

1 Adaptive Feature Swapping for Unsupervised Domain Adaptation

2

3

4

5

6 Anonymous Author(s)

7 Submission Id: 910*

8 ABSTRACT

9 The bottleneck of visual domain adaptation always lies in the learning
10 of domain invariant representations. In this paper, we present a
11 simple but effective technique named Adaptive Feature Swapping
12 for learning domain invariant features in Unsupervised Domain
13 Adaptation (UDA). Adaptive Feature Swapping aims to select
14 semantically irrelevant features from labeled source data and unlabeled
15 target data and swap these features with each other. Then the merged
16 representations are also utilized for training with prediction
17 consistency constraints. In this way, the model is encouraged
18 to learn representations that are robust to domain-specific information.
19 We develop two swapping strategies including channel
20 swapping and spatial swapping. The former encourages the model
21 to squeeze redundancy out of features and pay more attention to
22 semantic information. The latter motivates the model to be robust
23 to the background and focus on objects. We conduct experiments
24 on object recognition and semantic segmentation in UDA setting
25 and the results show that Adaptive Feature Swapping can promote
various existing UDA methods.

26 CCS CONCEPTS

- 27 • Computer systems organization → Embedded systems; Redundancy; Robotics; • Networks → Network reliability.

31 KEYWORDS

32 Visual Domain Adaptation, Domain Adaptive Semantic Segmentation,
33 Transfer Learning, Object Recognition

35 ACM Reference Format:

36 Anonymous Author(s). 2023. Adaptive Feature Swapping for Unsupervised
37 Domain Adaptation. In *Proceedings of Make sure to enter the correct conference
38 title from your rights confirmation email (Conference acronym 'XX)*. ACM,
39 New York, NY, USA, 10 pages. <https://doi.org/XXXXXXX.XXXXXXX>

40 1 INTRODUCTION

41 Domain adaptation improves a target task with insufficient annotations by knowledge transfer from a source domain with rich
42 annotations. To achieve reliable transfer, the domain discrepancy between the source and target domains should be mitigated in domain
43 adaptation. There are two main directions for mitigating domain
44 discrepancy including moment alignment [11, 19, 28, 30, 42, 56] and
45 adversarial training [8, 29, 40, 59], based on popular assumptions
46 such as Covariate Shift [41]. The former aims to minimize domain
47

48 Permission to make digital or hard copies of all or part of this work for personal or
49 classroom use is granted without fee provided that copies are not made or distributed
50 for profit or commercial advantage and that copies bear this notice and the full citation
51 on the first page. Copyrights for components of this work owned by others than ACM
52 must be honored. Abstracting with credit is permitted. To copy otherwise, or republish,
53 to post on servers or to redistribute to lists, requires prior specific permission and/or a
54 fee. Request permissions from permissions@acm.org.

55 Conference acronym 'XX, June 03–05, 2023, Woodstock, NY

56 © 2023 Association for Computing Machinery.

57 ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00

58 <https://doi.org/XXXXXXX.XXXXXXX>

59 discrepancy estimated by maximum mean discrepancy (MMD) [28],
60 correlation distance [42] or other distance metric [19, 30] calculated
61 on task-specific features. The latter direction of methods learn features
62 that are indistinguishable from an additional domain discriminator [8, 29] or utilizes adversarial learning between two classifiers
63 for aligning conditional feature distributions [21, 40, 59, 60].

64 Recently, domain-specific modeling and perturbation-tolerant
65 learning in domain adaptation have drawn considerable attention.
66 Modeling the domain-specific features enhances the learning
67 of domain invariant features [1, 4, 9] for UDA. For example,
68 Bousmalis et al. [1] proposed to learn domain-specific features
69 via time-consuming reconstruction at the pixel level and Cui et
70 al. [4] adopted a single fully connected layer whose responses
71 are restricted to vanish. By introducing auxiliary tasks like recon-
72 struction or vanishing bridge responses, these methods learn the
73 domain-specific features in implicit ways resulting in suboptimal
74 solutions with additional computation overhead.

75 Perturbation-tolerant learning imposes robustness over input
76 perturbations for enhancing the generalization ability of trained
77 models. Such robustness can be measured by the norm of the input-
78 output Jacobian of the network, and it correlates well with general-
79 ization [34]. To enhance the generalization in the target domain,
80 several methods [6, 39, 52] imposed consistency between the predic-
81 tions of two perturbed samples over a target sample with random
82 data augmentations. With expansion assumption, Wei et al. [50]
83 proved that enforcing locally consistent prediction provides accu-
84 racy guarantees on unlabeled target data for unsupervised domain
85 adaptation. Typical consistency constraints include L2 distance min-
86 imization [6], agreement maximization [39] and mutual information
87 maximization [52]. These methods usually require two complete
88 time-consuming forward computations of deep network. Besides,
89 the underlying information of samples from the two domains for
90 generating perturbation is not well explored. Perturbation based
91 on usual random augmentation may not capture the structure of
92 cross-domain data and restricts the generalization of trained model.

93 To alleviate aforementioned issues, we present Adaptive Feature
94 Swapping (AFS) to model the domain-specific features and acquire
95 sample perturbation more efficiently and effectively. The key in-
96 sight is to swap semantically irrelevant features between source
97 and target data, obtaining perturbed representations that maintain
98 the semantic information. By encouraging the consistency between
99 predictions of the representations after swapping and the original
100 ones, the model attains robustness that benefits better generaliza-
101 tion and more reliable knowledge transfer. Enforcing the model
102 to predict the ground truth label over the source representations
103 after swapping encourages the model to pay more attention to the
104 domain invariant features against the domain-specific ones. Be-
105 sides, AFS supports accuracy guarantees on unlabeled target data
106 as locally consistent prediction is ensured [50].

107 We develop two kinds of Adaptive Feature Swapping techniques
108 from spatial and channel perspectives to facilitate UDA. For Spatial
109

117 Feature Swapping, since the irrelevant information like the back-
 118 ground surrounding an object should not affect the prediction of
 119 the object, we attempt to swap the irrelevant features of source sam-
 120 ples with those of target samples. To enable reliable Spatial Feature
 121 Swapping, we extract the spatial attention to locating the irrelevant
 122 spatial features. We adopt activation-based attention [20] due to
 123 its computation efficiency, and the features with less attention are
 124 regarded as irrelevant ones. Initially, the model may mistakenly put
 125 more attention to some semantically irrelevant features. However,
 126 as the model is optimized during training, the attention is promoted
 127 to better locate the irrelevant features, which makes Spatial Feature
 128 Swapping more effective.

129 Sharing the same spirit of Spatial Feature Swapping, Channel
 130 Feature Swapping is conducted to squeeze the redundancy out of
 131 features, so the model tends to pay more attention to the semantic
 132 information. Similarly, the accumulation of feature responses of
 133 samples in a mini-batch is regarded as channel-wise attention, serv-
 134 ing to select irrelevant features for swapping. Since the two feature
 135 swapping strategies encourage the model to capture semantic infor-
 136 mation from different perspectives, they could cooperate well with
 137 each other which is verified in our experiments. Moreover, the rep-
 138 resentation with feature swapping randomly integrate feature from
 139 samples of other domain, exhibiting more diverse perturbations for
 140 domain adaptation task.

141 We evaluate our AFS on several benchmarks for object classifi-
 142 cation and semantic segmentation in UDA setting. Experimental
 143 results show that the proposed Adaptive Feature Swapping im-
 144 proves source-only model by a considerable margin. Besides, as a
 145 lightweight technique, Adaptive Feature Swapping can be easily
 146 plugged into various DA methods and boost their performances.

2 RELATED WORK

149 The unsupervised domain adaptation methods can be mainly di-
 150 vided into moment alignment based methods [19, 28, 30, 42, 56]
 151 and adversarial training based methods [8, 29, 40, 59]. Moment
 152 alignment based methods try to estimate the domain discrepancy
 153 related to distribution moments over deep representations. [45]
 154 firstly adopted the MMD [10] in the deep neural network and later
 155 DAN [28] utilized the multiple kernel variant of MMD over multi-
 156 ple layer representations for better moment matching. CAN [19]
 157 proposed to optimize a new metric that explicitly models the intra-
 158 class and the inter-class domain discrepancies. Enhanced Transport
 159 Distance (ETD) [24] builds an attention-aware transport distance,
 160 which can be viewed as the prediction feedback of the iteratively
 161 learned classifier, to measure the domain discrepancy.

162 Early Adversarial training based methods introduced a domain
 163 discriminator [8, 29] that is trained to distinguish source and target
 164 representations while the backbone network is trained to gener-
 165 ate representations to fool the domain discriminator. DANN [8]
 166 performed domain adversarial learning on features while [29] per-
 167 formed domain adversarial learning on the multilinear mapping
 168 of feature and classifier response. Recent methods [21, 40, 59, 60]
 169 utilized adversarial learning between two classifiers for aligning
 170 conditional feature distributions. In [59], the authors used two
 171 asymmetrical classifiers to estimate the conditional feature distri-
 172 butions with marginal loss. In [51], the authors propose to make
 173 the domain alignment proactively serve classification via feature

175 decomposition and alignment with the prior knowledge induced
 176 from the classification task.

177 Other methods focus on some characteristics of specific layers
 178 in a deep neural network for domain adaptation such as adaptively
 179 increasing the feature norms in AFN [53], maintaining both dis-
 180 crimination and diversity in BNM [3] and minimizing inter-class
 181 confusion of unlabeled target data in MCC [18]. Besides, Bousmalis
 182 et al. [1] developed a split model that models the shared domain in-
 183 variant features, and domain-specific features supporting effective
 184 reconstruction. Gong et al. [9] proposed a domain flow generation
 185 model to generate a continuous sequence of intermediate domains
 186 with various domain-specific information. Both methods require
 187 time-consuming reconstruction at the pixel level. Recently, Cui
 188 et al. [4] proposed to use a single fully connected layer, named
 189 gradually vanishing bridge (GVB) to capture domain-specific fea-
 190 tures. Dong et al. [5] developed a novel Knowledge Aggregation-
 191 induced Transferability Perception (KATP) module to distinguish
 192 transferable or untransferable knowledge across domains. He et
 193 al. [13] proposed Domain-specific Conditional Jigsaw Adaptation
 194 Network which simultaneously encourages the network to extract
 195 transferable and discriminative features. These methods require
 196 careful module design, while our method learns domain-specific
 197 features via simple feature swapping with prediction consistency
 198 constraints.

199 For the perturbation-tolerant learning, two random data aug-
 200 mentations are performed over a target sample and the consistency
 201 between the predictions of the two augmented samples is imposed
 202 in [6]. In [39], the authors proposed to minimize Min-Entropy Con-
 203 sensus loss that univocally selects a pseudo-label that maximizes
 204 the agreement between two perturbed versions of the same target
 205 sample. In [52], the authors explicitly maximized the mutual in-
 206 formation between the rotated image and the label. Our method
 207 differs from these methods by implementing the perturbation via
 208 feature swapping. The other related work is [49] which exchanged
 209 channel activations between the two convolutional layer activa-
 210 tions extracted from two modalities of a sample for better fusing the
 211 two modalities. [49] selected feature according to the magnitude of
 212 Batch-Normalization scaling factor. On the contrary, our method
 213 swaps both **channel and spatial** features from two **independent**
 214 **samples** according to carefully designed **aggregated attention**.

3 METHODOLOGY

215 We illustrate our method in object recognition task under unsuper-
 216 vised domain adaptation (UDA) setting for simplicity. Suppose that
 217 there are a source domain $D_s = \{(x_i^s, y_i^s)\}_{i=1}^{n_s}$ of n_s labeled samples
 218 and a target domain $D_t = \{(x_i^t)\}_{i=1}^{n_t}$ of n_t unlabeled samples. The
 219 two domains distribute differently that the input marginal distri-
 220 butions $p(x)$ vary. In UDA, both the two domains share the same
 221 K categories. The goal is to learn a domain invariant model that
 222 generalizes well in target domain.

223 We propose Adaptive Feature Swapping for learning domain
 224 invariant representations. As shown in Fig. 1, Adaptive Feature
 225 Swapping involves spatial swapping and channel swapping. Spa-
 226 tial swapping aims to select irrelevant features like background
 227 for source and target samples and swap these features with each
 228 other. Similarly, channel swapping selects irrelevant features along
 229

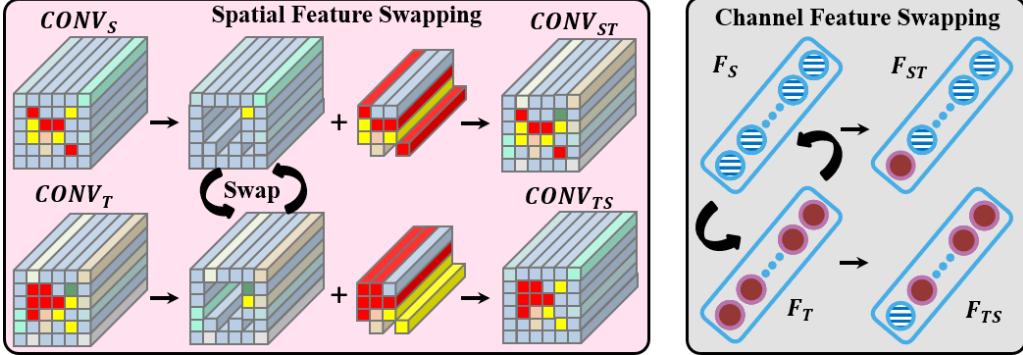


Figure 1: Illustration of Spatial Feature Swapping and Channel Feature Swapping. For the convolutional activations of a source (target) sample, the semantic features (in red and yellow) are kept and combined with the rest irrelevant features from convolutional activations of a target (source) sample, obtaining new convolutional activations. The Channel Feature Swapping exchanges some irrelevant elements between source and target features.

the channel and swaps these features between source and target samples. The prediction over the original feature and the one after swapping is enforced to be consistent. Both feature swapping strategies encourage the model to pay more attention to the semantic information that is invariant across domains.

3.1 Spatial Feature Swapping

We perform Spatial Feature Swapping on a specific convolutional layer activations. Suppose that for a batch of source samples and a batch of target samples both with size B , their convolutional layer activations are $CONV_S \in \mathcal{R}^{B \times C \times W \times H}$ and $CONV_T \in \mathcal{R}^{B \times C \times W \times H}$. To perform feature swapping, we need to determine which spatial feature is semantically irrelevant. Inspired by activation-based attention [20], we first map the batch of convolutional layer activations into an attention map and then select features with less attention for swapping. Concretely, the spatial attention $M_{sp}(A)$ for an activation tensor A is obtained via

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

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

To prevent a semantically relevant feature in $CONV_S$ to be swapped to $CONV_T$ that affects the semantic of swapped $CONV_{TS}$, we select features corresponding to the least r_1 percentage of aggregated attention ($M_{sp}(CONV_S) + M_{sp}(CONV_T)$) for swapping.

3.2 Channel Feature Swapping

In CNN-based classifier, the feature before the final classifier layer is usually of high dimension that involves redundancy. The network may tend to remember irrelevant features for recognition, especially for a small scale labeled source domain. This issue motivates us to adopt Channel Feature Swapping for encouraging 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 [17] or scaled transformation [7]. Among these solutions, GAP is simple yet effective and widely used in NIN [26], GoogLeNet [43], and ResNet [12]. 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 [53], MMD loss [28], MCC loss [18] or adversarial discriminator score [8].

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

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

349 where p_i^{st} is the classification response over representation of
 350 source sample x_i^s that is swapped with some features from target
 351 sample. And p_j^{ts} is the classification response over the representa-
 352 tion of target sample x_j^t that is swapped with some features from
 353 a source sample. p_j^t is the classification response of the original
 354 target sample x_j^t .
 355

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

362 The total loss of model learning with Adaptive Feature Swapping
 363 is

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

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

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

372 3.4 Theoretical Insight

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

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

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

382 In accordance with [50], we assume P , a distribution of unla-
 383 beled examples over input space X , is \mathcal{B} -separated with proba-
 384 bility $1 - \mu$ by G^* that $R_{\mathcal{B}}(G^*) \leq \mu$. Besides, let $\mathcal{M}(G_{pl}) = \{x :
 385 G_{pl}(x) \neq G^*(x)\}$ be the set of examples that mistakenly labeled
 386 by a pseudolabeler G_{pl} . The expansion assumption stated in [50]
 387 requires that $P_i(\mathcal{N}(V)) \leq \min(cP_i(V), 1)$ for all $V \subseteq X$ and all class
 388 i with $P_i(V) \leq a$, where P_i is the class-conditional distribution,
 389 $a = \max_i\{P_i(\mathcal{M}(G_{pl}))\}$ and a is required to be less than $1/3$ and
 390 $c > 3$.

391 Then it is proved in [50] that for any minimizer \widehat{G} of
 392 $\min_G \frac{c+1}{c-1} L_{0-1}(G, G_{pl}) + \frac{2c}{c-1} R_{\mathcal{B}}(G) - Err(G_{pl})$, we have

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

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

395 Therefore, our method could achieve the above bounded target
 396 error (Eqn. (8)) for unsupervised domain adaptation under ex-
 397 pansion assumption (this assumption is justified by Wei et al. [50] on
 398 real-world datasets with BigGAN). Specifically, during the train-
 399 ing process, the current model has relatively low $Err(G_{pl})$, which
 400 in turn lower the error rate $Err(\widehat{G})$ of the updated model. With
 401

402 stronger base model like MCC [18] and BNM [3], the $Err(G_{pl})$ is
 403 lower. Therefore, the trained model gets lower error rate on target
 404 domain, as verified experimentally in Section 4.

4 EXPERIMENTS

405 4.1 Datasets and Setup

406 We use two standard object classification benchmarks and a stan-
 407 dard semantic segmentation benchmark for unsupervised domain
 408 adaptation to evaluate the proposed methods. Office-Home [46] is
 409 a challenging benchmark with 65 categories in four domains: *Art*
 410 (*Ar*), *Clipart* (*Cl*), *Product* (*Pr*) and *Real-world* (*Rw*), which contains
 411 12 transfer tasks. VisDA-C [36] is a large-scale benchmark that con-
 412 tains the images from 12 categories of two very distinct domains,
 413 the synthetic domain, and the real-world domain. The synthetic
 414 domain contains 152,397 images and the latter contains 55,388 im-
 415 ages. Following standard protocol, we focus on the synthetic-to-real
 416 transfer task.

417 For the segmentation task, we evaluate our method by adapting
 418 the segmentation from game scenes, GTA5 [38] dataset, to real
 419 scenes, the Cityscapes [2] dataset. GTA5 contains 24,966 images
 420 with the resolution of 1914×1052. The Cityscapes dataset contains
 421 2,975 training images and 500 images for validation with the reso-
 422 lution of 2048×1024.

423 Following standard protocol, we use the training set of GTA5 as
 424 the labeled source domain and the unlabeled training images from
 425 Cityscapes as the target domain. We conduct evaluations on the
 426 validation set of Cityscapes and adopt the Intersection-over-Union
 427 (IoU) of each class and the mean-Intersection-over-Union (mIoU)
 428 as performance metrics. For the task GTA5→Cityscapes, we report
 429 the results on the common 19 classes.

430 4.2 Implementation Details

431 **Classification.** The methods are implemented based on PyTorch [35],
 432 and for fair comparisons, we use the ResNet-50/ResNet-101 [12]
 433 as the backbone of network, fixing the batch-size of 36 for all exper-
 434 iments. We adopt mini-batch SGD optimizer with momentum
 435 0.9 to train the model. For Office-Home, we use ResNet-50 as the
 436 backbone, and the hyperparameters r_1, r_2, λ_{cst} are set to 1/10, 1/15
 437 and 40, respectively. For VisDA-C we adopt ResNet-101 as backbone
 438 for fair comparisons. The hyperparameters r_1, r_2 and λ_{cst} are set
 439 to 1/4, 1/3 and 40, respectively. The Channel Feature Swapping is
 440 performed on the feature after mean pooling over the last residual
 441 block of ResNet [12]. The spatial Feature swapping is performed on
 442 the activations of the third residual block. We evaluate our approach
 443 by applying Adaptive Feature Swapping (AFS) to the source-only
 444 model and several mainstream DA methods including DANN [8],
 445 CDAN [29], BNM [3], MCC [18], and CST [27] based on their open-
 446 source codes. λ_{tran} for BNM, CDAN and MCC are kept the same as
 447 their original implementations. The reimplemented MCC is marked
 448 with † as MCC†. **Codes can be found in the Supplementary.**

449 **Segmentation.** We evaluate our approach by applying Channel
 450 Feature Swapping to the source-only model and the state-of-the-art
 451 domain adaptive segmentation method ProDA [57], DAFormer [15],
 452 and HRDA [16] based on their open-source codes. Our training is
 453 carried out on 2 TITAN RTX GPUs with PyTorch [35]. To train
 454 our models using SGD, the initial learning rate is set to 1e-4, and
 455

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

Method	plane	beybl	bus	car	horse	knife	mycle	persn	plant	sktb	train	trunk	Avg
DAN [28]	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
DRMEA [33]	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
DTA [22]	93.7	82.2	85.6	83.8	93.0	81.0	90.7	82.1	<u>95.1</u>	78.1	86.4	32.1	81.5
STAR [31]	95.0	84.0	<u>84.6</u>	73.0	91.6	91.8	85.9	78.4	94.4	84.7	<u>87.0</u>	42.2	82.7
RWOT [54]	95.1	80.3	<u>83.7</u>	90.0	92.4	68.0	92.5	82.2	87.9	78.4	90.4	68.2	84.0
ResNet [12] +AFS	67.7 95.1	27.4 59.3	50.0 83.9	61.7 72.9	69.5 <u>95.3</u>	13.7 66.2	85.9 <u>91.9</u>	11.5 63.4	64.4 93.6	34.4 62.3	84.2 84.1	19.2 17.0	49.1 73.8
DANN [8] +AFS	87.7 96.7	34.8 56.1	84.5 85.2	44.3 70.0	66.0 96.0	62.3 85.7	73.8 91.8	48.1 80.0	59.0 94.5	34.5 59.0	67.4 85.0	18.9 19.0	56.8 76.6
CDAN [29] +AFS	85.2 93.0	66.9 73.5	83.0 80.1	50.8 69.5	84.2 91.8	74.9 <u>97.0</u>	88.1 82.9	74.5 79.0	83.4 83.3	76.0 72.4	81.9 84.5	38.0 39.7	73.9 78.1
BNM [3] +AFS	95.6 <u>95.9</u>	86.7 <u>87.0</u>	77.4 78.1	48.9 53.6	92.0 93.7	83.3 90.5	76.2 78.8	76.3 78.8	88.0 87.4	64.2 84.9	84.7 83.4	46.3 47.0	76.6 79.9
MCC [18]† +AFS	91.5 93.1	82.6 78.5	75.9 81.1	61.5 79.0	91.4 93.9	92.9 94.2	81.6 88.1	79.6 85.1	85.8 91.0	88.4 85.6	69.5 85.2	51.8 41.6	80.2 83.1
CST+SAM [27]† +AFS	95.2 97.2	89.0 89.0	75.8 83.8	89.6 <u>87.5</u>	95.7 95.7	98.5 94.0	87.8 87.9	<u>84.1</u> 82.8	94.4 96.6	<u>90.9</u> 94.7	82.4 82.4	<u>54.5</u> 53.1	<u>86.5</u> 87.1

489 adjusted according to the ‘poly’ learning rate scheduler with a
 490 power of 0.9. The weight decay is set to 0.0005. These settings are the
 491 same as ProDA [57] for fair comparisons. In accordance with [15]
 492 and [16], we train our model with AdamW, a base learning rate of
 493 6e-5 for the encoder and 6e-4 for the decoder, a weight decay of 0.01,
 494 linear learning rate warmup with 1.5k iteration, and linear decay
 495 afterward. The data augmentation and other hyperparameters are
 496 kept the same as [15] and [16]. The hyperparameters r_2 , λ_{cst} are
 497 set to 1/15 and 15, respectively. The batch size of HRAD is set to
 498 1 rather than 2 as adopted in their open source code due to the
 499 limited computation resources (marked with †).

500 Note that ProDA [57] adopts pseudo label for training target
 501 samples, so the consistency loss on target domain is directly imple-
 502 mented as the cross entropy loss of prediction over representations
 503 with feature swapping and its corresponding pseudo labels. For
 504 DAFormer and HRDA, we swap features from the semantic de-
 505 coder and obtain consistency loss from both the source domain and
 506 target domain. Only Channel Feature Swapping is utilized since
 507 swapping spatial features changes their semantic in segmentation
 508 task. Besides, due to the limited computation resources, we only
 509 implemented ProDA for its first stage with a batch of 2 images
 510 (marked with †).

4.3 Classification Results

513 Results on VisDA-C are reported in Tab.1. We can observe that
 514 our method consistently improves the generalization ability of all
 515 baseline methods. Compared with source only model, integrating
 516 our AFS brings 24.7% accuracy improvement on average. The
 517 results show that our methods can directly improve the perfor-
 518 mance without any other tricks. As for CST+SAM [27], our simple
 519 model surpasses it by 0.6%. Note that input-consistency via ran-
 520 dom augmentation is involved in CST+SAM, leaving little room for
 521 our method for improvement. However, we can observe that our

522 method could still boost CST+SAM, showing its promising potential
 523 to boost existing unsupervised domain adaptation methods. Based
 524 on these experiments, we can infer that AFS can stably enhance
 525 the transfer ability of classifier.

526 Results on Office-Home are summarized in Tab. 2. Office-Home
 527 is a challenging dataset due to its large domain discrepancy. Utiliz-
 528 ing the proposed AFS into source only model brings 19.2% accuracy
 529 improvement on average. **Note that Source-only+AFS outper-**
 530 **forms SE [6], indicating the potential of AFS to generate more**
 531 **valuable perturbations compared with random augmenta-**
 532 **tion adopted in SE.** Compared with CST [27], our simple model
 533 surpasses it by 0.4% on average. Our AFS also boosts the CDAN,
 534 BNM and MCC baselines by considerable margins. Based on these
 535 promising results, we can infer that AFS can stably explore truly
 536 useful semantic information to better adapt the classifier.

537 **Ablation Study.** In Tab. 3, we validate the influence of the pro-
 538 posed spatial swapping (SS) and channel swapping (CS) in VisDA-C
 539 dataset. All the results reported in Tab. 3 show that our methods
 540 can stably improve the transferability of the model.

541 Firstly, we validate the effectiveness of using spatial swapping.
 542 Compared with baseline MCC, integrating spatial swapping (MCC+SS)
 543 shows a gain of 1.3% accuracy on average. Based on it, we can infer
 544 that spatial swapping helps the model to pay more attention to
 545 objects instead of the background. Secondly, the large improve-
 546 ment is achieved by adopting channel swapping into baseline MCC
 547 (MCC+CS) that encourages the model to focus on semantic infor-
 548 mation. Finally, we also utilize both two swapping strategies into
 549 baseline MCC and it performs better than MCC+SS and MCC+CS.
 550 Therefore, the two feature swapping strategies can benefit each
 551 other and help the model achieve better performance. **The compu-**
 552 **tation overhead analysis can be referred to the Supplemen-**
 553 **tary.**

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	ClAr	PrAr	RwCl	ArCl	PrCl	RwPr	ArPr	ClPr	RwRw	ArRw	ClRw	Pr Avg
DAN [28]	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
SE [6]	48.8	61.8	72.8	54.1	63.2	65.1	50.6	49.2	72.3	66.1	55.9	78.7	61.5
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 [59]	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 [24]	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 [55]	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
AAA [23]	56.7	78.3	82.1	66.4	<u>78.5</u>	79.4	<u>67.6</u>	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 [12]	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
DANN [8]	42.5	63.5	73.5	51.4	60.7	63.5	53.6	37.2	73.0	65.2	42.6	77.1	58.7
+AFS	48.6	69.2	75.4	56.0	64.7	68.2	56.4	42.8	74.5	67.5	49.1	79.7	62.7
CDAN [29]	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 [18]†	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	<u>59.0</u>	77.5	81.4	67.2	76.8	76.0	65.8	<u>57.3</u>	81.9	74.0	61.8	85.0	72.0
CST [27]†	<u>59.0</u>	<u>79.6</u>	83.4	68.4	77.1	<u>76.7</u>	68.9	56.4	83.0	75.3	<u>62.2</u>	<u>85.1</u>	<u>73.0</u>
+AFS	58.7	80.2	83.3	67.6	79.0	76.7	68.9	57.1	82.6	75.1	<u>65.5</u>	<u>85.7</u>	<u>73.4</u>

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	bcybl	bus	car	horse	knife	mcycle	persn	plant	sktb	train	trunk	Avg
MCC [18]†	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	<u>86.7</u>	<u>81.7</u>	<u>46.3</u>	81.5
MCC+CS	93.5	85.4	78.9	<u>70.6</u>	<u>93.4</u>	93.0	<u>86.0</u>	<u>81.6</u>	92.5	85.8	80.7	46.0	82.3
MCC+SS+CS	<u>93.1</u>	78.5	81.1	79.0	<u>93.9</u>	<u>94.2</u>	88.1	85.1	<u>91.0</u>	85.6	85.2	41.6	83.1

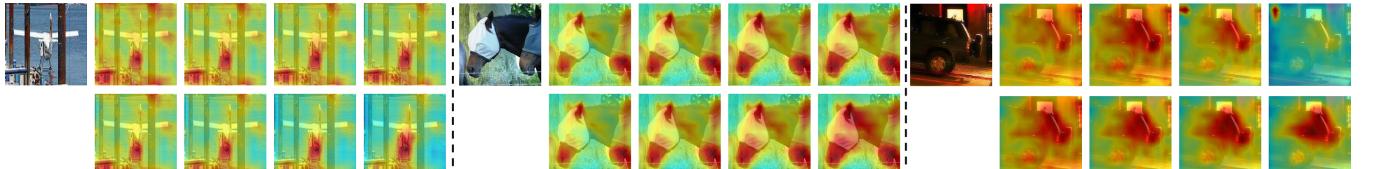


Figure 2: For each group of images, the 1st column are images from VisDA-C dataset. The 2nd to fifth columns show the attention maps extracted from the activations of third residual block of model trained with 1K, 3K, 5K, and 10K iterations, respectively. The model that produces attention maps of the 1st (2nd), 3rd (4th), and 5th (6th) rows is MCC† baseline (MCC+AFS). The model pays more attention to the regions in red than regions in blue.

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

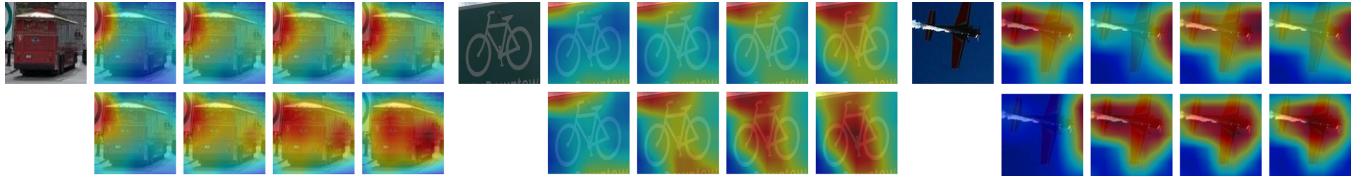


Figure 3: For each group of images, the 1st column are images from VisDA-C dataset. The 2nd to fifth columns show the attention maps extracted from the activations of the last residual block of model trained with 1K, 3K, 5K, and 10K iterations, respectively. The model that produces attention maps of the 1st (2nd), 3rd (4th), and 5th (6th) rows is MCC \dagger baseline (MCC+AFS). The model pays more attention to the regions in red than regions in blue.

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	S.walk	Build.	Wall	Fence	Pole	T.light	Sign	Veget.	Terrain	Sy	Person	Rider	Car	Truck	Bus	Train	M.bike	Bike	mIoU
AdaStruct [44]	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 [14]	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 [32]	87.0	27.1	79.6	27.3	23.3	28.3	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
ADVENT [47]	89.4	33.1	81.0	26.6	26.8	27.2	33.5	24.7	83.9	36.7	78.8	58.7	30.5	84.8	38.5	44.5	1.7	31.6	32.4	45.5
BDL [25]	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
MRKLD [62]	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 [58]	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 [61]	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 [57] \dagger	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 [15]	95.7	70.2	89.4	53.5	48.1	49.6	55.8	59.4	89.9	47.9	92.5	72.2	44.7	92.3	74.5	78.2	65.1	55.9	61.8	68.3
+AFS	<u>96.1</u>	<u>71.8</u>	89.4	53.8	44.5	49.9	56.6	62.5	<u>90.0</u>	51.0	90.9	71.7	45.0	92.3	72.3	<u>81.9</u>	<u>72.0</u>	56.6	63.6	<u>69.0</u>
HRDA [16] \dagger	95.5	69.1	<u>89.9</u>	<u>56.2</u>	<u>46.7</u>	<u>52.6</u>	58.4	<u>60.8</u>	89.6	45.4	<u>93.6</u>	74.2	25.9	<u>93.2</u>	82.6	80.3	66.6	<u>60.5</u>	66.8	68.8
+AFS	96.4	74.5	90.0	57.2	<u>45.2</u>	<u>54.5</u>	60.3	60.6	90.6	<u>49.3</u>	94.0	73.8	22.8	<u>93.5</u>	<u>81.0</u>	<u>84.3</u>	73.4	<u>62.6</u>	<u>67.3</u>	70.1

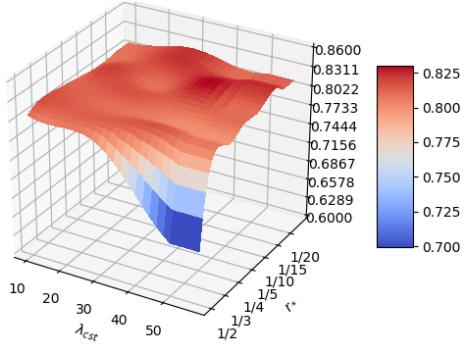


Figure 4: Hyper-parameter sensitivity analysis of AFS on VisDA-C based on MCC [18].

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 including learned representation visualization can be referred to the Supplementary.

Hyper-parameter Sensitivity For Classification Task. 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

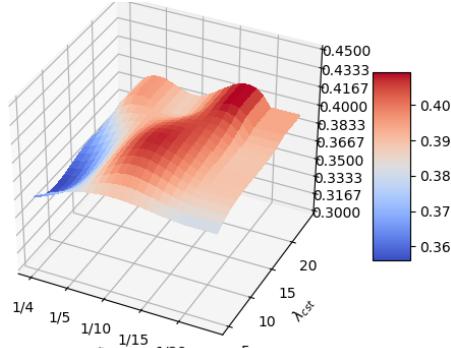


Figure 5: mIoUs of various Source only+AFS on GTA2Cityscape.

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. 4 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.

4.4 Segmentation Results

We report the quantitative evaluation of our method in Tab. 4. From the results, we can see that with Channel Feature Swapping,

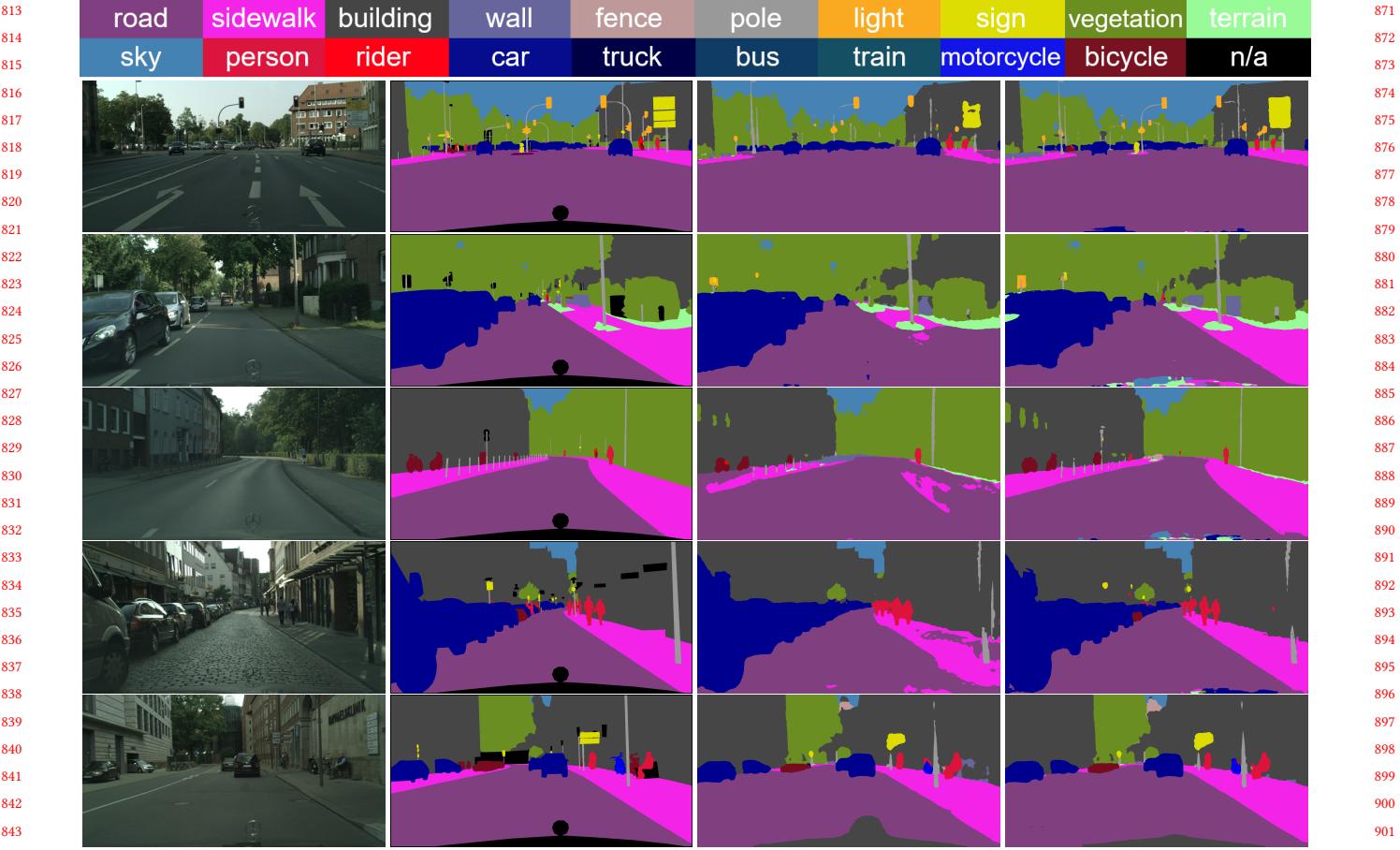


Figure 6: Qualitative results of validation images from Cityscapes where the source domain is GTA5 dataset. The 1st and 2nd columns are images and ground truth labels from validation set of Cityscape. The 3rd and 4th columns are segmentation results produced by DAFormer [15][†] and our method.

the source-only model is improved by 4.1% in mIoU. Besides, we can observe that applying our Channel Feature Swapping into the ProDA [57][†] with ResNet-101 improves the baseline ProDA[†] by 0.7% in mIoU. Channel Feature Swapping also boosts DAFormer [15] and HRDA [16] by 0.7% and 1.3% in mIoU, respectively. 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 provide visual comparison bewteen DAFormer [15] and DAFormer+AFS in Fig. 6. We also visually compare the baseline HRDA [16][†] and HRDA[†]+AFS in Fig. 6. As we can see from these figures, the predictions from our model appear less noisy. We accredit this to the enhanced generalization ability of utilizing Channel Feature Swapping. **More qualitative comparisons can be found in the Supplementary.**

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. 5. 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. It is promising to investiage Adaptive Feature Swapping in UDA under conditional distribution shift and domain generalization task and we leave these tasks for our future work.

REFERENCES

- [1] Konstantinos Bousmalis, George Trigeorgis, Nathan Silberman, Dilip Krishnan, and Dumitru Erhan. 2016. Domain separation networks. In *Advances in Neural Information Processing Systems*. 343–351.
- [2] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. 2016. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 3213–3223.
- [3] Shuhao Cui, Shuhui Wang, Junbao Zhuo, Liang Li, Qingming Huang, and Qi Tian. 2020. Towards discriminability and diversity: Batch nuclear-norm maximization under label insufficient situations. In *CVPR*. 3941–3950.
- [4] Shuhao Cui, Shuhui Wang, Junbao Zhuo, Chi Su, Qingming Huang, and Qi Tian. 2020. Gradually vanishing bridge for adversarial domain adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12455–12464.
- [5] Jiahua Dong, Yang Cong, Gan Sun, Zhen Fang, and Zhengming Ding. 2021. Where and how to transfer: knowledge aggregation-induced transferability perception for unsupervised domain adaptation. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2021).
- [6] Geoff French, Michal Mackiewicz, and Mark Fisher. 2018. Self-ensembling for visual domain adaptation. In *International Conference on Learning Representations*.
- [7] Jun Fu, Jing Liu, Huiji Tian, Yong Li, Yongjun Bao, Zhiwei Fang, and Hanqing Lu. 2019. Dual attention network for scene segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 3146–3154.
- [8] Yaroslav Ganin and Victor Lempitsky. 2015. Unsupervised domain adaptation by backpropagation. In *International conference on machine learning*. 1180–1189.
- [9] Rui Gong, Wen Li, Yuhua Chen, and Luc Van Gool. 2019. Dflow: Domain flow for adaptation and generalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2477–2486.
- [10] Arthur Gretton, Karsten Borgwardt, Malte Rasch, Bernhard Schölkopf, and Alex Smola. 2006. A kernel method for the two-sample-problem. *Advances in neural information processing systems* 19 (2006), 513–520.
- [11] Arthur Gretton, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola. 2012. A kernel two-sample test. *Journal of Machine Learning Research* 13, Mar (2012), 723–773.
- [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [13] Qi He, Zhaoquan Yuan, Xiao Wu, and Jun-Yan He. 2022. Domain-Specific Conditional Jigsaw Adaptation for Enhancing Transferability and Discriminability. In *Proceedings of the 30th ACM International Conference on Multimedia*. 6327–6336.
- [14] Judy Hoffman, Eric Tzeng, Taesung Park, Jun-Yan Zhu, Phillip Isola, Kate Saenko, Alexei Efros, and Trevor Darrell. 2018. Cycada: Cycle-consistent adversarial domain adaptation. In *International conference on machine learning*. PMLR, 1989–1998.
- [15] Lukas Hoyer, Dengxin Dai, and Luc Van Gool. 2021. DAFormer: Improving Network Architectures and Training Strategies for Domain-Adaptive Semantic Segmentation. *arXiv preprint arXiv:2111.14887* (2021).
- [16] Lukas Hoyer, Dengxin Dai, and Luc Van Gool. 2022. HRDA: Context-Aware High-Resolution Domain-Adaptive Semantic Segmentation. *arXiv preprint arXiv:2204.13132* (2022).
- [17] Jie Hu, Li Shen, and Gang Sun. 2018. Squeeze-and-excitation networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 7132–7141.
- [18] Ying Jin, Ximei Wang, Mingsheng Long, and Jianmin Wang. 2020. Minimum class confusion for versatile domain adaptation. In *ECCV*. 464–480.
- [19] Guoliang Kang, Lu Jiang, Yi Yang, and Alexander G Hauptmann. 2019. Contrastive adaptation network for unsupervised domain adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4893–4902.
- [20] Nikos Komodakis and Sergey Zagoruyko. 2017. Paying more attention to attention: improving the performance of convolutional neural networks via attention transfer. In *ICLR*.
- [21] Chen-Yu Lee, Tanmay Batra, Mohammad Haris Baig, and Daniel Ulbricht. 2019. Sliced wasserstein discrepancy for unsupervised domain adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 10285–10295.
- [22] Seungmin Lee, Dongwan Kim, Namil Kim, and Seong-Gyun Jeong. 2019. Drop to adapt: Learning discriminative features for unsupervised domain adaptation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 91–100.
- [23] Jingjing Li, Zhekai Du, Lei Zhu, Zhengming Ding, Ke Lu, and Heng Tao Shen. 2021. Divergence-agnostic Unsupervised Domain Adaptation by Adversarial Attacks. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2021).
- [24] Mengxue Li, Yi-Ming Zhai, You-Wei Luo, Peng-Fei Ge, and Chuan-Xian Ren. 2020. Enhanced transport distance for unsupervised domain adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 13936–13944.
- [25] Yunsheng Li, Lu Yuan, and Nuno Vasconcelos. 2019. Bidirectional learning for domain adaptation of semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6936–6945.
- [26] Min Lin, Qiang Chen, and Shuicheng Yan. 2014. Network In Network. In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14–16, 2014, Conference Track Proceedings*.
- [27] Hong Liu, Jianmin Wang, and Mingsheng Long. 2021. Cycle self-training for domain adaptation. *Advances in Neural Information Processing Systems* 34 (2021).
- [28] Mingsheng Long, Yue Cao, Jianmin Wang, and Michael Jordan. 2015. Learning transferable features with deep adaptation networks. In *International conference on machine learning*. 97–105.
- [29] Mingsheng Long, Zhangjie Cao, Jianmin Wang, and Michael I. Jordan. 2018. Conditional Adversarial Domain Adaptation. In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3–8, 2018, Montréal, Canada*. 1647–1657.
- [30] Mingsheng Long, Han Zhu, Jianmin Wang, and Michael I. Jordan. 2017. Deep transfer learning with joint adaptation networks. In *International conference on machine learning*. 2208–2217.
- [31] Zhihe Lu, Yongxin Yang, Xiatian Zhu, Cong Liu, Yi-Zhe Song, and Tao Xiang. 2020. Stochastic classifiers for unsupervised domain adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9111–9120.
- [32] Yawei Luo, Liang Zheng, Tao Guan, Junqing Yu, and Yi Yang. 2019. Taking a closer look at domain shift: Category-level adversaries for semantics consistent domain adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2507–2516.
- [33] You-Wei Luo, Chuan-Xian Ren, Pengfei Ge, Ke-Kun Huang, and Yu-Feng Yu. 2020. Unsupervised domain adaptation via discriminative manifold embedding and alignment. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 5029–5036.
- [34] Roman Novak, Yasaman Bahri, Daniel A Abolafia, Jeffrey Pennington, and Jascha Sohl-Dickstein. 2018. Sensitivity and Generalization in Neural Networks: an Empirical Study. In *International Conference on Learning Representations*.
- [35] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems* 32 (2019), 8026–8037.
- [36] Xingchao Peng, Ben Usman, Neela Kaushik, Judy Hoffman, Dequan Wang, and Kate Saenko. 2017. Visda: The visual domain adaptation challenge. *arXiv preprint arXiv:1710.06924* (2017).
- [37] Viraj Prabhu, Shivam Khare, Deeksha Kartik, and Judy Hoffman. 2021. Sentry: Selective entropy optimization via committee consistency for unsupervised domain adaptation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 8558–8567.
- [38] Stephan R Richter, Vibhav Vineet, Stefan Roth, and Vladlen Koltun. 2016. Playing for data: Ground truth from computer games. In *European conference on computer vision*. Springer, 102–118.
- [39] Subhankar Roy, Aliaksandr Siarohin, Enver Sangineto, Samuel Rota Bulo, Nicu Sebe, and Elisa Ricci. 2019. Unsupervised domain adaptation using feature-whitening and consensus loss. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 9471–9480.
- [40] Kuniaki Saito, Kohei Watanabe, Yoshitaka Ushiku, and Tatsuya Harada. 2018. Maximum classifier discrepancy for unsupervised domain adaptation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 3723–3732.
- [41] Hideyoshi Shimodaira. 2000. Improving predictive inference under covariate shift by weighting the log-likelihood function. *Journal of statistical planning and inference* 90, 2 (2000), 227–244.
- [42] Baochen Sun and Kate Saenko. 2016. Deep coral: Correlation alignment for deep domain adaptation. In *European Conference on Computer Vision*. Springer, 443–450.
- [43] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. 2015. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1–9.
- [44] Yi-Hsuan Tsai, Wei-Chih Hung, Samuel Schulter, Kihyuk Sohn, Ming-Hsuan Yang, and Manmohan Chandraker. 2018. Learning to adapt structured output space for semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 7472–7481.
- [45] Eric Tzeng, Judy Hoffman, Ning Zhang, Kate Saenko, and Trevor Darrell. 2014. Deep domain confusion: Maximizing for domain invariance. *arXiv preprint arXiv:1412.3474* (2014).
- [46] Hemanth Venkateswara, Jose Eusebio, Shayok Chakraborty, and Sethuraman Panchanathan. 2017. Deep hashing network for unsupervised domain adaptation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 5018–5027.
- [47] Tuan-Hung Vu, Himalaya Jain, Maxime Bucher, Matthieu Cord, and Patrick Pérez. 2019. Advent: Adversarial entropy minimization for domain adaptation in semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2517–2526.

929 930 931 932 933 934 935 936 937 938 939 940 941 942 943 944 945 946 947 948 949 950 951 952 953 954 955 956 957 958 959 960 961 962 963 964 965 966 967 968 969 970 971 972 973 974 975 976 977 978 979 980 981 982 983 984 985 986 987 988 989 990 991 992 993 994 995 996 997 998 999 1000 1001 1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012 1013 1014 1015 1016 1017 1018 1019 1020 1021 1022 1023 1024 1025 1026 1027 1028 1029 1030 1031 1032 1033 1034 1035 1036 1037 1038 1039 1040 1041 1042 1043 1044

- 1045 [48] Haoran Wang, Tong Shen, Wei Zhang, Ling-Yu Duan, and Tao Mei. 2020. Classes
1046 matter: A fine-grained adversarial approach to cross-domain semantic segmen-
1047 tation. In *European Conference on Computer Vision*. Springer, 642–659.
1048 [49] Yikai Wang, Wenbing Huang, Fuchun Sun, Tingyang Xu, Yu Rong, and Junzhou
1049 Huang. 2020. Deep multimodal fusion by channel exchanging. *Advances in
1050 Neural Information Processing Systems* 33 (2020).
1051 [50] Colin Wei, Kendrick Shen, Yining Chen, and Tengyu Ma. 2020. Theoretical
1052 Analysis of Self-Training with Deep Networks on Unlabeled Data. In *International
1053 Conference on Learning Representations*.
1054 [51] Guoqiang Wei, Cuiling Lan, Wenjun Zeng, Zhizheng Zhang, and Zhibo Chen.
1055 2021. ToAlign: Task-oriented Alignment for Unsupervised Domain Adaptation.
1056 *Advances in Neural Information Processing Systems* 34 (2021).
1057 [52] Liang Xiao, Jiaolong Xu, Dawei Zhao, Zhiyu Wang, Li Wang, Yiming Nie, and
1058 Bin Dai. 2021. Self-Supervised Domain Adaptation with Consistency Training. In
1059 *2020 25th International Conference on Pattern Recognition (ICPR)*. IEEE, 6874–6880.
1060 [53] Ruijia Xu, Guanbin Li, Jihan Yang, and Liang Lin. 2019. Larger norm more trans-
1061 ferable: An adaptive feature norm approach for unsupervised domain adaptation.
1062 In *ICCV*. 1426–1435.
1063 [54] Renjun Xu, Pelen Liu, Liyan Wang, Chao Chen, and Jindong Wang. 2020. Reliable
1064 weighted optimal transport for unsupervised domain adaptation. In *Proceedings
1065 of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4394–
1066 4403.
1067 [55] Zhongqi Yue, Qianru Sun, Xian-Sheng Hua, and Hanwang Zhang. 2021. Trans-
1068 porting causal mechanisms for unsupervised domain adaptation. In *Proceedings
1069 of the IEEE/CVF International Conference on Computer Vision*. 8599–8608.
- 1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
- [56] Werner Zellinger, Thomas Grubinger, Edwin Lughofer, Thomas Natschläger,
1103 and Susanne Saminger-Platz. 2017. Central Moment Discrepancy (CMD) for
1104 Domain-Invariant Representation Learning. In *5th International Conference on
1105 Learning Representations, ICLR 2017, Toulon, France, April 24–26, 2017, Conference
1106 Track Proceedings*. OpenReview.net.
1107 [57] Pan Zhang, Bo Zhang, Ting Zhang, Dong Chen, Yong Wang, and Fang Wen. 2021.
1108 Prototypical pseudo label denoising and target structure learning for domain
1109 adaptive semantic segmentation. In *Proceedings of the IEEE/CVF Conference on
1110 Computer Vision and Pattern Recognition*. 12414–12424.
1111 [58] Qiming Zhang, Jing Zhang, Wei Liu, and Dacheng Tao. 2019. Category Anchor-
1112 Guided Unsupervised Domain Adaptation for Semantic Segmentation. *Advances
1113 in Neural Information Processing Systems* 32 (2019), 435–445.
1114 [59] Yuchen Zhang, Tianle Liu, Mingsheng Long, and Michael Jordan. 2019. Bridging
1115 theory and algorithm for domain adaptation. In *International Conference on
1116 Machine Learning*. 7404–7413.
1117 [60] Yabin Zhang, Hui Tang, Kui Jia, and Mingkui Tan. 2019. Domain-Symmetric Net-
1118 works for Adversarial Domain Adaptation. In *Proceedings of the IEEE Conference
1119 on Computer Vision and Pattern Recognition*. 5031–5040.
1120 [61] Zhedong Zheng and Yi Yang. 2021. Rectifying pseudo label learning via un-
1121 certainty estimation for domain adaptive semantic segmentation. *International
1122 Journal of Computer Vision* 129, 4 (2021), 1106–1120.
1123 [62] Yang Zou, Zhiding Yu, Xiaofeng Liu, BVK Kumar, and Jinsong Wang. 2019.
1124 Confidence regularized self-training. In *Proceedings of the IEEE/CVF International
1125 Conference on Computer Vision*. 5982–5991.
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160