# ENHANCED SEPARABLE DISENTANGLEMENT FOR UNSUPERVISED DOMAIN ADAPTATION

*Youshan Zhang*     *Brian D. Davison*

Computer Science and Engineering, Lehigh University, Bethlehem, PA, USA
{yoz217, bdd3}@lehigh.edu

## ABSTRACT

Domain adaptation aims to mitigate the domain gap when transferring knowledge from an existing labeled domain to a new domain. However, existing disentanglement-based methods do not fully consider separation between domain-invariant and domain-specific features, which means the domain-invariant features are not discriminative. The reconstructed features are also not sufficiently used during training. In this paper, we propose a novel enhanced separable disentanglement (ESD) model. We first employ a disentangler to distill domain-invariant and domain-specific features. Then, we apply feature separation enhancement processes to minimize contamination between domain-invariant and domain-specific features. Finally, our model reconstructs complete feature vectors, which are used for further disentanglement during the training phase. Extensive experiments from three benchmark datasets outperform state-of-the-art methods, especially on challenging cross-domain tasks.

***Index Terms***— Unsupervised domain adaptation, Disentanglement, Domain discriminator

## 1. INTRODUCTION

Most existing machine learning models rely on large amounts of labeled training data to achieve high performance. Unfortunately, such a requirement cannot be met in many real-world applications. The number of labels is limited and manual annotation is expensive and time-consuming. Therefore, it is valuable to learn a model for a new domain from one with existing labeled samples. However, the difference between the two domains, termed the domain shift, can cause difficulty with direct use of models from one domain on another. Unsupervised domain adaptation (UDA) has emerged as a prominent method to address the domain shift.

Deep learning models have been widely used in UDA. Earlier methods rely on minimizing the discrepancy between the source and target distributions by proposing different loss functions, such as Maximum Mean Discrepancy (MMD) [1], CORrelation ALignment [2], Kullback-Leibler divergence [3]. To learn domain invariant features, adversarial domain adaptation methods aim to identify domain in-variant features by playing a min-max game between domain discriminator and feature extractor [4, 5]. However, most of these UDA methods do not consider constructing separable domain-specific features ($f_{ds}$) and domain-invariant features ($f_{di}$) to learn more discriminative representations.

Recently, disentanglement representation based methods can learn discriminative $f_{ds}$, which contains domain related information, and $f_{di}$, which contains intrinsic information related to different categories. Cross domain representation disentangler [6] bridged labeled source domain and unlabeled target domain via jointly learning cross-domain feature representation disentanglement and adaptation. Gonzalez et al. [7] proposed an image-to-image translation for representation disentangling based on GANs and cross-domain autoencoders. They separated the internal representation into three parts: shared part, which contains the domain invariant features for two domains, and two exclusive parts, which contain the domain-specific features. Peng et al. [8] minimized mutual information between $f_{ds}$ and $f_{di}$ to pursue implicit domain invariant features. Liu et al. [9] introduced a unified feature disentanglement network to learn a domain-invariant representation from multiple domains for image translation and manipulation. Gholami et al. [10] presented a multi-target domain adaptation information theoretic approach to find a shared latent space of all domains by simultaneously identifying the remaining private, domain-specific factors. However, these methods either cannot fully separate $f_{ds}$ and $f_{di}$ [6, 7, 9] or the reconstructed features are insufficiently used during training to facilitate the performance of a domain discriminator and a domain-invariant classifier [8, 10].

To alleviate these issues, we propose an enhanced separable disentanglement (ESD) model. Our contributions are:

- We propose a novel method for feature disentanglement representation learning: 1) teach a disentangler to distill domain-specific from domain-invariant features; 2) apply feature separation maximization processes to enhance the disentangler and improve the effectiveness of both kinds of features; 3) design a reconstructor to recover original feature prototypes which can be further re-utilized in steps 1) and 2) to improve performance.
- For feature separation maximization, we first propose a novel structural similarity loss to maximize the dissim-
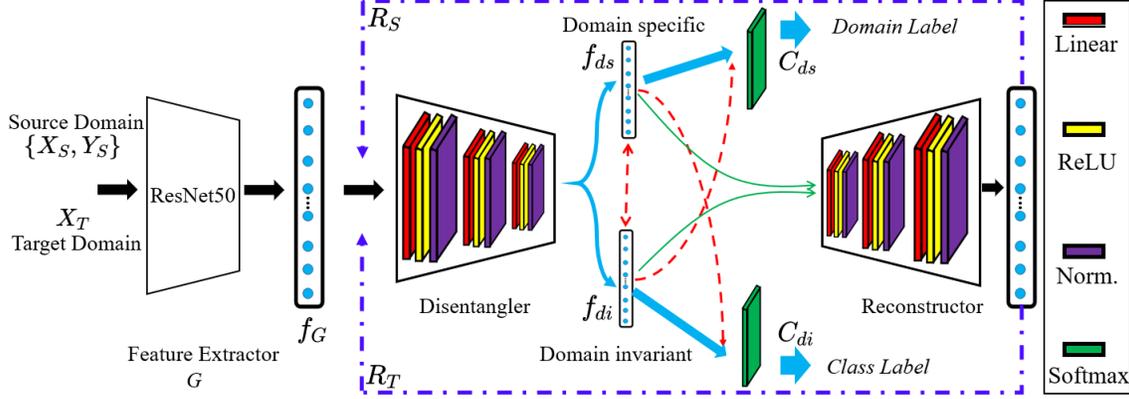
**Fig. 1**: An overview of ESD. We first employ ResNet50 as feature extractor $G$ to extract features of two domains. We then employ 1) a disentangler to distill domain-specific $f_{ds}$ and domain-invariant $f_{di}$ features and then train domain-specific classifier $C_{ds}$ and invariant classifier $C_{di}$, respectively (blue arrows); 2) feature separation processes to enhance the disentangler and minimize contamination between $f_{ds}$ and $f_{di}$ (red lines); 3) a reconstructor to recover original feature prototypes (green lines) and finally re-utilize the reconstructed features ($R_S$ and $R_T$) to improve $C_{ds}$ and $C_{di}$ (purple lines). Norm.: normalization.

ilarity between $f_{ds}$ and $f_{di}$, and then propose opposite binary cross-entropy loss and accurate loss to further remove contaminated information.
- The reconstructor recovers original feature prototypes using reconstructed features which are then used to update the learned classifier.

Extensive experiments on three benchmark datasets show that our ESD model outperforms state-of-the-art methods.

## 2. METHODS

**Problem and notation.** For UDA, given a source domain $\mathcal{D}_{\mathcal{S}} = \{X_S^i, Y_S^i\}_{i=1}^{n_s}$ of $n_s$ labeled samples in $K$ categories and a target domain $\mathcal{D}_{\mathcal{T}} = \{X_T^j\}_{j=1}^{n_t}$ without any labels ($Y_T$ for evaluation only), our ultimate goal is to learn a classifier under a feature extractor $G$, that produces lower generalization error in the target domain.

**Disentangler.** As shown in Fig. 1, feature extractor $G$ maps a labeled source domain and an unlabeled target domain into latent feature prototypes $f_G$. For a classical disentangler (D), it disentangles $f_G$ into domain-specific features $f_{ds}$, and domain-invariant features $f_{di}$. In the first step, for $f_{di}$, we aim to train a domain invariant classifier $C_{di}$ using the labeled source invariant features $f_{di}^S$ and make predictions for the target invariant features $f_{di}^T$ with typical cross-entropy loss as follows:

$$\mathcal{L}_{di} = \mathcal{L}_{\mathcal{C}}(C_{di}(f_{di}^S), Y_S), \quad (1)$$

where $\mathcal{L}_{\mathcal{C}}$ is cross-entropy loss. For $f_{ds}$, we aim to learn a domain discriminator classifier $C_{ds}$ using adversarial loss to distinguish the source domain-specific features $f_{ds}^S$ and the target domain-specific features $f_{ds}^T$ in the following equation.

$$\mathcal{L}_{ds} = \mathbb{E}_{x_s^i \sim f_{ds}^S} \log[C_{ds}(x_s^i)] + \mathbb{E}_{x_t^j \sim f_{ds}^T} \log[1 - C_{ds}(x_t^j)], \quad (2)$$

where domain labels of $f_{ds}^S$ and $f_{ds}^T$ are set as 1 and 0. However, a traditional disentangler cannot ensure that there is no contaminated information between $f_{di}$ and $f_{ds}$. We employ a feature separation enhancement step to alleviate this issue.

**Feature Separation Enhancement.** In step two, to guarantee the completed feature separation between $f_{di}$ and $f_{ds}$, we also employ three processes. First of all, we maximize the dissimilarity between $f_{di}$ and $f_{ds}$. Peng et al. [8] proposes to minimize the mutual information between them. However, there is no closed solution for minimizing mutual information, and the Monte Carlo sampling will lead to external computation. Instead, we directly maximize the dissimilarity between the disentangled features, which is equivalent to minimizing the similarity between $f_{di}$ and $f_{ds}$. We impose a batch-wise structural similarity loss.

$$\mathcal{L}_S = abs(\frac{(2\mu_{B_1}\mu_{B_2} + C_1)(2\sigma_{B_1 B_2} + C_2)}{(\mu_{B_1}^2 + \mu_{B_2}^2 + C_1)(\sigma_{B_1}^2 + \sigma_{B_2}^2 + C_2)}), \quad (3)$$

where $abs$ takes the absolute value, $B_1 \in f_{di}$ and $B_2 \in f_{ds}$ are batch-wise features, $\mu_{B_1}, \mu_{B_2}, \sigma_{B_1}, \sigma_{B_2}$, and $\sigma_{B_1 B_2}$ are mean, standard deviations of domain invariant and specific features batch, and cross-covariance for $(B_1, B_2)$. $C_1$ and $C_2$ are two variables to stabilize the division with weak denominator. This loss function is derived from structural similarity index measure (SSIM) [11]. It has the advantages of measuring luminance, contrast and structural difference between $B_1$ and $B_2$. Therefore, $\mathcal{L}_S$ has more capability of measuring the similarity between $f_{di}$ and $f_{ds}$. In addition, the range of the $\mathcal{L}_S$ is from 0 to 1, where 1 indicates high similarity between batch features, and 0 means they are not similar. During the training, we keep minimizing such a similarity, and thus contamination between $f_{di}$ and $f_{ds}$ is also minimized.

Secondly, to further ensure domain-specific features are fully segregated from domain-invariant features, we leverage

$f_{di}$ to fool the trained domain-specific classifier $C_{ds}$ using opposite binary cross-entropy loss in Eq. 4.

$$\mathcal{L}_O = \mathbb{E}_{x_s^i \sim f_{di}^S} \log[1 - C_{ds}(x_s^i)] + \mathbb{E}_{x_t^j \sim f_{di}^T} \log[C_{ds}(x_t^j)] \quad (4)$$

Differing from Eq. 2, the labels of $f_{di}^S$ and $f_{di}^T$ are set as 0 and 1, oppositely. Since we want the trained $C_{ds}$ to be unable to predict the correct domain labels of $f_{di}$, we effectively force $f_{ds}$ to contain no $f_{di}$ information.

Thirdly, to remove $f_{ds}$ from $f_{di}$, we design an accurate loss to force no correct labels can be predicted using trained domain invariant classifier $C_{di}$ as follows:

$$\mathcal{L}_A = \frac{|C_{di}(f_{ds}^S) \wedge C_{di}(f_{di}^S)|}{n_s} + \frac{|C_{di}(f_{ds}^T) \wedge C_{di}(f_{di}^T)|}{n_t}, \quad (5)$$

where $\wedge$ is the and operation, and $|\cdot|$ measures the length. The numerator measures the number of the same predictions of source and target $f_{ds}$ and $f_{di}$, respectively. Minimizing such an accurate loss, we can make sure that the predictions of $f_{ds}$ are different from $f_{di}$. Hence, we can force $f_{di}$ contains no $f_{ds}$ information.

Considering the above three processes, we can enhance the disentangler to minimize the contamination between $f_{di}$ and $f_{ds}$.

**Reconstructor.** In step three, to keep the information integrity of feature disentanglement, we train a reconstructor ($R$) to recover original feature prototypes ($f_G$) using the disentangled $f_{di}$ and $f_{ds}$. Let $R_S = R(f_{di}^S, f_{ds}^S)$ and $R_T = R(f_{di}^T, f_{ds}^T)$, the reconstruction loss is defined as:

$$\mathcal{L}_R = ||R_S - f_G^S||_2^2 + ||R_T - f_G^T||_2^2. \quad (6)$$

By this reconstructor R, we can guarantee that $R_S \sim f_G^S$ and $R_T \sim f_G^T$. Differing from previous work, we take advantage of reconstructed features to further improve the performance of $C_{di}$ and $C_{ds}$. As shown in purple lines in Fig. 1, we re-utilize the reconstructed features in the disentangler. At this stage, all parameters in step one and two are fixed, *i.e.*, $D, C_{di}$ and $C_{ds}$ are fixed. We directly apply the trained disentangler D to get the disentangled features $\hat{f}_{di}$ and $\hat{f}_{ds}$ using the reconstructed features $R(f_G)$. Therefore, we will repeat the minimization equations in step one and step two in Eq. 7.

$$\mathcal{L}_T = \hat{\mathcal{L}}_{di} + \hat{\mathcal{L}}_{ds} + \hat{\mathcal{L}}_S + \hat{\mathcal{L}}_O + \hat{\mathcal{L}}_A, \quad (7)$$

where in Eq. 1 to Eq. 5, we will directly compute the loss by replacing $f_{di}$ and $f_{ds}$ with $\hat{f}_{di}$ and $\hat{f}_{ds}$, and replacing $f_{di}^S, f_{di}^T, f_{ds}^S, f_{ds}^T$ with $\hat{f}_{di}^S, \hat{f}_{di}^T, \hat{f}_{ds}^S, \hat{f}_{ds}^T$, while all other components are the same as before.

**The overall training objective function.** The architecture of our proposed ESD model is shown in Fig. 1. Our model minimizes the following objective function.

$$\mathcal{L}(X_S, Y_S, X_T) = \arg\min((\mathcal{L}_{di} + \mathcal{L}_{ds}) + \alpha(\mathcal{L}_S + \mathcal{L}_O + \mathcal{L}_A) + \beta(\mathcal{L}_R + \mathcal{L}_T)) \quad (8)$$

where $(\cdot)$ represents loss functions in each step, $\alpha$ and $\beta$ are factors to balance the importance of steps two and three.

# 3. EXPERIMENTS

**Datasets.** We test our model using three public image datasets: Office-31, Office-Home and VisDA-2017. **Office-31** [12] has 4,110 images from three domains: Amazon (A), Webcam (W), and DSLR (D) in 31 classes. In experiments, A→W represents transferring knowledge from domain A to domain W. **Office-Home** [13] dataset contains 15,588 images from four domains: Art (Ar), Clipart (Cl), Product (Pr) and Real-World (Rw) in 65 classes. **VisDA-2017** [14] is a challenging dataset due to the big domain-shift between the synthetic images (152,397 images from VisDA) and the real images (55,388 images from COCO) in 12 classes. We evaluate our method on the setting of synthetic-to-real as the source-to-target domain and report accuracy of each category.

**Implementation details.** We implement our approach using PyTorch and extract features for the three datasets from finely tuned ResNet50 (Office-31, Office-Home) and ResNet101 (VisDA-2017) networks [15]. The 1,000 features are then extracted from the last fully connected layer for the source and target features. In the Disentangler $D$, the number of outputs of the first two Linear layers are 1000 and 512, and the output of the last Linear layer is the number of classes ($K$) in each dataset, while the output in the reconstruction layers is opposite ($K$, 512, and 1000). $C_1 = 0.01^2$ and $C_2 = 0.03^2$ as in [11]. The learning rate = 0.001, batch size = 32 and number of iterations = 100 with SGD optimizer. The balanced factors $\alpha = 0.3$ and $\beta = 0.1$.

**Comparison to state-of-the-art methods.** We compare the performance of our ESD model with 15 state-of-the-art methods. For a fair comparison, baseline results are directly reported from their original papers. From Tables 1 to 3, we can observe that the accuracy of the ESD model is ahead of all other methods in most tasks. For the Office-31 dataset, the average accuracy of ESD is 91.4%. It is superior to all other methods. ESD shows substantially better transferring ability than other methods in tasks W→A and D→A. In another difficult Office-Home dataset, which has more categories, more samples, and larger domain discrepancy than the Office-31 datasets, the average accuracy is 71.6%, which exceeds the SOTA methods. Our model has a particularly obvious improvement in the challenging VisDA-2017 dataset, which has many more samples and larger domain discrepancy than the other two datasets. Our model achieves the highest accuracy 89.2%, which is two percent higher than the best baseline CAN. By testing on three distinct datasets, we give evidence of ESD's broad applicability.

# 4. DISCUSSION

In all experiments, our method achieves the highest average accuracy. There are two prominent reasons for the success of our model. First of all, the proposed three novel loss functions ($\mathcal{L}_S, \mathcal{L}_O$ and $\mathcal{L}_A$) in the feature separation enhancement step are important to remove contaminated information be-

**Table 1**: Accuracy (%) on Office-Home dataset (based on ResNet50)

| Task | Ar→Cl | Ar→Pr | Ar→Rw | Cl→Ar | Cl→Pr | Cl→Rw | Pr→Ar | Pr→Cl | Pr→Rw | Rw→Ar | Rw→Cl | Rw→Pr | **Ave.** |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ResNet-50 [16] | 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 |
| DAN [17] | 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 [4] | 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 [18] | 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 |
| TAT [19] | 51.6 | 69.5 | 75.4 | 59.4 | 69.5 | 68.6 | 59.5 | 50.5 | 76.8 | 70.9 | 56.6 | 81.6 | 65.8 |
| ETD [20] | 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 |
| SymNets [21] | 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 |
| DCAN [22] | **54.5** | 75.7 | 81.2 | 67.4 | 74.0 | 76.3 | 67.4 | 52.7 | 80.6 | 74.1 | **59.1** | 83.5 | 70.5 |
| **ESD** | 53.2 | **75.9** | **82.0** | **68.4** | **79.3** | **79.4** | **69.2** | **54.8** | **81.9** | 74.6 | 56.2 | **83.8** | **71.6** |

**Table 2**: Accuracy (%) on VisDA-2017 dataset (based on ResNet101)

| Task | plane | bcycl | bus | car | horse | knife | mcycl | person | plant | sktbrd | train | truck | **Ave.** |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Source-only [16] | 55.1 | 53.3 | 61.9 | 59.1 | 80.6 | 17.9 | 79.7 | 31.2 | 81.0 | 26.5 | 73.5 | 8.5 | 52.4 |
| DANN [4] | 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 [17] | 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 |
| JAN [18] | 75.7 | 18.7 | 82.3 | 86.3 | 70.2 | 56.9 | 80.5 | 53.8 | 92.5 | 32.2 | 84.5 | 54.5 | 65.7 |
| MCD [23] | 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 |
| DADA [24] | 92.9 | 74.2 | 82.5 | 65.0 | 90.9 | 93.8 | 87.2 | 74.2 | 89.9 | 71.5 | 86.5 | 48.7 | 79.8 |
| STAR [25] | 95.0 | 84.0 | 84.6 | 73.0 | 91.6 | 91.8 | 85.9 | 78.4 | 94.4 | 84.7 | 87.0 | 42.2 | 82.7 |
| CAN [26] | **97.9** | 87.2 | 82.5 | 74.3 | **97.8** | 96.2 | 90.8 | 80.7 | 96.6 | 96.3 | 87.5 | 59.9 | 87.2 |
| **ESD** | 96.8 | **89.1** | **87.9** | **80.3** | 96.7 | **96.9** | **92.5** | **84.9** | **96.9** | **97.5** | **88.9** | **62.8** | **89.2** |

**Table 3**: Accuracy (%) on Office-31 (based on ResNet50)

| Task | A→W | A→D | W→A | W→D | D→A | D→W | **Ave.** |
|---|---|---|---|---|---|---|---|
| RTN [27] | 84.5 | 77.5 | 64.8 | 99.4 | 66.2 | 96.8 | 81.6 |
| ADDA [5] | 86.2 | 77.8 | 68.9 | 98.4 | 69.5 | 96.2 | 82.9 |
| JAN [18] | 85.4 | 84.7 | 70.0 | 99.8 | 68.6 | 97.4 | 84.3 |
| TAT [19] | 92.5 | 93.2 | 73.1 | **100** | 73.1 | **99.3** | 88.4 |
| TADA [28] | 94.3 | 91.6 | 73.0 | 99.8 | 72.9 | 98.7 | 88.4 |
| SymNets [21] | 90.8 | 93.9 | 72.5 | **100** | 74.6 | 98.8 | 88.4 |
| CAN [26] | 94.5 | 95.0 | 77.0 | 99.8 | 78.0 | 99.1 | 90.6 |
| **ESD** | **94.6** | **96.2** | **80.1** | 99.0 | **80.4** | 98.0 | **91.4** |



**Fig. 2**: Visualization of learned features using a 2D t-SNE view of task A→D in Office-31 dataset.

tween domain-specific and domain-invariant features, which improves the performance of $C_{di}$. Secondly, re-utilizing reconstructed features is also helpful in optimizing $C_{di}$ and $C_{ds}$ and leads to higher accuracy. We also observe that our model is suboptimal in some tasks (Rw→Cl in Office-Home and D→W in Office-31 dataset), so we cannot guarantee that our model always beats all other methods.

To show how different steps affect final performance, we also conduct an ablation study in Tab. 4 (step one is required). We find step two is more important than step three in improving final accuracy (90.8% v.s. 89.8%), but both contribute to final performance. To intuitively present adaptation performance and the effects of step two, we utilize t-SNE [29] to visualize the deep features of network activations in 2D space before and after adaptation (domain invariant/specific

**Table 4**: Ablation tests on Office-31 (minus steps II and III).

| Task | A→W | A→D | W→A | W→D | D→A | D→W | **Ave.** |
|---|---|---|---|---|---|---|---|
| ESD−II−III | 91.2 | 87.6 | 76.3 | 98.2 | 76.2 | 96.6 | 87.7 |
| ESD−II | 92.1 | 93.9 | 78.5 | 98.6 | 79.1 | 96.6 | 89.8 |
| ESD−III | 94.3 | 95.3 | 79.8 | 98.8 | 79.3 | 97.2 | 90.8 |
| **ESD** | **94.6** | **96.2** | **80.1** | **99.0** | **80.4** | **98.0** | **91.4** |

features) as shown in Fig. 2. Apparently, the distributions of domains A and D become more discriminative after adaptation ($f_{di}$), while many categories are mixed in the feature space before adaptation. In addition, the $f_{ds}$ can also distinguish the target domain from the source domain. Furthermore, the distributions of $f_{di}$ and $f_{ds}$ are different, which implies that contamination between them is minimized. This result indicates that ESD can learn more discriminative representations.

## 5. CONCLUSION

We have presented a novel enhanced separable disentanglement model to improve the separability between domain-specific and domain-invariant features. Specifically, we impose structural similarity loss, opposite binary cross-entropy loss, and accurate loss functions to minimize contamination among disentangled features. We further re-utilize the reconstructed features to improve the performance of domain discrimination and domain-invariant classification. Extensive experiments demonstrate that our ESD model outperforms state-of-the-art methods.

# 6. REFERENCES

[1] E. Tzeng, J. Hoffman, N. Zhang, K. Saenko, and T. Darrell, "Deep domain confusion: Maximizing for domain invariance," *arXiv preprint arXiv:1412.3474*, 2014.

[2] B. Sun and K. Saenko, "Deep CORAL: Correlation alignment for deep domain adaptation," in *European Conference on Computer Vision*. Springer, 2016, pp. 443–450.

[3] Z. Meng, J. Li, Y. Gong, and B. Juang, "Adversarial teacher-student learning for unsupervised domain adaptation," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 5949–5953.

[4] M. Ghifary, W. B. Kleijn, and M. Zhang, "Domain adaptive neural networks for object recognition," in *Pacific Rim International Conference on Artificial Intelligence*. Springer, 2014, pp. 898–904.

[5] E. Tzeng, J. Hoffman, K. Saenko, and T. Darrell, "Adversarial discriminative domain adaptation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 7167–7176.

[6] Y. Liu, Y. Yeh, T. Fu, S. Wang, W. Chiu, and Y. Frank Wang, "Detach and adapt: Learning cross-domain disentangled deep representation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 8867–8876.

[7] A. Gonzalez-Garcia, J. Van De Weijer, and Y. Bengio, "Image-to-image translation for cross-domain disentanglement," *Advances in Neural Information Processing Systems*, vol. 31, pp. 1287–1298, 2018.

[8] X. Peng, Z. Huang, X. Sun, and K. Saenko, "Domain agnostic learning with disentangled representations," in *International Conference on Machine Learning*, 2019, pp. 5102–5112.

[9] A. H. Liu, Y. Liu, Y. Yeh, and Y. F. Wang, "A unified feature disentangler for multi-domain image translation and manipulation," in *Advances in Neural Information Processing Systems*, 2018, pp. 2590–2599.

[10] B. Gholami, P. Sahu, O. Rudovic, K. Bousmalis, and V. Pavlovic, "Unsupervised multi-target domain adaptation: An information theoretic approach," *IEEE Transactions on Image Processing*, vol. 29, pp. 3993–4002, 2020.

[11] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE Transactions on Image Processing*, vol. 13, no. 4, pp. 600–612, 2004.

[12] K. Saenko, B. Kulis, M. Fritz, and T. Darrell, "Adapting visual category models to new domains," in *European Conference on Computer Vision*. Springer, 2010, pp. 213–226.

[13] H. Venkateswara, J. Eusebio, S. Chakraborty, and S. Panchanathan, "Deep hashing network for unsupervised domain adaptation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 5018–5027.

[14] X. Peng, B. Usman, N. Kaushik, J. Hoffman, D. Wang, and K. Saenko, "VisDA: The visual domain adaptation challenge," *arXiv preprint arXiv:1710.06924*, 2017.

[15] Y. Zhang and B. D. Davison, "Adversarial continuous learning in unsupervised domain adaptation.," in *ICPR Workshops (2)*, 2020, pp. 672–687.

[16] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.

[17] M. Long, Y. Cao, J. Wang, and M. I. Jordan, "Learning transferable features with deep adaptation networks," *arXiv preprint arXiv:1502.02791*, 2015.

[18] M. Long, H. Zhu, J. Wang, and M. I. Jordan, "Deep transfer learning with joint adaptation networks," in *Proceedings of the 34th International Conference on Machine Learning*. JMLR.org, 2017, vol. 70, pp. 2208–2217.

[19] M. Liu, H.and Long, J. Wang, and M. Jordan, "Transferable adversarial training: A general approach to adapting deep classifiers," in *International Conference on Machine Learning*, 2019, pp. 4013–4022.

[20] M. Li, Y. Zhai, Y. Luo, P. Ge, and C. Ren, "Enhanced transport distance for unsupervised domain adaptation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 13936–13944.

[21] Y. Zhang, H. Tang, K. Jia, and M. Tan, "Domain-symmetric networks for adversarial domain adaptation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5031–5040.

[22] S. Li, C. H. Liu, Q. Lin, B. Xie, Z. Ding, G. Huang, and J. Tang, "Domain conditioned adaptation network.," in *Proceedings of the AAAI Conference on Artificial Intelligence,*, 2020, pp. 11386–11393.

[23] K. Saito, K. Watanabe, Y. Ushiku, and T. Harada, "Maximum classifier discrepancy for unsupervised domain adaptation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3723–3732.

[24] H. Tang and K. Jia, "Discriminative adversarial domain adaptation," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020, vol. 34, pp. 5940–5947.

[25] Z. Lu, Y. Yang, X. Zhu, C. Liu, Y. Song, and T. Xiang, "Stochastic classifiers for unsupervised domain adaptation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 9111–9120.

[26] G. Kang, L. Jiang, Y. Yang, and A. G. Hauptmann, "Contrastive adaptation network for unsupervised domain adaptation," in *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4893–4902.

[27] M. Long, H. Zhu, J. Wang, and M. I. Jordan, "Unsupervised domain adaptation with residual transfer networks," in *Advances in Neural Information Processing Systems*, 2016, pp. 136–144.

[28] X. Wang, L. Li, W. Ye, M. Long, and J. Wang, "Transferable attention for domain adaptation," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2019, vol. 33, pp. 5345–5352.

[29] L. Van der Maaten and G. Hinton, "Visualizing data using t-SNE," *Journal of Machine Learning Research*, vol. 9, no. 11, 2008.