

TRANSDUCER ADAPTIVE ULTRASOUND VOLUME RECONSTRUCTION

Hengtao Guo¹, Sheng Xu², Bradford J. Wood², Pingkun Yan^{1*}

¹Department of Biomedical Engineering and Center for Biotechnology and Interdisciplinary Studies, Rensselaer Polytechnic Institute, Troy, NY 12180, USA

²National Institutes of Health, Center for Interventional Oncology, Radiology & Imaging Sciences, Bethesda, MD 20892, USA

ABSTRACT

Reconstructed 3D ultrasound volume provides more context information compared to a sequence of 2D scanning frames, which is desirable for various clinical applications such as ultrasound-guided prostate biopsy. Nevertheless, 3D volume reconstruction from freehand 2D scans is a very challenging problem, especially without the use of external tracking devices. Recent deep learning based methods demonstrate the potential of directly estimating inter-frame motion between consecutive ultrasound frames. However, such algorithms are specific to particular transducers and scanning trajectories associated with the training data, which may not be generalized to other image acquisition settings. In this paper, we tackle the data acquisition difference as a domain shift problem and propose a novel domain adaptation strategy to adapt deep learning algorithms to data acquired with different transducers. Specifically, feature extractors that generate transducer-invariant features from different datasets are trained by minimizing the discrepancy between deep features of paired samples in a latent space. Our results show that the proposed domain adaptation method can successfully align different feature distributions while preserving the transducer-specific information for universal freehand ultrasound volume reconstruction.

Index Terms— Ultrasound Volume Reconstruction, Deep Learning, Domain Adaptation

1. INTRODUCTION

Ultrasound (US) is a commonly used medical imaging modality in various clinical applications. US possesses many advantages, such as low cost, portable setup, and the capability of real-time imaging. Compared to a sequence of 2D US frames, a reconstructed 3D US image volume can provide richer context information, which is often highly desired. Thus, efficiently reconstructing 3D US volume is a critical component in many interventional tasks, such as magnetic resonance imaging (MRI) and US fusion guided prostate biopsy [1, 2, 3].

* indicates corresponding author.

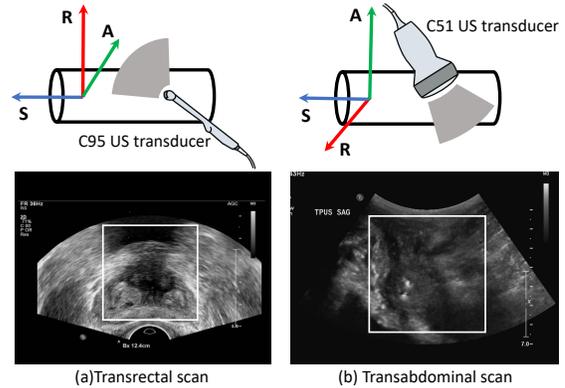


Fig. 1. (a) Abdominal and (b) transrectal scans use different ultrasound transducers along different motion trajectories. The cylinder in the first row represents patient body, and “RAS” indicates right, anterior and superior directions, respectively.

US volume registration from freehand ultrasound scans has traditionally implemented with tracking devices [4], either an optical or electromagnetic (EM) tracking system, to record the position and orientation of US transducer in 3D space. Sensorless freehand scans takes a step further by removing the requirement of tracking devices. The original method was supported by the speckle decorrelation algorithms [5], which estimates elevational distance between neighboring US images based on the speckle patterns correlation. Recent advances of deep learning (DL) methods have shown superior performance in automatic feature extraction. Prevost *et al.* [6] proposed to use convolutional neural network (CNN) to directly estimate the inter-frame motion between two 2D US frames for sensorless US volume reconstruction. In their latest work [7], two DL-reconstructed volumes from transversal and sagittal views are co-registered for a better reconstruction result. Our recent work [8] applies 3D CNN on a US video sub-sequence for better utilizing the temporal context information.

Although CNNs have achieved promising results in US volume reconstruction, these methods suffer from severe per-

formance degradation when applied to new datasets different from the training data. For example, as shown in Fig. 1, both transrectal and transabdominal scans can be used to facilitate prostate cancer diagnosis, but they have distinct motion trajectories and imaging properties. One network trained on transrectal scans cannot produce satisfactory volume reconstruction on transabdominal scans. Here we define the source domain (transrectal scans) as the dataset which serves as the training data of the CNN, and target domain (transabdominal scans) denotes the new dataset where the model is going to be applied to. Specifically, the domain shift is caused by difference between two datasets, leading to the model’s decreased performance. Our target is to efficiently transfer the model trained on source domain to the target domain given limited target labeled samples. Thus, we formulate reconstructing US volume different US transducers as a domain adaptation problem in this work.

Deep domain adaptation methods can be generally divided into two categories: adversarial-based and discrepancy-based [9]. The former methods propose to train a domain discriminator in an adversarial manner to enforce the feature vectors from both source and target domains to follow the same distribution [10, 11]. However, fooling the discriminator by generating mixed feature distributions does not help in our application. Our task is to make CNN accurately predict the relative position between two US frames. Merging the source and target feature distributions together without any high-level constraints contributes little to the task-specific feature learning: the adversarial strategy only pushes the target feature distributions close to that of the source, but does not enhance any specific feature learning that helps to accurately regress the inter-frame motion.

In the image registration field, Mahapatra and Ge [?] applied unsupervised domain adaptation for mono-modal medical image registration. An autoencoder was trained to extract latent feature vectors which are used for generating registered images through another generative adversarial network. Of note, the autoencoder was trained on chest X-ray images but the entire framework can be applied to registering images from other modalities such as brain MR images and retinal images. Another work by Zheng et al. [?] proposed a pairwise domain adaptation module to adapt the model from synthetic data to clinical data. Their primary assumption is that: if one X-ray image and one digitally reconstructed radiograph image were rendered from the same projection angle, the domain invariant feature extractor should extract consistent features from this real-synthetic image pair. In our work, on top of the discrepancy-based adaptation methods [12, 13], we propose a novel paired-sampling strategy and use a discrepancy loss to transfer task-specific feature learning from source domain to target domain. We hypothesize that if two US video sub-sequences acquired using different transducers have similar motion trajectories, they should be close to each other in the latent feature space. Our contributions are summarized as

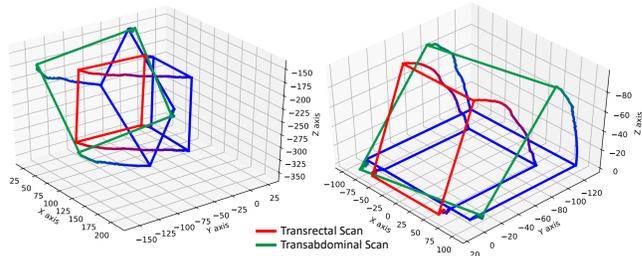


Fig. 2. US video sequence trajectories in 3D before (left) and after (right) alignment. Blue frame indicates the first frame of a video sequence. Red and green label the last frames of a transrectal scan and a transabdominal scan, respectively.

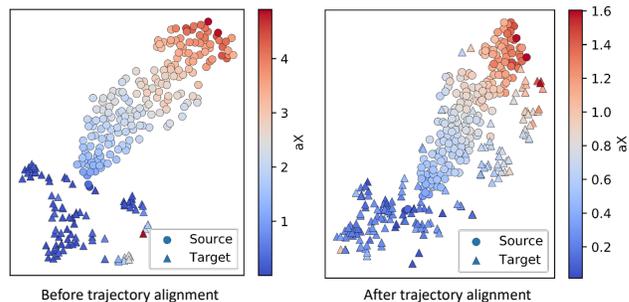


Fig. 3. For each US video, we compute a mean DOF vector throughout the sequence and use t-SNE [14] to project it into 2D space. The colorbar indicates the value of rotation α_X of each case, which is the most dominant motion direction. The trajectory alignment prevents the model’s performance being influenced by the distribution gap in label space.

follows:

1. We formulate our work on different US transducers as a domain adaptation problem. To the best of our knowledge, this reported work is the first to apply domain adaptation techniques to US volume reconstruction.
2. We propose a novel paired-sampling strategy with feature discrepancy minimization to facilitate model adaptation from the source to target domain. This strategy is specifically designed for registration-related domain adaptation problems.
3. Our results demonstrate that the proposed method can extract domain-invariant features while preserving task-specific feature learning.

2. MATERIALS AND METHODS

2.1. Transformation Space Alignment

A primary task in 3D ultrasound reconstruction is to obtain the relative spatial position of two or more consecutive US

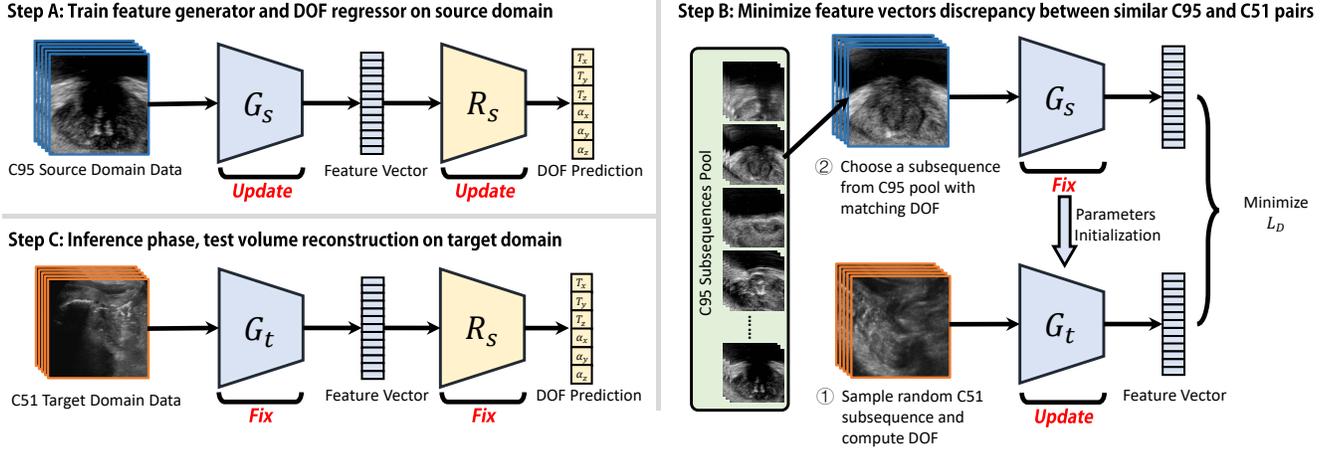


Fig. 4. Three steps of the proposed transducer adaptive ultrasound volume reconstruction (TAUVR) method, where transducer adaptation is achieved by minimizing the discrepancy between similar pairs sampled from both the source and target domains.

frames. Consider a small subsequence containing N consecutive frames as one sample unit, we can compute a relative transformation matrix and decompose it into 6 degrees of freedom (DOF) $Y = \{t_x, t_y, t_z, \alpha_x, \alpha_y, \alpha_z\}$, which contains the translations in millimeters and rotations in degrees. The network takes one video subsequence as the input for estimating the transformation parameters. We use each subsequence’s corresponding DOF vector as the groundtruth label [8] during the training process.

Since US transducers may have very different scanning trajectories for different applications (as in Fig. 2), this large motion difference will create label bias and can substantially impair the network performance. To alleviate this problem, we add a pre-processing step to roughly align the video sequence trajectory in 3D space. More precisely, we first scale the US videos to the same resolution and align the first frame of the video sequences to the same position. The sequence rotating center (transducer’s head) is overlapping with the origin $(0, 0, 0)$ of the 3D coordinate system. Thus, the label distributions of source domain and target domain are aligned together as in Fig. 3. Before the trajectory alignment, the source and target DOF label distributions are separated into two clusters; after the alignment, the label distributions are merged together, showing a smooth αX transition pattern. The trajectory alignment ensures that the model’s performance will not be impaired by the gap in label distributions.

2.2. Ultrasound Transducer Adaptation

We denote our source domain dataset (transrectal scans) as $\{X_s|Y_s\}$, where each image sample X_s represents a subsequence of $N = 5$ consecutive frames and its correspond label Y_s is a 6 DOF vector. In addition, we have another labeled but much smaller dataset on target domain (transabdominal scans) $\{X_t|Y_t\}$. Our proposed method for trans-

ducer adaptive ultrasound volume reconstruction (TAUVR) includes three consecutive steps as shown in Fig. 4.

Step A: A convolutional feature extractor G_s and a DOF regressor R_s are trained in the source domain in an end-to-end fashion [8]. The input to G_s is a $N \times W \times H$ subsequence tensor and the output is a 2048D feature vector. The DOF regressor is a linear layer that outputs 6 values for DOF regression. G_s and R_s are jointly trained by minimizing the mean squared error (MSE) loss between network’s output and groundtruth DOF labels.

Step B: In this step, we train a feature extractor G_t on target domain which produces both domain-invariant feature while preserves task-specific information. G_t is initialized with the parameters of G_s and shares the identical structure, and G_s ’s parameters are fixed in this step. We first create a source domain subsequence pool in which every transrectal video subsequence has a corresponding DOF label vector. During adaptation training, for every random target subsequence sample x_t , we compute its DOF vector y_t based on labeling information. Next, we search in the pool to find a source domain subsequence x_s that has the closest motion vector as y_t . With this paired subsequence serve as the input to corresponding networks, we yield a pair of latent feature vectors denoted as:

$$v_s = G_s(x_s), v_t = G_t(x_t) \quad (1)$$

G_t is trained by minimizing the discrepancy loss L_D , which is the L_2 norm between the two generators’ output feature vectors:

$$L_D = \frac{1}{P} \sum_{p=1}^P \|v_s^p - v_t^p\|_2 \quad (2)$$

where P denotes the total number of sampled pairs within one training epoch. The intuition of this paired sampling strategy

Method	Target Supervision	Source Domain (Transrectal Scans)								Target Domain (Abdominal Scans)							
		Distance Error (mm)				Final Drift (mm)				Distance Error (mm)				Final Drift (mm)			
		Min	Median	Max	Avg	Min	Median	Max	Avg	Min	Median	Max	Avg	Min	Median	Max	Avg
Source ADDA	No	-	-	-	-	-	-	-	-	7.31	14.35	21.59	15.73	11.87	28.15	43.66	26.73
Mixed Target	Fully	6.14	16.67	34.28	16.64	6.28	28.57	78.81	30.52	8.27	12.65	19.16	13.36	8.03	21.42	35.97	21.33
TAUVR	Weakly	-	-	-	-	-	-	-	-	7.31	10.02	20.68	12.67	6.87	22.02	32.13	20.34

Table 1. Performance of different methods on both source domain and target domain.

is to establish correspondence between source and target subsequences: when two subsequences from different domains have similar motion, we expect their extracted feature vectors to be close to each other in the latent space. This paired-sampling strategy takes rich information in the labeled source dataset as a reference to guide task-specific features learning in the target domain. Since the labels of target domain data are only used for sampling subsequence pairs while do not directly contribute to the loss function, we categorized our strategy as a weakly-supervised method.

Step C: The final step is also the inference testing phase on target domain data and does not involve any parameters update. The network used in this step is the concatenation of G_t from Step B and R_s from Step A. For a full-length US video sequence in the target domain test set, we use a sliding-window procedure to get the DOF motion vector prediction for every subsequence. By placing each frame into 3D space accordingly, a 3D US image volume can be reconstructed. The testing phase does not require any tracking devices and CNN estimates US frames relative position.

3. EXPERIMENTS

3.1. Settings

All the data utilized in this study are acquired by the Nation Institute of Health (NIH) from IRB-approved clinical trial:

Source domain contains 640 transrectal US video sequences, with each frame labeled a corresponding positioning matrix captured by EM-tracking device. An end-firing C95 transrectal ultrasound transducer captures axial images by steadily sweeping through the prostate from base to apex. The dataset is split into 500, 70 and 70 cases as training, validation and testing, respectively.

Target domain contains 12 transabdominal US video sequences acquired by C51 US transducer. 9 cases are used for training in Step B and the network’s parameters are saved after every epoch. 3 cases are used for testing in Step C.

Networks are trained for 300 epochs with batch size $K = 24$ using Adam optimizer [15]. Each US frame is cropped without exceeding the imaging field (white bounding box in Fig. 1) and then resized to 224×224 . The entire pipeline is implemented using the publicly available PyTorch library [16].

3.2. Results and Discussions

We present 4 baseline methods for comparison. As in Table 1, model on “Source” was trained on source domain and then directly tested on target domain; “Target” works in the opposite way; “Mixed” is trained on merged source and target domain using all available label for supervision; “ADDA” [10] uses unsupervised adversarial domain adaptation method to extract domain-invariant features. The proposed TAUVR achieved significantly lower average distance error and final drift comparing to both “Source” and “ADDA”. It is also comparable to the results of “Target” while the latter still have a huge domain shift problem between source and target domain because of the model’s overfitting to the transabdominal dataset.

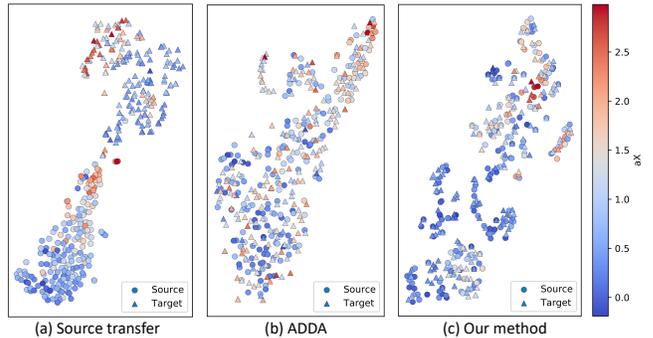


Fig. 5. t-SNE projections of the latent feature vectors obtained using (a) source domain model only, (b) ADDA, and (c) our proposed method TAUVR.

In addition, we evaluate the features quality through latent vector projections. As shown in Fig. 5, 2D tSNE projections of the extracted 2048D feature vectors are plotted, using the most dominant motion aX for color encoding.

On Fig. 5(a), points from source domain and target domain are roughly separated into two clusters, and within each cluster there exists a continuous changing pattern in aX encoding. This indicates that (1) the network trained on source domain exhibit an obvious domain gap on target data, and (2) the network, however, still preserves the task-specific information in feature vectors.

On Fig. 5(b), we observe that the distributions of two domains have been merged together through ADDA [10], as the adversarial training strategy tries to fool the domain dis-

criminator by generating “domain-invariant” features. However, since unsupervised learning poses no constraint on task-specific feature learning, the smooth color transition pattern disappears in the target domain (triangles), resulting in uninformative feature learning.

Our proposed method on Fig. 5(c) both merges the distributions of both domains and still keeps a gradual color transition in aX for each domain. These phenomena suggest that being benefited by the paired-sampling strategy, the network is extracting domain-invariant features while still preserves task-specific feature learning in both domains.

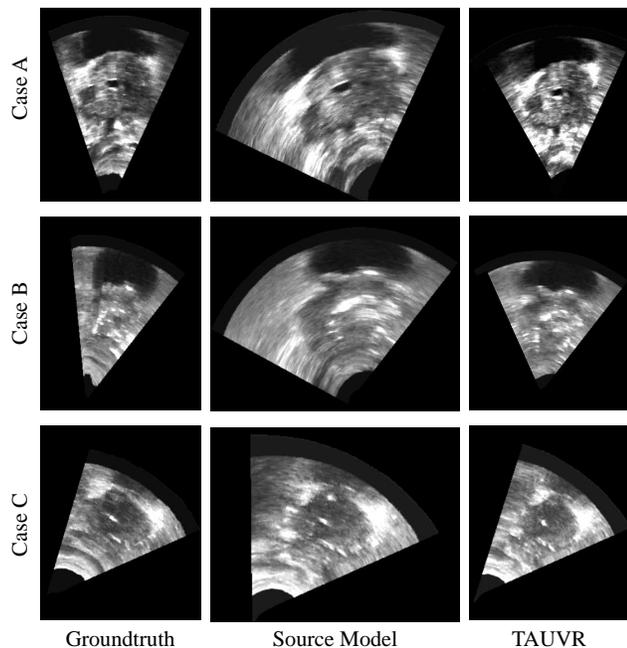


Fig. 6. Sagittal view of the reconstructed ultrasound volumes from 3 testing cases in the target domain (transabdominal scans). The model trained on source domain produces very deviated reconstruction result in the target domain (middle column). By applying the proposed TAUVR (right column), the reconstruction is much closer to the groundtruth volume.

3.3. Volume Reconstruction

We present the sagittal view of the reconstructed volumes for quality assessment in Fig. 6. All three test cases in the target domain (transabdominal scans) are presented by rows. From left to right, each column represents the reconstruction results from groundtruth labels, model trained only on source domain and our proposed TAUVR. As shown in the figure, by directly applying source model to the target data, the deep neural network may exhibit a over-fitting pattern that produce transducer trajectory prediction very close to that of the source domain. In other words, the trajectory prediction is deviated from the actual trajectory in the transabdominal scans. By

incorporating our pairwise domain adaptation methods, the third column (TAUVR) produces visually much closer volume reconstruction comparing with the groundtruth.

4. CONCLUSIONS

In this paper, we presented a novel pair-sampling strategy to enhance task-specific feature learning in target domain, using matched source domain samples as reference. The proposed transducer adaptive method (TAUVR) allows sensorless ultrasound volume reconstruction, yielding a network that is capable of extracting domain-invariant features and preserve task-specific feature learning. The proposed method achieves promising results on target domain while the performance does not degrade on source domain. A more detailed evaluation of the proposed method for additional datasets will be provided in a comprehensive future work.

5. COMPLIANCE WITH ETHICAL STANDARDS

This research study was conducted retrospectively using human subject data acquired through an IRB-approved clinical study, in accordance with the ethical standards of the institutional and/or national research committee of where the studies were conducted.

6. ACKNOWLEDGMENTS

All authors have no conflict of interest to report. This work was partially supported by National Institute of Biomedical Imaging and Bioengineering (NIBIB) of the National Institutes of Health (NIH) under awards R21EB028001 and R01EB027898, and through an NIH Bench-to-Bedside award made possible by the National Cancer Institute.

7. REFERENCES

- [1] Peter A Pinto, Paul H Chung, Ardeshir R Rastinehad, Angelo A Baccala, Jochen Kruecker, Compton J Benjamin, Sheng Xu, Pingkun Yan, Samuel Kadoury, Celene Chua, et al., “Magnetic resonance imaging/ultrasound fusion guided prostate biopsy improves cancer detection following transrectal ultrasound biopsy and correlates with multiparametric magnetic resonance imaging,” *The Journal of urology*, vol. 186, no. 4, pp. 1281–1285, 2011.
- [2] Grant Haskins, Jochen Kruecker, Uwe Kruger, Sheng Xu, Peter A Pinto, Brad J Wood, and Pingkun Yan, “Learning deep similarity metric for 3D MR-TRUS image registration,” *International journal of computer assisted radiology and surgery*, vol. 14, no. 3, pp. 417–425, 2019.

- [3] Hengtao Guo, Melanie Kruger, Sheng Xu, Bradford J Wood, and Pingkun Yan, "Deep adaptive registration of multi-modal prostate images," *Computerized Medical Imaging and Graphics*, vol. 84, pp. 101769, 2020.
- [4] Tiexiang Wen, Qingsong Zhu, Wenjian Qin, Ling Li, Fan Yang, Yaoqin Xie, and Jia Gu, "An accurate and effective fmm-based approach for freehand 3D ultrasound reconstruction," *Biomedical Signal Processing and Control*, vol. 8, no. 6, pp. 645–656, 2013.
- [5] Theresa A Tuthill, JF Krücker, J Brian Fowlkes, and Paul L Carson, "Automated three-dimensional us frame positioning computed from elevational speckle decorrelation," *Radiology*, vol. 209, no. 2, pp. 575–582, 1998.
- [6] Raphael Prevost, Mehrdad Salehi, Simon Jagoda, Navneet Kumar, Julian Sprung, Alexander Ladikos, Robert Bauer, Oliver Zettinig, and Wolfgang Wein, "3D freehand ultrasound without external tracking using deep learning," *Medical image analysis*, vol. 48, pp. 187–202, 2018.
- [7] Wolfgang Wein, Mattia Lupetti, Oliver Zettinig, Simon Jagoda, Mehrdad Salehi, Viktoria Markova, Dornoosh Zonoobi, and Raphael Prevost, "Three-dimensional thyroid assessment from untracked 2D ultrasound clips," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2020, pp. 514–523.
- [8] Hengtao Guo, Sheng Xu, Bradford Wood, and Pingkun Yan, "Sensorless freehand 3D ultrasound reconstruction via deep contextual learning," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2020, pp. 463–472.
- [9] Gabriela Csurka, "Domain adaptation for visual applications: A comprehensive survey," *arXiv preprint arXiv:1702.05374*, 2017.
- [10] Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell, "Adversarial discriminative domain adaptation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 7167–7176.
- [11] Judy Hoffman, Eric Tzeng, Taesung Park, Jun-Yan Zhu, Phillip Isola, Kate Saenko, Alexei Efros, and Trevor Darrell, "Cycada: Cycle-consistent adversarial domain adaptation," in *International conference on machine learning*. PMLR, 2018, pp. 1989–1998.
- [12] Kuniaki Saito, Kohei Watanabe, Yoshitaka Ushiku, and Tatsuya Harada, "Maximum classifier discrepancy for unsupervised domain adaptation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3723–3732.
- [13] Eric Tzeng, Judy Hoffman, Ning Zhang, Kate Saenko, and Trevor Darrell, "Deep domain confusion: Maximizing for domain invariance," *arXiv preprint arXiv:1412.3474*, 2014.
- [14] Laurens van der Maaten and Geoffrey Hinton, "Visualizing data using t-SNE," *Journal of Machine Learning Research*, vol. 9, no. 11, pp. 2579–2605, 2008.
- [15] Diederik P Kingma and Jimmy Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [16] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer, "Automatic differentiation in pytorch," in *NIPS 2017 Workshop Autodiff*, 2017.