

ECT-3DMedSAM: Efficient Cross Teaching Using Segment Anything Model for Semi-Supervised 3D Medical Image Segmentation

Zhewen Huang¹

Sara R. Guariglia²

Jiaqi Yang¹

Chia-Ling Tsai^{1,3}

ZHUANG5@GRADCENTER.CUNY.EDU

SARAROSE.X.GUARIGLIA@OPWDD.NY.GOV

JYANG2@GRADCENTER.CUNY.EDU

CHIALING.TSAI@QC.CUNY.EDU

¹ *Computer Science Department, The Graduate Center of the City University of NY, New York, NY*

² *Environmental Health & Neurodevelopment Laboratory, New York State Institute for Basic Research in Developmental Disabilities, Staten Island, NY*

³ *Computer Science Department, Queens College of the City University of NY, Flushing, NY*

Editors: Accepted for publication at MIDL 2026

Abstract

As precise manual annotation for medical imaging is both expert-intensive and costly, Semi-Supervised Medical Image Segmentation (SSMIS) provides a critical solution by leveraging large volumes of unlabeled data to achieve the high-performance segmentation necessary for anatomical structure analysis and disease diagnosis. Standard SSMIS models typically train specialized models with limited initialization, often failing to capture the complex semantic nuances of 3D anatomy. Foundation models offer superior generalization capabilities by leveraging large-scale pre-training but still struggle to adapt effectively when downstream annotations are limited. In this paper, we propose a novel cross-teaching framework tailored for the efficient adaptation of the 3D foundation model (MedSAM-2). We introduce a parameter-efficient design that shares frozen image and prompt encoders between two parallel, Low-Rank Adaptation (LoRA) learnable mask decoders. Furthermore, we replace the memory-intensive attention mechanism with a light-weight temporal propagation module for reducing the memory consumption while maintaining critical local volumetric coherence. Our model processes the same input volume through weakly and strongly augmentations to create a synergistic learning loop where the two decoders mutually supervise each other. We validate our method across three distinct datasets and modalities. Experimental results demonstrate that our framework effectively bridges the domain gap, achieving a 57.9% reduction in the average 95% Hausdorff Distance, substantially enhancing boundary precision for fine anatomical structures. Furthermore, our approach outperforms state-of-the-art baselines with a Dice score improvement of up to 2.8%, confirming its robustness and clinical reliability for volumetric segmentation.

Keywords: Semi-supervised learning, 3D image segmentation, Transformer, Segment Anything Model, Multi-modalities

1. Introduction

Medical image segmentation is a critical step in various clinical applications, ranging from anatomical structure analysis to disease diagnosis and treatment planning. However, training robust deep learning models typically requires large-scale datasets with precise pixel-level annotations, which are expensive and time-consuming to obtain due to the need for

expert knowledge (Ma et al., 2024). To mitigate the burden of manual annotation, Semi-Supervised Medical Image Segmentation (SSMIS) has emerged as a practical solution, enabling models to learn from limited labeled data alongside abundant unlabeled data.

Many influential SSMIS approaches rely on consistency learning, which encourages the model to produce invariant predictions under various perturbations. These methods can be broadly categorized based on the type of perturbation used. Input perturbation methods (Shu et al., 2023; Yang et al., 2022; Zou et al., 2021) encourage the model to maintain consistency across different perturbations of the same image input for the same model. Network perturbation methods (Liu et al., 2022; Chen et al., 2021) utilize and combine the strengths of different network architectures or initializations to generate diverse predictions for mutual supervision. Furthermore, hybrid frameworks (Luo et al., 2022; Yang et al., 2024; Chi et al., 2024) combine input and network perturbation, aiming at precise decision boundary in a larger variety of cases using mixed and multiple perturbations.

Recent progress in large-scale vision foundation models trained on broad and diverse data has opened new possibilities for reducing annotation costs in image analysis. The Segment Anything Model (SAM) (Kirillov et al., 2023) demonstrates remarkable zero-shot generalization capabilities. Following this success, recent works have introduced 3D-native foundation models in medical field (Wang et al., 2025; Ma et al., 2025a), which employ hierarchical vision transformers and memory attention modules to explicitly model volumetric dependencies. These pre-trained vision foundation models have demonstrated impressive segmentation performance and generalization capabilities for downstream tasks.

Despite these advancements, current semi-supervised methods face significant challenges when applied to complex 3D medical data with domain shifts. First, while hybrid frameworks combine multiple perturbations and have evolved to incorporate advanced transformer-based architectures, they typically rely on training encoders from limited initialization. Lacking the robust priors of foundation models, these methods remain susceptible to overfitting, particularly under extreme label scarcity or severe domain shifts. Second, directly leveraging large pretrained foundation models like MedSAM-2 (Ma et al., 2025a) presents a unique paradox. While their extensive prior knowledge provides strong generalization, it often leads to overconfidence in incorrect predictions on unseen domains, leading to over-segmentation and error accumulation that hinders the effective utilization of unlabeled data (Ma et al., 2025b). Third, direct application of traditional semi-supervised methods to modern foundation models is computationally prohibitive. These methods typically require maintaining multiple copies of the full network, which leads to massive memory consumption and slow training when applied to heavy backbones like MedSAM-2.

To address these problems and adapt medical foundation models to new tasks with sparse annotations, we propose a dual-stream, semi-supervised cross-teaching framework tailored for the efficient adaptation of MedSAM-2 to semi-supervised medical image segmentation. Unlike previous approaches that rely on architectural heterogeneity or computationally expensive model duplication, we introduce a parameter-efficient design that shares frozen image and prompt encoders between two parallel, learnable mask decoders.

This setup processes the same volumetric input through distinct augmentation strategies, a weakly augmented stream and a strongly augmented stream (Chen et al., 2021), while extracting unified features from the shared backbone. By decoupling the decoding process, we compel the two streams to learn robust, view-invariant representations through

mutual supervision. Rather than a unilateral teacher-student dynamic, our framework enforces a bidirectional consistency constraint where the predictions of the weak stream and the strong stream are jointly optimized to minimize discrepancy. This mutual alignment effectively mitigates foundation model overconfidence and prevents error accumulation. Our main contributions are summarized as follows:

- We propose a dual-stream cross-teaching framework, ECT-3DMedSAM, tailored for adapting foundation models to semi-supervised 3D medical segmentation. This approach effectively minimizes the overconfidence often observed in foundation models and bridges the domain gap between pre-training data and downstream medical tasks.
- We introduce a parameter-efficient fine-tuning strategy that integrates Low-Rank Adaptation (LoRA) with a simplified temporal propagation mechanism. We also freeze the massive image and prompt encoders and inject trainable LoRA layers solely into the mask decoder. In contrast to the naive dual-stream adaptations of MedSAM-2, our design reduces trainable parameters by 99% while effectively leveraging the pre-trained knowledge of foundation models, enabling the effective training of foundation models on limited labeled data.
- We conduct comprehensive experiments on three diverse public datasets in multi-modal scenarios. The results demonstrate that our method outperforms existing state-of-the-art semi-supervised learning methods and foundation models.

2. Methodology

2.1. Problem Setting

We address the task of semi-supervised volumetric medical image segmentation, where the goal is to leverage a small set of annotated data alongside a larger set of unannotated data to train a robust segmentation model. Let $\mathcal{D}_L = \{(X_i^{lab}, Y_i)\}_{i=1}^{N_L}$ denote the labeled dataset consisting of N_L volumes and their corresponding ground truth masks, and $\mathcal{D}_U = \{X_i^{unlab}\}_{i=1}^{N_U}$ denote the unlabeled dataset with N_U volumes, where typically $N_U \gg N_L$. Each input volume $X \in \mathbb{R}^{D \times H \times W}$ represents a 3D medical scan. For any input volume X (whether from \mathcal{D}_L or \mathcal{D}_U), we apply two distinct augmentations to generate a weak stream $X_{weak} = \mathcal{A}_{weak}(X)$ and a strong stream $X_{strong} = \mathcal{A}_{strong}(X)$. Let P as the volumetric prediction.

2.2. Overview

We propose a dual-stream semi-supervised cross-teaching framework, ECT-3DMedSAM, to adapt the pre-trained segmentation capabilities of MedSAM-2 to specific medical domains using limited labeled data. The overall framework is illustrated in Figure 1. To balance computational efficiency with effective consistency learning, we utilize a unified architectural design where both the labeled and unlabeled branches share the same frozen image and prompt encoders. Instead of training separate models, the extracted volumetric features are processed simultaneously by two parallel, LoRA (Hu et al., 2021) learnable mask decoders: one dedicated to the weakly-augmented stream and the other to the strongly-augmented

stream for learning robustness. These two decoder streams mutually supervise each other, enforcing the model to produce consistent predictions regardless of input perturbations.

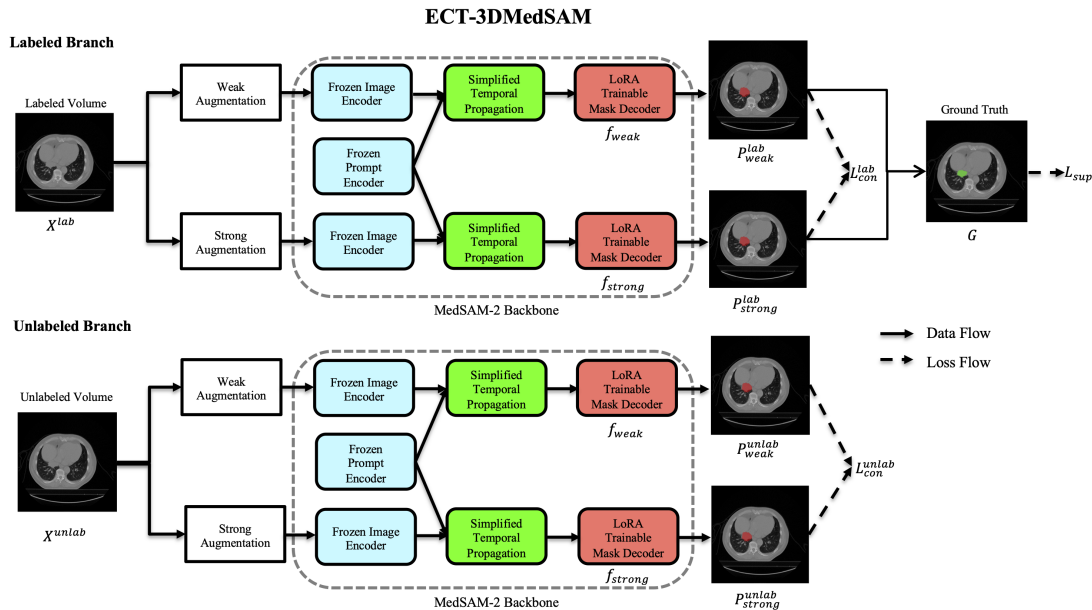


Figure 1: Overview of ECT-3DMedSAM. The architecture consists of two parallel data streams trained jointly: Labeled Branch (top) and Unlabeled Branch (bottom). Note that the weak decoder (f_{weak}) and the strong decoder (f_{strong}) are shared across both branches, which are mathematically identical. Both decoders receive gradient updates from the loss on the labeled branch and the unlabeled branch, enabling synergistic learning from both data sources.

2.3. Adapted MedSAM-2 Architecture

While MedSAM-2 demonstrates powerful segmentation capabilities, directly fine-tuning the entire architecture is computationally prohibitive and high data demanding when using two augmented streams with both labeled and unlabeled branches. Therefore, we introduce two critical architectural modifications to tailor it for efficient semi-supervised 3D segmentation.

Parameter-Efficient Fine-Tuning via LoRA. Instead of full fine-tuning, we adapt LoRA to the mask decoders to adapt the model to the target domain. We completely freeze the heavy image and prompt encoders to preserve the rich, pre-trained feature representations during training, effectively blocking the propagation of gradients that could degrade the foundation model’s generalization capabilities. We apply LoRA comprehensively across the mask decoder by a rank (r) of 16 and a scaling factor (α) of 32. Rather than limiting adaptation to the attention mechanism, we inject low-rank decomposition matrices into all linear and transposed convolutional layers. This ensures that the adaptation signal propagates through the entire decoding process from the transformer’s semantic processing to the

final spatial upscaling and mask generation. This strategy reduces the number of trainable parameters to less than 1M, 2% of the total model parameters, while allowing the decoder to learn task-specific segmentation boundaries.

Simplified Temporal Propagation. The original MedSAM-2 employs a sophisticated memory attention module with an 8-frame memory bank to handle long-range temporal occlusions in videos. However, in most of the tasks based on volumetric medical images, the spatial continuity between adjacent slices is actually more critical than long-range dependencies, which may introduce irrelevant noise from distant anatomical sections. To optimize computational efficiency and focus on local volumetric coherence, we replace the multi-frame memory bank with a simpler propagation mechanism. Specifically, the mask prediction M_t at slice t is conditioned solely on the image embedding E_t and the prediction from the immediate previous slice M_{t-1} . This light-weight propagation further reduces memory and computing resource consumption.

2.4. Consistency Learning and Loss Functions

Supervised Loss (\mathcal{L}_{sup}). For the labeled subset, we apply supervision to both the weakly and strongly augmented predictions. To prioritize geometric overlap in the semi-supervised setting, we employ a composite loss function with a 1:20 weighting ratio between Focal and Dice loss:

$$\mathcal{L}_{sup} = \sum_{s \in \{weak, strong\}} (1.0 \cdot \mathcal{L}_{Focal}(P_s^{lab}, G) + 20.0 \cdot \mathcal{L}_{Dice}(P_s^{lab}, G)) \quad (1)$$

where P_s represents the prediction from stream s , and G is the ground truth.

Consistency Losses (\mathcal{L}_{con}). We enforce consistency between the weak and strong streams for both labeled and unlabeled data. This bidirectional constraint ensures that the model learns robust features invariant to the intensity of augmentation. We utilize Mean Squared Error (MSE) to minimize the discrepancy:

$$\mathcal{L}_{con}^{lab} = \|P_{weak}^{lab} - P_{strong}^{lab}\|_2^2, \quad \mathcal{L}_{con}^{unlab} = \|P_{weak}^{unlab} - P_{strong}^{unlab}\|_2^2 \quad (2)$$

Total Loss. The final objective function integrates these components, with both consistency terms weighted by a time-dependent ramp-up factor $\lambda(t)$ to ensure stable convergence:

$$\mathcal{L}_{total} = \mathcal{L}_{sup} + \lambda(t) \cdot (\mathcal{L}_{con}^{lab} + \mathcal{L}_{con}^{unlab}) \quad (3)$$

3. Experiments

3.1. Datasets and Metrics

In our experiments, we use three volumetric datasets across different modalities (CT and MRI) to evaluate performance. For the semi-supervised learning setting, we partition the training set such that only 20% of the cases are used as labeled data, while the remaining cases serve as unlabeled data.

NSCLC-Radiomics Dataset (Aerts et al., 2014) comprises CT scans from 422 patients with non-small cell lung cancer (NSCLC). In our experimental setup, we utilize the first 200 cases for training and the remaining 222 cases for testing.

The PROMISE12 Dataset (Litjens et al., 2014) consists of T2-weighted MRI scans from patients with various prostatic diseases acquired at multiple locations to introduce domain variability. We use 30 cases for training and 7 cases for testing. Detailed information of this dataset is provided in the appendix.

The M&Ms Dataset (Campello et al., 2021) contains Cardiac Magnetic Resonance (CMR) images acquired from four distinct scanner vendors (Siemens, Philips, GE, and Canon) across six clinical centers. The dataset provides ground truth masks strictly for the Left Ventricle (LV), Right Ventricle (RV), and Left Ventricle Myocardium (MYO) at the End-Diastolic (ED) and End-Systolic (ES) phases. We utilize a split of 256 cases for training and 64 cases for testing.

The evaluation metrics include the Dice Similarity Coefficient (DSC), Jaccard index, 95% Hausdorff Distance (95HD), and Average Surface Distance (ASD).

3.2. Implementation Details

Our framework is implemented using PyTorch and trained on NVIDIA L40S GPUs. The model is trained for 50 epochs using the Adam optimizer with an initial learning rate of 1×10^{-4} and a weight decay of 1×10^{-4} . We employ a polynomial learning rate scheduler with a power of 0.9 to decay the learning rate during training.

The total batch size is set to 4, composed of a balanced mix of 2 labeled and 2 unlabeled volumes per iteration. The consistency loss weight λ is set to 0.1, utilizing a Gaussian ramp-up function over the beginning epochs to stabilize early training. All data augmentations are applied at the volumetric level rather than slice-by-slice. For weak data augmentations, random rotation and flipping are used. For strong data augmentations, in addition to the same geometric transfer in weak data augmentations, color jitter, blur and cutout are used. All input volumes are resized to 512×512 spatial resolution for processing. For each class, one prompt point is randomly generated from the mask in the middle frame in labeled branch. The unlabeled branch uses the center point for prompts, and it shares identical encoder and decoder parameters with the labeled branch. In the test, we use the mask decoder trained with weak augmentation volumes. All the baselines are finetuned using the training set for fair comparison.

3.3. Comparison with State-of-the-Art Methods

Table 1 demonstrates the performance comparison between our model, SSMIS approaches, and state-of-the-art 3D medical foundation models on the NSCLC dataset for cancer segmentation. As observed, our proposed method achieves the best DSC of 72.31%. Notably, our method shows superior boundary precision, decreasing 63% in 95HD and 64.5% in ASD compared to MedSAM-2, proving its efficacy in delineating precise margins.

Table 2 illustrates the performance comparison on the PROMISE12 dataset for prostate segmentation across three distinct institutions. Our method shows the most balanced performance across the three institutions, achieving the highest average DSC of 74.17%. It also outperforms all comparison methods on the challenging UCL domain with a DSC of 79.66%, indicating superior generalization capability. Furthermore, our method achieves best in boundary metrics, reducing the 95HD from 37.34 to 17.68 and the ASD from 11.92

Table 1: Comparison of different methods on NSCLC dataset (the best performance is marked as **bold**, and the second-best is underlined). **Avg.** means the average performance of all volumes in the dataset. **# of Params** is the number of trainable parameters.

Method	# of Params	DSC (%) ↑	Jaccard (%) ↑	95HD (mm) ↓	ASD (mm) ↓
		Avg.			
UNet (Ronneberger et al., 2015)	31.05M	9.39	5.22	228.93	136.14
nnUNet (Isensee et al., 2021)	126.82M	34.48	24.65	156.84	64.46
BCP (Bai et al., 2023)	18.92M	25.40	18.38	120.86	90.77
SAM-Med3D (Wang et al., 2025)	100.51M	30.85	21.51	45.49	20.21
MedSAM-2 (Ma et al., 2025a)	38.96M	<u>71.00</u>	59.36	<u>17.81</u>	<u>4.93</u>
ECT-3DMedSAM (Ours)	0.94M	72.31	<u>57.55</u>	6.59	1.75

to 4.69 compared to MedSAM-2, confirming its reliability in edges in multi-domain medical image segmentation.

Table 2: Comparison of different methods on PROMISE12 dataset

Method	Segmentation			DSC ↑	Jaccard ↑	95HD ↓	ASD ↓
	BIDMC	HK	UCL	Avg.			
UNet	66.71	54.41	5.86	37.12	27.08	113.22	50.37
nnUNet	86.98	80.29	24.54	58.31	48.02	53.70	19.71
BCP	<u>72.58</u>	39.45	28.87	44.38	33.92	45.10	17.08
SAM-Med3D	26.82	45.22	26.14	31.78	20.02	98.86	34.18
MedSAM-2	45.93	90.96	<u>77.15</u>	<u>72.18</u>	63.57	<u>37.34</u>	<u>11.92</u>
ECT-3DMedSAM (Ours)	55.88	<u>84.24</u>	79.66	74.17	<u>60.68</u>	17.68	4.69

Table 3 presents the quantitative results on the M&Ms dataset, which poses significant challenges due to domain shifts across four different scanner vendors. First, we observe that our proposed model adapts MedSAM-2 to the specific domain, doubling the segmentation accuracy of MedSAM-2 to an average DSC of 30.21% and demonstrating the effectiveness of our cross-teaching framework. Second, our model achieves the best performance in boundary metrics, reducing the 95HD to 18.40 and ASD to 5.73. This demonstrates that our method produces topologically coherent segmentation with precise boundaries. Finally, although our model underperforms BCP a little in overlap metrics, our method shows superior robustness on the most challenging domain, Vendor C. While nnUNet and BCP drop to 27.52% and 16.93% DSC respectively on Vendor C, our method achieves 35.08%, indicating a stronger capability to generalize to unseen or difficult domains compared to models trained from limited initialization.

Table 3: Comparison of different methods on M&Ms dataset

Method	(RV / MYO / LV Segmentation) DSC ↑				DSC ↑	Jaccard ↑	95HD ↓	ASD ↓
	Vendor A	Vendor B	Vendor C	Vendor D	Avg.			
UNet	29.94 / 20.18 / 34.60	<u>23.66</u> / 11.22 / 33.58	18.07 / 8.60 / 14.46	<u>56.19</u> / 46.00 / <u>57.93</u>	29.66	21.53	91.59	82.74
nnUNet	47.10 / 44.16 / 65.94	24.45 / 24.25 / 43.30	10.50 / 23.74 / <u>48.33</u>	22.88 / 19.34 / 27.21	36.67	25.65	<u>43.70</u>	<u>14.69</u>
BCP	<u>37.93</u> / <u>39.44</u> / 55.59	18.72 / 40.34 / <u>55.96</u>	14.93 / 16.16 / 19.71	25.33 / <u>28.33</u> / 19.16	<u>34.37</u>	<u>24.47</u>	88.57	69.54
SAM-Med3D	7.26 / 3.42 / 7.40	1.58 / 1.64 / 3.62	8.74 / 1.16 / 12.61	5.94 / 3.60 / 8.04	4.77	2.69	76.70	30.30
MedSAM-2	3.18 / 19.41 / 0.16	17.34 / 16.34 / 3.90	<u>22.37</u> / 30.53 / 12.25	73.27 / 17.64 / 3.41	15.45	12.54	44.82	17.37
ECT-3DMedSAM (Ours)	9.63 / 4.84 / 74.28	19.12 / 6.26 / 64.20	28.40 / 11.33 / 65.51	14.49 / 5.34 / 62.49	30.21	21.81	18.40	5.73

Figure 2 and Figure 3 visualizes the superiority of our model. We observe distinct failure modes among the baselines. Standard CNN-based methods like UNet struggle with volumetric integrity and localization, particularly in unseen domains such as the UCL and HK domains on PROMISE12 dataset. In contrast, nnUNet demonstrates strong baseline performance. However, it still lacks the capacity to adapt to the severe domain shift as a model trained from limited initialization. Additionally, foundation models tend to exhibit overconfidence. They often ignore anatomical margins, generating topologically implausible shapes that leak into adjacent tissues.

Our model successfully bridges these gaps. By making use of the pre-trained knowledge of foundation models and actively utilizing the 80% unlabeled data through our cross-teaching, our model effectively adapts to the challenging domains where nnUNet and other CNN-based methods fail. The dual-stream consistency constraints filter out the overconfident errors typical of foundation models, producing spatially coherent segmentations that are robust in difficult domain scenarios.

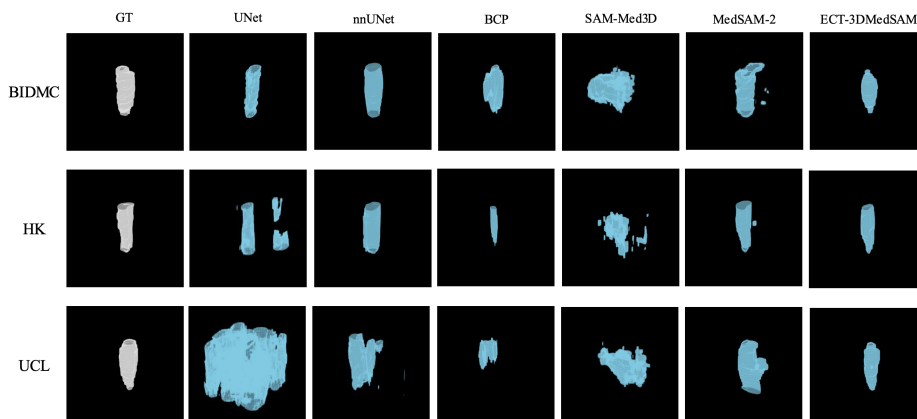


Figure 2: 3D Qualitative comparison of segmentation results on each institution in the PROMISE12 dataset: BIDMC (top row), HK (middle row), and UCL (bottom row). The columns display the Ground Truth (GT) and predictions from five state-of-the-art methods compared to our proposed ECT-3DMedSAM.

3.4. Ablation Studies

We conduct ablation studies in order to verify the effectiveness of each module in our method. Table 4 shows the effectiveness of each module in our model. The baseline MedSAM-2 baseline suffers from significant boundary errors due to overconfidence. Incorporating the dual-stream mechanism substantially refines these boundaries, halving the 95HD from 41.40 to 20.21. Furthermore, the addition of the unlabeled data stream proves critical for generalization, increasing the DSC to 73.96%. Exchanging sophisticated memory attention based propagation approach with light-weight temporal propagation achieves the best results with a peak DSC of 74.17% and the lowest ASD of 4.69.

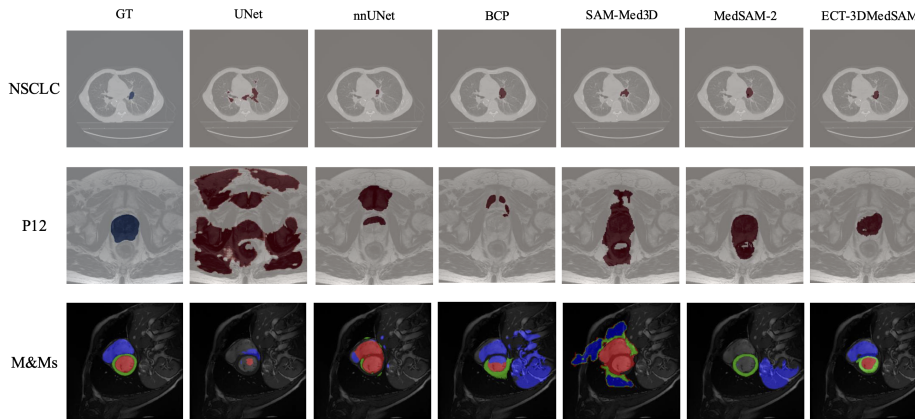


Figure 3: Qualitative comparison of segmentation results across three datasets: NSCLC (top row), PROMISE12 (middle row), and M&Ms (bottom row). For the M&Ms dataset, anatomical structures are color-coded: LV (Red), RV (Blue), and MYO (Green).

Table 4: Ablation experiments on the PROMISE12 dataset. **Base** is MedSAM-2 with frozen image and prompt encoders and the LoRA trainable decoder. **DS** is adding dual-stream framework. **CT** is adding the unlabeled branch for cross-teaching. **SP** is using the simplified temporal propagation.

Base	DS	CT	SP	DSC \uparrow	Jaccard \uparrow	95HD \downarrow	ASD \downarrow
\checkmark				64.97	54.81	41.40	13.09
\checkmark	\checkmark			71.57	58.58	20.21	<u>5.05</u>
\checkmark	\checkmark	\checkmark		<u>73.96</u>	<u>60.50</u>	<u>18.89</u>	5.22
\checkmark	\checkmark	\checkmark	\checkmark	74.17	60.68	17.48	4.69

We also evaluate the performance of different weight ratios between Dice and Focal loss in \mathcal{L}_{sup} . As demonstrated in Table 5, when the Focal loss weight is too high, the optimization lacks strong global structural constraints. The model struggles to enforce volumetric coherence, resulting in suboptimal performance. When the Dice loss weight is too high, the model fails to resolve subtle boundary details, leading to the slight degradation of the boundary metrics. The weight of 20.0 effectively balances the gradient magnitudes between the Dice loss and Focal loss, ensuring the model optimizes for global structure without sacrificing the ability to refine hard boundaries.

As observed in the qualitative analysis, standard CNN-based methods struggle with precise localization, while foundation models suffer from overconfidence on unseen domains. Theoretically, both failure modes contribute to an elevated false positive rate. To validate this, we quantified the pixel-level False Positive Rate (FPR, over-segmentation) and False Negative Rate (FNR, under-segmentation) on the NSCLC and PROMISE12 datasets. As presented in Table 6, our method achieves the lowest FPR, confirming its efficacy in sup-

Table 5: Ablation study of different Dice and Focal loss ratios on the PROMISE12 dataset. **DFR** is the weight ratio between Dice and Focal loss.

DFR	DSC \uparrow	Jaccard \uparrow	95HD \downarrow	ASD \downarrow
1:20	51.13	36.26	23.66	7.14
1:1	65.51	52.79	22.05	6.07
10:1	70.48	58.30	19.84	5.45
20:1	<u>74.17</u>	<u>60.68</u>	17.48	4.69
40:1	74.63	60.72	<u>18.12</u>	<u>4.75</u>

pressing both failure localization and overconfident leakage. Furthermore, it maintains a highly competitive FNR compared to CNN-based baselines.

Table 6: False negative rates (**FNR**) and false positive rates (**FPR**) of different methods on the NSCLC and PROMISE12 datasets. Both rates are calculated by dividing the total foreground in ground truth.

Method	NSCLC		PROMISE12	
	FPR	FNR	FPR	FNR
UNet	49.21	0.52	3.48	0.45
nnUNet	14.51	0.40	0.80	<u>0.30</u>
BCP	10.86	0.62	0.79	0.52
SAM-Med3D	27.76	0.50	1.59	0.60
MedSAM-2	<u>1.18</u>	0.14	<u>0.47</u>	0.21
ECT-3DMedSAM (Ours)	0.47	<u>0.18</u>	0.18	0.33

4. Conclusion

In this paper, we proposed a MedSAM-2 based dual-stream cross-teaching framework to address the challenges of semi-supervised 3D medical image segmentation under limited supervision. By synergizing the pre-trained medical segmentation capabilities of the MedSAM-2 foundation model with a consistency-based cross-teaching paradigm, our approach effectively mitigates the overconfidence inherent in large-scale foundation models while preventing overfitting to the limited source domain.

Unlike methods that rely on computationally expensive full fine-tuning or heavy memory banks, we also introduced a parameter-efficient architecture sharing frozen image and prompt encoders between two parallel, LoRA-adapted mask decoders. Furthermore, by replacing the standard memory attention module with a streamlined simple propagation, we reduced the memory consumption while maintaining the volumetric coherence essential for medical scans. These architectural optimizations enable the effective training of foundation models on limited labeled data, overcoming the resource constraints that typically hinder their adaptation to downstream medical tasks.

Our extensive experiments across diverse modalities demonstrate that our method consistently outperforms both conventional SSMIS architectures and 3D medical foundation models, especially in boundary metrics. These results confirm that identifying robust, consistency-invariant features via dual-stream cross-teaching is a powerful strategy for adapting foundation models to complex, multi-modalities medical semi-supervised segmentation tasks.

While our dual-stream consistency effectively suppresses the overconfidence common in foundation models, it can lead to over-conservative predictions when the domain shift is severe or feature ambiguity is high. In these scenarios, the consistency loss forces the model to converge towards the background class to minimize disagreement. Future work will focus on mitigating this over-suppression by incorporating uncertainty-aware consistency mechanisms. Specifically, we aim to explore adaptive weighting strategies (Luo et al., 2021) that can dynamically relax consistency constraints in high-uncertainty regions, thereby improving sensitivity in challenging, low-contrast domains.

Acknowledgments

This work was supported by PSC-CUNY Award 65406-00 53 and New York State Institute for Basic Research (Office for People with Developmental Disabilities).

References

- Hugo J. W. L. Aerts et al. Decoding tumour phenotype by noninvasive imaging using a quantitative radiomics approach. *Nature Communications*, 5(1):4006, June 2014. ISSN 2041-1723. doi: 10.1038/ncomms5006. URL <https://www.nature.com/articles/ncomms5006>. Publisher: Nature Publishing Group.
- Yunhao Bai, Duowen Chen, Qingli Li, Wei Shen, and Yan Wang. Bidirectional Copy-Paste for Semi-Supervised Medical Image Segmentation. In *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11514–11524, Vancouver, BC, Canada, June 2023. IEEE. ISBN 979-8-3503-0129-8. doi: 10.1109/CVPR52729.2023.01108. URL <https://ieeexplore.ieee.org/document/10204177/>.
- Víctor M. Campello et al. Multi-Centre, Multi-Vendor and Multi-Disease Cardiac Segmentation: The M&Ms Challenge. *IEEE Transactions on Medical Imaging*, 40(12):3543–3554, December 2021. ISSN 1558-254X. doi: 10.1109/TMI.2021.3090082. URL <https://ieeexplore.ieee.org/document/9458279>.
- Xiaokang Chen, Yuhui Yuan, Gang Zeng, and Jingdong Wang. Semi-Supervised Semantic Segmentation with Cross Pseudo Supervision. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2613–2622, Nashville, TN, USA, June 2021. IEEE. ISBN 978-1-6654-4509-2. doi: 10.1109/CVPR46437.2021.00264. URL <https://ieeexplore.ieee.org/document/9577639/>.
- Hanyang Chi, Jian Pang, Bingfeng Zhang, and Weifeng Liu. Adaptive Bidirectional Displacement for Semi-Supervised Medical Image Segmentation. In *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4070–4080, Seattle, WA, USA, June 2024. IEEE. ISBN 979-8-3503-5300-6. doi: 10.1109/CVPR52733.2024.00390. URL <https://ieeexplore.ieee.org/document/10658193/>.
- Edward J. Hu et al. LoRA: Low-Rank Adaptation of Large Language Models, October 2021. URL <http://arxiv.org/abs/2106.09685>. arXiv:2106.09685 [cs].
- Fabian Isensee, Paul F. Jaeger, Simon A. A. Kohl, Jens Petersen, and Klaus H. Maier-Hein. nnU-Net: a self-configuring method for deep learning-based biomedical image segmentation. *Nature Methods*, 18(2):203–211, February 2021. ISSN 1548-7105. doi: 10.1038/s41592-020-01008-z. URL <https://www.nature.com/articles/s41592-020-01008-z>.
- Alexander Kirillov et al. Segment Anything. In *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 3992–4003, Paris, France, October 2023. IEEE. ISBN 979-8-3503-0718-4. doi: 10.1109/ICCV51070.2023.00371. URL <https://ieeexplore.ieee.org/document/10378323/>.

- Geert Litjens et al. Evaluation of prostate segmentation algorithms for MRI: The PROMISE12 challenge. *Medical Image Analysis*, 18(2):359–373, February 2014. ISSN 1361-8415. doi: 10.1016/j.media.2013.12.002. URL <https://www.sciencedirect.com/science/article/pii/S1361841513001734>.
- Yuyuan Liu et al. Perturbed and Strict Mean Teachers for Semi-supervised Semantic Segmentation. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4248–4257, New Orleans, LA, USA, June 2022. IEEE. ISBN 978-1-6654-6946-3. doi: 10.1109/CVPR52688.2022.00422. URL <https://ieeexplore.ieee.org/document/9879201/>.
- Xiangde Luo, Minhao Hu, Tao Song, Guotai Wang, and Shaoting Zhang. Semi-Supervised Medical Image Segmentation via Cross Teaching between CNN and Transformer. In *Proceedings of The 5th International Conference on Medical Imaging with Deep Learning*, pages 820–833. PMLR, December 2022. URL <https://proceedings.mlr.press/v172/luo22b.html>. ISSN: 2640-3498.
- Xiangde Luo et al. Efficient Semi-supervised Gross Target Volume of Nasopharyngeal Carcinoma Segmentation via Uncertainty Rectified Pyramid Consistency. In *Medical Image Computing and Computer Assisted Intervention – MICCAI 2021*, pages 318–329, Cham, 2021. Springer International Publishing. ISBN 978-3-030-87196-3. doi: 10.1007/978-3-030-87196-3_30.
- Jun Ma et al. Segment anything in medical images. *Nature Communications*, 15(1):654, January 2024. ISSN 2041-1723. doi: 10.1038/s41467-024-44824-z. URL <https://www.nature.com/articles/s41467-024-44824-z>. Publisher: Nature Publishing Group.
- Jun Ma et al. MedSAM2: Segment Anything in 3D Medical Images and Videos, April 2025a. URL <http://arxiv.org/abs/2504.03600>. arXiv:2504.03600 [eess].
- Qinghe Ma et al. Steady Progress Beats Stagnation: Mutual Aid of Foundation and Conventional Models in Mixed Domain Semi-Supervised Medical Image Segmentation. In *2025 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5175–5185, June 2025b. doi: 10.1109/CVPR52734.2025.00488. URL <https://ieeexplore.ieee.org/document/11094453>. ISSN: 2575-7075.
- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-Net: Convolutional Networks for Biomedical Image Segmentation. In *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*, pages 234–241, Cham, 2015. Springer International Publishing. ISBN 978-3-319-24574-4. doi: 10.1007/978-3-319-24574-4_28.
- Yucheng Shu, Hengbo Li, Bin Xiao, Xiuli Bi, and Weisheng Li. Cross-Mix Monitoring for Medical Image Segmentation With Limited Supervision. *Trans. Multi.*, 25:1700–1712, January 2023. ISSN 1520-9210. doi: 10.1109/TMM.2022.3154159. URL <https://doi.org/10.1109/TMM.2022.3154159>.
- Haoyu Wang et al. SAM-Med3D: A Vision Foundation Model for General-Purpose Segmentation on Volumetric Medical Images. *IEEE Transactions on Neural Networks and*

Learning Systems, 36(10):17599–17612, October 2025. ISSN 2162-2388. doi: 10.1109/TNNLS.2025.3586694. URL <https://ieeexplore.ieee.org/document/11105528>.

Lihe Yang, Wei Zhuo, Lei Qi, Yinghuan Shi, and Yang Gao. ST++: Make Self-training Work Better for Semi-supervised Semantic Segmentation. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4258–4267, New Orleans, LA, USA, June 2022. IEEE. ISBN 978-1-6654-6946-3. doi: 10.1109/CVPR52688.2022.00423. URL <https://ieeexplore.ieee.org/document/9880151/>.

Yuan Yang, Lin Zhang, and Lei Ren. Semi-supervised medical image segmentation via cross teaching between MobileNet and MobileViT. *Image Vision Comput.*, 150(C), October 2024. ISSN 0262-8856. doi: 10.1016/j.imavis.2024.105196. URL <https://doi.org/10.1016/j.imavis.2024.105196>.

Yuliang Zou et al. Pseudoseg: Designing Pseudo Labels for Semantic Segmentation. In *International Conference on Learning Representations*, 2021.

Appendix A. Detailed information of the PROMISE12 Dataset

Table 7 shows detailed information of the PROMISE12 dataset used in the experiment.

Table 7: Detailed information of the PROMISE12 Dataset

Institution	Case number	Field strength (T)	Endorectal coil	Manufactor
HK	12	1.5	Endorectal	Siemens
BIDMC	12	3	Endorectal	GE
UCL	13	1.5 and 3	No	Siemens