MAPSeg: Unified Unsupervised Domain Adaptation for Heterogeneous Medical Image Segmentation Based on 3D Masked Autoencoding and Pseudo-Labeling

Xuzhe Zhang^{1†}, Yuhao Wu^{2†}, Elsa Angelini^{1,3}, Ang Li⁴, Jia Guo¹, Jerod M. Rasmussen⁵,
Thomas G. O'Connor⁶, Pathik D. Wadhwa⁵, Andrea Parolin Jackowski⁷,
Hai Li², Jonathan Posner², Andrew F. Laine^{1‡}, Yun Wang^{2,8‡}

¹Columbia University ²Duke University ³Télécom Paris, LTCI, Institut Polytechnique de Paris

⁴University of Maryland, College Park ⁵University of California, Irvine ⁶University of Rochester

⁷Universidade Federal de São Paulo ⁸Emory University

Abstract

Robust segmentation is critical for deriving quantitative measures from large-scale, multi-center, and longitudinal medical scans. Manually annotating medical scans, however, is expensive and labor-intensive and may not always be available in every domain. Unsupervised domain adaptation (UDA) is a well-studied technique that alleviates this label-scarcity problem by leveraging available labels from another domain. In this study, we introduce Masked Autoencoding and Pseudo-Labeling Segmentation (MAPSeg), a unified UDA framework with great versatility and superior performance for heterogeneous and volumetric medical image segmentation. To the best of our knowledge, this is the first study that systematically reviews and develops a framework to tackle four different domain shifts in medical image segmentation. More importantly, MAPSeg is the first framework that can be applied to centralized, federated, and test-time UDA while maintaining comparable performance. We compare MAPSeg with previous state-of-the-art methods on a private infant brain MRI dataset and a public cardiac CT-MRI dataset, and MAPSeg outperforms others by a large margin (10.5 Dice improvement on the private MRI dataset and 5.7 on the public CT-MRI dataset). MAPSeg poses great practical value and can be applied to real-world problems. GitHub: https://github.com/XuzheZ/MAPSeg/.

1. Introduction

Quantitative measures from medical scans serve as biomarkers for various types of medical research and clinical practice. For instance, neurodevelopmental studies

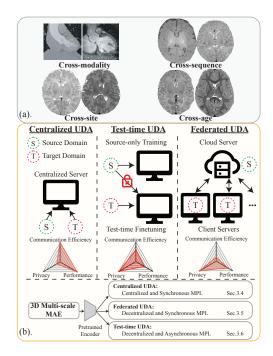


Figure 1. (a). Illustrations of four different domain shifts in medical images. (b). Overview of different UDA settings and how MAPSeg can fit into different scenarios.

utilize metrics such as brain volume and cortex thickness/surface area from infant brain magnetic resonance imaging (MRI) to investigate the early brain development and neurodevelopmental disorders [2, 11, 24, 62]. Therefore, robust segmentation of medical images acquired from large-scale, multi-center, and longitudinal studies is desired, yet often challenged by the domain shifts across different imaging techniques and even within a single modality (Fig.1a). For example, computed tomography (CT) and MRI provide markedly different signals for the same structure (e.g., cardiac regions, Fig.1a). MRI, a widely adopted radiation-free imaging technique, bears various types of in-

[†]Co-first authors.

[‡]Co-senior supervising authors.

herent heterogeneity, including cross-sequence (*e.g.*, distinct contrasts for the same tissue in T1/T2 sequences) and cross-site (*e.g.*, contrast of the same tissue in the same sequence varies with acquisition scanner and setup). Moreover, subject-dependent physiological changes also lead to domain shift. For example, contrasts of white matter and grey matter vary while the human brain undergoes significant growth and expansion within both cortical and subcortical regions during early postnatal years [20], which contributes to the cross-age domain shift (Fig.1a).

The prevalent heterogeneities in medical images lead to suboptimal performance when deep neural networks trained in one source domain are applied to another target domain. To address this challenge, we introduce a *unified* unsupervised domain adaptation (UDA) framework for volumetric and heterogeneous medical image segmentation, named Masked Autoencoding and Pseudo-Labeling Segmentation (MAPSeg). To the best of our knowledge, MAPSeg is the first framework that can be used in centralized, federated, and test-time UDA for volumetric medical image segmentation while maintaining comparable performance. This versatility is particularly advantageous in the field of medical image segmentation, where data sharing is restricted and annotations are expensive. While centralized UDA delivers the best performance in most cases, the strict requirement of co-located data limits its application in multiinstitutional studies due to regulations such as the Health Insurance Portability and Accountability Act (HIPAA) and EU General Data Protection Regulation (GDPR) [38, 68]. MAPSeg circumvents this restriction with federated and test-time adaptation, enabling clinical and research collaboration across different medical centers. In contrast, some previous studies, despite showing promising results in one scenario, may become infeasible or suffer significant performance drop in others due to the requirement for co-located data or synchronous adaptation.

In addition, we conduct extensive experiments on a private infant brain MRI dataset, which includes expertprovided annotations, to evaluate MAPSeg on crosssequence, cross-site, and cross-age adaptation tasks. MAPSeg is also compared with previously reported stateof-the-art (SOTA) results on a public cardiac CT → MRI segmentation task. MAPSeg consistently outperforms previous SOTA methods by a large margin (10.5 Dice improvement on the private MRI dataset and 5.7 on the public CT-MRI dataset in the centralized UDA setting). While previous studies have separately explored one of the abovementioned domain shifts [12, 64, 79], they may not generalize to others. For example, cross-age domain shift is mainly composed of changes in brain size and contrast, and methods based on image-translation fail to handle it as they also change the size when translating data from target domain to source domain, leading to segmentation errors. We systematically evaluate MAPSeg across various domain shifts and imaging modalities, demonstrating its consistent and generalizable effectiveness.

Moreover, in all three UDA settings, MAPSeg does not rely on any target labels for model validation and selection. On the contrary, some previous studies on cardiac CT → MRI segmentation [3, 4] validate and select the best model using labeled target data, which may not be readily available in real-world problems. We demonstrate that MAPSeg surpasses the previous SOTA results without using any target label for validation, and the performance drop between using and without using target label is minor (0.9 mean Dice). This further justifies its practical value in real-world medical image segmentation tasks. The contributions of this study are multi-fold:

- 1. We propose MAPSeg, a unified UDA framework capable of handling various domain shifts in medical image segmentation.
- MAPSeg is suitable for universal UDA scenarios, suggesting its versatility and practical value for real-world problems.
- 3. MAPSeg is extensively evaluated on both private and public datasets, outperforming previous SOTA methods by a large margin. We conduct detailed ablation studies to investigate the impact of each component of MAPSeg.

2. Related Work

2.1. Masked Image Modeling

Masked image modeling (MIM) represents a category of methods that learn representations from corrupted or incomplete images [7, 13, 28, 56], and can naturally serve as a pretext task for self-supervised learning. For example, masked autoencoder (MAE) trains an encoder by reconstructing missing regions from a masked image input and has demonstrated improved generalization and performance in downstream tasks [26, 39, 41, 66, 74–76]. MAPSeg heavily relies on MIM, leveraging MAE and masked pseudo-labeling (MPL), to achieve versatile UDA.

2.2. Pseudo-Labeling

Pseudo-labeling facilitates learning from limited or imperfect data and is prevalent in semi- and self-supervised learning [35, 43, 58]. Consistency regularization is widely used in pseudo-label learning [42, 60], which is a scheme that forces the model to output consistent prediction for inputs with different degrees of perturbation (*e.g.*, weakly- and strongly-augmented images). Mean Teacher [67], a teacher-student framework that generates pseudo-labels from the teacher model (which is a temporal ensembling of the student model), is also a common strategy. In this work, we utilize the teacher model to generate pseudo labels based on complete images and guide the learning of student model

on masked images.

2.3. Unsupervised Domain Adaptation

Discrepancy minimization, adversarial learning, and pseudo-labeling are the three main directions explored in UDA. Previous studies have explored minimizing the discrepancy between source and target domains within different spaces, such as input [3, 29, 79], feature [14– 16, 18, 27, 50], and output spaces [36, 71], and they sometimes overlap with approaches base on adversarial learning as the supervisory signal to align two distributions may come from statistical distance metrics [21, 49] or a discriminator model [15, 29, 71]. Meanwhile, self-training with pseudo-label is also a prevalent technique [81, 82, 85] and has shown significant improvement on natural image segmentation [32–34]. Hover et al. [34] proposed masked image consistency as a plug-in to improve previous UDA baselines. In contrast, MAPSeg leverages the synergy between MAE and MPL, and employs MPL as a standalone component for various scenarios.

In this study, we exploit the vanilla pseudo-labeling with three straightforward yet crucial measures to stabilize the training. We hypothesize that random masking is an ideal strong perturbation for consistency regularization in pseudo-labeling, and the model pretrained via MAE can be efficiently adapted to infer semantics of missing regions from visible patches. This hypothesis is justified in Sec. 4.3. In addition, we leverage the anatomical distribution prior in medical images and make predictions jointly based on local and global contexts, which also help mitigate the pseudo-label drifts. We demonstrate the superior performance and versatility of MAPSeg in different UDA scenarios in the following sections.

2.4. Federated Learning

Federated Learning (FL) is a distributed learning paradigm that aims to train models on decentralized data [54]. FL has attracted great attention in the research community in the last few years and numerous works have focused on the key challenges raised by FL such as data/system heterogeneity [65] and communication/computation efficiency [44]. By virtue of keeping privacy-sensitive medical data local, FL has been adopted for various medical image analysis tasks [23]. Sheller et al. [61] pioneered FL for brain tumor segmentation on multimodal brain scans in a multiinstitutional collaboration and showed its promising performance compared to centralized training. Yang et al. [77] proposed a federated semi-supervised learning framework for COVID-19 detection that relaxed the requirement for all clients to have access to ground truth annotations. Fed-Mix [72] further alleviated the necessity for all clients to possess dense pixel-level labels, allowing users with weak bounding-box labels or even image-level class labels to collaboratively train a segmentation model. In contrast, MAPSeg assumes all clients have completely unlabeled data when extended to federated UDA scenario. Mushtaq et al. [55] proposed a Federated Alternate Training (FAT) scheme that leverages both labeled and unlabeled data silos. It employs mixup [80] and pseudo-labeling to enable self-supervised learning on the unlabeled participants. MAPSeg, on the other hand, adopts masked pseudolabeling and global-local feature collaboration for adapting to unlabeled target domains. Yao et al. [78] introduced the federated multi-target domain adaptation problem and a solution termed *DualAdapt*. It decouples the local-classifier adaptation with client-side self-supervised learning from the feature alignment via server-side mixup and adversarial training. MAPSeg addresses the same federated multitarget UDA problem, and we compare our results to those of FAT and DualAdapt in Sec. 4.3.

2.5. Test-Time UDA

While federated UDA eases the constraint of centralized data, its learning paradigm still requires synchronous learning across server and clients. Test-time UDA [5, 10, 19, 27, 40, 47] assumes the unavailability of source-domain data when adapting to target domains. This assumption significantly limits the applicability of methods based on image translation, adversarial learning, and feature distribution alignment which require simultaneous access to both source and target data. Gandelsman et al. [19] explored using MAE retraining during test-time to improve classification without employing pseudo labeling. Chen et al. [5] proposed using prototype and uncertainty estimation for denoised pseudo labeling of 2D fundus images. Karani et al. [40] designed a 2D denoising autoencoder to refine pseudo labels. He et al. [27] employed AE during testtime to align source and target feature distributions by minimizing AE reconstruction loss. We demonstrate that, with slight performance drop on source domain, MAPSeg can be extended to test-time UDA with comparable performance to that of centralized UDA on target domain (Sec. 4.3).

3. Methods

3.1. Preliminary

In this section, we introduce each component of MAPSeg (Fig.2) and how MAPSeg can serve as a unified solution to centralized, federated, and test-time UDA (Fig.1b). We deploy MAPSeg for domain adaptative 3D segmentation of heterogeneous medical images and it consists of three components: (1) 3D masked multi-scale autoencoding for self-supervised pre-training, (2) 3D masked pseudo-labeling for domain adaptive self-training, and (3) global-local feature collaboration to fuse global and local contexts for the final segmentation task. The hybrid cross-entropy and Dice loss

(Eq.1) is often adopted for regular supervised segmentation training, and we employ it as the basic component of the objective functions for MAPSeg:

$$\mathcal{L}_{seg}(\hat{y}, y) = -\frac{1}{n} \sum_{i} \sum_{j} y_{i,j} \log(\hat{y}_{i,j}) - \frac{2 \sum y \hat{y} + \epsilon}{\sum y + \sum \hat{y} + \epsilon}$$
(1)

where n denotes the number of pixels, $y_{i,j}$ and $\hat{y}_{i,j}$ represent the ground truth label and predicted probability for the ith pixel to belong to the jth class, and ϵ is used to prevent zero-division.

In the following sections, notations are defined as: x and y indicate the original image and label of the randomly sampled local patch; X and Y refer to downsampled global scan and label; the subscripts s and t refer to the source and target domains, respectively; the superscript M indicates the image is masked $(e.g., x_t^M)$ refers to a masked local patch from the target domain).

3.2. 3D Multi-Scale Masked Autoencoder (MAE)

In this study, we propose a 3D variant of MAE using a 3D CNN backbone (Fig.2a). The detailed configuration can be found in Appendix Sec. 1.1. Training is jointly performed on two image sources with identical size (96^3 voxels): local patches x randomly sampled from the volumetric scan, and the whole scan downsampled to the same size, denoted as X. Both x and X are masked before feeding into the MAE: x is divided into non-overlapping 3D sub-patches with size 8³, of which 70% are masked out randomly based on a uniform distribution (Fig.2a); The same procedure is applied to X with patch size 4^3 since it contains a larger field-ofview (FOV). The masked versions of x and X are denoted as x^M and X^M , respectively. We train the MAE encoder and decoder to reconstruct x/X based on x^M/X^M using mean squared error on the masked-out regions as the objective function.

3.3. 3D Masked Pseudo-Labeling (MPL)

MPL uses a teacher-student framework which is a standard strategy in semi-/self-supervised learning [22, 67] to provide stable pseudo labels on an unlabeled target domain during training. After MAE pre-training, we keep the MAE encoder g and append a segmentation decoder h to build the segmentation model $f=h\circ g$ (Fig.2b-c). Given an input image x_s and label y_s from the source domain and an input image x_t from the target domain, the teacher model f_θ takes as input the target image x_t and generates pseudo labels $f_\theta(x_t)$, with gradient detached. The student model f_ϕ is then optimized by minimizing the segmentation loss between the predictions of x_t^M/x_s^M and $f_\theta(x_t)/y_s$, which can be formulated as:

$$\mathcal{L}_{MPL} = \mathcal{L}_{Seg}(f_{\phi}(x_t^M), f_{\theta}(x_t)) + \beta \mathcal{L}_{Seg}(f_{\phi}(x_s^M), y_s)$$
(2)

where β is the weight of source prediction and set as 0.5. The teacher model's parameters θ are then updated during training via exponential moving average (EMA) based on the student model's parameters ϕ [67].

$$\theta_{t+1} \leftarrow \alpha \theta_t + (1 - \alpha)\phi_t,$$
 (3)

where t and t+1 indicate training iterations and α is the EMA update weight. For model initialized from the large-scale MAE pretraining, we set α as 0.999 during the first 1,000 steps and 0.9999 afterwards. For model pretrained on small-scale source and target datasets (e.g., only dozens of scans), we set α as 0.99 during the first 1,000 steps, 0.999 during the next 2,000 steps, and 0.9999 for the remaining training. The teacher model f_{θ} is initialized with student model's parameters ϕ after some warm-up training (e.g., 1,000 iterations) on the source-domain data.

3.4. 3D Global-Local Collaboration (GLC)

Directly applying MPL for UDA segmentation with large domain shift (e.g., cross-modality/sequence) may lead to unreliable pseudo-label and disrupt the training. Therefore, we design a GLC module (Fig.2c) to improve pseudolabeling by leveraging the spatial global-local contextual relations induced by the inherent anatomical distribution prior in medical images. With the image encoder pretrained to extract image features at both local and global levels during multi-scale MAE, we take advantage of the global-local contextual relations by concatenating local and global semantic features in the latent space and make prediction based on the fused features. We differ from previous study [8] by only applying GLC on the output of the encoder g instead of all layers to save computation cost and employing a different regularization to prevent segmentation decoder from predicting solely based on local features.

In GLC, a binary mask M is used to indicate the corresponding location of the local patch x inside the downsampled global volume X. The encoder g takes as input xand X and generates the local latent feature $\chi_{loc} = g(x)$ as well as cropped and resized global latent feature χ_{alo} = $upsample(M \odot g(X)), \text{ where } \odot \text{ indicates cropping } g(X)$ based on M followed by upsampling to match the spatial size of χ_{loc} . Therefore, segmenting a local patch x can be rewritten as $f(x) = h(\chi_{loc} \oplus \chi_{glo})$, where \oplus is the concatenation along channel dimension (Fig.2c). In addition, f is also trained on downsampled global volume X with $\mathcal{L}_{Seq}(f(X),Y)$), in which the global latent feature g(X)is duplicated and $f(X) = h(g(X) \oplus g(X))$, to prevent model from solely relying on local semantic features and encourage the encoder to extract meaningful semantic features from both local and global levels.

We also add a regularization term between the χ_{loc} and χ_{glo} to maintain their similarity following [8]. Instead of the \mathcal{L}_2 regularization used in [8], we maximize the cosine

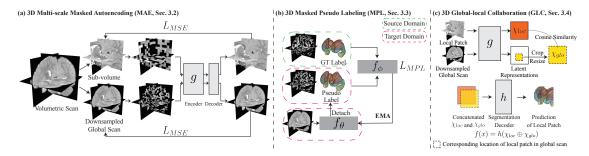


Figure 2. Components of the proposed MAPSeg framework. (a) 3D multi-scale masked autoencoding. (b) 3D masked pseudo labeling in source and target domains. (c) 3D Global-local collaboration.

similarity between the χ_{loc} and χ_{qlo} as

$$\mathcal{L}_{cos}(x, X) = 1 - \frac{\chi_{loc} \cdot \chi_{glo}}{\max(\|\chi_{loc}\|_2, \|\chi_{glo}\|_2, \epsilon)}$$
(4)

where ϵ is used to prevent zero-division. The loss function for GLC calculated on the source data is formulated as:

$$\mathcal{L}_{GLC}^{S} = \gamma(\mathcal{L}_{Seg}(f_{\phi}(X_s), Y_s) + \mathcal{L}_{Seg}(f_{\phi}(X_s^M), Y_s)) + \delta(\mathcal{L}_{cos}(x_s, X_s) + \mathcal{L}_{cos}(x_s^M, X_s^M))$$
(5)

where γ and δ are the weights of the auxiliary global loss and cosine similarity, and set as $\gamma=0.05$ and $\delta=0.025$ in our experiments. Similarly, the GLC loss is also calculated on the target data based on pseudo-label $f_{\theta}(X_t)$ and formulated as:

$$\mathcal{L}_{GLC}^{T} = 2\gamma \mathcal{L}_{Seg}(f_{\phi}(X_{t}^{M}), f_{\theta}(X_{t})) + 2\delta \mathcal{L}_{cos}(x_{t}^{M}, X_{t}^{M})$$
(6)

Therefore, the overall loss function of GLC is:

$$\mathcal{L}_{GLC} = \mathcal{L}_{GLC}^S + \mathcal{L}_{GLC}^T \tag{7}$$

With the regular fully-supervised segmentation loss on source data $\mathcal{L}_{FSS} = \beta \mathcal{L}_{Seg}(f_{\phi}(x_s), y_s)$, where β is defined as in Eq.2, the overall objective function \mathcal{L} for centralized UDA is formulated as:

$$\mathcal{L} = \mathcal{L}_{FSS} + \mathcal{L}_{MPL} + \mathcal{L}_{GLC} \tag{8}$$

It is clear that Eq.8 requires centralized and synchronous access to source and target data. In the section 3.5 and 3.6, we demonstrate how MAPSeg can be adapted to federated (decentralized and synchronous access to data) and test-time (decentralized and asynchronous access to data) UDA scenarios.

3.5. Extension to Federated UDA

In reality, labeled source-domain data and unlabeled target-domain data are often collected at different sites. We consider a practical scenario where a server (*e.g.* a major hospital) hosts potentially large amount of both labeled and unlabeled scans, and distributed clients (*e.g.* clinics or imaging sites) possess only unlabeled images. This is an underexplored scenario as FL typically assumes either fully or

partially labeled data from all clients. We extend MAPSeg to solve this federated multi-target UDA problem according to the details in Algorithm 1 of Appendix Sec. 1.2. Specifically, the server updates the student model f_{ϕ} by minimizing the loss for the labeled source-domain data D_S :

$$\mathcal{L}_{s} = \beta(\mathcal{L}_{seg}(f_{\phi}(x_{s}), y_{s}) + \mathcal{L}_{seg}(f_{\phi}(x_{s}^{M}), y_{s})) + \gamma(\mathcal{L}_{seg}(f_{\phi}(X_{s}), Y_{s}) + \mathcal{L}_{seg}(f_{\phi}(X_{s}^{M}), Y_{s})) + \delta(\mathcal{L}_{cos}(x_{s}, X_{s}) + \mathcal{L}_{cos}(x_{s}^{M}, X_{s}^{M}))$$
(9)

The clients update the student model f_{ϕ} by minimizing the loss for its own unlabeled target-domain data D_T^k :

$$\mathcal{L}_{u} = \beta(\mathcal{L}_{seg}(f_{\phi}(x_{t}^{M}), f_{\theta}(x_{t})) + \mathcal{L}_{seg}(f_{\phi}(x_{t}), f_{\theta}(x_{t})))$$

$$+ \gamma(\mathcal{L}_{seg}(f_{\phi}(X_{t}^{M}), f_{\theta}(X_{t})) + \mathcal{L}_{seg}(f_{\phi}(X_{t}), f_{\theta}(X_{t})))$$

$$+ \delta(\mathcal{L}_{cos}(x_{t}, X_{t}) + \mathcal{L}_{cos}(x_{t}^{M}, X_{t}^{M}))$$
(10)

Comparing to the centralized UDA loss (Eq.8), we decompose it into two components: fully supervised loss for server training (Eq.9) and self-supervised loss for client updates (Eq.10), which avoids the need for centralized data. After each local update, each client sends the EMA teacher model parameters θ to the server for aggregation following typical federated averaging [54].

3.6. Extension to Test-time UDA

Test-time UDA often involves two separate stages of training, including the source-only training at one center and the target-only finetuning at another site. In the federated UDA setting, Eq.9 and Eq.10 are jointly used to update the server model through synchronous federated averaging after each round. We can further ease the constraint of synchronous communication between source and target sites by training f_{ϕ} on the source data using Eq.9 for some (e.g. 1,000) warm-up steps before distributing the model parameters ϕ to the target site for initializing the teacher model f_{θ} . On the target site, f_{θ} provides stable pseudo-labels to guide the self-supervised training with Eq.10 and is updated by the EMA of ϕ following Eq.3. We find that in this asynchronous setting MAPSeg still performs well on the target-

domain data, albeit with a minor performance tradeoff on the source-domain data (see Tab.3).

3.7. Implementation Details

Model architecture and implementation. We implement the encoder backbone g using 3D-ResNet-like CNN. The segmentation decoder h is adapted from DeepLabV3 [6]. The framework is implemented using PyTorch. More details of the model and the training procedure are provided in Appendix Sec. 1.1 and Sec. 1.2.

Selecting the best model. For choosing the best model during training, some studies choose to train for fixed iterations and use the last checkpoint. On the other hand, some of the previous UDA studies [3, 4] face a dilemma in selecting the best model during training by validating against a holdout portion of target-domain labels, which is unrealistic as UDA assumes full absence of target labels. We demonstrate that MPL not only provides an efficient pathway to domain adaptative segmentation but also serves as an indicator of how well the model is being adapted to the target domain. We validate the model after each epoch and the best model is selected based on the score: $Score = Dice_{Src} - 0.5 \times$ $\overline{\mathcal{L}_{Seg}}(f_{\phi}(x_t^M), f_{\theta}(x_t)),$ where $Dice_{Src}$ is the Dice score on source-domain validation set and $\overline{\mathcal{L}_{Seg}}(f_{\phi}(x_t^M), f_{\theta}(x_t))$ is the mean of $\mathcal{L}_{Seg}(f_{\phi}(x_t^M), f_{\theta}(x_t))$ during the last training epoch. From Eq.1, it is clear that $\lim_{\hat{y}\to y} \mathcal{L}_{seq}(\hat{y},y) =$ -1, therefore, Score has an upper bound of 1.5. We demonstrate in Tab.4 that the difference between validation using target labels versus Score is acceptable (81.2 vs. 80.3). Even without accessing target labels for validation, MAPSeg still surpasses the previous SOTA results that use target labels for validation. It is worth noting that we only use target labels for validation in Tab.4 for a fair comparison with previously reported results; other results presented use Score for validation by default. For federated and test-time UDA, $Score = -\overline{\mathcal{L}_{Seq}}(f_{\phi}(x_t^M), f_{\theta}(x_t)).$

4. Experiments and Results

4.1. Datasets

Brain MRI Datasets. We include 2,421 (1,163 T1w) brain MRI scans acquired from newborn to toddler in this study. Among them, 2,306 are unannotated scans dedicated for the 3D multi-scale MAE pretraining. These MRI scans are acquired from multiple sites with different sequence parametrization and scanner types. All scans are preprocessed with skull stripping [30] and bias-field correction [69]. These MRI brain scans were acquired worldwide, and detailed descriptions can be found in Appendix Sec. 1.4.

To evaluate cross-sequence/site/age UDA segmentation for seven subcortical regions (*i.e.*, hippocampus (HC), amygdala (AD), caudate (CD), putamen (PT), pallidum

Table 1. Performance of centralized UDA on brain MRI segmentation.

		Cw	oc Coa					
Cross-Sequence Dice(%)↑								
Method	НС	AD	CD	PT	PD	TM	AB	Avg
AdvEnt[71]	56.7	52.7	66.7	66.1	61.8	74.1	40.1	59.8
DAFormer[32]	40.5	53.3	62.2	64.7	45.9	61.8	39.9	52.6
HRDA[33]	42.6	37.7	66.5	71.9	0.0	67.6	0.3	40.9
MIC[34]	40.3	47.0	72.5	52.9	0.0	62.1	0.0	39.3
DAR-UNet[79]	61.3	65.2	76.7	75.8	68.1	82.0	48.4	68.2
MAPSeg (Ours)	70.3	73.2	81.4	83.9	76.5	89.6	69.2	77.7
Cross-Site								
Method				Dice	(%)↑			
Menion	HC	AD	CD	PT	PD	TM	AB	Avg
AdvEnt[71]	27.1	6.7	21.0	23.1	12.5	36.0	20.5	21.0
DAFormer[32]	40.0	45.8	75.3	70.0	68.4	64.0	51.3	59.3
HRDA[33]	30.9	44.3	80.8	79.8	66.4	83.0	53.4	62.7
MIC[34]	48.1	36.2	67.7	82.8	69.5	66.8	52.3	60.5
DAR-UNet[79]	51.9	43.6	69.8	55.2	55.5	81.2	45.8	57.6
MAPSeg (Ours)	70.0	53.5	85.6	85.4	67.9	88.1	61.4	73.1
		-	Cross-A	ge				
				-	(%)↑			
Method	HC	AD	CD	PT	PD	TM	AB	Avg
AdvEnt[71]	58.7	54.1	44.0	63.8	56.9	78.0	30.9	55.2
DAFormer[32]	30.2	65.7	72.7	55.8	38.4	88.8	57.3	58.4
HRDA[33]	48.6	66.6	81.9	67.7	35.7	74.1	56.0	61.5
MIC[34]	61.3	66.0	80.9	73.4	44.3	76.1	51.0	64.7
DAR-UNet[79]	58.8	56.3	64.4	64.5	53.6	82.6	28.6	58.8
MAPSeg (Ours)	75.8	76.7	83.1	71.4	58.2	90.7	70.1	75.2

(PD), thalamus (TM), and accumbens (AB)), our analysis include manual segmentation of 115 scans. They comprise independent subjects from the BCP cohort (*BCP50*) with private expert segmentation for both T1w and T2w scans (acquired from 0 to 24 months postnatal age); 5 newborn scans from the ECHO cohort (*ECHO5*) with private expert segmentation; and 10 newborn scans from the M-CRIB project (*MCRIB10*) with publicly available segmentation [1].

Cardiac CT-MRI Dataset. Following the previous studies [3, 4], we include 40 independent scans (20 CT and 20 MRI) of cardiac regions from Multi-Modality Whole Heart Segmentation (MMWHS) Challenge 2017 dataset [53, 83, 84] with ground truth labels of ascending aorta (AA), left atrium blood cavity (LAC), left ventricle blood cavity (LVC), and myocardium of the left ventricle (MYO). Similarly, we apply bias-field correction to the MRI scans.

4.2. Dataset Partition

Pretraining. For multi-scale MAE pretraining on brain MRI scans, we have four models pretrained on different amounts of data to investigate the influence of pretraining data size. The model pretrained on large-scale data takes advantage of all 2,306 unannotated scans introduced in Sec. 4.1. Since there is no overlapping with the annotated scans, the pretrained model can be directly applied to all downstream UDA tasks (*i.e.*, cross-site/age/sequence). We also pretrain the model solely relying on source and target training data of each task.

For multi-scale MAE pretraining on cardiac CT-MRI

scans, the model is only pretrained on training scans of source (16 CT scans) and target (16 MRI scans) domains, following the partition adopted by previous studies.

Cross-Sequence UDA segmentation of brain. The model is trained on T1w MRI scans (source domain) and tested on T2w MRI scans (target domain). The *BCP50* dataset is randomly split into two non-overlapping subsets of 25 subjects per each. The model is trained on T1w scans of the first group (source domain 18 scans for training and 7 for validation) and T2w scans of the second group (target domain 15 for training and 10 for testing). The best validation model is then applied to the T2w testing scans.

Cross-Site UDA segmentation of brain. The model is trained on a single site (*BCP50*, source domain) and tested on two other sites (*MCRIB10* and *ECHO5*, target domains). Utilizing 50 T2w MRI scans from BCP as the source domain, we randomly select 40 scans for training and 10 for validation. Six scans from *MCRIB10* and three scans from *ECHO5* are used for UDA training, and remaining scans are used for testing.

Cross-Age UDA segmentation of brain. We also conduct experiments in cross-age segmentation using longitudinal scans from *BCP50*. We set the 24 T2w MRI scans of 12-24 month-old infant as the source domain and 14 T2w MRI scans of 0-6 month-old infants as the target domain. For the source domain, 19 scans are randomly sampled for training and remaining 5 scans are used for validation. For the target domain, 8 scans are used for UDA training and 6 scans are used for testing.

Cross-Modality UDA segmentation of cardiac. For the cardiac scans, for a fair comparison, we follow the same partition employed by the previous studies. We set CT as the source domain and MRI as the target domain, and use 16 CT scans and 16 MRI scans for training, 4 CT scans for validation, and the remaining 4 MRI scans for testing.

4.3. Results

Centralized Domain Adaptation. To assess MAPSeg's performance in different UDA tasks for infant brain MRI segmentation, we compare it with methods utilizing adversarial entropy minimization [71], image translation [79], and pseudo-labeling [32-34]. The results are reported in Tab.1. MAPSeg consistently outperforms its counterparts across all tasks. DAR-UNet ranks second in the cross-sequence task but shows degraded performance in others, partially due to translation error (details in Appendix). Among pseudo-labeling approaches, HRDA and MIC achieve the second best performance in cross-site and cross-age tasks, respectively. However, they fail to segment pallidum and accumbens in the cross-sequence task. A major challenge here is the small size of subcortical regions (accounting for approximately 2% of overall voxels) and significant class imbalance (e.g., thalamus comprises about

Table 2. Performance of federated UDA on brain MRI segmentation.

Method	Dice(%) ↑				
Wichiod	Cross-Sequence	Cross-Site	Cross-Age		
FAT[55]	27.6	63.8	69.0		
DualAdapt[78]	28.4	66.1	54.8		
Fed-MAPSeg (ours)	69.9	73.6	71.0		

Table 3. Comparison between centralized and test-time UDA on brain MRI segmentation. Performance of source domain are reported on source validation set.

Task	Centralized UDA		Test-time UDA		Λ		
lask	Source	Target	Source	Target	Δ_{Source}	Δ_{Target}	
X-seq	84.0	77.7	79.2	75.9	-4.8	-1.8	
X-age	85.8	75.2	84.2	72.9	-1.6	-2.3	
X-site	85.7	73.1	79.9	70.3	-5.8	-2.8	

0.8% of overall voxels, while accumbens accounts for only 0.03%). This imbalance poses a significant challenge for previous pseudo-labeling methods. Additional visualizations and discussions are available in Appendix Sec. 1.7.

Federated Domain Adaptation. To evaluate our framework in the federated domain adaptation setting, we designate the labeled source-domain dataset as the server dataset and the unlabeled target-domain datasets as the client datasets. In the cross-sequence setting, the 25 T1w scans of the first group are considered as the server dataset, and the 25 T2w scans of the second group are split roughly equally into three disjoint client datasets. In the cross-site setting, the *BCP50* is considered as the server dataset, and the *ECHO5* and *MCRIB10* naturally serve as two different client datasets. In the cross-age setting, we treat the scans from the first age group as the server dataset, and split the scans from the second age group equally into two client datasets.

We compare our Fed-MAPSeg with two other related work, FAT [55] and DualAdapt [78]. To our best knowledge, there is no direct comparison from the literature that addresses this challenging federated multi-target unsupervised domain adaptation for 3D medical image segmentation. FAT [55] proposes an alternating training scheme between the labeled and unlabeled data silos and adopts a mixup approach to augment the unlabeled input data for self-supervised learning with pseudo-labels. Adapt [78] considers a similar single-source to multi-target unsupervised domain adaptation setting, except that it only reports segmentation performance for 2D image datasets such as the DomainNet [57] and CrossCity [9]. Implementation details for our Fed-MAPSeg as well as the baselines are included in Appendix Sec. 1.3. We report our results in Tab.2. Fed-MAPSeg not only outperforms the two baselines by a large margin (esp. in the the cross-sequence setting), it also maintains a fairly close performance compared to the centralized UDA.

Table 4. Performance of centralized UDA on cardiact $CT \rightarrow MRI$ segmentation. Underline indicates the target labels are not used for validation.

Cardiac $CT \rightarrow MRI$ segmentation							
Method		Dice(%)↑					
Method	AA	LAC	LVC	MYO	Avg		
PnP-AdaNet[15]	43.7	47.0	77.7	48.6	54.3		
SIFA-V1[3]	67.0	60.7	75.1	45.8	62.1		
SIFA-V2[4]	65.3	62.3	78.9	47.3	63.4		
DAFormer[32]	75.2	59.4	72.0	57.1	65.9		
MPSCL[48]	62.8	76.1	80.5	55.1	68.6		
MA-UDA[37]	71.0	67.4	77.5	57.1	68.7		
SE-ASA[18]	68.3	74.6	81.0	55.9	69.9		
FSUDA-V1[45]	62.4	72.1	81.2	66.5	70.6		
PUFT[14]	69.3	77.4	83.0	63.6	73.3		
SDUDA[12]	72.8	79.3	82.3	64.7	74.8		
FSUDA-V2[46]	72.5	78.6	82.6	68.4	75.5		
MAPSeg (Ours)	78.5	81.8	92.1	68.8	80.3		
wiar seg (Ours)	78.2	81.8	92.9	72.0	81.2		

Table 5. Ablation studies of MAPSeg components on cross-sequence brain MRI segmentation.

Co	mponer	Performance	
MAE	GLC	MPL	Dice(%) ↑
			31.6
\checkmark			51.3
	✓		53.0
		✓	39.5
	✓	✓	59.0
\checkmark	✓		71.3
\checkmark		✓	75.3
\checkmark	✓	✓	77.7

Test-Time Domain Adaptation. We further extend MAPSeg to Test-time UDA, and the results for different tasks are reported in Tab.3. With decentralized data and asynchronous training, MAPSeg still performs very well in all tasks, with performance drop smaller than 3% in the target domain. However, we observe a slightly more performance degradation in the source domain (Tab.3), particularly in cross-sequence and cross-site tasks, suggesting that the model suffers from forgetting of the source domain knowledge during test-time UDA.

Cross-Modality Segmentation of Cardiac. To evaluate the generalizability of MAPSeg, we further conduct experiment for cross-modality cardiac segmentation and the results are reported in Tab.4. MAPSeg surpasses all previously reported results. Results of MRI \rightarrow CT segmentation can be found in Appendix Sec. 1.5.

Ablation Studies. To further investigate each component of MAPSeg, we conduct ablation studies focusing on MAE, GLC, MPL, masking ratio, masking patch size of local patch, and pretraining data size in the context of cross-sequence segmentation. From Tab.5, it is clear that directly applying MPL only brings a minor improvement, suggesting using MPL alone suffers from pseudo-label drifts. By incorporating GLC to leverage global-local contexts, MPL



Figure 3. Ablation studies on masking ratio, patch size, and pretrain data. Experiments on masking ratio and patch size are conducted on cross-sequence task.

yields better results. MAE pretraining significantly boosts the performance from using MPL alone (39.5 to 75.3), justifying MAE and MPL are complementary parts in MAPSeg. Combining MAE, MPL, and GLC together yields the optimal performance.

The impact of masking ratio and local patch size is reported in Fig.3. The masking ratio and patch size remain the same in MAE and MPL. The results indicate that MAPSeg is more sensitive to patch size. A patch size of 4 or 16 decreases the performance significantly. For the masking ratio, MAPSeg achieves optimal performance when 70% of the regions are masked out. Additionally, we evaluate model's performance using only source and target training data (< 50 scans) for MAE pretraining, much fewer than the large-scale pretraining (> 2,000 scans). This suggests that, even with dozens of scans involved in MAE, MAPSeg still delivers comparable performance. Another benefit of largescale pretraining is its immediate applicability to new target domains; the pretrained encoders can be directly employed for MPL, bypassing the need for training from scratch. Additional analyses about sensitivity to other hyperparameters can be found in Appendix Sec. 1.6.

5. Conclusions

In this paper, we introduce the MAPSeg framework as a unified UDA framework that works on centralized, federated, and test-time UDA scenarios. We evaluate it under multiple domain shift and adaptation settings, and it outperforms all the baselines in all scenarios. We conduct extensive ablation study to demonstrate the effectiveness of each component.

6. Acknowledgements

This work was supported by NIH grants R00HD103912 (Y.W.), R01HL121270 (R.G.B. & A.F.L.), R01MH121070 (J.P. & A.P.J.), and NSF grant CNS-2112562 (H.L.), as well as by Duke Science and Technology (Y.W. & H.L.).

References

- [1] Bonnie Alexander, Andrea L Murray, Wai Yen Loh, Lillian G Matthews, Chris Adamson, Richard Beare, Jian Chen, Claire E Kelly, Sandra Rees, Simon K Warfield, et al. A new neonatal cortical and subcortical brain atlas: the melbourne children's regional infant brain (m-crib) atlas. *Neuroimage*, 147:841–851, 2017. 6, 3
- [2] Danielle A Baribeau, Annie Dupuis, Tara A Paton, Christopher Hammill, Stephen W Scherer, Russell J Schachar, Paul D Arnold, Peter Szatmari, Rob Nicolson, Stelios Georgiades, et al. Structural neuroimaging correlates of social deficits are similar in autism spectrum disorder and attention-deficit/hyperactivity disorder: analysis from the pond network. Translational psychiatry, 9(1):72, 2019. 1
- [3] Cheng Chen, Qi Dou, Hao Chen, Jing Qin, and Pheng-Ann Heng. 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, 2019. 2, 3, 6, 8
- [4] Cheng Chen, Qi Dou, Hao Chen, Jing Qin, and Pheng Ann Heng. 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. 2, 6, 8, 3
- [5] Cheng Chen, Quande Liu, Yueming Jin, Qi Dou, and Pheng-Ann Heng. Source-free domain adaptive fundus image segmentation with denoised pseudo-labeling. In Medical Image Computing and Computer Assisted Intervention–MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part V 24, pages 225–235. Springer, 2021. 3
- [6] Liang-Chieh Chen, George Papandreou, Florian Schroff, and Hartwig Adam. Rethinking atrous convolution for semantic image segmentation. arXiv preprint arXiv:1706.05587, 2017. 6, 1
- [7] Mark Chen, Alec Radford, Rewon Child, Jeffrey Wu, Heewoo Jun, David Luan, and Ilya Sutskever. Generative pretraining from pixels. In *Proceedings of the 37th Interna*tional Conference on Machine Learning, pages 1691–1703. PMLR, 2020. 2
- [8] Wuyang Chen, Ziyu Jiang, Zhangyang Wang, Kexin Cui, and Xiaoning Qian. Collaborative global-local networks for memory-efficient segmentation of ultra-high resolution images. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2019. 4
- [9] Y. Chen, W. Chen, Y. Chen, B. Tsai, Y. Wang, and M. Sun. No more discrimination: Cross city adaptation of road scene segmenters. In 2017 IEEE International Conference on Computer Vision (ICCV), pages 2011–2020, Los Alamitos, CA, USA, 2017. IEEE Computer Society. 7
- [10] Boris Chidlovskii, Stephane Clinchant, and Gabriela Csurka. Domain adaptation in the absence of source domain data. In Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, page 451–460, New York, NY, USA, 2016. Association for Computing Machinery. 3

- [11] K Guadalupe Cruz, Yi Ning Leow, Nhat Minh Le, Elie Adam, Rafiq Huda, and Mriganka Sur. Cortical-subcortical interactions in goal-directed behavior. *Physiological reviews*, 103(1):347–389, 2023. 1
- [12] Zhiming Cui, Changjian Li, Zhixu Du, Nenglun Chen, Guodong Wei, Runnan Chen, Lei Yang, Dinggang Shen, and Wenping Wang. Structure-driven unsupervised domain adaptation for cross-modality cardiac segmentation. *IEEE Transactions on Medical Imaging*, 40(12):3604–3616, 2021. 2, 8, 3
- [13] Carl Doersch, Abhinav Gupta, and Alexei A. Efros. Unsupervised visual representation learning by context prediction. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2015. 2
- [14] Shunjie Dong, Zixuan Pan, Yu Fu, Dongwei Xu, Kuangyu Shi, Qianqian Yang, Yiyu Shi, and Cheng Zhuo. Partial unbalanced feature transport for cross-modality cardiac image segmentation. *IEEE Transactions on Medical Imaging*, 42 (6):1758–1773, 2023. 3, 8
- [15] Qi Dou, Cheng Ouyang, Cheng Chen, Hao Chen, Ben Glocker, Xiahai Zhuang, and Pheng-Ann Heng. Pnp-adanet: Plug-and-play adversarial domain adaptation network at unpaired cross-modality cardiac segmentation. *IEEE Access*, 7:99065–99076, 2019. 3, 8
- [16] Liang Du, Jingang Tan, Hongye Yang, Jianfeng Feng, Xiangyang Xue, Qibao Zheng, Xiaoqing Ye, and Xiaolin Zhang. Ssf-dan: Separated semantic feature based domain adaptation network for semantic segmentation. In Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV), 2019. 3
- [17] A David Edwards, Daniel Rueckert, Stephen M Smith, Samy Abo Seada, Amir Alansary, Jennifer Almalbis, Joanna Allsop, Jesper Andersson, Tomoki Arichi, Sophie Arulkumaran, et al. The developing human connectome project neonatal data release. Frontiers in neuroscience, 16, 2022.
- [18] Wei Feng, Lie Ju, Lin Wang, Kaimin Song, Xin Zhao, and Zongyuan Ge. Unsupervised domain adaptation for medical image segmentation by selective entropy constraints and adaptive semantic alignment. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(1):623–631, 2023. 3, 8
- [19] Yossi Gandelsman, Yu Sun, Xinlei Chen, and Alexei Efros. Test-time training with masked autoencoders. In Advances in Neural Information Processing Systems, pages 29374– 29385. Curran Associates, Inc., 2022. 3
- [20] John H Gilmore, Rebecca C Knickmeyer, and Wei Gao. Imaging structural and functional brain development in early childhood. *Nature Reviews Neuroscience*, 19(3):123–137, 2018. 2
- [21] Yves Grandvalet and Yoshua Bengio. Semi-supervised learning by entropy minimization. In Advances in Neural Information Processing Systems. MIT Press, 2004. 3
- [22] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent-a new approach

- to self-supervised learning. Advances in neural information processing systems, 33:21271–21284, 2020. 4
- [23] Hao Guan and Mingxia Liu. Federated learning for medical image analysis: A survey. CoRR, abs/2306.05980, 2023. 3
- [24] Heather Cody Hazlett, Hongbin Gu, Brent C Munsell, Sun Hyung Kim, Martin Styner, Jason J Wolff, Jed T Elison, Meghan R Swanson, Hongtu Zhu, Kelly N Botteron, et al. Early brain development in infants at high risk for autism spectrum disorder. *Nature*, 542(7641):348–351, 2017.
- [25] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [26] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16000–16009, 2022. 2, 1
- [27] Yufan He, Aaron Carass, Lianrui Zuo, Blake E Dewey, and Jerry L Prince. Autoencoder based self-supervised test-time adaptation for medical image analysis. *Medical image analysis*, 72:102136, 2021. 3
- [28] Olivier Henaff. Data-efficient image recognition with contrastive predictive coding. In *Proceedings of the 37th International Conference on Machine Learning*, pages 4182–4192. PMLR, 2020. 2
- [29] 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 *Proceedings of the 35th International Conference on Machine Learning*, pages 1989–1998. PMLR, 2018. 3
- [30] Andrew Hoopes, Jocelyn S Mora, Adrian V Dalca, Bruce Fischl, and Malte Hoffmann. Synthstrip: skull-stripping for any brain image. *NeuroImage*, 260:119474, 2022. 6
- [31] Brittany R Howell, Martin A Styner, Wei Gao, Pew-Thian Yap, Li Wang, Kristine Baluyot, Essa Yacoub, Geng Chen, Taylor Potts, Andrew Salzwedel, et al. The unc/umn baby connectome project (bcp): An overview of the study design and protocol development. *NeuroImage*, 185:891–905, 2019. 3
- [32] Lukas Hoyer, Dengxin Dai, and Luc Van Gool. Daformer: Improving network architectures and training strategies for domain-adaptive semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9924–9935, 2022. 3, 6, 7, 8
- [33] Lukas Hoyer, Dengxin Dai, and Luc Van Gool. Hrda: Context-aware high-resolution domain-adaptive semantic segmentation. In *Computer Vision – ECCV 2022*, pages 372–391, Cham, 2022. Springer Nature Switzerland. 6
- [34] Lukas Hoyer, Dengxin Dai, Haoran Wang, and Luc Van Gool. Mic: Masked image consistency for context-enhanced domain adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11721–11732, 2023. 3, 6, 7
- [35] Zijian Hu, Zhengyu Yang, Xuefeng Hu, and Ram Nevatia. Simple: Similar pseudo label exploitation for semisupervised classification. In *Proceedings of the IEEE/CVF*

- Conference on Computer Vision and Pattern Recognition (CVPR), pages 15099–15108, 2021. 2
- [36] Jiaxing Huang, Shijian Lu, Dayan Guan, and Xiaobing Zhang. Contextual-relation consistent domain adaptation for semantic segmentation. In *Computer Vision – ECCV 2020*, pages 705–722, Cham, 2020. Springer International Publishing. 3
- [37] Wen Ji and Albert C. S. Chung. Unsupervised domain adaptation for medical image segmentation using transformer with meta attention. *IEEE Transactions on Medical Imaging*, pages 1–1, 2023. 8, 3
- [38] Peter Kairouz, H. Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, Rafael G. L. D'Oliveira, Hubert Eichner, Salim El Rouayheb, David Evans, Josh Gardner, Zachary Garrett, Adrià Gascón, Badih Ghazi, Phillip B. Gibbons, Marco Gruteser, Zaid Harchaoui, Chaoyang He, Lie He, Zhouyuan Huo, Ben Hutchinson, Justin Hsu, Martin Jaggi, Tara Javidi, Gauri Joshi, Mikhail Khodak, Jakub Konecný, Aleksandra Korolova, Farinaz Koushanfar, Sanmi Koyejo, Tancrède Lepoint, Yang Liu, Prateek Mittal, Mehryar Mohri, Richard Nock, Ayfer Özgür, Rasmus Pagh, Hang Qi, Daniel Ramage, Ramesh Raskar, Mariana Raykova, Dawn Song, Weikang Song, Sebastian U. Stich, Ziteng Sun, Ananda Theertha Suresh, Florian Tramèr, Praneeth Vepakomma, Jianyu Wang, Li Xiong, Zheng Xu, Qiang Yang, Felix X. Yu, Han Yu, and Sen Zhao. Advances and open problems in federated learning. Foundations and Trends® in Machine Learning, 14(1-2):1-210, 2021. 2
- [39] Ioannis Kakogeorgiou, Spyros Gidaris, Bill Psomas, Yannis Avrithis, Andrei Bursuc, Konstantinos Karantzalos, and Nikos Komodakis. What to hide from your students: Attention-guided masked image modeling. In *Computer Vision ECCV 2022*, pages 300–318, Cham, 2022. Springer Nature Switzerland. 2
- [40] Neerav Karani, Ertunc Erdil, Krishna Chaitanya, and Ender Konukoglu. Test-time adaptable neural networks for robust medical image segmentation. *Medical Image Analysis*, 68: 101907, 2021. 3
- [41] Xiangwen Kong and Xiangyu Zhang. Understanding masked image modeling via learning occlusion invariant feature. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6241–6251, 2023. 2
- [42] Samuli Laine and Timo Aila. Temporal ensembling for semisupervised learning. arXiv preprint arXiv:1610.02242, 2016.
- [43] Dong-Hyun Lee et al. Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks. In Workshop on challenges in representation learning, ICML, page 896. Atlanta, 2013. 2
- [44] Ang Li, Jingwei Sun, Xiao Zeng, Mi Zhang, Hai Li, and Yiran Chen. Fedmask: Joint computation and communication-efficient personalized federated learning via heterogeneous masking. In *Proceedings of the 19th ACM Conference on Embedded Networked Sensor Systems*, page 42–55, New

- York, NY, USA, 2021. Association for Computing Machinery. 3
- [45] Shaolei Liu, Siqi Yin, Linhao Qu, and Manning Wang. Reducing domain gap in frequency and spatial domain for cross-modality domain adaptation on medical image segmentation. *Proceedings of the AAAI Conference on Artificial Intelligence*, 37(2):1719–1727, 2023. 8, 3
- [46] Shaolei Liu, Siqi Yin, Linhao Qu, Manning Wang, and Zhijian Song. A structure-aware framework of unsupervised cross-modality domain adaptation via frequency and spatial knowledge distillation. *IEEE Transactions on Medical Imag*ing, pages 1–1, 2023. 8, 3
- [47] Yuang Liu, Wei Zhang, and Jun Wang. Source-free domain adaptation for semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1215–1224, 2021. 3
- [48] Zhizhe Liu, Zhenfeng Zhu, Shuai Zheng, Yang Liu, Jiayu Zhou, and Yao Zhao. Margin preserving self-paced contrastive learning towards domain adaptation for medical image segmentation. *IEEE Journal of Biomedical and Health Informatics*, 26(2):638–647, 2022. 8, 3
- [49] Mingsheng Long, Yue Cao, Jianmin Wang, and Michael Jordan. Learning transferable features with deep adaptation networks. In *Proceedings of the 32nd International Conference on Machine Learning*, pages 97–105, Lille, France, 2015. PMLR. 3
- [50] Mingsheng Long, Han Zhu, Jianmin Wang, and Michael I Jordan. Unsupervised domain adaptation with residual transfer networks. In *Advances in Neural Information Processing Systems*. Curran Associates, Inc., 2016. 3
- [51] Ilya Loshchilov and Frank Hutter. SGDR: Stochastic gradient descent with warm restarts. In *International Conference* on Learning Representations, 2017. 2
- [52] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019. 2
- [53] Xinzhe Luo and Xiahai Zhuang. X-metric: An n-dimensional information-theoretic framework for groupwise registration and deep combined computing. IEEE Transactions on Pattern Analysis and Machine Intelligence, 45(7): 9206–9224, 2023. 6
- [54] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communicationefficient learning of deep networks from decentralized data. In Artificial intelligence and statistics, pages 1273–1282. PMLR, 2017. 3, 5
- [55] Erum Mushtaq, Yavuz Faruk Bakman, Jie Ding, and Salman Avestimehr. Federated alternate training (fat): Leveraging unannotated data silos in federated segmentation for medical imaging. In 20th IEEE International Symposium on Biomedical Imaging, ISBI 2023, Cartagena, Colombia, April 18-21, 2023, pages 1–5. IEEE, 2023. 3, 7
- [56] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A. Efros. Context encoders: Feature learning by inpainting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 2

- [57] Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment matching for multi-source domain adaptation. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1406–1415, 2019. 7
- [58] Andra Petrovai and Sergiu Nedevschi. Exploiting pseudo labels in a self-supervised learning framework for improved monocular depth estimation. In *Proceedings of* the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 1578–1588, 2022. 2
- [59] Fernando Pérez-García, Rachel Sparks, and Sébastien Ourselin. Torchio: A python library for efficient loading, preprocessing, augmentation and patch-based sampling of medical images in deep learning. Computer Methods and Programs in Biomedicine, 208:106236, 2021.
- [60] Mehdi Sajjadi, Mehran Javanmardi, and Tolga Tasdizen. Regularization with stochastic transformations and perturbations for deep semi-supervised learning. In Advances in Neural Information Processing Systems. Curran Associates, Inc., 2016. 2
- [61] Micah J. Sheller, G. Anthony Reina, Brandon Edwards, Jason Martin, and Spyridon Bakas. Multi-institutional deep learning modeling without sharing patient data: A feasibility study on brain tumor segmentation. In *Brainlesion: Glioma, Multiple Sclerosis, Stroke and Traumatic Brain Injuries*, pages 92–104, Cham, 2019. Springer International Publishing. 3
- [62] Mark D Shen, Meghan R Swanson, Jason J Wolff, Jed T Elison, Jessica B Girault, Sun Hyung Kim, Rachel G Smith, Michael M Graves, Leigh Anne H Weisenfeld, Lisa Flake, et al. Subcortical brain development in autism and fragile x syndrome: evidence for dynamic, age-and disorder-specific trajectories in infancy. *American Journal of Psychiatry*, 179 (8):562–572, 2022.
- [63] Carole H. Sudre, M. Jorge Cardoso, and Sebastien Ourselin. Longitudinal segmentation of age-related white matter hyperintensities. *Medical Image Analysis*, 38:50–64, 2017.
- [64] Yue Sun, Kun Gao, Zhengwang Wu, Guannan Li, Xiaopeng Zong, Zhihao Lei, Ying Wei, Jun Ma, Xiaoping Yang, Xue Feng, Li Zhao, Trung Le Phan, Jitae Shin, Tao Zhong, Yu Zhang, Lequan Yu, Caizi Li, Ramesh Basnet, M. Omair Ahmad, M. N. S. Swamy, Wenao Ma, Qi Dou, Toan Duc Bui, Camilo Bermudez Noguera, Bennett Landman, Ian H. Gotlib, Kathryn L. Humphreys, Sarah Shultz, Longchuan Li, Sijie Niu, Weili Lin, Valerie Jewells, Dinggang Shen, Gang Li, and Li Wang. Multi-site infant brain segmentation algorithms: The iseg-2019 challenge. IEEE Transactions on Medical Imaging, 40(5):1363–1376, 2021. 2
- [65] Minxue Tang, Xuefei Ning, Yitu Wang, Jingwei Sun, Yu Wang, Hai Li, and Yiran Chen. Fedcor: Correlation-based active client selection strategy for heterogeneous federated learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10102–10111, 2022. 3
- [66] Yucheng Tang, Dong Yang, Wenqi Li, Holger R. Roth, Bennett Landman, Daguang Xu, Vishwesh Nath, and Ali Hatamizadeh. Self-supervised pre-training of swin transformers for 3d medical image analysis. In *Proceedings of*

- the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 20730–20740, 2022. 2
- [67] Antti Tarvainen and Harri Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. In Advances in Neural Information Processing Systems. Curran Associates, Inc., 2017. 2, 4
- [68] Nguyen Binh Truong, Kai Sun, Siyao Wang, Florian Guitton, and Yike Guo. Privacy preservation in federated learning: Insights from the GDPR perspective. *CoRR*, abs/2011.05411, 2020. 2
- [69] Nicholas J Tustison, Brian B Avants, Philip A Cook, Yuanjie Zheng, Alexander Egan, Paul A Yushkevich, and James C Gee. N4itk: improved n3 bias correction. *IEEE transactions* on medical imaging, 29(6):1310–1320, 2010. 6
- [70] K. Van Leemput, F. Maes, D. Vandermeulen, and P. Suetens. Automated model-based tissue classification of mr images of the brain. *IEEE Transactions on Medical Imaging*, 18(10): 897–908, 1999.
- [71] T. Vu, H. Jain, M. Bucher, M. Cord, and P. Perez. Advent: Adversarial entropy minimization for domain adaptation in semantic segmentation. In 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 2512–2521, Los Alamitos, CA, USA, 2019. IEEE Computer Society. 3, 6, 7
- [72] Jeffry Wicaksana, Zengqiang Yan, Dong Zhang, Xijie Huang, Huimin Wu, Xin Yang, and Kwang-Ting Cheng. Fedmix: Mixed supervised federated learning for medical image segmentation. *IEEE Transactions on Medical Imag*ing, 42(7):1955–1968, 2023. 3
- [73] Adrian Wolny, Lorenzo Cerrone, Athul Vijayan, Rachele Tofanelli, Amaya Vilches Barro, Marion Louveaux, Christian Wenzl, Sören Strauss, David Wilson-Sánchez, Rena Lymbouridou, et al. Accurate and versatile 3d segmentation of plant tissues at cellular resolution. *Elife*, 9:e57613, 2020. 1
- [74] Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Jianmin Bao, Zhuliang Yao, Qi Dai, and Han Hu. Simmin: A simple framework for masked image modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9653–9663, 2022. 2
- [75] Zhenda Xie, Zigang Geng, Jingcheng Hu, Zheng Zhang, Han Hu, and Yue Cao. Revealing the dark secrets of masked image modeling. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 14475–14485, 2023.
- [76] Zhenda Xie, Zheng Zhang, Yue Cao, Yutong Lin, Yixuan Wei, Qi Dai, and Han Hu. On data scaling in masked image modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10365–10374, 2023.
- [77] Dong Yang, Ziyue Xu, Wenqi Li, Andriy Myronenko, Holger R. Roth, Stephanie Harmon, Sheng Xu, Baris Turkbey, Evrim Turkbey, Xiaosong Wang, Wentao Zhu, Gianpaolo Carrafiello, Francesca Patella, Maurizio Cariati, Hirofumi Obinata, Hitoshi Mori, Kaku Tamura, Peng An, Bradford J. Wood, and Daguang Xu. Federated semi-supervised learning for covid region segmentation in chest ct using multi-

- national data from china, italy, japan. *Medical Image Analysis*, 70:101992, 2021. 3
- [78] Chun-Han Yao, Boqing Gong, Hang Qi, Yin Cui, Yukun Zhu, and Ming-Hsuan Yang. Federated multi-target domain adaptation. In 2022 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV), pages 1081–1090, 2022. 3, 7
- [79] Kai Yao, Zixian Su, Kaizhu Huang, Xi Yang, Jie Sun, Amir Hussain, and Frans Coenen. A novel 3d unsupervised domain adaptation framework for cross-modality medical image segmentation. *IEEE Journal of Biomedical and Health Informatics*, 26(10):4976–4986, 2022. 2, 3, 6, 7
- [80] Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *International Conference on Learning Representa*tions, 2018. 3
- [81] Pan Zhang, Bo Zhang, Ting Zhang, Dong Chen, Yong Wang, and Fang Wen. Prototypical pseudo label denoising and target structure learning for domain adaptive semantic segmentation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 12414–12424, 2021. 3
- [82] Zhedong Zheng and Yi Yang. Rectifying pseudo label learning via uncertainty estimation for domain adaptive semantic segmentation. *International Journal of Computer Vision*, 129(4):1106–1120, 2021. 3
- [83] Xiahai Zhuang. Multivariate mixture model for myocardial segmentation combining multi-source images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41 (12):2933–2946, 2019. 6
- [84] Xiahai Zhuang and Juan Shen. Multi-scale patch and multi-modality atlases for whole heart segmentation of mri. *Medical Image Analysis*, 31:77–87, 2016. 6
- [85] Yang Zou, Zhiding Yu, Xiaofeng Liu, B.V.K. Vijaya Kumar, and Jinsong Wang. Confidence regularized self-training. In Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV), 2019. 3

MAPSeg: Unified Unsupervised Domain Adaptation for Heterogeneous Medical Image Segmentation Based on 3D Masked Autoencoding and Pseudo-Labeling

Supplementary Material

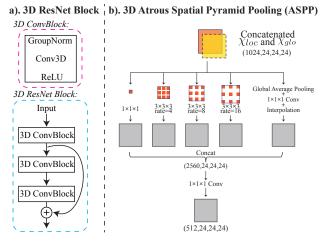
1. Appendix

1.1. Model Architecture

MAPSeg is implemented using PyTorch. Detailed configurations of model and training can be found below.

3D Multi-Scale Masked Autoencoder (MAE). We implement the 3D MAE using 3D ResNet Blocks [25, 73] instead of Vision Transformers, different from the previous study [26], due to the constraint of GPU memory. The encoder consists of eight 3D ResNet Blocks. The 3D ResNet Block is depicted in Suppl.Fig.1a. Following the previous study [26], we adopt an asymmetric design by employing a lightweight decoder (Suppl.Tab.1).

3D Global-Local Collaboration (GLC). The segmentation backbone (Suppl.Tab.1) consists of the pretrained encoder and a segmentation decoder that is adapted from DeepLabV3 [6]. In the decoding path, we take advantage of the Atrous Spatial Pyramid Pooling (ASPP), which employs dilated convolution at multiple scales and provides access to larger FOV (Suppl.Fig.1b). After feature extraction, the GLC module fuses the local and global features and forms a latent representation with a dimension of 1024, which is then fed into the ASPP layer. During training, each local sub-volume with size of $96 \times 96 \times 96$ is randomly sampled from global scan. During inference, the final output is formed by sliding window with stride of 80 across entire volumetric scan.



Suppl.Fig. 1. Illustrations of 3D ResNet Block and 3D Atrous Spatial Pyramid Pooling (ASPP) layer.

Encoder								
Layer Name	Input Size	Output Size	Architecture					
enc_res1	(1,96,96,96)	(512,24,24,24)	$\begin{bmatrix} 4 \times 4 \times 4, 512 \\ 3 \times 3 \times 3, 512 \\ 3 \times 3 \times 3, 512 \end{bmatrix} \times 1$					
enc_res2.x	(512,24,24,24)	(512,24,24,24)	$\begin{bmatrix} 3 \times 3 \times 3, 512 \\ 3 \times 3 \times 3, 512 \\ 3 \times 3 \times 3, 512 \end{bmatrix} \times 7$					
MAE Decoder								
Layer name	Input size	Output size	Architecture					
trans_conv1	(512,24,24,24)	(32,96,96,96)	4×4×4, 32, stride 4					
dec_res1	(32,96,96,96)	(16,96,96,96)	$\begin{bmatrix} 3 \times 3 \times 3, 16 \\ 3 \times 3 \times 3, 16 \\ 3 \times 3 \times 3, 16 \end{bmatrix} \times 1$					
final_recon	(16,96,96,96)	(1,96,96,96)	3×3×3, 1, stride 1					
	Seg	mentation Decoder						
Layer name	Input size	Output size	Architecture					
ASPP	(1024,24,24,24)	(512,24,24,24)	Suppl. Fig.1b					
trans_conv2	(512,24,24,24)	(64,96,96,96)	4×4×4, 64, stride 4					
seg_head	(64,96,96,96)	(cls_num,96,96,96)	$\begin{bmatrix} 3 \times 3 \times 3, 64 \\ 1 \times 1 \times 1, \text{ cls_num} \end{bmatrix} \times 1$					

Suppl.Tab. 1. Architectures of different components of MAPSeg. Building blocks ([kernal size, output channels]) are shown in brackets, with the number of blocks stacked. Downsampling is performed by the first block of enc_res1 with a stride of 4.

1.2. Training Recipe

MAE Pretraining. For the MAE Pretraining, we follow the training configurations listed in Suppl.Tab.2. Each minibatch contains a pair of randomly sampled local patch x and downsampled global scan X. The masking patch in Suppl.Tab.2 only applies to x and is always half-sized for X because of the larger FOV. For example, in the ablation study of masking patch size, a masking patch of 16 to x indicates a masking patch of 8 to x. We implement the augmentation using TorchIO [59]. During the MAE stage, we employ random 3D affine transformation, with isotropic scaling 75-150% and rotation [-40°, 40°].

Centralized UDA. For the centralized UDA on brain MRI segmentation tasks, detailed training configuration can be found in Suppl.Tab.3. Similarly, each mini-batch contains a pair of x and X from the source domain and another pair from the target domain (four $96 \times 96 \times 96$ patches). During warmup epochs, the model is only trained on source domain. We utilize *Score* to select the best model and the patience is set as 50 epochs. For the target domain, we design a similar random 3D affine transformation, with scaling 70-130% and rotation $[-30^{\circ}, 30^{\circ}]$. A stronger augmentation

config	value			
masking patch	8×8×8			
masking ratio	70%			
optimizer	AdamW [52]			
learning rate	$2e^{-4}$			
weight decay	0.05			
optim. momentum	$\beta_1, \beta_2 = 0.9, 0.95$			
lr scheduler	cosine annealing [51]			
ii schedulei	T_{max} =20, min_lr=1 e^{-6}			
total epochs	300			
annealing epochs	last 100			
batch size	4			
iters/epoch	500			
aug. prob.	0.35			
augmentation	random affine			

Suppl.Tab. 2. MAE Pretraining Configurations

config	value
masking patch	8×8×8
masking ratio	70%
optimizer	AdamW [52]
learning rate	$1e^{-4}$
weight decay	0.01
optim. momentum	$\beta_1, \beta_2 = 0.9, 0.999$
lr scheduler	cosine annealing warm restart [51]
ii schedulei	T_0 =10, T_{mult} =2, min_lr=1 e^{-8}
total epochs	100
warmup epochs	first 10
annealing epochs	all
early stop	50
batch size	1
iters/epoch	100
aug. prob.	0.35
	random affine
source aug.	random bias field
	random gamma trans.
target aug.	random affine

Suppl.Tab. 3. Centralized UDA configurations for brain MRI segmentation.

strategy is applied to the source domain, consisting of random affine (scaling 70-140% and rotation [-30°, 30°]), random bias field [63, 70], and random gamma transformation ($\gamma \in [e^{-0.4}, e^{0.4}]$). For the centralized UDA on public cardiac CT \rightarrow MRI segmentation, we use the same configuration except for training epochs of 150 and warmup epochs of 50. For MRI \rightarrow CT cardiac segmentation, we use a less aggressive augmentation strategy because MRI is noisier than CT. We set the scaling ratio to 85-115% and rotation to [-15°, 15°] for both source and target domains, and exclude random bias field and gamma transformation. The warmup epoch is set as 70.

Federated UDA. For the federated UDA tasks, we follow the procedure detailed in Algorithm 1. We initialize the encoder of the global model f_ϕ with the encoder pretrained on the large-scale data mentioned in Sec.4.3. We set the global FL round R=100. We set both the server and client update steps to 1 epoch with batch size of 1. Training configuration inherits mostly from that of the centralized UDA, except a global cosine annealing learning rate schedule is adopted to decay the learning rate from $1e^{-4}$ to $1e^{-6}$ over the course of the FL rounds.

Test-Time UDA. For the test-time UDA tasks, we follow

Algorithm 1 Federated MAPSeg (Fed-MAPSeg)

Require: Source domain dataset $D_S = \{(x_s, y_s)\}$ and target domain datasets $D_T^k = \{(x_t^k)\}$ for each client k, pretrained global model f_ϕ , number of FL round R, number of server update steps T_s , number of client update steps T_t

- 1: **for** $r = 1, 2, \cdots, R$ **do**
- 2: Initialize server EMA teacher model: $\theta \leftarrow \phi$
- 3: **for** $t = 1, 2, \cdot, T_s$ **do**
- 4: Sample patches (x_s,y_s) from D_S and generate downsampled global volume and masked inputs X_s, X_s^M, x_s^M
- 5: Update f_{ϕ} on server by minimizing \mathcal{L}_s (Eq.9)
- 6: Update server EMA teacher model parameter θ with (Eq.3)
- 7: end for
- 8: Server broadcast θ to clients
- 9: **for** each client k in parallel **do**
- 10: $\phi_k \leftarrow \theta, \theta_k \leftarrow \theta$
- 11: **for** $t = 1, 2, \dots, T_t$ **do**
- 12: Sample patches x_t^k from D_T^k and generate downsampled global volume and masked inputs $X_t^k, (X_t^k)^M, (x_s^k)^M$
- 13: Generate pseudolabels for unmasked inputs x_t^k and X_t^k using the teacher model f_{θ_k} : $f_{\theta_k}(x_t^k)$ and $f_{\theta_k}(X_t^k)$
- 14: Update f_{ϕ_k} by minimizing \mathcal{L}_u (Eq.10)
- 15: Update client EMA teacher model parameter with (Eq.3)
- 16: end for
- 17: Upload θ_k to server
- 18: end for
- 19: The server aggregates θ_k from clients:

$$\bar{\theta} \leftarrow \sum_{k} \frac{|D_T^k|}{\sum_{k} |D_T^k|} \theta_k$$

- 20: Update server model parameters $\phi \leftarrow \bar{\theta}$
- 21: **end for**

the same configuration as listed in Suppl.Tab.3. The difference is that the model can only access source domain data (image and label) during warmup epochs and can only access target domain data (image only) after that, while centralized UDA has synchronous access to both source and target domain data throughout the whole training process.

1.3. Implementation of Comparing Methods

For other comparing methods in centralized UDA, we adapt their official implementations. For DAFormer, HRDA, and MIC, we modify the ground truth labels to make them denser, as we observe that the original sparse annotations cause trouble for those methods. Specifically, we crop the scans to include only brain regions. In addition to having foreground classes of 7 subcortical regions (which account for approximately 2% of overall voxels), we assign another foreground class to the remaining brain regions. Therefore, there are 9 classes for DAFormer, HRDA, and MIC, 8 foreground and 1 background classes. This modification significantly improves the results. For the FL baselines FAT [55] and DualAdapt [78], since there is no public official implementation available, we implement both methods following the description in the original papers and finetune thoroughly. We use the same network backbone initialized with the same pretrained encoder and training configuration (FL rounds, global learning rate schedule, local update steps, batch size, etc.) as Fed-MAPSeg whenever possible.

1.4. Dataset Description

We include a diverse collection of 2,421 brain MRI scans from several international projects, each with its unique focus on infant brain development. From the Developing Human Connectome Project (dHCP) V1.0.2 data release [17] in the UK, we incorporate 983 scans (426 T1weighted, T1w), acquired shortly after birth. The Baby Connectome Project (BCP) [31] in the USA contributes 892 scans (519 T1w), featuring longitudinal data. Additionally, from the Environmental Influences on Child Health Outcomes (ECHO) project, also in the USA, we have 433 scans (218 T1w) from newborn infants. The 'Maternal Adversity, Inflammation, and Neurodevelopment' (Healthy Minds) project from Brazil, conducted at Hospital São Paulo - Federal University of São Paulo (UNIFESP), adds 103 T2-weighted (T2w) MRI scans, acquired shortly after birth and available in the National Institute of Mental Health Data Archive (collection ID 3811). Lastly, the Melbourne Children's Regional Infant Brain (M-CRIB) project [1] from Australia provides 10 additional T2w scans. All studies involved have received Institutional Review Board (IRB) approvals. MAPSeg takes normalized scans as inputs. During training, the intensity of each volumetric scan is clipped at a percentile randomly drawn from a uniform distribution $\mathcal{U}(99, 100)$, then normalized to 0-1. During inference, the intensity clip is fixed at 99.5%. The top 0.5% intensity is clipped as 1 to cope with outlier pixels (hyperintensities) that are usual in MRI.

1.5. Results of MRI \rightarrow CT cardiac segmentation

The performance of MAPSeg on the public cardiac MRI \rightarrow CT segmentation is reported in Suppl.Tab.4. Similarly, we use the same dataset partition as previous studies. MAPSeg consistently outperforms other baseline methods, although the performance gap is smaller than CT \rightarrow MRI.

Suppl.Tab. 4. Results of cardiact MRI→CT segmentation.

Cardiac $CT \rightarrow MRI$ segmentation							
Method	Dice(%) ↑						
Method	AA	LAC	LVC	MYO	Avg		
PnP-AdaNet[15]	74.0	68.9	61.9	50.8	63.9		
SIFA-V1[3]	81.1	76.4	75.7	58.7	73.0		
SIFA-V2[4]	81.3	79.5	73.8	61.6	74.1		
DAFormer[32]	85.5	88.2	74.5	60.2	77.1		
MPSCL[48]	90.3	87.1	86.5	72.5	84.1		
MA-UDA[37]	90.8	88.7	77.6	67.4	81.1		
SE-ASA[18]	83.8	85.2	82.9	71.7	80.9		
FSUDA-V1[45]	86.4	86.9	84.8	81.8	85.0		
PUFT[14]	88.1	88.5	87.5	74.1	84.6		
SDUDA[12]	87.9	88.1	88.4	78.7	85.8		
FSUDA-V2[46]	88.2	88.9	85.2	82.2	86.1		
MAPSeg (Ours)	93.3	<u>87.3</u>	<u>89.1</u>	78.9	<u>87.1</u>		



Suppl.Fig. 2. Downstream cross-sequence centralized UDA performance vs. MAE pretraining iterations.

1.6. Additional Analysis

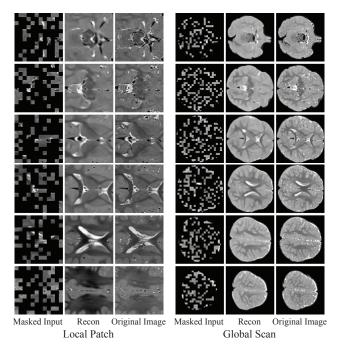
Influence of MAE Pretraining on UDA Results. We conduct an additional analysis to investigate the relationship between MAE training steps and downstream UDA performance. The experiments are conducted on cross-sequence brain MRI segmentation (Suppl.Fig.2). We observe significant improvement in UDA performance at the first 75,000 MAE training steps, which then gradually saturates. We choose 150,000 MAE training steps as the benefits of further training diminish.

Sensitivity to hyperparameters. We conduct additional experiments on cross-sequence brain MRI segmentation to investigate the sensitivity of MAPSeg to hyperparameters (Suppl.Tab.5). Specifically, we investigate the step size (α) of EMA update as well as weights of loss terms (γ and δ). When one parameter is varying, other parameters remain unchanged. We notice that the performance is relatively stable across a wide range of hyperparameters. Since we did not tune the hyperparameters extensively during development, the default parameters may not represent the optimal setting.

https://www.developingconnectome.org/datarelease/data-release-user-quide/

Suppl.Tab. 5. Influence of hyperparameters on results, bold indicates used parameters.

α	0.999/0.9999	0.99/0.999	0.99	0.999	0.9999
Dice (%)	77.73	74.00	74.26	74.74	78.06
$\overline{\gamma}$	0.05	0.5	0.1	0.01	0.005
Dice (%)	77.73	77.22	77.97	77.98	77.99
δ	0.025	0.25	0.1	0.01	0.0025
Dice (%)	77.73	76.74	78.08	77.82	78.57

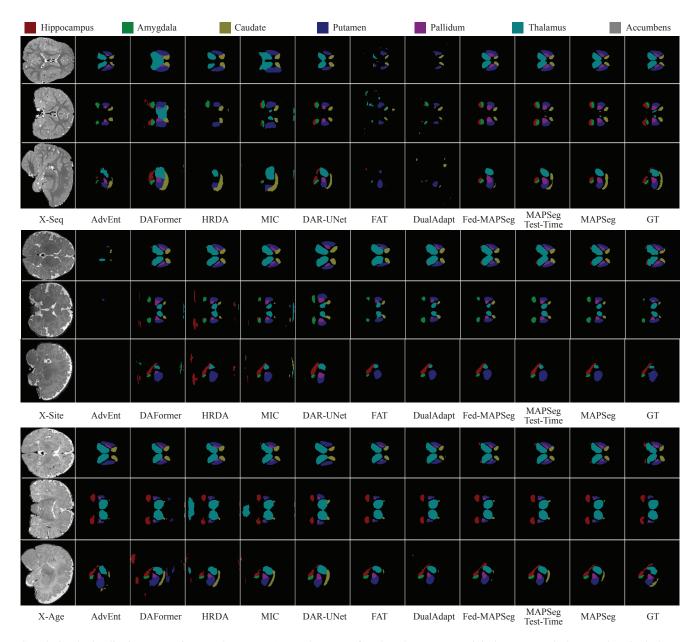


Suppl.Fig. 3. A randomly sampled T2w scan in cross-sequence task. MAE parameters is same as in Suppl.Tab.2

1.7. Visualization

MAE. Some visualizations of MAE results (axial slices) are provided in Suppl.Fig.3.

UDA Results. We provide qualitative comparisons of different methods on cross-sequence (X-Seq), cross-site (X-Site), and cross-age (X-Age) brain MRI segmentation tasks in Suppl.Fig.4. MAPSeg consistently provides accurate segmentation in different UDA settings. It is worth noting that, despite the second best performance in cross-sequence, DAR-UNet tends to oversegment on cross-site and crossage tasks, partially because of translation errors. On crossite and cross-age tasks, despite DAFormer, HRDA, and MIC generate reasonably good segmentation inside the subcortical regions, they exhibit extensive false positives outside the subcortical regions, leading to suboptimal overall Dice score.



Suppl.Fig. 4. Qualitative comparisons. Three rows (top to bottom) of each task represent axial plane, coronal plane, and sagittal plane, respectively.