

Supplementary File

Toward Improving Robustness and Accuracy in Unsupervised Domain Adaption

Anonymous submission

Authors' Response

Hyperparameter Analysis

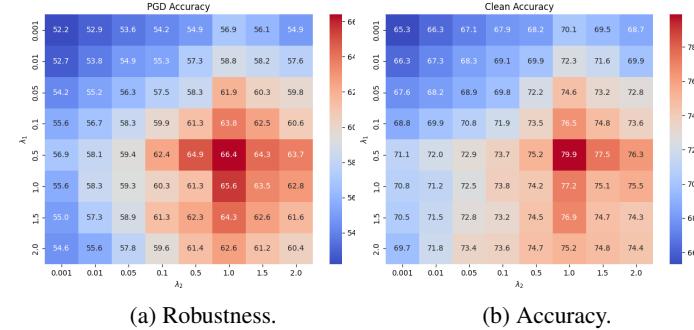
There are two main categories of hyper-parameters in our method: the loss weights λ_1 , λ_2 , and the confidence threshold (T). We evaluate the impact of loss weights (λ_1 , λ_2) on robustness (under a PGD20 attack with $\epsilon = 8/255$) and accuracy for the VISDA (Real→Synthetic) dataset. As shown in Figure 1, setting $\lambda_1 = 0.5$ and $\lambda_2 = 1.0$ achieves a balanced trade-off, improving both accuracy and robustness across clean and adversarial examples. Additionally, we examine the effect of varying confidence threshold values for pseudo-label incorporation during training on the VISDA (Real→Synthetic) and PACS (Photo→Clipart) datasets. Figure 2 indicates that both robustness and accuracy remain stable when the confidence threshold is set between 0.5 and 0.8. Setting the threshold too low (e.g., 0.1) degrades overall performance, while a higher threshold restricts the amount of data incorporated, which can also limit performance gains

Black-box Robustness

Table 1 illustrate the robustness of different methods against black-box attacks on the VISDA and PACS (Photo→Sketch) datasets. Using a naturally trained DANN model with ResNet-18 as the substitute, we generate black-box adversarial examples for target models using Momentum Iterative Fast Gradient Sign Method (Dong et al. 2018) (MIFGSM), a highly transferable attack. This attack is configured with $\epsilon = 8/255$ and 10 iterations to rigorously test the various methods resilience. The naturally trained DANN method shows the lowest resistance to black-box attacks compared to other training methods. Our CAM + SPLR approach not only achieves superior clean accuracy but also consistently outperforms others in black-box robustness across both datasets.

Visual analysis using ScoreCam and AblationCam

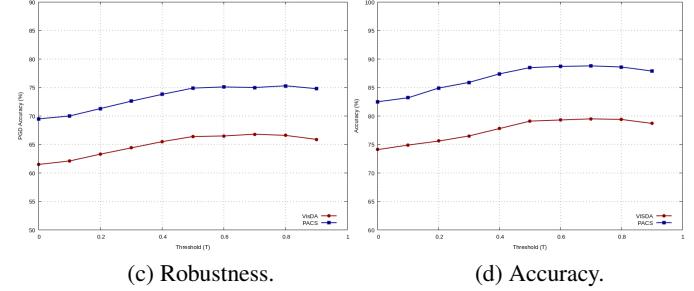
To evaluate the effectiveness of the proposed method (CAM+SPLR), we employ two visualization techniques namely **AblationCAM** and **ScoreCAM** in Figure 3. These visualization methods allow us to gain insights into the model's interpretability by highlighting image regions that most influence the model's predictions, thus offering a window into its decision-making process. In both adversar-



(a) Robustness.

(b) Accuracy.

Figure 1: Parameter sensitivity analysis on λ_1 and λ_2 on robustness and accuracy for VISDA(Real→Synthetic) dataset.



(c) Robustness.

(d) Accuracy.

Figure 2: Effect of varying confidence thresholds on robustness and accuracy for VISDA (Real→Synthetic) and PACS (Photo→Clipart) datasets.

ial and clean visual ablations, we observe that both images activate the same regions, whether using AblationCAM or ScoreCAM. AblationCAM (Ramaswamy et al. 2020) highlights key regions by selectively removing parts of the model's layers and measuring the impact on class scores. It identifies important areas in the input image by evaluating how much the prediction score drops when certain features are masked. While ScoreCAM (Wang et al. 2020) generates activation maps by applying perturbations to different regions of the input image and measuring their contribution to the predicted class score.

Visual Comparison

Figure 4 illustrates a visual ablation study highlighting the effectiveness of our method in capturing semantically relevant regions across both clean and adversarial images through consistent attention mapping. The Score-CAM

Table 1: An illustration of black-box robustness comparison across the VisDA and PACS (Photo→Sketch) datasets.

Source → Target	Syn → Real		Real → Syn		Ph → Sk	
	Clean	MI-FGSM	Clean	MI-FGSM	Clean	MI-FGSM
DANN(Ganin et al. 2016)	67.5	52.9	78.5	65.9	74.5	60.5
UDA+AT	49.6	44.5	52.8	48.1	69.5	64.8
Trades(Zhang et al. 2019)	51.7	48.3	57.4	53.3	72.3	68.2
Mart(Wang et al. 2019)	50.8	46.9	56.0	53.5	71.1	66.8
ARTUDA(Yang et al. 2021a)	54.9	52.1	59.5	56.6	67.2	63.8
SRoUDA(Zhu et al. 2023)	51.3	47.6	61.6	58.3	66.5	63.7
DART(Wang et al. 2024)	65.6	62.4	73.8	69.8	79.7	76.1
Ours	69.5	67.3	79.1	76.9	83.5	80.6

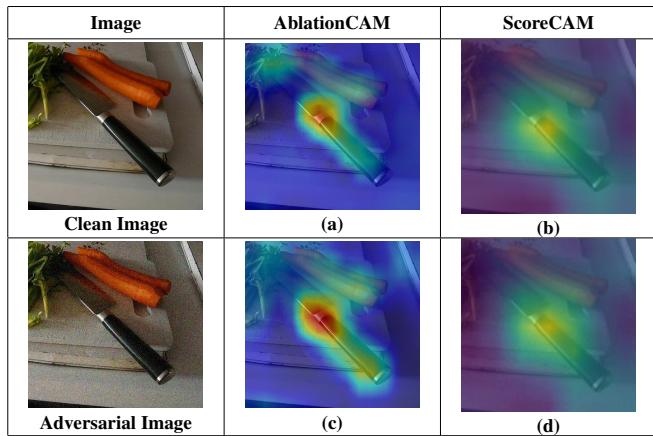


Figure 3: Effectiveness of CAM+SPLR in visualizing key regions for clean and adversarial images using AblationCAM and ScoreCAM. The first column presents the clean image and its adversarial counterpart. The second and third columns show the respective AblationCAM and ScoreCAM maps, highlighting the critical areas influencing the model’s predictions. Notably, the method consistently predicts the “Knife” label accurately across both clean and adversarial image.

maps produced by our approach demonstrate a focused attention on meaningful areas integral to accurate classification. In contrast to Trades (Zhang et al. 2019), our approach preserves semantic consistency across clean and adversarial inputs by aligning attention maps, ensuring the model consistently attends to the same regions of interest despite perturbations. This is particularly apparent in the similarity of the Score-CAM maps for both clean and adversarial cases in our method, indicating a strong alignment in feature representations. This alignment reduces the model’s vulnerability to adversarial distractions, keeping the attention on essential object regions even under adversarial influence.

Performance Comparison

Table 2 demonstrates a robustness comparison between our method and TRADES under both white-box and black-box attacks on the VISDA (Re→Syn) and PACS (Photo→Clipart) datasets. The results clearly show that our method consistently outperforms TRADES across both attack scenarios. It is due to its emphasis on aligning key se-

mantic regions across clean and adversarial inputs. By incorporating consistent attention mapping with self-pseudo label refinement (SPLR), our approach not only enhances focus on key semantic regions but also continuously refines pseudo labels. This dual strategy promotes stable attention across layers and improves overall robustness and accuracy in unsupervised domain adaptation.

Table 2: Comparison of robustness between our method and TRADES across different adversarial attacks (FGSM, PGD10, PGD20, and Black-box (MI-FGSM)) using the VISDA and PACS datasets.

Dataset	Method	Clean	FGSM	PGD10	PGD20	MI-FGSM
VISDA (Re→Syn)	DANN	78.5	19.5	1.2	0.5	65.9
	Trades	57.4	52.3	47.8	45.6	53.3
	Our	79.1	74.7	70.8	66.4	76.9
PACS (Ph→Cl)	DANN	80.2	21.5	2.1	0.7	69.9
	Trades	80.9	75.4	72.4	66.5	76.3
	Our	88.5	84.7	80.5	74.9	85.1

Supplementary

Contributions

- We introduce a novel self-training method called Consistent Attention Mapping with Self Pseudo Label Refinement (CAM+SPLR), designed to enhance both the robustness and accuracy of UDA models. This method leverages adversarial target data generated from pseudo labels to strengthen robustness through a self-training paradigm in two step gradient descent process. Simultaneously, it encourages consistency between the attention maps of clean examples and their adversarial counterparts, while progressively refining the pseudo labels.
- We propose the Consistent Attention Mapping (CAM) method to prevent the model from concentrating on less informative regions that may be influenced by adversarial perturbations or noisy pseudo-labels. During training of the TargetNet model, CAM ensures that attention maps remain consistent between clean target data processed by the frozen Anchor model and their adversarial counterparts processed by the TargetNet model. By focusing on semantically relevant key areas, CAM enhances the learning of more discriminative features.
- We further introduce the Self Pseudo Label Refinement (SPLR) method to prevent the model from overfitting caused by inevitable noisy pseudo labels. To achieve this, we progressively refine the pseudo labels during the training of the TargetNet model by incorporating feedback from the labeled source data . This refinement occurs during the second step of the gradient descent process, ensuring that the model remains accurate and resilient as it updates.
- We achieve improvement in robustness and gain in standard accuracy across multiple datasets (Office-Home (Wang et al. 2021), PACS (Li et al. 2017), VisDA (Peng et al. 2017)) compared to state-of-the-art methods (DART (Wang et al. 2024), SRoUDA (Zhu et al. 2023), and ARTUDA (Yang et al. 2021a)). Specifically, we observed remarkable average robustness gain at $\epsilon =$

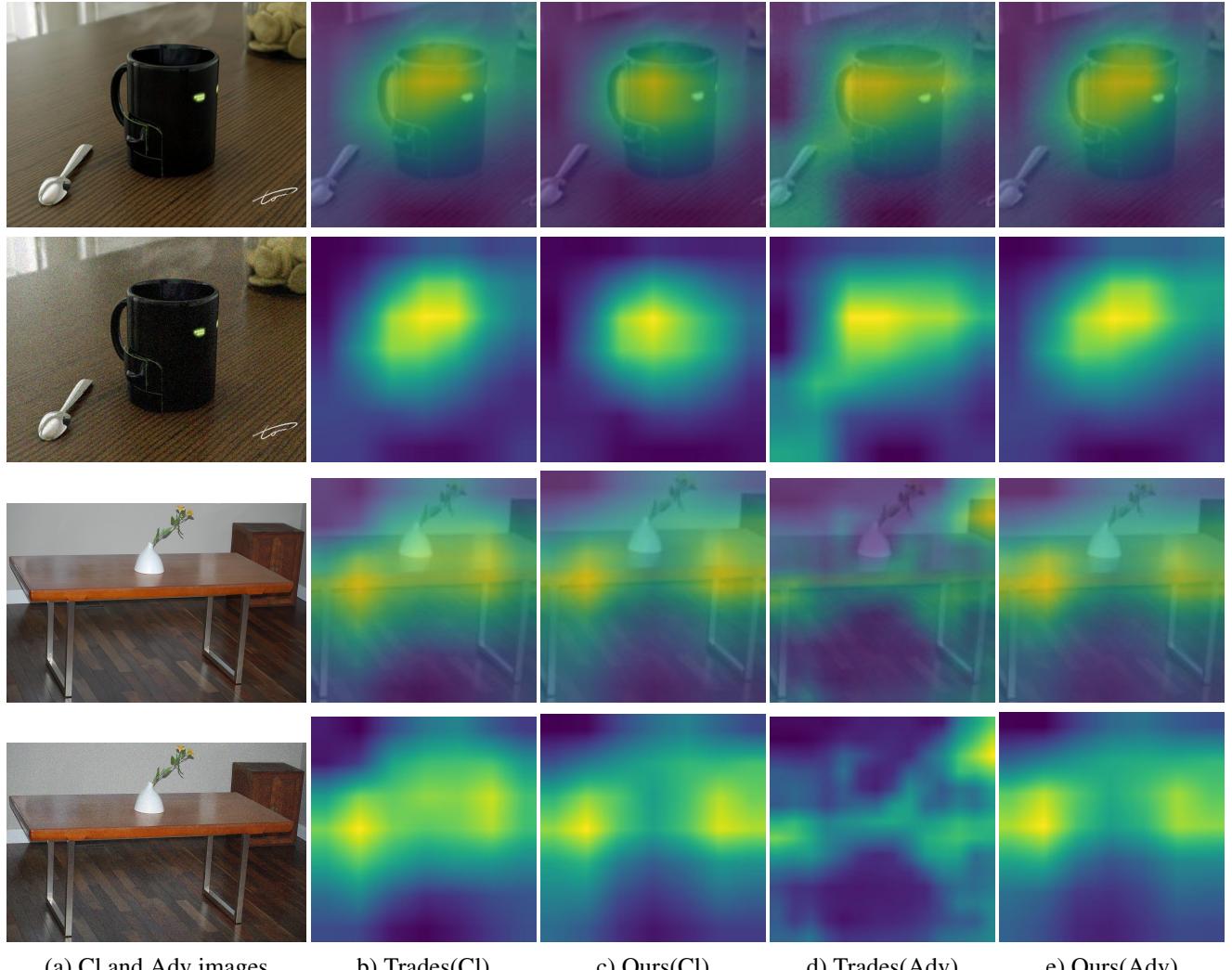


Figure 4: Visual ablation comparing our method with TRADES in identifying key regions for clean and adversarial images using Score-CAM. The first column displays the original clean image and its adversarial counterpart in each row. Subsequent columns present the Score-CAM maps for both our method and TRADES, under both clean and adversarial conditions, showcasing the regions that influence the model’s predictions. Cl and Adv represent Clean and Adversarial images, respectively.

2/255 of 5.2%, 4.9%, and 10.2% on the OfficeHome, VisDA, and PACS datasets, respectively. Additionally, average accuracy improved by 0.9%, 8.1%, and 6% over the UDA baseline DANN (Ganin et al. 2016) on the OfficeHome, VisDA, and PACS datasets, respectively.

CAM+SPLR algorithm

The training steps of the proposed method, as defined in Algorithm 1.

Related Work

Unsupervised Domain Adaption

Unsupervised domain adaption is a scenario in which the rich knowledge of the labelled source domain is transferred to the unlabelled target domain to perform different machine learning tasks. There exists various approaches (Ganin et al. 2016; Long et al. 2018, 2017; Saito et al. 2018; Sun and Saenko 2016; Xu et al. 2019; Choi et al. 2019; Yang et al. 2021b), which minimize the distributional differences between the source domain and target domain to perform UDA. DANN (Ganin et al. 2016) is one of the important approaches of UDA that utilize the GAN (Goodfellow et al. 2020) method to learn the domain invariant feature between the source domain and target domain via discriminator and minimize the distributional gaps. CDAN (Long et al. 2018) extends the concept of DANN by adding class conditions to learn more discriminative domain invariant features. Similarly, the authors in (Long et al. 2017), proposed a UDA approach namely, JAN that learns features by aligning the joint distribution of domain-specific layers of source and target domains. Next, the MCD approach (Saito et al. 2018) attempts to align source and target distributions by utilizing task-specific decision boundaries. The authors in (Sun and Saenko 2016) develop a technique, CORAL, to lean non-linear transformation that aligns correlations of activation layers in deep neural networks. Apart from only minimizing the domain gaps, the prior approaches (Xu et al. 2019; Choi et al. 2019; Yang et al. 2021b) emphasised the use of pseudo labels, generated using the source model, for training the model in the target domain to reduce the domain divergence. These UDA approaches primarily focus on improving the feature alignment from the labelled source domain to the unlabelled target domain to perform specific machine learning tasks; however, they do not consider improving the adversarial robustness of the models.

Adversarial Robustness

Deep neural networks find applications in various domains, including autonomous driving vehicles, recognition systems, and security-related applications (malware, intrusion, spam detections *etc.*) Despite their higher applicability in threat detection and classification, deep networks are vulnerable to attacks. It introduces the challenges for training adversarially robust neural networks, enhancing reliability while dealing with maliciously manipulated inputs. Such adversarial attacks can readily perturb the trained weight of the constructed classifier due to the high memorization capability of the deep networks (Arpit et al. 2017). One of

the popular attacks, known as Fast Gradient Sign Methods (FGSM) (Szegedy et al. 2013; Goodfellow, Shlens, and Szegedy 2014), creates the adversarial example by a one-step gradient ascent across the model loss surface. Similarly, Projected Gradient Descent (PGD) (Madry et al. 2017) is a multi-step or iterative perturbation generation method for adversarial examples, a classical and effective method to generate the perturbations. Other attacks, such as Moment Iterative FGSM (MI-FGSM) (Dong et al. 2018) and Multiplicative adversarial examples(multiadv) (Lo and Patel 2021) are more frequent in neural networks. Similarly, the PGT-AT (Madry et al. 2017) method employs max-min optimization to generate the adversarial examples and train the model with these examples only.

Adversarial Robustness of Unsupervised Domain adaption

Different methods have been proposed in the existing literature to enhance the robustness of deep learning models. However, only a few methods are proposed to enhance the robustness in unsupervised domain adaption (Awais et al. 2021; Wang et al. 2024; Lo and Patel 2022; Zhu et al. 2023; Yang et al. 2021a). The authors in (Awais et al. 2021) utilize the adversarial pre-trained Imagenet model to improve the robustness of the unsupervised domain adaption. Though the mechanism provides some degree of robustness, but assumes a pre-trained model is somewhere impractical in real scenarios. Similarly, ARTUDA (Lo and Patel 2022) also propose a self-supervised method to achieve the robustness of unsupervised domain adaption in white-box attacks. It uses an additional regularizer and UDA model losses to minimize the distance between target and adversarial logits. Further, ARTUDA utilizes self-supervised signals to generate adversarial examples. SRoUDA (Zhu et al. 2023) introduces a meta-learning-based adversarial training method to improve the robustness. It utilizes a pre-trained UDA model to generate pseudo labels for target domains to generate adversarial images. Afterwards, adversarial training is performed on the target model to enhance robustness by fine-tuning the pseudo-label predictor. Finally, DART (Wang et al. 2024) also used pseudo labels to generate the adversarial example and re-train the UDA model by utilizing the source domain and adversarial target domain by considering the joint loss along with classifier and discriminator loss.

Experiments

Robustness Evaluation in Feature Space

This section utilizes t-SNE visualization to evaluate the robustness of the proposed method in the feature space of adversarial and clean examples of the target data. In feature space, adding a small perturbation in the clean image results in large changes in the feature space; thus, the model predicts the wrong class corresponding to the perturb images. To study this, we evaluate our method in feature space for the VisDA dataset (*Real* → *Syn*) to determine the robustness of the TargetNet model. We choose features from the last layer of the feature extractor (*i.e.*, ResNet-50) for clean and adversarial examples of the target data.

Algorithm 1: Proposed Method Algorithm.

Input: Source and target domain datasets: $\mathcal{D}^s = \{x_s^i, y_s^i\}_{i=1}^n$ and $\mathcal{D}^t = \{x_t^j\}_{j=1}^m$, Pre-trained UDA model F_p , Anchor model F_a , TargetNet model F_t , Batch size B , Learning rate lr , Training epoch $epoch_{max}$, Threshold confidence T , λ_1, λ_2 Hyperparameters;

Output: Adversarial trained target model F_t ;

- 1 Pre-training UDA model F_p using: $\min(\mathcal{L}_{CE}(F_s(x_s), y_s) + \omega \mathcal{L}_{dd}(x_s, x_t))$;
- 2 Initialize F_t and F_a by copying parameters from F_p ;
- 3 Anchor model F_a is frozen in training process;
- 4 **for** $i = 1$ to $epoch_{max}$ **do**
- 5 Sampling a random mini-batch B from \mathcal{D}^t and \mathcal{D}^s ;
- 6 Compute the hard pseudo label y_t for unlabeled target data, include $y_t \geq TD^t$ using F_t ;
- 7 Generate the target adversarial image x_t using pseudo labels y_t ;
- 8 Train F_t using Adversarial Target data \hat{x}_t, y_t ;
- 9 Compute attention map from F_a and F_t using x_t and \hat{x}_t respectively;
- 10 Compute the loss $L_1(\theta_M)$ gradient with the target pseudo labels;
- 11 Update the model F_t by $\theta_M' = \theta_M - \eta_1 \cdot \nabla L_1(\theta_M)$;
- 12 Compute the new loss $L_2(\theta_M')$ and gradient with the target pseudo labels and labelled source data (x_s, y_s) ;
- 13 Update the model F_t by $\theta_M'' = \theta_M' - \eta_2 \cdot \nabla L_2(\theta_M')$;
- 14 **return** TargetNet model F_t ;

Table 3: An illustration of comparison on PGD 20 attack at $\epsilon = 2/255$ using VisDA dataset.

Source→Target	Syn→Re		Re→Syn		Avg Accuracy	
Method	Clean	PGD	Clean	PGD	Clean	PGD
DANN(Ganin et al. 2016)	67.5	0.3	78.5	0.5	73.0	0.4
UDA+AT	69.6	58.3	85.7	82.0	77.6	70.1
UDA+Trades(Zhang et al. 2019)	68.1	57.9	85.1	81.5	76.6	69.7
UDA+Mart(Wang et al. 2019)	64.8	58.1	82.1	83.9	73.4	71.0
ARTUDA (Yang et al. 2021a)	45.2	32.5	72.5	62.6	58.8	47.5
SRoUDA(Zhu et al. 2023)	48.2	33.4	81.2	72.9	64.7	53.1
DART(Wang et al. 2024)	69.5	58.0	87.3	85.3	78.4	71.6
Ours	72.8	65.9	89.5	87.1	81.1	76.5

From Figure 5, While UDA+AT reveals a significant distribution gap between clean and adversarial data, and the SRoUDA method reduces this gap to some extent, our proposed method (CAM+SPLR) effectively align the clean and adversarial examples within the target data.

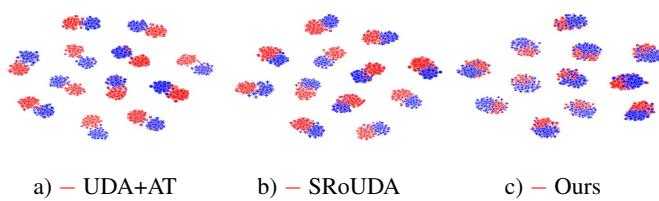


Figure 5: The t-SNE visualization of extracted features from model trained with UDA+AT, SRoUDA, and Ours) on the **Real→ Synthetic** source target domain, respectively. The blue symbols represent clean target data, while the red symbols denote the adversarial examples of the target data. Our method demonstrates a remarkable overlap between clean and adversarial examples in the feature space, as shown in (c).

Comparison results on PGD20 at $\epsilon = 2/255$

Table 3 shows the results for $\epsilon = 2/255$ using various methods on the VisDA dataset.

Additional Visual results

Figure 6 show the attention maps defense against White-box attacks on the VisDA dataset.

References

- Arpit, D.; Jastrzebski, S.; Ballas, N.; Krueger, D.; Bengio, E.; Kanwal, M. S.; Maharaj, T.; Fischer, A.; Courville, A.; Bengio, Y.; et al. 2017. A Closer Look at Memorization in Deep Networks. In *International conference on machine learning*, 233–242. PMLR.
- Awais, M.; Zhou, F.; Xu, H.; Hong, L.; Luo, P.; Bae, S.-H.; and Li, Z. 2021. Adversarial robustness for unsupervised domain adaptation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 8568–8577.
- Choi, J.; Jeong, M.; Kim, T.; and Kim, C. 2019. Pseudo-Labeling Curriculum for Unsupervised Domain Adaptation. arXiv:1908.00262.
- Dong, Y.; Liao, F.; Pang, T.; Su, H.; Zhu, J.; Hu, X.; and Li, J. 2018. Boosting adversarial attacks with momentum. In *Proceedings of*

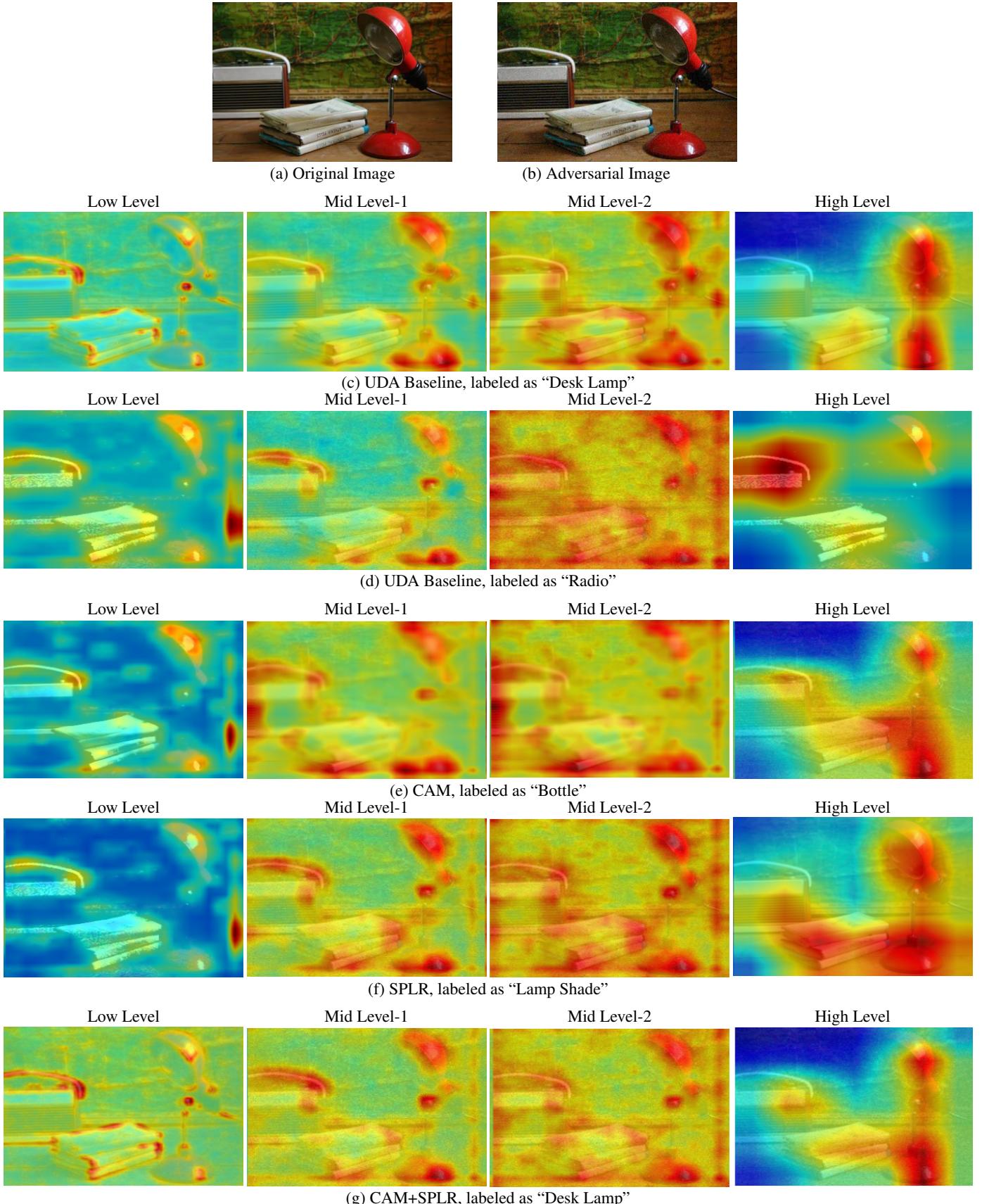


Figure 6: Attention maps for defense against White-box PGD20 attacks ($\epsilon = 8/255$) on VisDA Dataset. (a) Clean image, and (b) the corresponding adversarial image. (c) and (d) are attention maps of Clean and Adversarial images, and ((e) and (f)) display the attention maps for CAM and SPLR method individually applied to the adversarial images, while (g) attention map for the combined CAM+SPLR method.

- the IEEE conference on computer vision and pattern recognition*, 9185–9193.
- Ganin, Y.; Ustinova, E.; Ajakan, H.; Germain, P.; Larochelle, H.; Laviolette, F.; March, M.; and Lempitsky, V. 2016. Domain-Adversarial Training of Neural Networks. *Journal of Machine Learning Research*, 17(59): 1–35.
- Goodfellow, I.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A.; and Bengio, Y. 2020. Generative adversarial networks. *Commun. ACM*, 63(11): 139–144.
- Goodfellow, I. J.; Shlens, J.; and Szegedy, C. 2014. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*.
- Li, D.; Yang, Y.; Song, Y.-Z.; and Hospedales, T. M. 2017. Deeper, broader and artier domain generalization. In *Proceedings of the IEEE International Conference on Computer Vision*, 5542–5550.
- Lo, S.-Y.; and Patel, V. 2022. Exploring adversarially robust training for unsupervised domain adaptation. In *Proceedings of the Asian Conference on Computer Vision*, 4093–4109.
- Lo, S.-Y.; and Patel, V. M. 2021. Multav: Multiplicative adversarial videos. In *2021 17th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS)*, 1–6.
- Long, M.; CAO, Z.; Wang, J.; and Jordan, M. I. 2018. Conditional Adversarial Domain Adaptation. In Bengio, S.; Wallach, H.; Larochelle, H.; Grauman, K.; Cesa-Bianchi, N.; and Garnett, R., eds., *Advances in Neural Information Processing Systems*, volume 31, 1–11.
- Long, M.; Zhu, H.; Wang, J.; and Jordan, M. I. 2017. Deep Transfer Learning with Joint Adaptation Networks. In Precup, D.; and Teh, Y. W., eds., *Proceedings of the 34th International Conference on Machine Learning*, volume 70, 2208–2217.
- Madry, A.; Makelov, A.; Schmidt, L.; Tsipras, D.; and Vladu, A. 2017. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*.
- Peng, X.; Usman, B.; Kaushik, N.; Hoffman, J.; Wang, D.; and Saenko, K. 2017. Visda: The visual domain adaptation challenge. *arXiv preprint arXiv:1710.06924*.
- Ramaswamy, H. G.; et al. 2020. Ablation-CAM: Visual explanations for deep convolutional networks via gradient-free localization. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 983–991.
- Saito, K.; Watanabe, K.; Ushiku, Y.; and Harada, T. 2018. Maximum Classifier Discrepancy for Unsupervised Domain Adaptation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Sun, B.; and Saenko, K. 2016. Deep CORAL: Correlation Alignment for Deep Domain Adaptation. In *Computer Vision – ECCV 2016 Workshops*, 443–450. Cham: Springer International Publishing.
- Szegedy, C.; Zaremba, W.; Sutskever, I.; Bruna, J.; Erhan, D.; Goodfellow, I.; and Fergus, R. 2013. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*.
- Wang, H.; Wang, Z.; Du, M.; Yang, F.; Zhang, Z.; Ding, S.; Mardziel, P.; and Hu, X. 2020. Score-CAM: Score-weighted visual explanations for convolutional neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, 24–25.
- Wang, J.; Tian, K.; Ding, D.; Yang, G.; and Li, X. 2021. Unsupervised Domain Expansion for Visual Categorization. *ACM Transactions on Multimedia Computing Communications and Applications (TOMM)*. In press.
- Wang, Y.; Hazimeh, H.; Ponomareva, N.; Kurakin, A.; Hammoud, I.; and Arora, R. 2024. DART: A Principled Approach to Adversarially Robust Unsupervised Domain Adaptation. *arXiv preprint arXiv:2402.11120*.
- Wang, Y.; Zou, D.; Yi, J.; Bailey, J.; Ma, X.; and Gu, Q. 2019. Improving adversarial robustness requires revisiting misclassified examples. In *International conference on learning representations*, 1–14.
- Xu, R.; Li, G.; Yang, J.; and Lin, L. 2019. Larger Norm More Transferable: An Adaptive Feature Norm Approach for Unsupervised Domain Adaptation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 1426–1435.
- Yang, J.; Li, C.; An, W.; Ma, H.; Guo, Y.; Rong, Y.; Zhao, P.; and Huang, J. 2021a. Exploring robustness of unsupervised domain adaptation in semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 9194–9203.
- Yang, J.; Shi, S.; Wang, Z.; Li, H.; and Qi, X. 2021b. ST3D: Self-Training for Unsupervised Domain Adaptation on 3D Object Detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 10368–10378.
- Zhang, H.; Yu, Y.; Jiao, J.; Xing, E.; El Ghaoui, L.; and Jordan, M. 2019. Theoretically principled trade-off between robustness and accuracy. In *International conference on machine learning*, 7472–7482. PMLR.
- Zhu, W.; Yin, J.-L.; Chen, B.-H.; and Liu, X. 2023. SRoUDA: meta self-training for robust unsupervised domain adaptation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, 3852–3860.