	81.8	25.9	75.9	57.3	26.2	76.3	29.8	32.1	7.2	29.5	32.5	41.4	
CyCADA [12]	86.7	35.6	80.1	19.8	17.5	38.0	39.9	41.5	82.7	27.9	73.6	64.9	19.0	65.0	12.0	28.6	4.5	31.1	42.0	42.7	
CLAN [31]	87.0	27.1	79.6	27.3	23.3	28.5	35.5	24.2	83.6	27.4	74.2	58.6	28.0	76.2	33.1	36.7	6.7	31.9	31.4	43.2	
APODA [55]	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 [44]	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 [47]	89.4	33.1	81.0	26.6	26.8	27.2	33.5	23.5	84.7	38.9	36.7	78.8	58.7	30.5	84.8	38.5	44.5	1.7	31.6	32.4	45.5
BDL [22]	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	
FADA [48]	91.0	50.6	86.0	43.4	29.8	36.8	43.4	25.0	86.8	38.3	87.4	64.0	38.0	85.2	31.6	46.1	6.5	25.4	37.1	50.1	
CBST [63]	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 [64]	91.0	55.4	80.0	33.7	21.4	37.3	32.9	24.5	85.0	34.1	80.8	57.7	24.6	84.1	27.8	30.1	26.9	26.0	42.3	47.1	
CAG-UDA [59]	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	
Seg-U [62]	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	
Source-only	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	
+AFS	76.4	21.1	74.3	19.6	23.6	33.3	36.2	26.3	82.5	22.9	72.7	61.2	32.2	73.5	32.9	28.0	0.8	28.8	27.2	40.7	
ProDA [58]†	87.6	55.3	77.7	41.9	35.3	41.6	43.9	48.2	86.6	40.7	81.0	65.6	22.9	86.8	38.0	49.8	0.0	41.9	50.7	52.4	
+AFS	90.5	55.8	80.8	41.8	36.0	41.9	45.2	48.6	86.7	41.1	81.0	66.2	22.8	87.5	39.9	49.6	0.0	42.6	51.6	53.1	
DAFormer [13]	<u>95.7</u>	<u>70.2</u>	<u>89.4</u>	<u>53.5</u>	<u>48.1</u>	<u>49.6</u>	<u>55.8</u>	<u>59.4</u>	<u>89.9</u>	<u>47.9</u>	<u>92.5</u>	<u>72.2</u>	<u>44.7</u>	<u>92.3</u>	<u>74.5</u>	<u>78.2</u>	<u>65.1</u>	<u>55.9</u>	<u>61.8</u>	<u>68.3</u>	
+AFS	96.1	71.8	89.4	53.8	<u>44.5</u>	49.9	56.6	62.5	90.0	51.0	<u>90.9</u>	<u>71.7</u>	<u>45.0</u>	<u>92.3</u>	<u>72.3</u>	<u>81.9</u>	<u>72.0</u>	56.6	63.6	69.0	

We report the quantitative evaluation of our method in Tab. 4. From the results, we can see that with Channel Feature Swapping, the source-only model is improved by 4.1% in mIoU. Besides, we can observe that applying our Channel Feature Swapping into the ProDA [58]† with ResNet-101 improves the baseline ProDA† by 0.7% in mIoU. Channel Feature Swapping also boosts DAFormer [13] by 0.7% in mIoU. For the challenging domain adaptive segmentation task, the 0.7% improvement in mIoU is not marginal, which validates the effectiveness of Channel Feature Swapping in enhancing the transfer ability of the model.

We visually compare the baseline DAFormer [13] and DAFormer+AFS in Fig. 7. As we can see from Fig. 7, the predictions from our model appear less

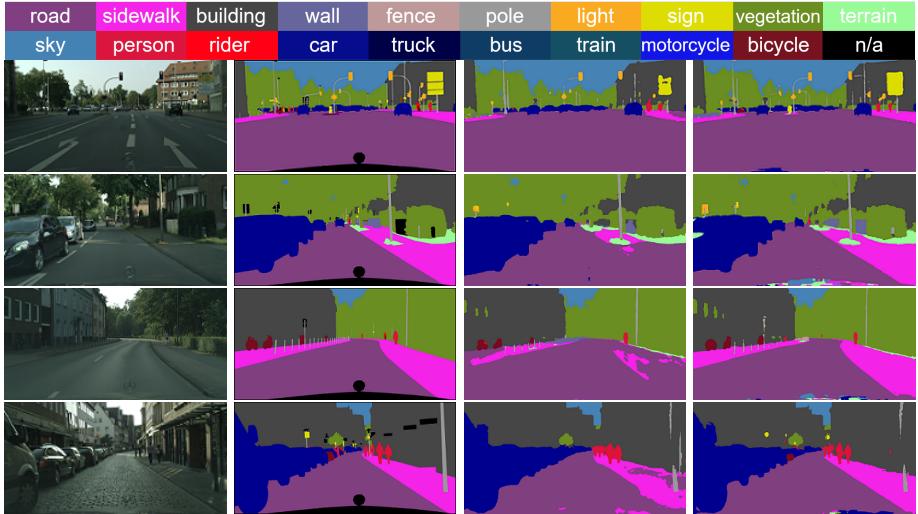


Fig. 7. Qualitative results of validation images from Cityscapes where the source domain is GTA5 dataset. The 1st and 2nd columns are images and ground truth labels from validation set of Cityscape. The 3rd and 4th columns are segmentation results produced by DAFormer [13]† and our method.

noisy, like the “sign” in the first row and the sidewalk for all rows. We credit this to the enhanced generalization ability of utilizing Channel Feature Swapping. More qualitative comparisons can be found in the Supplementary.

Hyper-parameter Sensitivity We conduct experiments on semantic segmentation task w.r.t to r_2 and λ_{cst} , since Spatial Feature Swapping is not adopted for it destroys the semantic for segmentation task. We vary $r_2 \in \{1/4, 1/5, 1/10, 1/15, 1/20\}$ and $\lambda_{cst} \in \{5, 10, 15, 20\}$, and the results are shown in Fig. 6. We can observe that Source only+AFS is not that sensitive to r_2 and λ_{cst} , and can achieve competitive results under a wide range of hyper-parameter values. Intuitively, the redundancy of features for segmentation tasks is less since it requires much more information for dense prediction. The swapping ratio should be smaller compared with classification.

5 Conclusion

Aiming to model domain-specific features and achieve prediction robustness toward sample perturbation in an efficient manner, we propose Adaptive Feature Swapping for learning domain invariant representations in Unsupervised Domain Adaptation (UDA). Specifically, we utilize attention mechanisms to select semantically irrelevant features from labeled source data and unlabeled target data, and swap these features from each other. Then consistency between predictions of original representations and the representations after swapping is enforced such that the model is insensitive to the irrelevant features. We develop two swapping strategies including channel swapping and spatial swapping and extensive experiments on object recognition and semantic segmentation in UDA setting validate the effectiveness of feature swapping.

630 References

- 631 1. Bousmalis, K., Trigeorgis, G., Silberman, N., Krishnan, D., Erhan, D.: Domain
632 separation networks. In: Advances in Neural Information Processing Systems. pp.
633 343–351 (2016)
- 634 2. Cordts, M., Omran, M., Ramos, S., Rehfeld, T., Enzweiler, M., Benenson, R.,
635 Franke, U., Roth, S., Schiele, B.: The cityscapes dataset for semantic urban scene
636 understanding. In: Proceedings of the IEEE conference on computer vision and
637 pattern recognition. pp. 3213–3223 (2016)
- 638 3. Cui, S., Wang, S., Zhuo, J., Li, L., Huang, Q., Tian, Q.: Towards discriminability
639 and diversity: Batch nuclear-norm maximization under label insufficient situations.
640 In: CVPR. pp. 3941–3950 (2020)
- 641 4. Cui, S., Wang, S., Zhuo, J., Su, C., Huang, Q., Tian, Q.: Gradually vanishing bridge
642 for adversarial domain adaptation. In: Proceedings of the IEEE/CVF Conference
643 on Computer Vision and Pattern Recognition. pp. 12455–12464 (2020)
- 644 5. French, G., Mackiewicz, M., Fisher, M.: Self-ensembling for visual domain adap-
645 tation. In: International Conference on Learning Representations (2018)
- 646 6. Ganin, Y., Lempitsky, V.: Unsupervised domain adaptation by backpropagation.
647 In: International conference on machine learning. pp. 1180–1189 (2015)
- 648 7. Gong, R., Li, W., Chen, Y., Gool, L.V.: Dlow: Domain flow for adaptation and
649 generalization. In: Proceedings of the IEEE/CVF Conference on Computer Vision
650 and Pattern Recognition. pp. 2477–2486 (2019)
- 651 8. Grandvalet, Y., Bengio, Y., et al.: Semi-supervised learning by entropy minimiza-
652 tion. CAP **367**, 281–296 (2005)
- 653 9. Gretton, A., Borgwardt, K., Rasch, M., Schölkopf, B., Smola, A.: A kernel method
654 for the two-sample-problem. Advances in neural information processing systems
655 **19**, 513–520 (2006)
- 656 10. Gretton, A., Borgwardt, K.M., Rasch, M.J., Schölkopf, B., Smola, A.: A kernel
657 two-sample test. Journal of Machine Learning Research **13**(Mar), 723–773 (2012)
- 658 11. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recogni-
659 tion. In: Proceedings of the IEEE conference on computer vision and pattern recogni-
660 tion. pp. 770–778 (2016)
- 661 12. Hoffman, J., Tzeng, E., Park, T., Zhu, J.Y., Isola, P., Saenko, K., Efros, A., Dar-
662 rell, T.: Cycada: Cycle-consistent adversarial domain adaptation. In: Interna-
663 tional conference on machine learning. pp. 1989–1998. PMLR (2018)
- 664 13. Hoyer, L., Dai, D., Van Gool, L.: Daformer: Improving network architectures
665 and training strategies for domain-adaptive semantic segmentation. arXiv preprint
666 arXiv:2111.14887 (2021)
- 667 14. Hu, J., Shen, L., Sun, G.: Squeeze-and-excitation networks. In: Proceedings of the
668 IEEE conference on computer vision and pattern recognition. pp. 7132–7141 (2018)
- 669 15. Jin, Y., Wang, X., Long, M., Wang, J.: Minimum class confusion for versatile
670 domain adaptation. In: ECCV. pp. 464–480 (2020)
- 671 16. Kang, G., Jiang, L., Yang, Y., Hauptmann, A.G.: Contrastive adaptation network
672 for unsupervised domain adaptation. In: Proceedings of the IEEE/CVF Conference
673 on Computer Vision and Pattern Recognition. pp. 4893–4902 (2019)
- 674 17. Komodakis, N., Zagoruyko, S.: Paying more attention to attention: improving the
675 performance of convolutional neural networks via attention transfer. In: ICLR
676 (2017)
- 677 18. Lee, C.Y., Batra, T., Baig, M.H., Ulbricht, D.: Sliced wasserstein discrepancy for
678 unsupervised domain adaptation. In: Proceedings of the IEEE Conference on Com-
679 puter Vision and Pattern Recognition. pp. 10285–10295 (2019)

- 675 19. Lee, S., Kim, D., Kim, N., Jeong, S.G.: Drop to adapt: Learning discriminative
676 features for unsupervised domain adaptation. In: Proceedings of the IEEE/CVF
677 International Conference on Computer Vision. pp. 91–100 (2019) 678
- 678 20. Li, J., Du, Z., Zhu, L., Ding, Z., Lu, K., Shen, H.T.: Divergence-agnostic unsupervised
679 domain adaptation by adversarial attacks. *IEEE Transactions on Pattern
680 Analysis and Machine Intelligence* (2021) 681
- 681 21. Li, M., Zhai, Y.M., Luo, Y.W., Ge, P.F., Ren, C.X.: Enhanced transport distance
682 for unsupervised domain adaptation. In: Proceedings of the IEEE/CVF Conference
683 on Computer Vision and Pattern Recognition. pp. 13936–13944 (2020) 684
- 683 22. Li, Y., Yuan, L., Vasconcelos, N.: Bidirectional learning for domain adaptation of
684 semantic segmentation. In: Proceedings of the IEEE/CVF Conference on Computer
685 Vision and Pattern Recognition. pp. 6936–6945 (2019) 686
- 686 23. Liang, J., Hu, D., Feng, J.: Do we really need to access the source data? source hy-
687 pothesis transfer for unsupervised domain adaptation. In: International Conference
688 on Machine Learning. pp. 6028–6039 (2020) 689
- 689 24. Liang, J., Hu, D., Feng, J.: Domain adaptation with auxiliary target domain-
690 oriented classifier. In: Proceedings of the IEEE/CVF Conference on Computer
691 Vision and Pattern Recognition. pp. 16632–16642 (2021) 692
- 692 25. Lin, M., Chen, Q., Yan, S.: Network in network. In: 2nd International Conference
693 on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14–16, 2014,
694 Conference Track Proceedings (2014) 695
- 694 26. Liu, H., Wang, J., Long, M.: Cycle self-training for domain adaptation. *Advances
695 in Neural Information Processing Systems* **34** (2021) 696
- 696 27. Long, M., Cao, Y., Wang, J., Jordan, M.: Learning transferable features with deep
697 adaptation networks. In: International conference on machine learning. pp. 97–105
698 (2015) 699
- 699 28. Long, M., Cao, Z., Wang, J., Jordan, M.I.: Conditional adversarial domain adap-
700 tation. In: Advances in Neural Information Processing Systems 31: Annual Con-
701 ference on Neural Information Processing Systems 2018, NeurIPS 2018, December
702 3–8, 2018, Montréal, Canada. pp. 1647–1657 (2018) 703
- 703 29. Long, M., Zhu, H., Wang, J., Jordan, M.I.: Deep transfer learning with joint adap-
704 tation networks. In: International conference on machine learning. pp. 2208–2217
705 (2017) 706
- 705 30. Lu, Z., Yang, Y., Zhu, X., Liu, C., Song, Y.Z., Xiang, T.: Stochastic classifiers for
706 unsupervised domain adaptation. In: Proceedings of the IEEE/CVF Conference
707 on Computer Vision and Pattern Recognition. pp. 9111–9120 (2020) 708
- 708 31. Luo, Y., Zheng, L., Guan, T., Yu, J., Yang, Y.: Taking a closer look at domain
709 shift: Category-level adversaries for semantics consistent domain adaptation. In:
710 Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recog-
711 nition. pp. 2507–2516 (2019) 712
- 712 32. Luo, Y.W., Ren, C.X., Ge, P., Huang, K.K., Yu, Y.F.: Unsupervised domain adap-
713 tation via discriminative manifold embedding and alignment. In: Proceedings of
714 the AAAI Conference on Artificial Intelligence. vol. 34, pp. 5029–5036 (2020) 715
- 715 33. Van der Maaten, L., Hinton, G.: Visualizing data using t-sne. *Journal of machine
716 learning research* **9**(11) (2008) 717
- 717 34. Novak, R., Bahri, Y., Abolafia, D.A., Pennington, J., Sohl-Dickstein, J.: Sensi-
718 tivity and generalization in neural networks: an empirical study. In: International
719 Conference on Learning Representations (2018) 720
- 719 35. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen,
720 T., Lin, Z., Gimelshein, N., Antiga, L., et al.: Pytorch: An imperative style, high-
721 performance deep learning library. *Neural computation* **32**, 1–53 (2020) 722

- 720 performance deep learning library. *Advances in neural information processing systems* **32**, 8026–8037 (2019)
- 721
- 722 36. Peng, X., Usman, B., Kaushik, N., Hoffman, J., Wang, D., Saenko, K.: Visda: The
723 visual domain adaptation challenge. arXiv preprint arXiv:1710.06924 (2017)
- 724 37. Prabhu, V., Khare, S., Kartik, D., Hoffman, J.: Sentry: Selective entropy optimization
725 via committee consistency for unsupervised domain adaptation. In: Proceedings of the
726 IEEE/CVF International Conference on Computer Vision. pp. 8558–
727 8567 (2021)
- 728 38. Richter, S.R., Vineet, V., Roth, S., Koltun, V.: Playing for data: Ground truth
729 from computer games. In: European conference on computer vision. pp. 102–118.
730 Springer (2016)
- 731 39. Roy, S., Siarohin, A., Sangineto, E., Bulo, S.R., Sebe, N., Ricci, E.: Unsupervised
732 domain adaptation using feature-whitening and consensus loss. In: Proceedings of the
733 IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp.
734 9471–9480 (2019)
- 735 40. Saito, K., Watanabe, K., Ushiku, Y., Harada, T.: Maximum classifier discrepancy
736 for unsupervised domain adaptation. In: Proceedings of the IEEE conference on
737 computer vision and pattern recognition. pp. 3723–3732 (2018)
- 738 41. Sun, B., Saenko, K.: Deep coral: Correlation alignment for deep domain adaptation.
739 In: European Conference on Computer Vision. pp. 443–450. Springer (2016)
- 740 42. Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D.,
741 Vanhoucke, V., Rabinovich, A.: Going deeper with convolutions. In: Proceedings
742 of the IEEE conference on computer vision and pattern recognition. pp. 1–9 (2015)
- 743 43. Tsai, Y.H., Hung, W.C., Schulter, S., Sohn, K., Yang, M.H., Chandraker, M.:
744 Learning to adapt structured output space for semantic segmentation. In: Proceedings
745 of the IEEE conference on computer vision and pattern recognition. pp. 7472–7481 (2018)
- 746 44. Tsai, Y.H., Sohn, K., Schulter, S., Chandraker, M.: Domain adaptation for struc-
747 tured output via discriminative patch representations. In: Proceedings of the
748 IEEE/CVF International Conference on Computer Vision. pp. 1456–1465 (2019)
- 749 45. Tzeng, E., Hoffman, J., Zhang, N., Saenko, K., Darrell, T.: Deep domain confusion:
750 Maximizing for domain invariance. arXiv preprint arXiv:1412.3474 (2014)
- 751 46. Venkateswara, H., Eusebio, J., Chakraborty, S., Panchanathan, S.: Deep hashing
752 network for unsupervised domain adaptation. In: Proceedings of the IEEE confer-
753 ence on computer vision and pattern recognition. pp. 5018–5027 (2017)
- 754 47. Vu, T.H., Jain, H., Bucher, M., Cord, M., Pérez, P.: Advent: Adversarial entropy
755 minimization for domain adaptation in semantic segmentation. In: Proceedings
756 of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp.
757 2517–2526 (2019)
- 758 48. Wang, H., Shen, T., Zhang, W., Duan, L.Y., Mei, T.: Classes matter: A fine-
759 grained adversarial approach to cross-domain semantic segmentation. In: European
760 Conference on Computer Vision. pp. 642–659. Springer (2020)
- 761 49. Wang, Y., Huang, W., Sun, F., Xu, T., Rong, Y., Huang, J.: Deep multimodal
762 fusion by channel exchanging. *Advances in Neural Information Processing Systems*
763 **33** (2020)
- 764 50. Wei, C., Shen, K., Chen, Y., Ma, T.: Theoretical analysis of self-training with deep
765 networks on unlabeled data. In: International Conference on Learning Representa-
766 tions (2020)
- 767 51. Wei, G., Lan, C., Zeng, W., Zhang, Z., Chen, Z.: Toalign: Task-oriented alignment
768 for unsupervised domain adaptation. *Advances in Neural Information Processing
769 Systems* **34** (2021)

- 765 52. Xiao, L., Xu, J., Zhao, D., Wang, Z., Wang, L., Nie, Y., Dai, B.: Self-supervised do-
766 main adaptation with consistency training. In: 2020 25th International Conference
767 on Pattern Recognition (ICPR). pp. 6874–6880. IEEE (2021)
768 53. Xu, R., Liu, P., Wang, L., Chen, C., Wang, J.: Reliable weighted optimal transport
769 for unsupervised domain adaptation. In: Proceedings of the IEEE/CVF Conference
770 on Computer Vision and Pattern Recognition. pp. 4394–4403 (2020)
771 54. Xu, R., Li, G., Yang, J., Lin, L.: Larger norm more transferable: An adaptive
772 feature norm approach for unsupervised domain adaptation. In: ICCV. pp. 1426–
773 1435 (2019)
774 55. Yang, J., Xu, R., Li, R., Qi, X., Shen, X., Li, G., Lin, L.: An adversarial perturba-
775 tion oriented domain adaptation approach for semantic segmentation. In: Proceed-
776 ings of the AAAI Conference on Artificial Intelligence. vol. 34, pp. 12613–12620
777 (2020)
778 56. Yue, Z., Sun, Q., Hua, X.S., Zhang, H.: Transporting causal mechanisms for un-
779 supervised domain adaptation. In: Proceedings of the IEEE/CVF International
780 Conference on Computer Vision. pp. 8599–8608 (2021)
781 57. Zellinger, W., Grubinger, T., Lughofe, E., Natschläger, T., Saminger-Platz, S.:
782 Central moment discrepancy (CMD) for domain-invariant representation learning.
783 In: 5th International Conference on Learning Representations, ICLR 2017, Toulon,
784 France, April 24–26, 2017, Conference Track Proceedings. OpenReview.net (2017)
785 58. Zhang, P., Zhang, B., Zhang, T., Chen, D., Wang, Y., Wen, F.: Prototypical pseudo
786 label denoising and target structure learning for domain adaptive semantic seg-
787 mentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and
788 Pattern Recognition. pp. 12414–12424 (2021)
789 59. Zhang, Q., Zhang, J., Liu, W., Tao, D.: Category anchor-guided unsupervised
790 domain adaptation for semantic segmentation. Advances in Neural Information
791 Processing Systems **32**, 435–445 (2019)
792 60. Zhang, Y., Tang, H., Jia, K., Tan, M.: Domain-symmetric networks for adversarial
793 domain adaptation. In: Proceedings of the IEEE Conference on Computer Vision
794 and Pattern Recognition. pp. 5031–5040 (2019)
795 61. Zhang, Y., Liu, T., Long, M., Jordan, M.: Bridging theory and algorithm for do-
796 main adaptation. In: International Conference on Machine Learning. pp. 7404–7413
797 (2019)
798 62. Zheng, Z., Yang, Y.: Rectifying pseudo label learning via uncertainty estimation
799 for domain adaptive semantic segmentation. International Journal of Computer
800 Vision **129**(4), 1106–1120 (2021)
801 63. Zou, Y., Yu, Z., Kumar, B., Wang, J.: Unsupervised domain adaptation for seman-
802 tic segmentation via class-balanced self-training. In: Proceedings of the European
803 conference on computer vision (ECCV). pp. 289–305 (2018)
804 64. Zou, Y., Yu, Z., Liu, X., Kumar, B., Wang, J.: Confidence regularized self-training.
805 In: Proceedings of the IEEE/CVF International Conference on Computer Vision.
806 pp. 5982–5991 (2019)