# MT-UDA: Towards Unsupervised Cross-Modality Medical Image Segmentation with Limited Source Labels

Ziyuan Zhao[1,2], Kaixin Xu[2], Shumeng Li[1,2], Zeng Zeng[2], and Cuntai Guan[1]

[1] Nanyang Technological University, Singapore
[2] Institute for Infocomm Research, A*STAR, Singapore

**Abstract.** The success of deep convolutional neural networks (DCNNs) benefits from high volumes of annotated data. However, annotating medical images is laborious, expensive, and requires human expertise, which induces the label scarcity problem. Especially When encountering the domain shift, the problem becomes more serious. Although deep unsupervised domain adaptation (UDA) can leverage well-established source domain annotations and abundant target domain data to facilitate cross-modality image segmentation and also mitigate the label paucity problem on the target domain, the conventional UDA methods suffer from severe performance degradation when source domain annotations are scarce. In this paper, we explore a challenging UDA setting - limited source domain annotations. We aim to investigate how to efficiently leverage unlabeled data from the source and target domains with limited source annotations for cross-modality image segmentation. To achieve this, we propose a new label-efficient UDA framework, termed MT-UDA, in which the student model trained with limited source labels learns from unlabeled data of both domains by two teacher models respectively in a semi-supervised manner. More specifically, the student model not only distills the intra-domain semantic knowledge by encouraging prediction consistency but also exploits the inter-domain anatomical information by enforcing structural consistency. Consequently, the student model can effectively integrate the underlying knowledge beneath available data resources to mitigate the impact of source label scarcity and yield improved cross-modality segmentation performance. We evaluate our method on MM-WHS 2017 dataset and demonstrate that our approach outperforms the state-of-the-art methods by a large margin under the source-label scarcity scenario.

**Keywords:** Segmentation · Unsupervised Domain Adaptation · Semi-supervised Learning · Self-ensembling

## 1 Introduction

Deep convolutional neural networks (DCNNs) have obtained promising performance on medical image segmentation tasks [17,18], which further promotes the development of automated medical image analysis. DCNNs are data-hungry and

require large amounts of well-annotated data, however, in real-world clinical settings, medical image annotations are pricey and labor-intensive, which require extensive domain knowledge from biomedical experts. This leads to that scarce annotations are available for training DCNNs, *i.e.*, label scarcity.

To alleviate the burden on human annotation, plenty of methods beyond supervised learning have been proposed for improving label efficiency on medical imaging [19], including self-supervised learning [27], semi-supervised learning [2,29] and disentangled representation learning [3]. In recent years, semi-supervised learning (SSL) methods based on the self-ensembling strategy [14,20,7] have received much attention in medical image analysis, achieving state-of-the-art results in many SSL benchmarks. For instance, Laine and Aila [14] propose Temporal Ensembling to enforce the consistent outputs of the network-in-training across different epochs. Tarvainen and Valpola [20] build the mean teacher (MT) model based on the exponential moving average (EMA) of the weights of the student network, forcing the prediction consistency and further boosting the model performance. Subsequently, many studies have been devoted to leveraging abundant unlabeled data based on MT to mitigate the paucity-of-label problem in biomedical image segmentation [8,15,26]. These methods, however, are presented for label-efficient learning on a single partially labeled dataset, failing to use cross-domain information well when using multi-domain datasets.

On the other hand, given the various imaging modalities with different physical principles, such as CT and MR, the domain shift problem is severe in cross-modality image segmentation, resulting in significantly reduced performance when applying well-trained DCNNs on one domain (*e.g.*, MR) to another domain (*e.g.*, CT), especially in the absence of target labels. To tackle this serious issue, much research has been devoted to investigating unsupervised domain adaptation (UDA) for minimizing the discrepancy between the source and target domains, consequently boosting the generalization ability on the target domain for cross-modality medical image segmentation [11,25]. Inspired by the great success of generative adversarial networks (GANs) on image-to-image translation [1,12,13], many approaches have been developed with adversarial learning from different perspectives for domain alignment, including image-level adaptation [4,28], feature-level adaptation [10,9,22] and their mixtures [5,6]. For example, Chen *et al.* design a synergistic image and feature adaptation model [6], which achieves the state-of-the-art performance in UDA for cross-modality medical image segmentation. Despite the success of adversarial learning in UDA, these methods heavily rely on abundant source labels, which become sub-optimal when only limited source labels are available in clinical deployment.

These motivate us to advocate studying a practical, challenging, and different UDA setting from the past, where only limited source labels are accessible. In this paper, we investigate the feasibility of integrating SSL into UDA under source label scarcity and propose a novel label-efficient UDA framework for cross-modality medical image segmentation. We first present a dual cycle alignment module (DCAM) to bridge the appearance gap across domains, synthe-
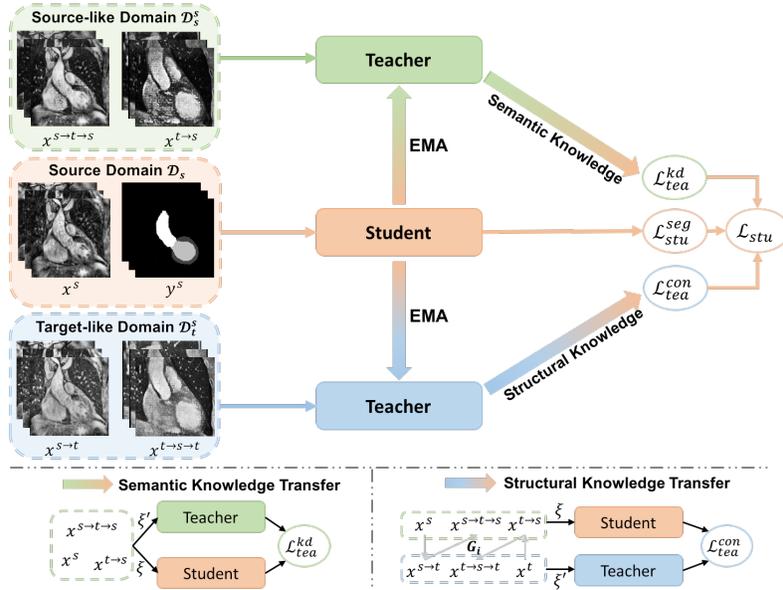
**Fig. 1.** Overall framework of our proposed MT-UDA. The student model learns from labeled source samples $\mathcal{D}_s^l$ by the $\mathcal{L}_{stu}^{seg}$ loss, and distills the intra-domain semantic knowledge and inter-domain anatomical information from *source-like domain* and *target-like domain* by $\mathcal{L}_{tea}^{kd}$ and $\mathcal{L}_{tea}^{con}$, simultaneously.

sizing *source-like domain* images and *target-like domain* images via adversarial learning [12]. We further develop an MT framework [20] for UDA, named MT-UDA, to exploit the knowledge from both intermediate domains. In MT-UDA, the student model distills the intra-domain semantic knowledge by encouraging the prediction consistency of the source domain and exploits the inter-domain anatomical information by enforcing the structural consistency across domains. We evaluate the proposed MT-UDA on a public multi-modality cardiac image segmentation dataset, MM-WHS 2017, and demonstrate that our method outperforms the state-of-the-art methods by a lot under the challenging UDA scenario.

## 2   Methodology

Let $\mathcal{D}_s^l = \{(\mathbf{x}_i^s, y_i^s)\}_{i=1}^N$ and $\mathcal{D}_s^u = \{(\mathbf{x}_i^s)\}_{i=N+1}^M$ denote the labeled samples and unlabeled samples from source domain (*e.g.*, MR), respectively. In conventional UDA setting, abundant labeled source data $\mathcal{D}_s^l$ is given, *i.e.*, $N = M$. Differently, in our setting, only limited labeled source data $\mathcal{D}_s^l$ is used for UDA, *i.e.*, $N \ll M$, which is more practical and challenging. We aim to exploit $\mathcal{D}_s^l$, $\mathcal{D}_s^u$ and unlabeled samples $\mathcal{D}_t = \{(\mathbf{x}_i^t)\}_{i=1}^P$ from target domain (*e.g.*, CT) for UDA to improve the model performance on the target domain. The overview of the

4      Zhao et al.

proposed method is presented in Fig. 1. Firstly, two sets of synthetic images, *i.e.*, *source-like domain* $\mathcal{D}_s^s$ and *target-like domain*, $\mathcal{D}_t^s$ are generated with the proposed dual cycle alignment module to alleviate the notorious domain discrepancy in appearance (see Fig. 2). To leverage the knowledge beneath real images $\mathcal{D}_s$, $\mathcal{D}_t$ and synthetic ones $\mathcal{D}_s^s$, $\mathcal{D}_t^s$, we propose an MT framework for label-efficient UDA, named MT-UDA, in which, the student model explore the knowledge beneath *source-like domain* and *target-like domain* through two teacher models simultaneously for comprehensive integration.

### 2.1  Dual Cycle Alignment Module

To reduce the semantic gap across domains, we generate synthetic samples for two domains using generative adversarial networks [12]. We design a dual cycle alignment module (DCAM) based on CycleGANs [30] to narrow the domain shift bidirectionally, as demonstrated in Fig. 2. To be specific, the target generator $G_t$ aims to transform source domain inputs to target domain distribution, *i.e.*, $G_t(x^s) = x^{s\to t}$, whereas the discriminator $D_t$ aims to differentiate whether the images are fake target images $x^{s\to t}$ or real ones $x^t$. Similarly, with $x^t$, $G_s$ aims to generate $x^{t\to s}$, while $D_s$ aims to classify the transferred images $x^{t\to s}$ and the original images $x^s$. In CycleGAN, a reverse generator is employed to impose a cycle consistency between source domain images $x^s$ and reconstructed images $x^{s\to t\to s}$. It is noted that both the reverse generator and the source generator $G_s$ aim to generate source-like images, therefore, we share the weights between them. In similar fashion, we refactor $G_t$ to generate $x^{t\to s\to t}$. Different from CycleGAN, we further force the discriminator $D_s$ to differentiate source images $x^s$, synthetic source images $x^{t\to s}$ or reconstructed source images $x^{s\to t\to s}$ in order to bridge the domain gap better, Similarly, we construct a powerful discriminator $D_t$. Finally, we can obtain two newly-augmented intermediate domains, *i.e.*, *source-like domain* $\mathcal{D}_s^s = \{x^{t\to s}, x^{s\to t\to s}\}$ and *target-like domain* $\mathcal{D}_t^s = \{x^{s\to t}, x^{t\to s\to t}\}$.
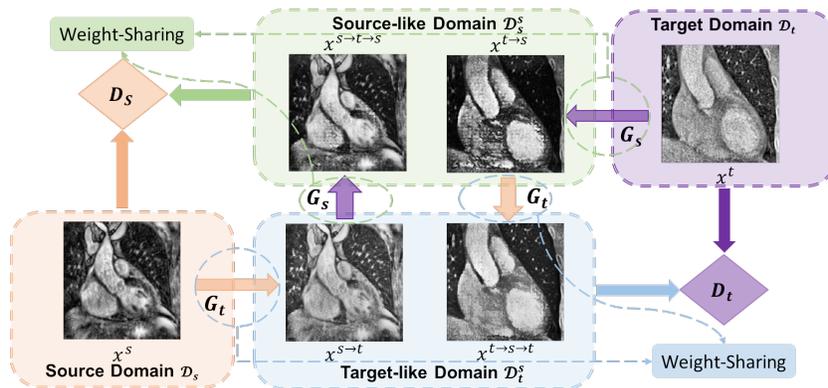


**Fig. 2.** Overall framework of Dual Cycle Alignment Module (DCAM).

## 2.2   Semantic Knowledge Transfer

Following image-level adaptation by DCAM, *source-like domain* images $\mathcal{D}_s^s$ and source domain images $\mathcal{D}_s$ maintain a similar visual appearance, allowing us to leverage the knowledge beneath $\mathcal{D}_s^s$ to improve the segmentation performance on $\mathcal{D}_s$ under label scarcity. As shown in Fig. 1, we follow the mean teacher (MT) paradigm and adopt the same architecture for the student and teacher models based on self-ensembling [21]. Specifically, the teacher model $f_{\theta'}$ at training step $t$ is updated with the exponential moving average (EMA) weights of the student model $f_\theta$, *i.e.*, $\theta'_t = \alpha \theta'_{t-1} + (1-\alpha)\theta_t$, where $\alpha$ is the EMA decay rate that reflects the influence level of the current student model parameters. Given different perturbations (*e.g.*, noises $\xi$ and $\xi'$) to the inputs of teacher and student models, we expect their predictions to be consistent by minimizing the difference between them with a mean square error (MSE) loss $\mathcal{L}_{tea}^{kd}$ as

$$\mathcal{L}_{tea}^{kd} = \frac{1}{N}\sum_{i=1}^{N}\left\|f\left(x_i;\theta'_t,\xi'\right) - f\left(x_i;\theta_t,\xi\right)\right\|^2 , \tag{1}$$

where $f(\cdot)$ is the segmentation network. $f\left(x_i;\theta_t,\xi\right)$ and $f\left(x_i;\theta_t,\xi'\right)$ represent the outputs of the student model and the teacher model, respectively.

## 2.3   Structural Knowledge Transfer

Despite distinct differences like image appearance across domains, the transformed images obtained from generators should have the same structural information as the original ones. In other words, source domain image $x_i^s$ and its synthesis target-like image $x_i^{s \to t}$ should have the same segmentation masks, *i.e.*, $y^s = y_i^{s \to t}$. In this regard, We propose a teacher model for keeping structural consistency between predictions of source images and corresponding synergistic target images, *i.e.*, $f(x;\theta,\xi) = f(G_i[x];\theta',\xi')$, where $x$ are source (-like) domain images, and $G_i$ is generator $G_t$ or reverse generator $G_s$. Transferring structural knowledge across domains not only regularizes the student model for semi-supervised learning, but also helps increase adaptation performance at the feature level. Instead of the conventional consistency loss, *e.g.*, MSE loss [16], we exploit the structural information based on weighted self-information [23,24], and calculate the structural consistency loss $\mathcal{L}_{tea}^{con}$ between the teacher and student networks as

$$\mathcal{L}_{tea}^{con} = \frac{1}{N}\sum_{i=1}^{N}\frac{1}{H \times W}\sum_{v=1}^{\mathcal{V}}\left\|\mathbf{I}_{i,v}^s - \mathbf{I}_{i,v}^t\right\|^2 \tag{2}$$

where $\mathcal{V} = \{1, 2, \ldots, H \times W\}$, $\mathbf{I}_{i,v}^s = -\mathbf{p}_{i,v}^s \circ \log \mathbf{p}_{i,v}^s$ is the weighted self-information of the predicted label at $v$-th pixel of $i$-th input from the student network, and similarly $\mathbf{I}_{i,v}^t$ is that from the teacher network. The notation $\circ$ is Hadamard product and $\log$ is the logarithmic expression using base 2.

### 2.4   MT-UDA Framework

With the supervision of corresponding labels $y^s$, the student model is trained by the supervised loss $\mathcal{L}_{stu}^{seg}$ as

$$\mathcal{L}_{stu}^{seg} = \frac{1}{2}\left[\mathcal{L}_{ce}\left(y^s, p_{stu}^s\right) + \mathcal{L}_{dice}\left(y^s, p_{stu}^s\right)\right], \tag{3}$$

where $\mathcal{L}_{ce}$ and $\mathcal{L}_{dice}$ are cross-entropy loss and dice loss, respectively, and $p_{stu}^s$ is the predictions of the student model on source labeled images $x^s$. Based on the above discussion, we integrate Eq. 1, Eq. 2 and Eq. 3, and the training objective for the student model is formulated as

$$\mathcal{L}_{stu} = \mathcal{L}_{stu}^{seg} + \lambda_{kd}\mathcal{L}_{tea}^{kd} + \lambda_{con}\mathcal{L}_{tea}^{con}, \tag{4}$$

where $\lambda_{kd}$ and $\lambda_{con}$ are the trade-off parameters with the associated losses. With the MT-UDA framework, we can distill the knowledge from *source-like domain* and *target-like domain* together for more accurate cross-modality image segmentation.

## 3   Experiments and Results

**Dataset and pre-processing.** We evaluated our method on the Multi-Modality Whole Heart Segmentation (MM-WHS) 2017 dataset, consisting of unpaired 20 MR and 20 CT volumes with ground truth masks. We employed MR as source domain and CT as target domain. Following general UDA setting as in [5], each modality was first randomly split with 16 scans for training and 4 scans for testing. To validate the performance under the source-label scarcity scenario, we randomly selected 4 annotated MR scans for training in comparison experiments. For data pre-processing, following previous work [9], we cropped all the coronal slices into centering at the heart region after resampling with unit spacing. Four cardiac substructures, *i.e.*, ascending aorta (AA), left atrium blood cavity (LAC), left ventricle blood cavity (LVC), and myocardium of the left ventricle (MYO) were selected for segmentation.

**Implementation details.** We followed [30] to optimize the proposed dual cycle alignment module for generating *source-like domain* images $\mathcal{D}_s^s$ and *target-like domain* images $\mathcal{D}_t^s$. Similar to [28], we verified our model on the transformed source-like images $x^{t\to s}$ instead of target domain images $x^t$, since our model was trained on source domain under source label scarcity. We implemented U-Net [18] as our network backbone for both student and teacher models in MT-UDA. We trained the framework for a total of 150 iterations and used Adam optimizer with the initial learning rate of $1 \times 10^{-4}$, momentum of 0.9, learning rate warm up over the first 20 iterations, and cosine decay of the learning rate with the SGD optimizer. Following [20], the EMA decay rate $\alpha$ was set to 0.999 for two teacher models, and hyperparameters $\lambda_{con}$ and $\lambda_{kd}$ were ramped up individually with

**Table 1.** Comparison results of different methods. suffix -4 or -16 after method names stand for the number of labelled source scans used for training.

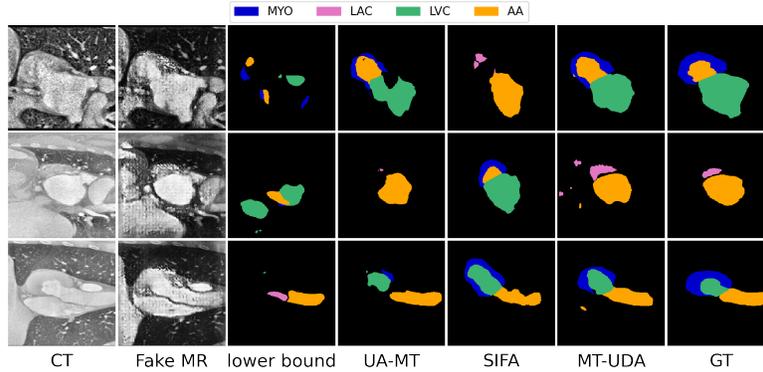| Method | | AA | LAC | LVC | MYO | Avg | AA | LAC | LVC | MYO | Avg |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | Dice ↑ | | | | | ASD ↓ | | |
| W/o Adaptation - 4 | | 5.6 | 17.8 | 12.1 | 5.5 | 10.3 | 36.9 | 24.6 | 38.5 | 35.6 | 33.9 |
| UDA-16 | PnP-AdaNet [9] | 74 | 68.9 | 61.9 | 50.8 | 63.9 | 12.8 | 6.3 | 17.4 | 14.7 | 12.8 |
| | SIFA-v1 [5] | 81.1 | 76.4 | 75.7 | 58.7 | 73 | 10.6 | 7.4 | 6.7 | 7.8 | 8.1 |
| | SIFA-v2 [6] | 81.3 | 79.5 | 73.8 | 61.6 | 74.1 | 7.9 | 6.2 | 5.5 | 8.5 | 7 |
| UDA-4 | DCAM | 19.3 | 28.1 | 34.1 | 6.4 | 22 | 32.5 | 21.8 | 17.7 | 22.8 | 23.7 |
| | SIFA-v2 [6] | 50.5 | 59.6 | 31.9 | 28.9 | 42.7 | 8.8 | 7.3 | 15.8 | 13.2 | 11.3 |
| SSL-4 | MT [20] | 3.6 | 26.8 | 14.5 | 4.6 | 12.4 | 34.5 | 22.7 | **5.7** | 17.6 | 20.1 |
| | UA-MT [26] | 20.1 | 40.5 | 2.5 | 11.3 | 18.6 | 40.1 | 23.3 | 43.2 | 20.9 | 31.9 |
| UDA +SSL-4 | DCAM+MT [20] | 35.3 | 31.6 | 48.4 | 11.2 | 31.6 | 39.9 | 39.8 | 10.5 | 14.6 | 23.7 |
| | DCAM+UA-MT [26] | 61.3 | 59.7 | 46.5 | 19.2 | 46.7 | 5.6 | 8.3 | 8.2 | 10.6 | 8.2 |
| | MT-UDA (Ours) | **72.7** | **71.4** | **60.7** | **41.7** | **61.6** | **5.3** | **5.7** | 6.7 | **6.1** | **5.9** |



**Fig. 3.** Visualization of segmentation results generated by different methods.

the sigmoid-shaped function $\lambda(t) = 0.01 \cdot e^{\left(-5(1-t/t_{\max})^2\right)}$, where $t$ and $t_{max}$ were the current and the last step, respectively. Data augmentation such as random rotation was applied in all the experiments for a fair comparison. We evaluate different methods on Dice score and average surface distance (ASD) with the largest 3D connected component of each substructure.

**Comparison with other methods.** We compare our methods with the state-of-the-art UDA methods in cardiac segmentation, *i.e.,* Pnp-AdaNet [10] and SIFA [5,6], as well as two recent popular SSL approaches, including MT [20] and UA-MT [26]. In Table 1, we list the results of PnP-AdaNet and SIFA with 16 labeled source scans in cardiac segmentation. Since SIFA-v2 [6] obtains the best segmentation performance on each substructure, we further train SIFA-v2 on 4 labeled MR scans to simulate the source-label scarcity scenario. It is observed
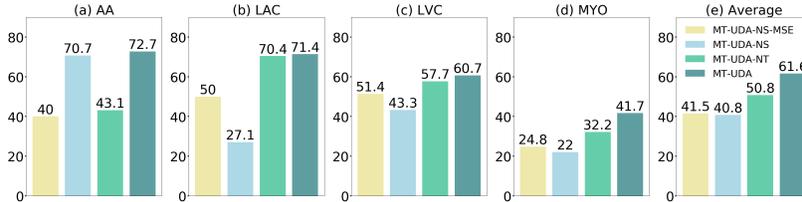
**Fig. 4.** Ablation results (Dice[%]) on different components.

that SIFA-v2 obtains severely degraded performance on target domain when using 4 labeled source domain scans, which can be attributed to the source label scarcity. We also directly test the U-Net trained on 4 labeled MR scans from the source domain as our lower bound, referred as W/o Adaptation-4. By taking advantage of image-to-image translation *i.e.,* DCAM, a great improvement can be achieved when testing W/o Adaptation on fake MR images $x^{t \to s}$, but it is still not optimal with the average dice of merely 22% across the substructures. It is worth noting that MT and UA-MT can help improve the segmentation performance on target domain by leveraging unlabeled source domain images. Along with image appearance alignment, MT and UA-MT can achieve promising improvement on cross-modality segmentation, which demonstrates the feasibility of integrating SSL into UDA for label-efficient UDA. By simultaneously exploiting all available data sources, the proposed MT-UDA obtains the best segmentation results with the average dice of 61.6%, outperforming SIFA-v2 (4 training MR scans) by a large margin and achieving comparable performance with the state-of-the-art methods, but only requires 1/4 source labels. We further visualize the segmentation results on testing data of different methods including the best methods of UDA and SSL, *i.e.* SIFA-v2 and DCAM+UA-MT in Fig. 3. It is observed that our method can generate more reliable masks with fewer false positives than other methods.

**Ablation studies of our method.** To evaluate the effectiveness of different components of MT-UDA, we conduct ablation experiments on various variants. Specifically, we remove one of the teacher models, separately, *i.e.*, W/o semantic knowledge transfer (MT-UDA-NS) and W/o structural knowledge transfer (MT-UDA-NT). We further implement the MSE loss in MT-UDA-NS to evaluate the efficacy of the structural loss, *i.e.*, MT-UDA-NS-MSE. Fig 4 demonstrates the ablation results of different substitutes. We can see that both types of knowledge transfer can benefit the model performance on unsupervised cross-domain segmentation. In comparison with the MSE loss in structural knowledge transfer, the proposed loss based on the weighted self-information can better improve the segmentation performance on some substructures such as AA, benefiting from the structural consistency across domains.

## 4    Conclusion

In this work, we present a novel label-efficient UDA framework, MT-UDA, which integrates SSL into UDA for cross-modality medical image segmentation under source label scarcity. By bridging both source and target domains to intermediate domains through knowledge transfer, the student model can leverage intra-domain semantic knowledge and exploit inter-domain structural knowledge concurrently, thereby mitigating both the domain discrepancy and source label scarcity. We evaluate the proposed MT-UDA on MM-WHS 2017 dataset, and demonstrate that our method outperforms the state-of-the-art UDA methods by a lot under the challenging source-label scarcity scenario.

## References

1. Arjovsky, M., Chintala, S., Bottou, L.: Wasserstein generative adversarial networks. In: International conference on machine learning. pp. 214–223. PMLR (2017)
2. Bai, W., Oktay, O., Sinclair, M., Suzuki, H., Rajchl, M., Tarroni, G., Glocker, B., King, A., Matthews, P.M., Rueckert, D.: Semi-supervised learning for network-based cardiac mr image segmentation. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 253–260. Springer (2017)
3. Chartsias, A., Joyce, T., Papanastasiou, G., Semple, S., Williams, M., Newby, D.E., Dharmakumar, R., Tsaftaris, S.A.: Disentangled representation learning in cardiac image analysis. Medical image analysis **58**, 101535 (2019)
4. Chen, C., Dou, Q., Chen, H., Heng, P.A.: Semantic-aware generative adversarial nets for unsupervised domain adaptation in chest x-ray segmentation. In: International workshop on machine learning in medical imaging. pp. 143–151. Springer (2018)
5. Chen, C., Dou, Q., Chen, H., Qin, J., Heng, P.A.: Synergistic image and feature adaptation: Towards cross-modality domain adaptation for medical image segmentation. Proceedings of the AAAI Conference on Artificial Intelligence **33**(01), 865–872 (Jul 2019)
6. Chen, C., Dou, Q., Chen, H., Qin, J., Heng, P.A.: Unsupervised bidirectional cross-modality adaptation via deeply synergistic image and feature alignment for medical image segmentation. IEEE transactions on medical imaging **39**(7), 2494–2505 (2020)
7. Cheplygina, V., de Bruijne, M., Pluim, J.P.: Not-so-supervised: a survey of semi-supervised, multi-instance, and transfer learning in medical image analysis. Medical image analysis **54**, 280–296 (2019)
8. Cui, W., Liu, Y., Li, Y., Guo, M., Li, Y., Li, X., Wang, T., Zeng, X., Ye, C.: Semi-supervised brain lesion segmentation with an adapted mean teacher model. In: International Conference on Information Processing in Medical Imaging. pp. 554–565. Springer (2019)

9. Dou, Q., Ouyang, C., Chen, C., Chen, H., Glocker, B., Zhuang, X., Heng, P.: Pnp-adanet: Plug-and-play adversarial domain adaptation network at unpaired cross-modality cardiac segmentation. IEEE Access **7**, 99065–99076 (2019)

10. Dou, Q., Ouyang, C., Chen, C., Chen, H., Heng, P.A.: Unsupervised cross-modality domain adaptation of convnets for biomedical image segmentations with adversarial loss. In: Proceedings of the 27th International Joint Conference on Artificial Intelligence. p. 691–697. IJCAI'18, AAAI Press (2018)

11. Ganin, Y., Lempitsky, V.: Unsupervised domain adaptation by backpropagation. In: International conference on machine learning. pp. 1180–1189. PMLR (2015)

12. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. In: Ghahramani, Z., Welling, M., Cortes, C., Lawrence, N., Weinberger, K.Q. (eds.) Advances in Neural Information Processing Systems. vol. 27. Curran Associates, Inc. (2014)

13. Hoffman, J., Tzeng, E., Park, T., Zhu, J.Y., Isola, P., Saenko, K., Efros, A., Darrell, T.: Cycada: Cycle-consistent adversarial domain adaptation. In: International conference on machine learning. pp. 1989–1998. PMLR (2018)

14. Laine, S., Aila, T.: Temporal ensembling for semi-supervised learning. In: 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings. OpenReview.net (2017)

15. Li, K., Wang, S., Yu, L., Heng, P.A.: Dual-teacher: Integrating intra-domain and inter-domain teachers for annotation-efficient cardiac segmentation. In: Martel, A.L., Abolmaesumi, P., Stoyanov, D., Mateus, D., Zuluaga, M.A., Zhou, S.K., Racoceanu, D., Joskowicz, L. (eds.) Medical Image Computing and Computer Assisted Intervention – MICCAI 2020. pp. 418–427. Springer International Publishing, Cham (2020)

16. Li, X., Yu, L., Chen, H., Fu, C.W., Xing, L., Heng, P.A.: Transformation-consistent self-ensembling model for semisupervised medical image segmentation. IEEE Transactions on Neural Networks and Learning Systems (2020)

17. Long, J., Shelhamer, E., Darrell, T.: Fully convolutional networks for semantic segmentation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 3431–3440 (2015)

18. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: International Conference on Medical image computing and computer-assisted intervention. pp. 234–241. Springer (2015)

19. Tajbakhsh, N., Jeyaseelan, L., Li, Q., Chiang, J.N., Wu, Z., Ding, X.: Embracing imperfect datasets: A review of deep learning solutions for medical image segmentation. Medical Image Analysis **63**, 101693 (2020)

20. Tarvainen, A., Valpola, H.: Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. In: Proceedings of the 31st International Conference on Neural Information Processing Systems. p. 1195–1204. NIPS'17, Curran Associates Inc., Red Hook, NY, USA (2017)

21. Tarvainen, A., Valpola, H.: Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. arXiv preprint arXiv:1703.01780 (2017)

22. Tzeng, E., Hoffman, J., Saenko, K., Darrell, T.: Adversarial discriminative domain adaptation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 7167–7176 (2017)

23. Vu, T.H., Jain, H., Bucher, M., Cord, M., Pérez, P.: Advent: Adversarial entropy minimization for domain adaptation in semantic segmentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 2517–2526 (2019)

24. Vu, T.H., Jain, H., Bucher, M., Cord, M., Pérez, P.: Dada: Depth-aware domain adaptation in semantic segmentation. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 7364–7373 (2019)
25. Yang, J., Dvornek, N.C., Zhang, F., Chapiro, J., Lin, M., Duncan, J.S.: Unsupervised domain adaptation via disentangled representations: Application to cross-modality liver segmentation. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 255–263. Springer (2019)
26. Yu, L., Wang, S., Li, X., Fu, C.W., Heng, P.A.: Uncertainty-aware self-ensembling model for semi-supervised 3d left atrium segmentation. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 605–613. Springer (2019)
27. Zeng, Z., Xulei, Y., Qiyun, Y., Meng, Y., Le, Z.: Sese-net: Self-supervised deep learning for segmentation. Pattern Recognition Letters **128**, 23–29 (2019)
28. Zhang, Y., Miao, S., Mansi, T., Liao, R.: Task driven generative modeling for unsupervised domain adaptation: Application to x-ray image segmentation. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 599–607. Springer (2018)
29. Zhao, Z., Zeng, Z., Xu, K., Chen, C., Guan, C.: Dsal: Deeply supervised active learning from strong and weak labelers for biomedical image segmentation. IEEE Journal of Biomedical and Health Informatics (2021)
30. Zhu, J.Y., Park, T., Isola, P., Efros, A.A.: Unpaired image-to-image translation using cycle-consistent adversarial networks. In: Proceedings of the IEEE international conference on computer vision. pp. 2223–2232 (2017)