

Segmenting Medical Images with Limited Data

Zhaoshan Liu^{a,*}, Qiuji Lv^{a,b,*}, Chau Hung Lee^c, Lei Shen^{a,**}

^a*Department of Mechanical Engineering, National University of Singapore, 9 Engineering Drive 1, Singapore, 117575, Singapore*

^b*School of Intelligent Systems Engineering, Sun Yat-sen University, No.66, Gongchang Road, Guangming District, 518107, China*

^c*Department of Radiology, Tan Tock Seng Hospital, 11 Jalan Tan Tock Seng, Singapore, 308433, Singapore*

Abstract

While computer vision has proven valuable for medical image segmentation, its application faces challenges such as limited dataset sizes and the complexity of effectively leveraging unlabeled images. To address these challenges, we present a novel semi-supervised, consistency-based approach termed the data-efficient medical segmenter (DEMS). The DEMS features an encoder-decoder architecture and incorporates the developed online automatic augments (OAA) and residual robustness enhancement (RRE) blocks. The OAA augments input data with various image transformations, thereby diversifying the dataset to improve the generalization ability. The RRE enriches feature diversity and introduces perturbations to create varied inputs for different decoders, thereby providing enhanced variability. Moreover, we introduce a sensitive loss to further enhance consistency across different decoders and stabilize the training process. Extensive experimental results on both our own and three public datasets affirm the effectiveness of DEMS. Under extreme data shortage scenarios, our DEMS achieves 16.85% and 10.37% improvement in dice score compared with the U-Net and top-performed state-of-the-art method, respectively. Given its superior data efficiency, DEMS could present significant advancements in medical segmentation under small data regimes. The project homepage can be accessed at <https://github.com/NUS-Tim/DEMS>.

Keywords: Semi-Supervised Learning, Data Augmentation, Medical Image Segmentation, Medical Ultrasound

*Equal contribution

**Corresponding author

Email addresses: e0575844@u.nus.edu (Zhaoshan Liu), lvqj5@mail2.sysu.edu.cn (Qiuji Lv), chau_hung_lee@ttsh.com.sg (Chau Hung Lee), mpeshel@nus.edu.sg (Lei Shen)

1. Introduction

Recently, there has been a rapid advancement in artificial intelligence [1–3], with computer vision emerging as a particularly prominent area of research. It has witnessed remarkable progress in a variety of applications [4–7], including facial recognition, autonomous navigation, and medical image segmentation. Medical image segmentation leverages the neural network model to generate high-accuracy mask prediction and has achieved significant success in various scenarios [8–10]. Although model assistance greatly enhances medical diagnosis, its usage faces unique challenges. In particular, training the neural network often requires large labeled datasets, yet collecting medical images is not as straightforward as collecting natural images [11]. This can be attributed to several factors, including cost and privacy concerns [12]. First, acquiring medical images is time-consuming and requires specialized equipment and expert intervention. Additionally, in some scenarios, there may not be a sufficient number of patients available for data collection. Second, scanned images require additional labeling, which further increases the time cost. Finally, patient privacy [13] must be considered, making a significant portion of datasets publicly unavailable. In this context, the number of accessible images may be limited, and a substantial amount of unlabeled images may remain underutilized for model training.

To facilitate effective model training with limited data and bolster generalizability, a multitude of data augmentation (DA) methods have been presented. Prevalent DA methods can be categorized into three types, including conventional DA, generative adversarial network (GAN)-based DA, and automatic DA. Conventional DA leverages various image transformations such as flipping, rotation, contrast adjustment, and scaling [14] to generate augmented variants. This approach is the most commonly implemented [15–17], yet the augmentation pipeline design heavily relies on experience. The GAN-based DA leverages GAN [18] to synthesize artificial images and enlarge datasets [19–21]. This method, while more versatile, presents challenges such as variable synthetic quality, extensive training times [22, 23], and strong dependence of model performance on dataset size [22]. Automatic DA [24] integrates multiple conventional DA transformations to enhance data diversity [25–27]. However, its application faces various challenges, including extensive computational demands, complex stage design, and potential model confounding [28].

Semi-supervised learning (SSL) has been proposed to effectively utilize unlabeled images [29] and has garnered significant attention in the field of medical image segmentation. Its applications predominantly rely on various intrinsic principles such as consistency regularization [30], pseudo-labeling [31], and prior knowledge [32], with consistency regularization being the most prevalent. The consistency-based method encourages the model to generate consistent outputs when taking inputs that undergo various operations or perturbations, thereby enhancing its generalization ability. A significant number of consistency-based methods adopt the mean teacher (MT) architecture [33] as their foundation [34–36]. The MT-based meth-

ods comprise a teacher subnet and a student subnet, where the parameters of the teacher network are updated from the student network using the exponential moving average algorithm. However, challenges arise with the MT-based approaches, such as the constraint on teacher network performance and limited model variability [37]. In addition to the MT-based architecture, various innovative methods demonstrate remarkable performance, such as those based on encoder-decoder architecture [38–40] and GAN architecture [41]. With sophisticated network and loss design, these approaches have demonstrated outstanding performance, while challenges such as variability restriction and stability limitation [42] persist.

To bolster model performance with limited data comprising unlabeled images, we propose a data-efficient medical segmenter (DEMS). The DEMS employs an encoder-decoder architecture and comprises our developed online automatic augments (OAA) and residual robustness enhancement (RRE) blocks. The OAA diversifies visual input with varying transformations, thereby providing enhanced data diversity to bolster generalization ability. The RRE block enriches feature diversity and introduces perturbations to generate varied inputs for different decoders and thus offers greater variability. Additionally, we introduce a novel sensitivity loss to enhance consistency across different decoders and stabilize model training. We perform extensive experiments on our own and three public datasets and the results show that DEMS outperforms state-of-the-art (SOTA) methods. Furthermore, DEMS exhibits greater performance leadership under severe data regimes, thereby highlighting its exceptional data efficiency. To sum up, our main contributions are:

- We propose a data-efficient medical segmenter (DEMS), a novel semi-supervised segmentation approach that effectively utilizes limited and unlabeled data and features superior performance in small data scenarios.
- We introduce an online automatic augments (OAA) to diversify data through various image transformations to enhance generalization ability. We propose a residual robustness enhancement (RRE) block to enrich feature diversity and inject perturbations to produce varied inputs for different decoders, therefore enhancing the variability. Furthermore, we introduce a sensitivity loss to improve consistency across varying decoders and stabilize training.
- Extensive results from both our own and three public datasets underscore the superiority of DEMS. Notably, DEMS shows increasing performance advantages in extreme data shortages, emphasizing its extraordinary efficiency in severe data shortages.

The rest of this paper is organized as follows. Section 2 meticulously reviews the related works, focusing on DA and SSL in medical image segmentation. Existing challenges and potential improvements are then summarized. In Section 3, we provide a detailed illustration of the proposed approach, including the introduction of DEMS, OAA, RRE block and connection structure, and loss function. Section 4

presents the datasets, experimental setup, and evaluation metrics. In Section 5, we illustrate the experimental results, conduct comprehensive analysis and visualization, and perform intensive ablation experiments. We conclude our work and provide insightful future perspectives in Section 6.

2. Related Work

2.1. Data Augmentation

Numerous researchers leverage conventional DA as a component in their proposed methods. For instance, Li et al. [15] presented a HAL-IA method to reduce the annotation cost, in which Gaussian noise, random flip, and random crop were leveraged for DA. Yu and colleagues [43] developed a UNesT network for local communication, in which various DA transformations such as random flip, rotation, and intensity change were leveraged to train the renal segmentation model. Li et al. [44] introduced a GFUNet that integrates Fourier transform with U-Net and leverages varying transformations including random histogram matching, rotation, and shifting to perform DA. Zhao and colleagues [16] developed a novel knowledge distillation framework, in which DA transformations including random scaling, random rotation, and random elastic deformation were utilized. Furthermore, Isensee et al. [17] presented a nnU-Net with a preset DA pipeline that sequentially applies transformations such as rotation, scaling, Gaussian noise, Gaussian blur, etc. The conventional DA is intuitive and effective while the design of the augmentation pipeline predominantly relies on experience.

Leveraging GAN to synthesize artificial images and enlarge datasets for DA presents an effective alternative. It generates images at the pixel level, thus being considered more versatile compared with the conventional DA. In 2022, Chai et al. [19] developed a DPGAN consisting of three variational auto-encoder GANs to synthesize artificial images and labels. Moreover, an extra discriminator was leveraged to promote the image reality and relationship across images and latent vectors. Li and colleagues [20] proposed a semi-supervised framework to capture the joint image-label distribution and synthesize chest X-ray and liver computed tomography images. Kugelman et al. [21] leveraged Style-GAN [45] to create image-mask pairs for optical coherence tomography images, in which the image and mask are generated on adjacent channels. Pandey and colleagues [46] developed a two-stage GAN approach, in which the mask is initially synthesized, followed by the image leveraging the generated mask. Besides, Iqbal et al. [47] presented a data expansion network in 2023 to synthesize images solely and leveraged a trained convolutional neural network (CNN) to label the generated images. The GAN-based DA faces varied challenges such as varying synthesizing quality, extensive training consumption, and high data appetite.

The automatic DA has garnered significant attention for its extraordinary augmentation diversity. Various approaches adapt reinforcement learning to perform

automatic DA. For example, Yang et al. [25] leveraged the validation accuracy to update the recurrent neural network controller. The AADG framework presented by Lyu and colleagues [26] introduces a novel proxy task, in which Sinkhorn distance was utilized to maximize the diversity across various augmented domains. Qin et al. [27] proposed a joint-learning strategy that combines Dueling DQN [48] and segmentation modules to search for maximum performance improvement. Xu and colleagues [49] devised a differentiable approach to update the parameters using stochastic relaxation and the Monte Carlo method. Besides the methods based on reinforcement learning, the MedAugment [50] developed in 2023 augments each input image with different transformations sampled from two groups of transformations. Zhao et al. [51] presented a novel automatic DA method, in which spatial transform and appearance transform were leveraged to synthesize labeled examples. Additionally, Eaton and colleagues [52] utilized the online mix-up [53] algorithm to enhance the performance of brain glioma segmentation. These methods arise with various challenges including high computation costs, considerable stage complexity, and potential model confounding.

2.2. Semi-Supervised Segmentation

The MT-based approaches have been widely adopted in the field of semi-supervised medical image segmentation. Zhang et al. [34] proposed an uncertainty-guided mutual consistency learning framework with two branches for pixel-wise classification and level set function regression, respectively. Lei and colleagues [35] presented an ASE-Net comprising both discriminator and segmentation networks. The adversarial consistency training approach leverages two discriminators to extract prior relationships, and the dynamic convolution-based bidirectional attention component is utilized to adaptively adjust the network weights. Xu et al. [36] introduced a shadow-consistent SSL approach featuring shadow augmentation and shadow dropout mechanisms. The shadow augmentation mechanism enhances the samples by integrating simulated shadow artifacts, and the shadow dropout mechanism compels the network to identify the boundary using shadow-free pixels. Lyu and colleagues [54] developed an AFAM-Net, which utilizes a reconstruction task to capture anatomical information and an adaptive feature aggregation strategy to transfer and filter features. The unsupervised consistency loss is formulated based on the reconstruction and segmentation tasks. Additionally, Yu et al. [55] presented an uncertainty-aware framework in which the teacher model generates target outputs together with the uncertainty of each prediction through Monte Carlo sampling. Although the MT-based methods have demonstrated promising performance in various scenarios, various challenges exist. Firstly, the performance of the teacher network can be constrained. Secondly, the shared structure restricts the model variability.

In addition to the MT-based methods, Tang et al. [38] developed an MGCC model that adopts the encoder-decoder architecture and incorporates the multi-scale attention gate for feature enhancement. Wu and colleagues [39] proposed an MC-Net consisting of an encoder and two decoders and calculated the prediction discrep-

ancy between decoder outputs to promote mutual consistency. Subsequently, the MC-Net+ with further intra-model diversity was presented [40], in which an additional decoder with the nearest interpolating operation was incorporated. The ASS-GAN presented by Zhai et al. [41] leverages two generators and a discriminator to perform adversarial learning, in which the predicted masks of a generator were utilized to guide the other. Additionally, Wang and colleagues [37] developed a mutual correction framework, in which a contrastive difference review module was leveraged to replace the potential moving average. Moreover, they introduced a novel rectification loss based on the ratio of potential mispredicted area loss to area size. Luo et al. [56] introduced a CTCT that employs both U-Net and Swin-UNet [57] to perform cross-teaching between the CNN and transformer. These approaches can face restricted variability and limited training stability.

2.3. Challenges and Improvements

The primary challenges in achieving high-accuracy medical image segmentation lie in leveraging limited and unlabeled data more effectively through DA and SSL. To leverage limited data, current methods encompass conventional DA, GAN-based DA, and automatic DA. However, each of these approaches faces unique challenges. The design of the conventional DA pipeline heavily relies on empirical knowledge. The GAN-based DA may experience inconsistent synthesis quality, high computational costs, and significant data requirements. The automatic DA can lead to considerable computational overhead, procedural complexity, and potential model confusion. As an alternative, we propose a training-free OAA approach that involves minimal computational costs and eliminates stage complexity. This approach augments input data with tailored transformations during training in an online augmentation manner to enhance generalizability. In terms of utilizing unlabeled data, existing architectures are primarily categorized into MT-based and miscellaneous methods. However, methods based on various architectures encounter varying challenges. MT-based approaches are limited by constrained teacher network performance and model variability. Miscellaneous methods can face limited variability and reduced stability. To this end, we propose a meticulously designed encoder-decoder architecture, incorporating the innovative RRE blocks and sensitivity loss. The RRE block diversifies features and introduces perturbations to yield varied decoder inputs, therefore enhancing variability. The sensitivity loss enhances the consistency across varying decoders and stabilizes the training process. We combine a multi-term loss function with a warming-up function to enhance training stability, especially during the early training stages.

3. Methods

3.1. Data-Efficient Medical Segmenter

The architecture of DEMS is depicted in Fig. 1. The DEMS utilizes an encoder-decoder architecture and incorporates the developed OAA and RRE blocks. It should

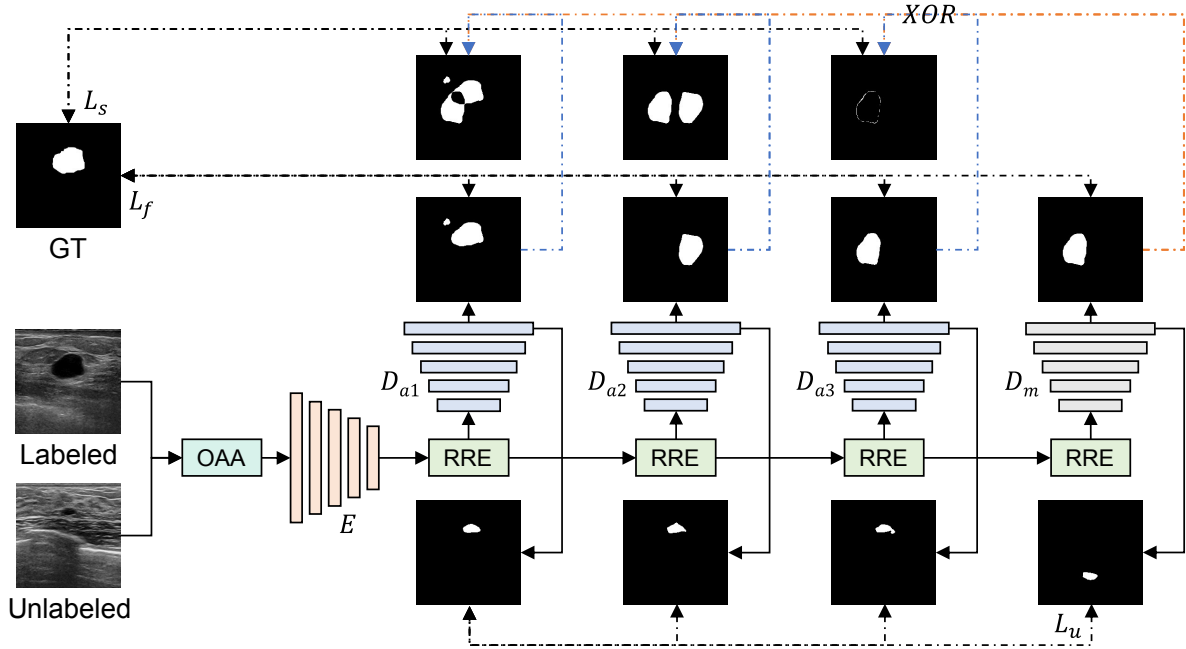


Figure 1: Detailed architecture of DEMS. The DEMS consists of an encoder E , a main decoder D_m , and three auxiliary decoders D_{a1} , D_{a2} , and D_{a3} . It comprises proposed online automatic augementer (OAA) and residual robustness enhancement (RRE) blocks. The OAA augments the visual input with varying image transformations to diversify the data to enhance generalization ability. The RRE block enriches the feature diversity and introduces perturbations to create varied decoder inputs, thereby bolstering the variability. The loss function comprises fusion loss L_f , sensitivity loss L_s , and unsupervised loss L_u . The introduced sensitivity loss further enhances the consistency across various decoders and stabilizes model training. The L_s is formulated based on the XOR operation across the main and auxiliary decoder pairs. The XOR operation is depicted with an overlay of blue and orange dash-dot lines. GT stands for the ground truth.

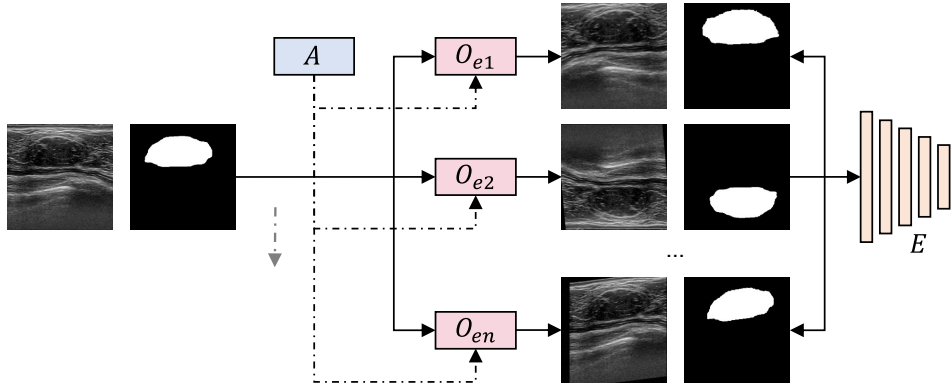


Figure 2: Workflow of the proposed OAA. Each input image and its corresponding mask undergo diverse DA transformations O sampled from augmentation spaces A at each new epoch e . The encoder receives diverse inputs in successive epochs as training proceeds, thereby enhancing the generalization capability. The gray dash-dot line depicts the training progress.

be noted that the auxiliary decoders are used exclusively during the training phase. The structure of the encoder E , main decoder D_m , and auxiliary decoders D_{a1} , D_{a2} , and D_{a3} follows that of the U-Net [58]. The OAA augments input data with a variety of image transformations, thereby diversifying the dataset to improve the generalization capability. The RRE block diversifies features and introduces perturbations to generate varying decoder inputs, thereby providing enhanced variability. The loss function L consists of fusion loss L_f and sensitivity loss L_s for labeled images, and unsupervised loss L_u for unlabeled images. We propose the sensitive loss to further enhance the consistency across various decoders and stabilize model training.

3.2. Online Automatic Augmenter

We revisit the fundamental aspects of MedAugment to establish the basis for our proposed improvements. MedAugment is an offline automatic DA approach that expands the input dataset into a larger dataset during a separate preparatory stage ahead of model training. It comprises $N = 4$ augmentation branches along with an additional branch. Each augmentation branch performs $M = [2, 3]$ transformations on each input. The transformations are sampled from the pixel augmentation space A_p and spatial augmentation space A_s using one of the sub-strategies in the sampling strategy Π . The strategy comprises four sub-strategies, sampling $[1, 0, 1, 0]$ and $[2, 3, 1, 2]$ transformations from A_p and A_s , respectively. The sampled transformations for each branch are subsequently shuffled. The maximum magnitude of each transformation M_A and the augment probability P_A are controlled based on the augmentation level L , and the applied magnitude is uniformly sampled within the maximum boundary. The additional branch is maintained to preserve the source visual information.

We introduce the OAA based on MedAugment to perform online automatic DA

Algorithm 1 Pseudocode of OAA.

Require: Pixel augmentation space $A_p = [\text{brightness, contrast, posterize, sharpness, Gaussian blur, Gaussian noise}]$, spatial augmentation space $A_s = [\text{rotate, horizontal flip, vertical flip, scale, x-axis translate, y-axis translate, x-axis shear, y-axis shear}]$, number of transformation $M = [2, 3]$, sampling strategy $\Pi = [\pi_1, \pi_2, \pi_3, \pi_4]$, augmentation level $L = 5$, maximum transformation magnitude M_A , transformation probability $P_A = 0.2L$, input image and mask pair P_{IM} ;

Ensure: Output pair P'_{IM} ;

- 1: Sample π from Π
 - 2: Sample M transformations $O = [o_1, \dots, o_M]$ using π from A_p, A_s
 - 3: Shuffle O
 - 4: **for all** o **do**
 - 5: Calculate M_A, P_A using L
 - 6: Uniformly sample magnitude $m_A \in M_A$
 - 7: **end for**
 - 8: $P'_{IM} = OP_{IM}$
 - 9: Out P'_{IM}
-

and diversify the input data to enhance the generalization capability. We consolidate multiple branches into a single main branch to perform one-to-one augmentation. Following this setup, each input data undergoes augmentation by the transformations sampled using one of the sub-strategies from Π . We uniformly distribute the probability of employing each sub-strategy. The additional branch is excluded as it does not contribute to the one-to-one augmentation setup. Given the limited number of transformations sampled from A_p , and considering that transformations in A_s do not compromise the effectiveness of medical images, we adopt sampling with replacement to offer enriched data diversity to enhance generalization ability. This modification offers a more comprehensive set of transformation combinations, making the combinations with identical transformations available for sampling. We illustrate the pseudocode of OAA in Algorithm 1 and it demonstrates the process of augmenting an input image and its corresponding mask within a single epoch. It should be noted that several of the transformations do not hold magnitude. In Fig. 2, we demonstrate the workflow of OAA, illustrating how an example input data is augmented throughout the training process. The input image and mask undergo diverse transformations at each successive epoch, thereby ensuring remarkable generalization ability. We emphasize that the OAA can be seamlessly integrated into the established DA pipeline [59] with one line of code, and utilizing it is as straightforward as leveraging arbitrary conventional DA transformations.

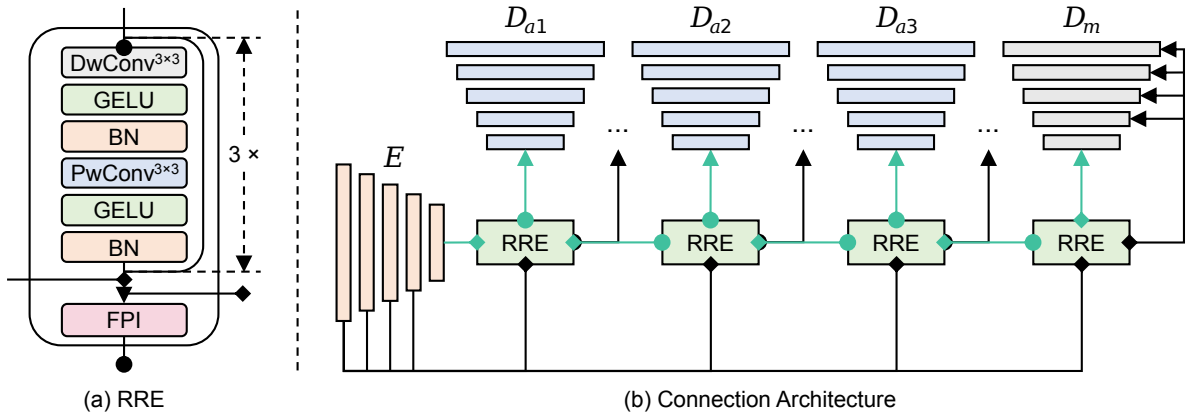


Figure 3: Detailed structure of the RRE block and the connection structure between the encoder and varying decoders. The RRE block features two distinct input-output pairs denoted with rhombus and circle by the shapes of the starting and ending arrows. It mainly encompasses residual connection, depthwise convolution (DwConv), pointwise convolution (PwConv), and feature perturbation injection (FPI) block. We denote the output of the encoder at varying blocks with f_1, f_2, f_3, f_4 , and f_5 . We represent the streams starting at f_1, f_2, f_3 , and f_4 with black arrows, and streams starting at f_5 with teal arrows. For clarity and conciseness, the skip connections across the encoder and main decoder are exclusively depicted. BN denotes the batch normalization layer.

3.3. Block and Connection Structure

We develop the RRE block to diversify features and introduce perturbations to produce varied decoder inputs to enhance variability. We illustrate the structure of the RRE block in Fig. 3a. The RRE block features two distinct input-output pairs symbolized by rhombus and circle, respectively. The primary components of the RRE block include residual connection [60], depthwise convolution (DwConv), pointwise convolution (PwConv), and feature perturbation injection (FPI) block. The residual connection helps train the deep network by mitigating the vanishing gradient problem, thus ensuring a smoother flow of information. The configuration of DwConv and PwConv endows the model with powerful representational capabilities while maintaining computational efficiency. The FPI block injects diverse perturbations encompassed from feature noise [30], feature dropout [30], and dropout [61].

The connection structure between the encoder and decoders is demonstrated in Fig. 3b. The encoder comprises five convolution blocks, with four max pooling layers inserted succeeding each of the first four blocks. Each convolution block consists of two 3×3 convolution layers, followed by a batch normalization layer and the GELU activation function. The decoder is structured into four upsampling stages, each consisting of an upsampling block, a feature concatenation layer, and a convolution block. The upsampling block comprises an upsampling layer, followed by a 3×3 convolution layer, a batch normalization layer, and the GELU activation function.

Given the encoder output at varying blocks f_1, f_2, f_3, f_4 , and f_5 , the first RRE block receives inputs from all blocks through two distinct input ports. Subsequent blocks take both f_1, f_2, f_3, f_4 , and outputs from preceding blocks as inputs. Additionally, the output ports associated with the skip connections between the main decoder and auxiliary decoders differ. The meticulous connection design ensures that each decoder input undergoes a distinct number of convolution operations, diverse perturbations, or both.

3.4. Loss Function

The fusion loss L_f fuses binary cross entropy loss L_{BCE} and dice loss L_{DSC} to ensure robust segmentation performance across various scales. Primarily, the L_{BCE} focuses on pixel-level accuracy, penalizing deviations of predicted pixel values from the actual labels. In contrast, the L_{DSC} emphasizes object-level accuracy by measuring the overlap between the prediction and the ground truth (GT). Given the GT y , decoder prediction \hat{y} , and pixel index i , the L_f for each decoder can be calculated using Eq. 1:

$$\begin{aligned} L_f &= 0.5 \cdot L_{BCE}(\hat{y}, y) + L_{DSC}(\hat{y}, y) \\ &= -\frac{1}{2N} \sum_{i=1}^N [y_i \cdot \log(\hat{y}_i) + (1 - y_i) \cdot \log(1 - \hat{y}_i)] + 1 - \frac{2 \cdot \sum_{i=1}^N \hat{y}_i y_i}{\sum_{i=1}^N \hat{y}_i + \sum_{i=1}^N y_i} \end{aligned} \quad (1)$$

We present a novel sensitivity loss L_s to further bolster the consistency across various decoders and enhance training stability. The formulation of L_s is based on an intuitive principle in binary segmentation where a mismatch in predictions across two different decoders suggests an error in one of them. To this end, we term the area predicted by different decoders with various binarized results as the sensitivity area and leverage the sensitive loss to encourage the decoder pairs to reduce the area size. Given the output of the main and arbitrary auxiliary decoder \hat{y}_m , and \hat{y}_a , and the binarization threshold T of 0.5, the L_s for each $\hat{y}_m - \hat{y}_a$ pair can be formulated based on XOR operation using Eq. 2:

$$L_s = \frac{1}{N} \sum_{i=1}^N (\hat{y}_{m,i} > T) \oplus (\hat{y}_{a,i} > T) \quad (2)$$

The unsupervised loss L_u is meticulously designed to harness the information embedded within unlabeled images. The implementation of L_u bears resemblance to that of the sensitivity loss, in which the essence is to boost the consistency across different decoders. We utilize the mean squared error loss to compute the L_u , and the computation for each $\hat{y}_m - \hat{y}_a$ pair is formulated as shown in Eq. 3:

$$L_u = \sum_{i=1}^N (\hat{y}_{m,i} - \hat{y}_{a,i})^2 \quad (3)$$

The loss functions delineated in Eq. 1, Eq. 2, and Eq. 3 compute the loss for each decoder or decoder pair, and the complete loss terms can be formulated based on their mean. Specifically, the \bar{L}_f is derived from the average losses for four individual decoders, while \bar{L}_s and \bar{L}_u are calculated based on the mean losses of three distinct decoder pairs. To reduce the risk of sensitivity and unsupervised losses destabilizing the training, especially in the initial stages, we progressively increase the weight of these two terms as training progresses. Specifically, we employ the Gaussian warming-up function λ [55] to adjust the weight based on the training step t and maximum training step t_{max} using Eq. 4:

$$\lambda(t, t_{max}) = \begin{cases} 1 & t_{max} = 0 \\ \exp\left(-5\left(1 - \frac{t}{t_{max}}\right)^2\right) & 0 \leq t \leq t_{max} \end{cases} \quad (4)$$

Considering the convergence and stability of the training process, we comprehensively integrate the three loss terms along with the Gaussian warming-up function to formulate the overall loss function L . The overall loss function embodies the cumulative attributes of each loss term and can be expressed using Eq. 5:

$$L = \bar{L}_f + \lambda(\bar{L}_s + \bar{L}_u) \quad (5)$$

4. Experiments

4.1. Datasets

We utilize both our own and three public datasets for performance evaluation. We choose ultrasound (US) datasets as quintessential examples due to their inherent challenges. Specifically, the US images not only suffer from speckle noise and low quality [62] but also face significant data shortages. Conducting experiments under these challenging conditions is crucial for producing more convincing and robust results. The designated datasets include the own stomach nasogastric tube (SNGT), public breast ultrasound (BUS) [63, 64], breast US images (BUSI) [65], and digital database thyroid image (DDTI) [66, 67].

The SNGT dataset consists of 221 US images with stomach and feeding tubes collected from two aspects. First, we conducted a retrospective search through the institutional imaging database PACS from 1 January 2020 to 31 December 2022 to retrieve existing stomach images, either with or without the feeding tube. Second, we prospectively recruited patients with in-situ feeding tubes for US scanning following the approval from the ethics review board and the acquisition of patient consent. The US scans were performed using standard commercial US machines from manufacturers including General Electric, Siemens, and Toshiba. A curvilinear probe operating at a frequency of 2-8 MHz was utilized to capture the images. The captured images were then cropped and four objects including the liver, stomach, tube, and pancreas were annotated using the open-source tool Labelme [68]. In this

study, we employ the SNGT dataset for tube binary segmentation based on our loss function design.

The BUS dataset is collected from the UDIAT Diagnostic Centre of the Parc Taulí Corporation utilizing the Siemens ACUSON Sequoia C512 system. It is composed of 163 images, including 110 images featuring benign lesions and 53 with cancerous masses. The average image resolution of the BUS dataset is 760×570 . The BUSI dataset is collected from 600 female patients aged between 25 and 75 years using the LOGIQ E9 and LOGIQ E9 Agile US systems. It encompasses 780 images, of which 437, 210, and 133 are benign, malignant, and normal images, respectively. The average image resolution of the BUSI dataset is approximately 500×500 . In this study, we leverage benign and malignant images for image segmentation. The DDTI dataset comprises 637 B-mode thyroid US images and includes a variety of lesions such as thyroiditis, cystic nodules, adenomas, and thyroid cancers. The US images were curated at the IDIME US Department in Colombia, with patient selection based on the TI-RADS description.

4.2. Setup and Evaluation

The datasets are divided into training and validation subsets in a ratio of 7:3. For semi-supervised methods, we utilize 20% and 40% of labeled images for training. Images and masks are preprocessed to a resolution of 224×224 . We employ SGD as the optimizer with a base learning rate of 0.01. The momentum equals 0.9 and the weight decay is set to 0.0001. The learning rate is updated using the cosine annealing schedule [69]. The implemented loss function is shown in Eq. 5, in which the λ is updated every 150 iterations. The batch size is set to 8 and the maximum training iterations equals 20000. We perform our experiments using AMD Ryzen 5965WX and NVIDIA RTX 4090. We report the mean and standard deviation of three independent runs with varied seeds in percentage. Model performance is evaluated using five evaluation metrics including dice score (DSC), intersection over union (IoU), sensitivity (SEN), precision (PRE), and pixel accuracy (PA). Metric calculations are conducted using true positives (TP), true negatives (TN), false positives (FP), and false negatives (FN), as detailed from Eq. 6 to Eq. 10:

$$DSC = \frac{2 \times TP}{2 \times TP + FP + FN} \quad (6)$$

$$IoU = \frac{TP}{TP + FP + FN} \quad (7)$$

$$SEN = \frac{TP}{TP + FN} \quad (8)$$

$$PRE = \frac{TP}{TP + FP} \quad (9)$$

$$PA = \frac{TP + TN}{TP + TN + FP + FN} \quad (10)$$

5. Results

5.1. Metric Comparison

Table 1: Performance of DEMS and SOTA methods on the SNGT dataset using 20% and 40% labeled data. We report the mean and standard deviation of three independent runs in percentage. The best results are highlighted in bold.

Method	Venue	Labeled	Unlabeled	DSC	IoU	SEN	PRE	PA
U-Net [58]	MICCAI	30 (20%)	0	37.74±2.50	27.78±1.17	44.84±6.33	40.87±1.83	99.09±0.10
U-Net	MICCAI	61 (40%)	0	47.47±0.92	36.96±1.05	47.10±1.20	57.46±2.25	99.31±0.02
CCT [30]	CVPR	30 (20%)	124 (80%)	34.27±1.32	25.73±0.87	35.35±3.89	42.08±1.61	99.22±0.08
UA-MT [55]	MICCAI			33.60±0.61	25.65±0.62	34.57±1.74	41.87±1.56	99.25±0.03
MC-Net+ [40]	MEDIA			38.68±0.61	29.69±0.57	38.69±0.92	46.99±2.36	99.29±0.02
MGCC [38]	ArXiv			44.22±0.01	33.62±0.64	50.69±2.06	47.41±1.96	99.18±0.06
CPS [70]	CVPR			33.66±1.62	25.50±0.74	34.97±3.85	39.48±0.56	99.23±0.04
CTCT [56]	PMLR			41.83±1.08	31.82±0.66	44.28±2.37	50.64±1.17	99.32±0.00
ICT [71]	NN			34.13±0.98	26.52±0.73	33.61±1.32	41.28±0.50	99.28±0.02
R-Drop [72]	NIPS			32.38±1.14	24.79±0.92	32.71±1.88	39.62±0.64	99.27±0.01
URPC [73]	MICCAI			35.42±0.90	27.68±0.72	35.19±2.72	42.92±1.95	99.28±0.02
DEMS (Ours)	-			54.59±0.44	43.39±0.54	56.84±0.64	59.60±2.41	99.32±0.03
CCT [30]	CVPR	61 (40%)	93 (60%)	42.81±0.86	32.69±0.69	45.92±1.69	49.69±1.23	99.29±0.02
UA-MT [55]	MICCAI			46.34±1.38	35.88±0.77	48.38±2.51	54.29±1.04	99.32±0.02
MC-Net+ [40]	MEDIA			48.30±0.84	37.67±0.71	49.67±0.84	57.95±2.23	99.36±0.02
MGCC [38]	ArXiv			53.76±0.80	42.74±0.70	58.18±1.05	57.45±2.23	99.34±0.02
CPS [70]	CVPR			42.50±0.99	32.80±0.71	43.00±1.61	53.84±2.68	99.33±0.02
CTCT [56]	PMLR			46.01±1.09	35.78±0.68	48.11±1.25	54.69±2.16	99.34±0.01
ICT [71]	NN			46.47±1.69	36.20±1.18	48.48±2.19	55.17±2.15	99.34±0.01
R-Drop [72]	NIPS			38.68±1.04	29.79±0.75	38.10±0.89	49.51±1.63	99.32±0.01
URPC [73]	MICCAI			44.09±0.45	34.58±0.57	45.46±0.75	51.65±0.74	99.35±0.01
DEMS (Ours)	-			59.66±0.21	48.21±0.62	61.48±1.01	64.52±0.46	99.37±0.02

We show the performance comparison between DEMS and SOTA methods on the SNGT dataset in Tab. 1. As shown, the DEMS outperforms SOTA methods to a large extent on the SNGT dataset, followed by the MGCC. The DEMS achieves a DSC of 54.59% with merely 30 labeled images for training. This result represents an improvement of 16.85% and 10.37% over U-Net and MGCC, respectively. When the number of labeled images increases to 61, the DEMS reaches a DSC of 59.66%, outperforming the U-Net and MGCC with 12.19% and 5.90% DSC leadership, respectively. It should be noted that the uniformly high PA across all methods does not necessarily indicate superior performance but may be influenced by the small size of the object being segmented. In the SNGT dataset, the overall size of the tube is significantly smaller than the background, allowing models to achieve high PA even if categorizing all pixels as background. We illustrate the performance of various methods on the BUS, BUSI, and DDTI datasets in Tab. 2, Tab. 3, and Tab. 4, respectively. From the observations, it is evident that the DEMS outperforms the SOTA methods at varying degrees. As shown in Tab. 2, the DEMS reaches a DSC of 76.14% and 83.03% using 20% and 40% labeled images for training, respectively. Conversely,

Table 2: Performance of DEMS and SOTA methods on the BUS dataset using 20% and 40% labeled data.

Method	Venue	Labeled	Unlabeled	DSC	IoU	SEN	PRE	PA
U-Net [58]	MICCAI	22 (20%)	0	61.68±1.60	51.24±1.13	62.42±3.26	71.60±3.82	96.96±0.35
U-Net	MICCAI	45 (40%)	0	74.02±1.76	65.04±1.34	72.54±3.13	84.23±0.90	97.95±0.03
CCT [30]	CVPR	22 (20%)	92 (80%)	61.91±1.17	52.62±0.81	58.64±0.87	77.89±1.43	96.97±0.07
UA-MT [55]	MICCAI			62.09±1.14	52.77±0.75	61.69±2.38	73.63±2.28	97.08±0.06
MC-Net+ [40]	MEDIA			65.92±0.55	56.56±0.54	64.26±1.78	76.62±1.65	97.27±0.09
MGCC [38]	ArXiv			72.35±0.88	60.72±0.83	73.05±3.90	78.64±3.17	97.64±0.03
CPS [70]	CVPR			55.42±0.63	45.74±0.57	53.39±0.90	70.38±2.01	96.67±0.06
CTCT [56]	PMLR			63.86±1.77	54.95±1.85	62.50±1.47	76.67±2.73	97.21±0.07
ICT [71]	NN			63.82±0.74	54.44±0.99	62.91±1.03	76.98±2.22	97.23±0.08
R-Drop [72]	NIPS			54.90±1.25	46.76±0.80	52.16±0.50	72.54±5.88	96.72±0.15
URPC [73]	MICCAI			59.74±0.95	51.13±1.02	57.38±0.37	72.35±2.37	96.87±0.09
DEMS (Ours)	-			76.14±0.03	65.59±0.04	81.12±0.84	77.80±0.85	97.69±0.01
CCT [30]	CVPR	45 (40%)	69 (60%)	76.15±0.52	67.48±0.58	74.74±0.57	85.17±1.39	97.93±0.08
UA-MT [55]	MICCAI			73.96±1.02	65.87±0.81	72.11±1.07	84.53±0.68	97.98±0.13
MC-Net+ [40]	MEDIA			76.10±0.66	67.67±0.62	73.50±0.93	86.66±2.78	98.08±0.07
MGCC [38]	ArXiv			82.99±1.05	73.98±0.99	81.84±3.22	88.96±2.10	98.33±0.08
CPS [70]	CVPR			69.83±0.74	61.81±0.68	67.05±1.17	83.15±1.51	97.74±0.06
CTCT [56]	PMLR			75.11±0.98	66.85±0.74	72.12±1.24	85.50±1.58	98.14±0.03
ICT [71]	NN			74.89±0.37	66.63±0.66	73.06±0.78	86.14±2.30	98.02±0.06
R-Drop [72]	NIPS			65.59±0.66	57.77±0.79	62.15±1.18	80.17±3.20	97.30±0.12
URPC [73]	MICCAI			71.22±0.49	62.52±0.57	67.11±0.73	83.26±1.15	97.76±0.01
DEMS (Ours)	-			83.03±1.42	74.13±1.66	85.11±0.61	84.72±1.21	98.47±0.05

CCT and MGCC outperform DEMS on PRE by margins of 0.09% and 4.24%, respectively. Regarding the metrics illustrated in Tab. 3, the DEMS achieves a DSC of 76.01% using 90 labeled images for training. The second highest DSC is achieved by MGCC with the highest PRE of 79.06%. When the number of labeled images equals 180, the DEMS realizes the highest SEN of 81.04% while slightly underperforming MGCC with a DSC lag of 0.35%. The MGCC also presents a marginally higher IoU, PRE, and PA. For the results indicated in Tab. 4, the DEMS reaches the highest DSC of 73.90% and 77.78% when trained using 20% and 40% labeled images, respectively. Nevertheless, MGCC reaches the highest PRE of 78.99% when trained using 40% labeled images, achieving a 0.87% leadership compared with DEMS.

5.2. Analysis and Visualization

We present the predicted masks across DEMS and SOTA methods trained using 40% labeled images on different datasets in Fig. 4. The illustration demonstrates that the masks predicted by DEMS show the most closely resemble compared with the GT. Specifically, the DEMS can predict the contour and the size of the objects most accurately. In the first and second columns that feature small object predictions, most methods such as ICT misidentify non-existent areas and fail to highlight the correct regions. For regular object predictions, although most methods can accurately identify the majority of regions, certain approaches such as CTCT misidentify a large number of areas and can overlook the correct ones. To further illustrate

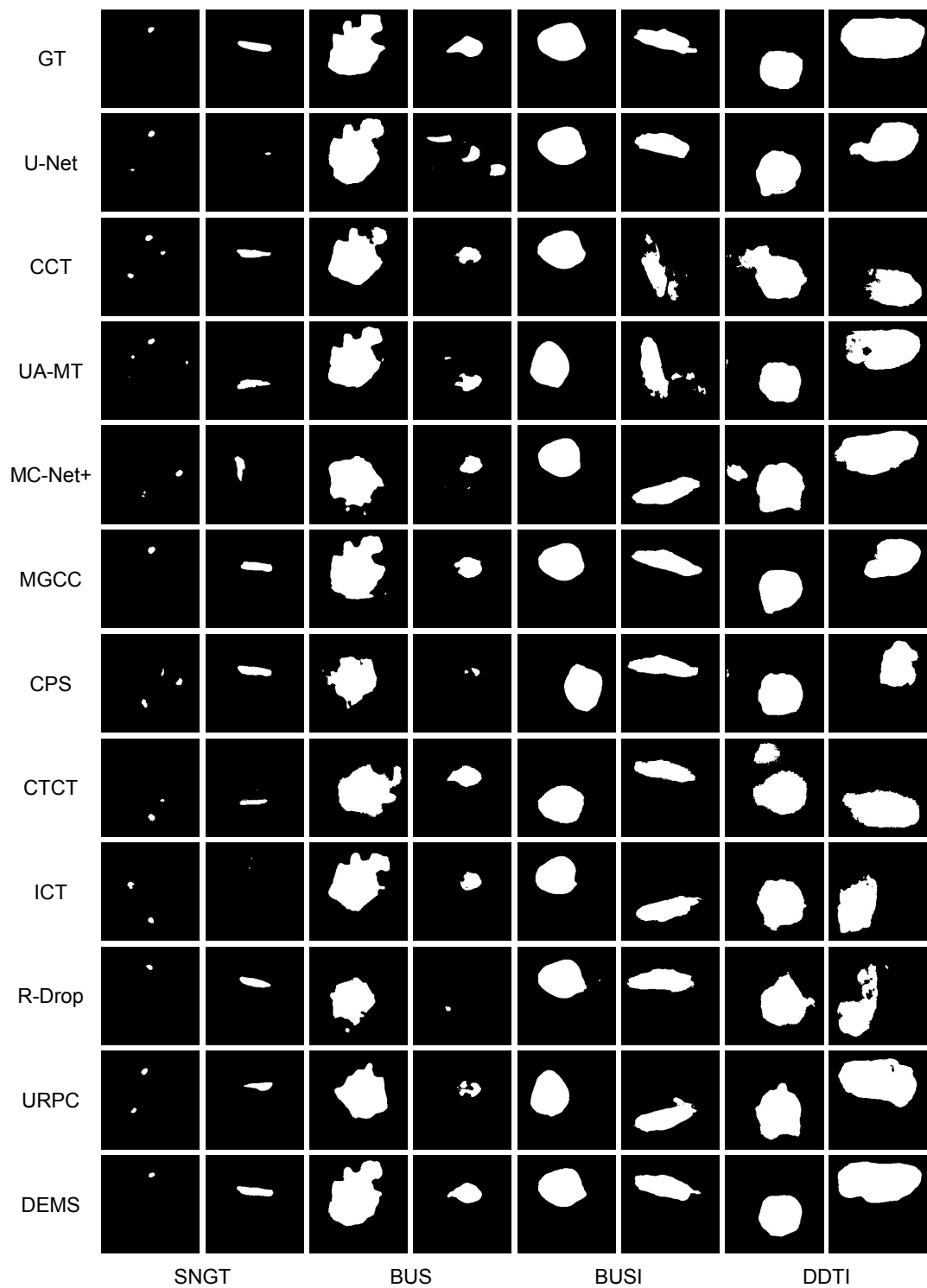


Figure 4: Predicted masks across DEMS and SOTA methods on the four datasets using 40% labeled data.

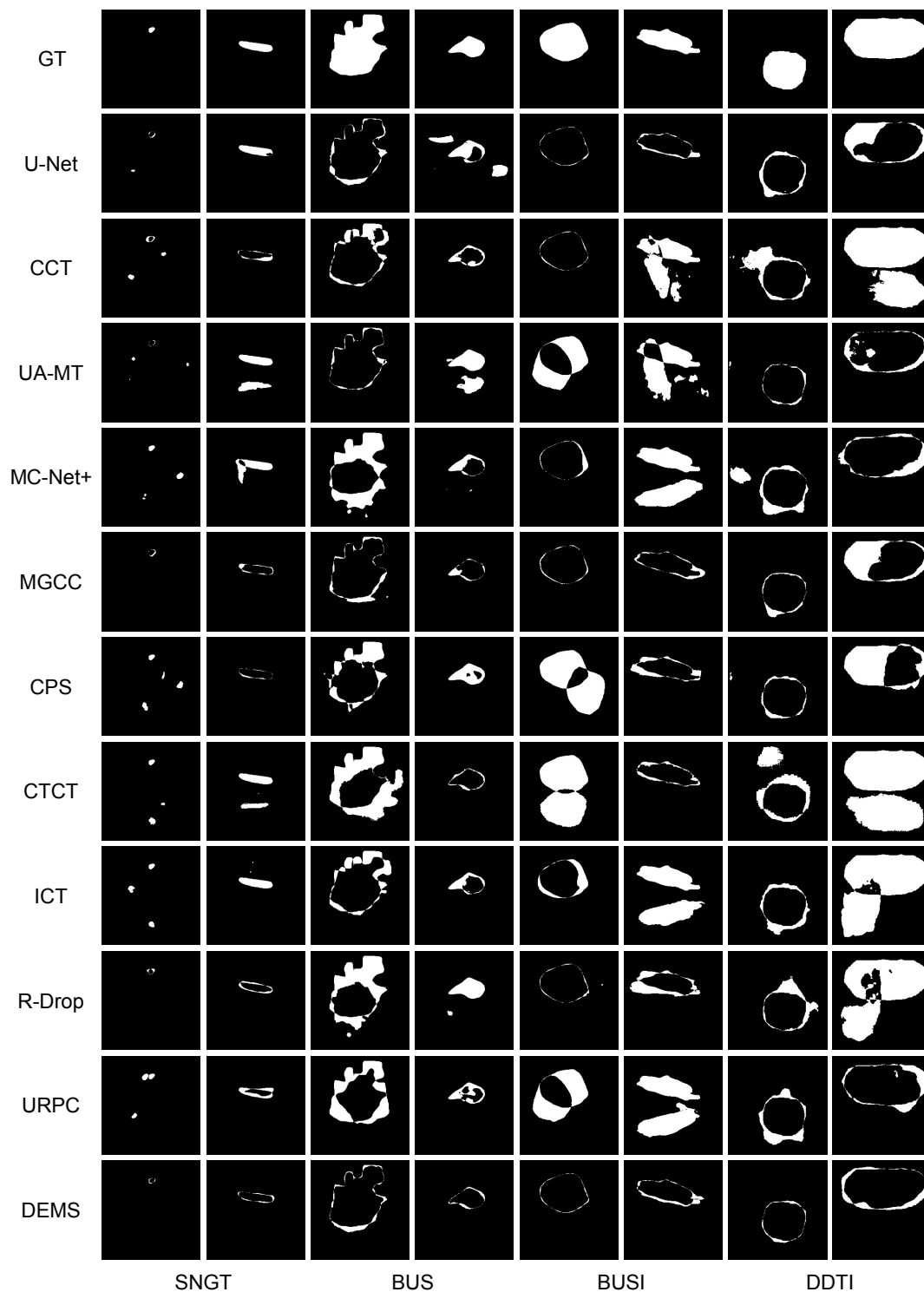


Figure 5: XOR outputs of predicted and GT masks across DEMS and SOTA methods on the four datasets using 40% labeled data.

Table 3: Performance of DEMS and SOTA methods on the BUSI dataset using 20% and 40% labeled data.

Method	Venue	Labeled	Unlabeled	DSC	IoU	SEN	PRE	PA
U-Net [58]	MICCAI	90 (20%)	0	70.49±0.82	60.73±0.70	74.32±1.36	75.07±1.01	93.45±0.03
U-Net	MICCAI	180 (40%)	0	76.18±0.61	66.56±0.60	78.08±2.22	80.24±2.10	94.81±0.10
CCT [30]	CVPR	90 (20%)	362 (80%)	70.04±0.67	60.82±0.68	70.48±0.81	77.30±0.45	93.98±0.09
UA-MT [55]	MICCAI			69.09±0.43	59.68±0.53	70.69±0.71	76.14±1.36	93.72±0.17
MC-Net+ [40]	MEDIA			69.78±1.05	60.71±0.67	69.89±2.63	77.99±0.65	93.94±0.09
MGCC [38]	ArXiv			75.10±0.70	65.95±0.95	76.71±1.38	79.06±0.72	94.40±0.25
CPS [70]	CVPR			67.59±0.73	58.77±0.65	68.14±0.30	75.84±1.44	93.68±0.07
CTCT [56]	PMLR			69.66±0.76	60.82±0.62	69.94±1.23	77.97±0.94	93.98±0.07
ICT [71]	NN			69.31±0.19	59.69±0.55	69.91±0.54	77.45±0.94	93.84±0.17
R-Drop [72]	NIPS			64.71±0.77	55.99±0.81	64.91±0.61	75.61±2.00	93.29±0.09
URPC [73]	MICCAI			68.91±0.97	59.92±0.78	68.49±1.34	77.54±0.64	94.02±0.00
DEMS (Ours)	-			76.01±0.54	66.84±0.55	77.76±0.26	79.51±1.12	94.47±0.22
CCT [30]	CVPR	180 (40%)	272 (60%)	72.09±0.63	62.79±0.68	72.95±1.90	78.29±0.83	94.38±0.03
UA-MT [55]	MICCAI			72.50±0.65	62.91±0.82	74.50±0.51	77.29±1.51	94.32±0.23
MC-Net+ [40]	MEDIA			73.01±0.71	63.88±0.68	75.43±0.50	77.66±1.11	94.55±0.22
MGCC [38]	ArXiv			79.52±0.66	70.71±0.65	79.89±0.86	83.47±0.79	95.20±0.20
CPS [70]	CVPR			70.85±0.57	61.69±0.59	70.72±1.04	78.10±0.38	94.18±0.25
CTCT [56]	PMLR			72.13±0.84	62.92±0.82	74.01±1.75	77.06±0.52	94.41±0.14
ICT [71]	NN			72.20±0.53	62.78±0.60	74.75±2.20	77.04±2.29	94.15±0.38
R-Drop [72]	NIPS			72.06±0.88	62.68±0.68	73.38±2.10	77.97±1.18	94.49±0.14
URPC [73]	MICCAI			71.33±0.53	61.82±0.67	74.06±1.49	75.26±0.25	94.11±0.07
DEMS (Ours)	-			79.17±0.19	70.09±0.26	81.04±1.69	82.52±1.70	95.06±0.38

the advantages of DEMS, we demonstrate the XOR outputs between prediction and GT masks across various methods on the four datasets in Fig. 5. Observations reveal that DEMS consistently produces the smallest XOR areas in various images, maintaining high prediction accuracy even for objects with complex geometries depicted in the third column. In addition to DEMS, the MGCC also shows desirable performance in most cases except for the examples in the last column. As for the remaining methods, they may accurately predict one or several of the images but often underperform on others. The superior prediction accuracy across varying objects demonstrates the desirable feature capture ability of the developed DEMS.

To thoroughly investigate the performance of DEMS, we additionally train DEMS using 60% and 80% labeled images and compare its performance with the U-Net trained with 100% labeled images. We term the performance of the U-Net in this setup as the upper bound. Observing Fig. 6, it becomes clear that DEMS consistently outperforms the U-Net upper bound across all datasets, achieving the highest DSC of over 68%, 87%, 82%, and 80% on the SNGT, BUS, BUSI, and DDTI datasets, respectively. Additionally, two patterns are observed. Firstly, DEMS can surpass the upper bound by utilizing a smaller percentage of labeled images on smaller datasets. Specifically, on the smaller SNGT and BUS datasets, the DEMS surpasses the upper bound using merely 40% labeled images. For the relatively larger BUSI and DDTI datasets, the DEMS exceeds the U-Net upper bound using 80% and 60% labeled images for training, respectively. Moreover, the performance advantage is

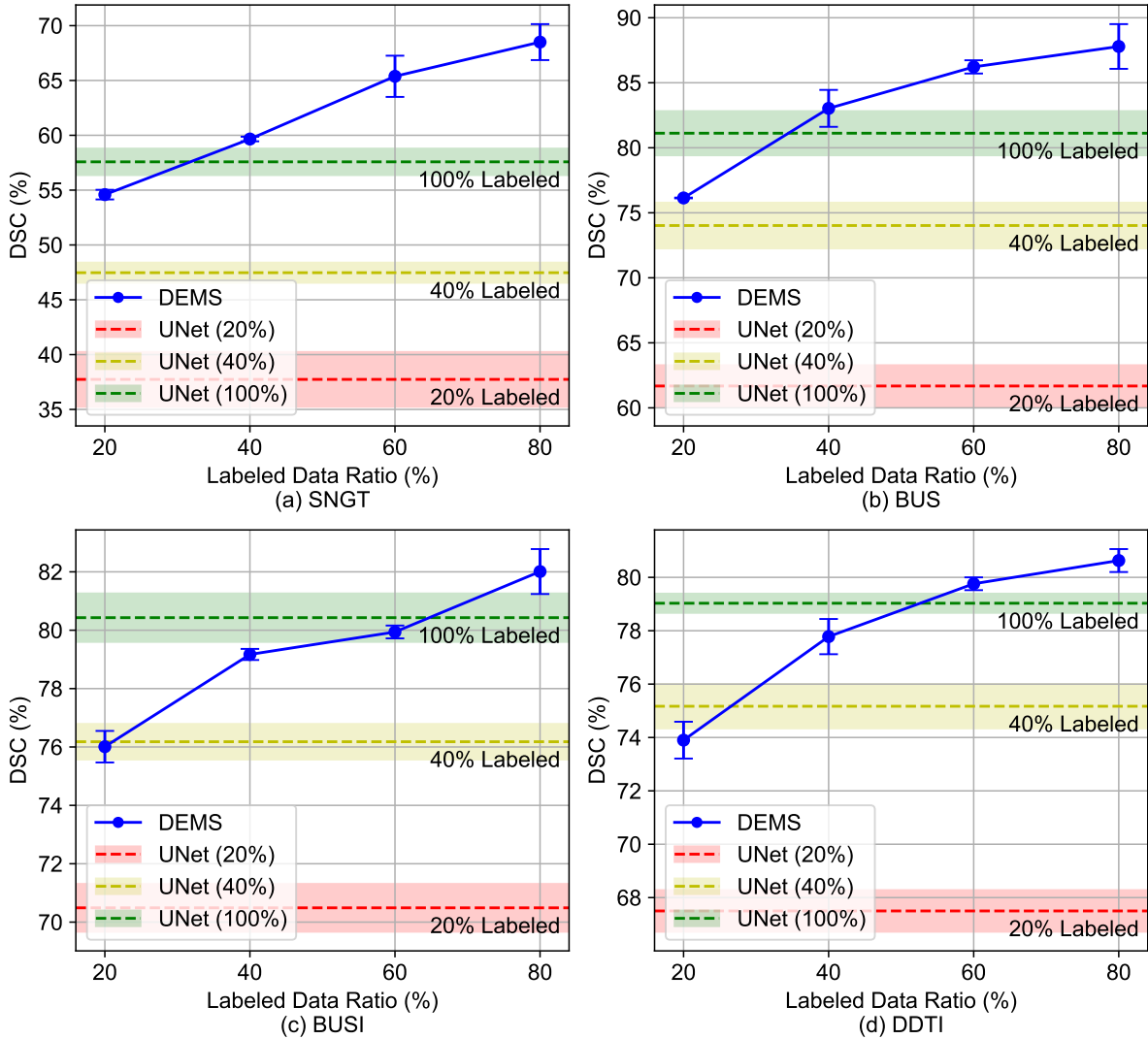


Figure 6: Boundary performance among DEMS and U-Net on the four datasets. The error bar depicts the standard deviation.

Table 4: Performance of DEMS and SOTA methods on the DDTI dataset using 20% and 40% labeled data.

Method	Venue	Labeled	Unlabeled	DSC	IoU	SEN	PRE	PA
U-Net [58]	MICCAI	89 (20%)	0	67.50±0.78	54.96±0.99	76.25±1.49	68.94±1.27	92.79±0.25
U-Net	MICCAI	178 (40%)	0	75.17±0.83	63.70±0.65	81.38±3.24	75.35±1.73	94.29±0.12
CCT [30]	CVPR	89 (20%)	356 (80%)	68.62±0.72	56.71±0.58	77.54±1.58	69.35±1.93	92.75±0.34
UA-MT [55]	MICCAI			68.72±0.56	56.72±0.55	76.60±0.16	69.50±1.20	92.99±0.11
MC-Net+ [40]	MEDIA			69.31±0.58	57.60±0.56	75.90±2.00	71.64±1.17	93.22±0.11
MGCC [38]	ArXiv			67.71±0.77	54.65±0.58	73.99±2.90	71.25±1.48	92.47±0.19
CPS [70]	CVPR			67.50±0.65	55.81±0.65	72.81±1.18	71.10±1.69	92.99±0.27
CTCT [56]	PMLR			69.66±0.80	57.86±0.68	77.25±1.35	70.90±0.80	93.18±0.07
ICT [71]	NN			69.33±0.78	57.85±0.73	76.01±0.69	70.72±1.68	93.23±0.24
R-Drop [72]	NIPS			65.93±0.79	53.86±0.78	72.82±4.61	69.79±4.57	92.56±0.51
URPC [73]	MICCAI			68.73±0.41	56.80±0.68	74.85±1.58	71.45±1.10	93.09±0.15
DEMS (Ours)	-			73.90±0.69	61.86±0.79	80.22±0.83	74.72±0.65	93.87±0.21
CCT [30]	CVPR	178 (40%)	267 (60%)	72.07±0.79	60.86±0.72	79.98±1.94	71.99±0.27	93.83±0.06
UA-MT [55]	MICCAI			71.77±0.59	60.88±0.68	77.06±0.52	74.12±1.10	94.06±0.22
MC-Net+ [40]	MEDIA			73.67±0.69	62.75±0.57	78.60±1.05	76.18±0.48	94.35±0.11
MGCC [38]	ArXiv			76.53±0.72	65.89±0.86	79.73±0.67	78.99±0.73	94.84±0.15
CPS [70]	CVPR			72.02±0.88	60.87±0.89	76.80±1.67	74.70±0.28	94.10±0.11
CTCT [56]	PMLR			73.60±0.58	62.64±0.63	79.45±0.98	74.61±0.54	94.41±0.17
ICT [71]	NN			73.09±0.51	61.84±0.69	79.90±0.53	73.82±0.63	93.97±0.19
R-Drop [72]	NIPS			68.99±1.09	57.95±0.78	72.33±1.71	74.16±0.80	93.89±0.21
URPC [73]	MICCAI			72.72±0.59	61.89±0.75	78.61±0.67	74.16±2.06	94.20±0.24
DEMS (Ours)	-			77.78±0.66	67.22±0.43	83.48±0.46	78.12±0.81	95.03±0.05

more pronounced on the smaller datasets compared to larger ones. Taking training with 80% labeled images as an example, the DEMS achieves a DSC superiority of approximately 11% and 7% over the U-Net upper bound on the smaller SNGT and BUS datasets. However, this superiority diminishes to about 2% on the relatively larger BUSI and DDTI datasets. Extensive comparisons across various datasets reveal that DEMS can utilize fewer labeled images to achieve greater performance improvements, showcasing its superior data efficiency.

Given that MGCC marginally outperforms DEMS, as noted in Tab. 3, we conduct extensive supplementary experiments using 60% and 80% labeled images to assess the performance across various methods more comprehensively. We show the performance of different methods under the complete range of labeled percentages across four datasets in Fig. 7. Through observation, it is evident that the DEMS achieves the highest performance relative to other SOTA methods as the number of labeled images increases. We observe prominent leadership on the SNGT dataset and moderate leadership on the BUS, BUSI, and DDTI datasets. Interestingly, enhanced DSC performance is sometimes observed with fewer labeled images. For instance, the MC-Net+ reaches a higher DSC using 60% labeled images compared with 80% labeled images on the SNGT dataset. This phenomenon can be attributed to the limited model stability due to noise susceptibility or overfitting under severe data shortages. This inference aligns with the observations that improved performance with fewer labeled data occurs more frequently on the smaller SNGT and

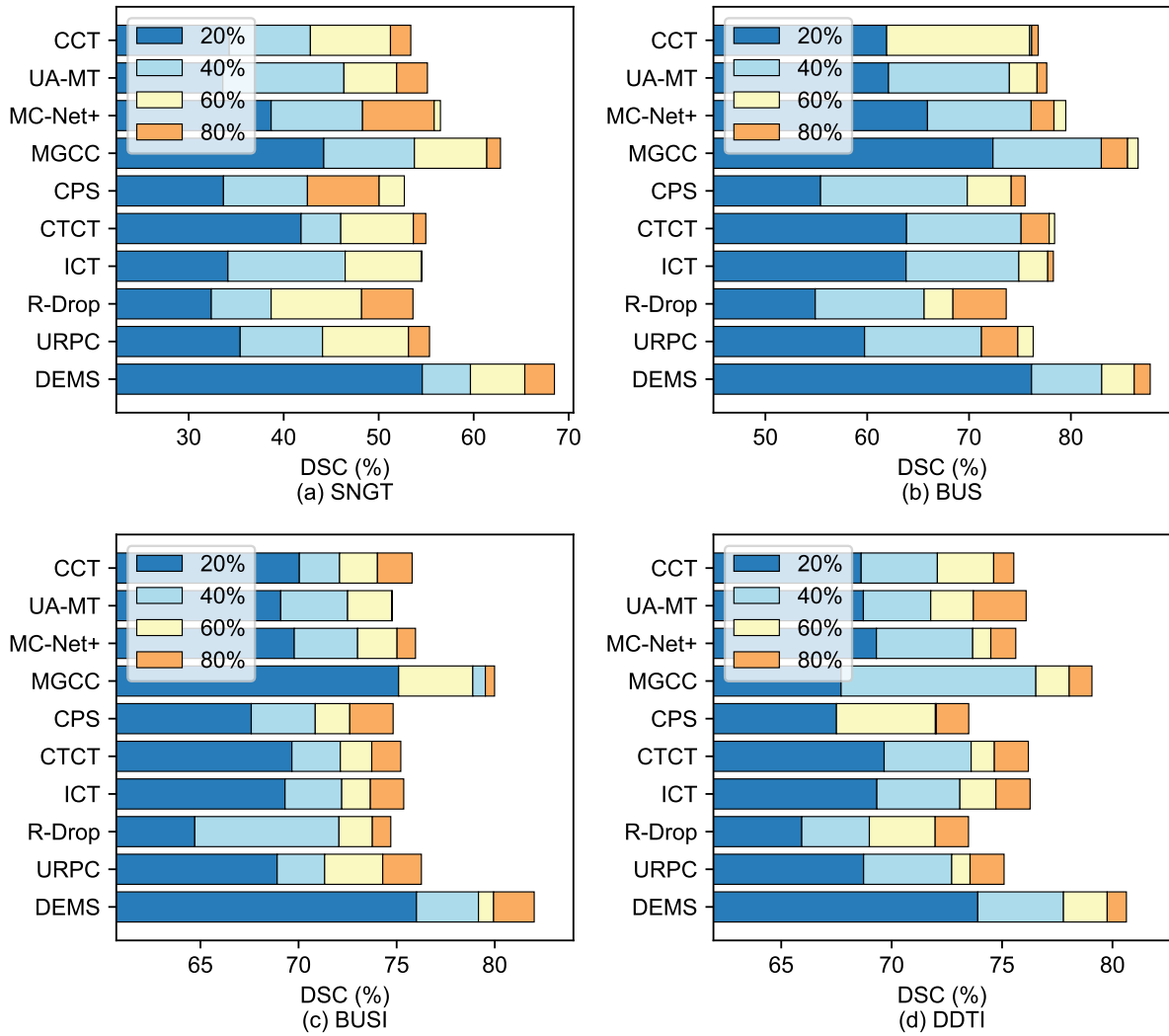


Figure 7: Performance of DEMS and SOTA methods under the complete range of labeled percentages on the four datasets.

BUS datasets than on the larger BUSI and DDTI datasets. Furthermore, the results depict that performance improvements are more modest on the larger datasets as the percentage of labeled images increases. This outcome is expected, as the model may capture relatively abundant features using a low percentage of labeled images on larger datasets.

We employ the Bland-Altman plot to visualize the prediction consistency for each image in Fig. 8. We count the object pixels in the predicted and GT masks, scale each by dividing by 224^2 , and visualize each prediction-GT pair. The analysis clearly shows the predictions of DEMS demonstrate superior consistency in comparison with the GT. The mean differences are notably minimal, hovering close to 0%, with the maximal and minimal results of 1.05% and 0.22% for the BUSI and SNGT datasets, respectively. The standard deviation lies within an ideal range and varies following the actual size of various objects. Specifically, the highest and lowest standard deviations of 5.41% and 0.63% are achieved on the BUSI and SNGT datasets, respectively. Furthermore, most points fall within the 95% limits of agreement, indicating relatively high consistency between the predicted and GT masks. It should be noted that several outliers are observed in each dataset, indicating the presence of inaccuracies in the prediction of these masks. This observation meets our expectations, as datasets containing hundreds of images might include patterns that are evident in the validation subset but absent or underrepresented in the training subset. Additionally, the datasets exhibit no uniform trend in differences, suggesting an absence of a strict systematic pattern in the observed discrepancies. An exception to this is the SNGT dataset, in which the difference increases with the escalation of the mean pixel ratio. This can be attributed to the fact that most of the tubes in the SNGT dataset are substantially small, thereby compelling the model to underpredict the size of the larger tubes. This assumption is consistent with the observation that the outliers predominantly fall along the negative axis. In this scenario, the inconsistency between predictions and GT is predominantly influenced by the dataset-specific attributes rather than the inherent constraints of DEMS.

We conduct cross-dataset evaluation experiments to assess the generalization capability of the proposed DEMS. The DEMS is trained on the BUS dataset and evaluated on the BUSI dataset for inter-dataset experimentation, with the results being compared to those from intra-dataset experiments. As shown in Fig. 9, the DEMS exhibits superior generalization ability and achieves relatively minor performance decreases across various metrics. When trained with 20% labeled images, DEMS achieves comparable performance in both inter-dataset and intra-dataset settings. This observation may be attributed to the fact that insufficient labeled data does not adequately highlight significant differences between the two datasets. With an increase in the number of labeled images, the performance gap widens and the largest gap is observed when training the DEMS with 40% labeled data. In this scenario, we observe a performance decrease of around 9%, 9%, 13%, and 1% in DSC, IoU, PRE, and PA, respectively. Conversely, the SEN shows an approximate 1% improvement in performance. As the number of labeled training images

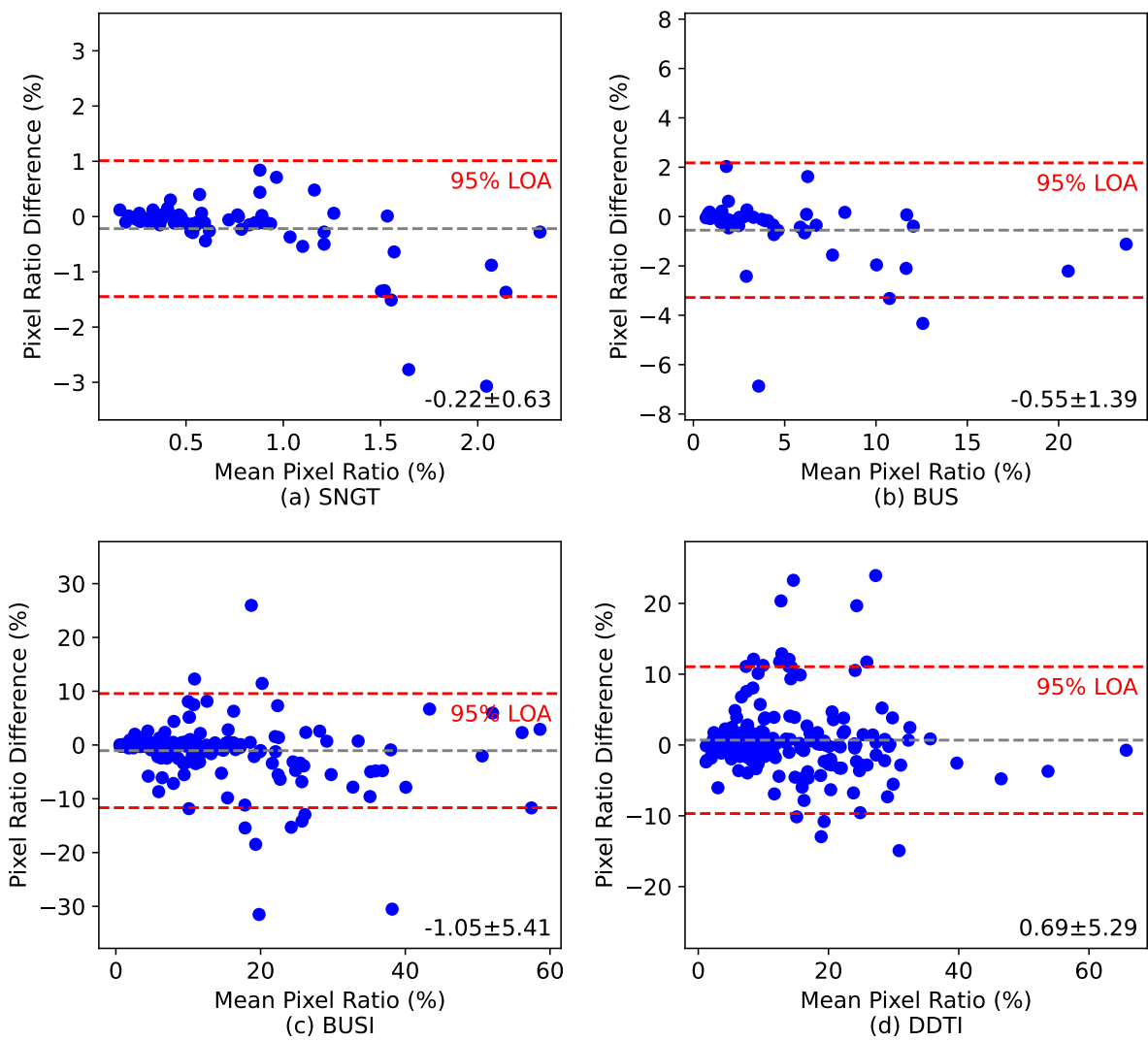


Figure 8: Bland-Altman analysis of DEMS on the four datasets using 80% labeled images. LOA denotes the 95% limits of agreement.

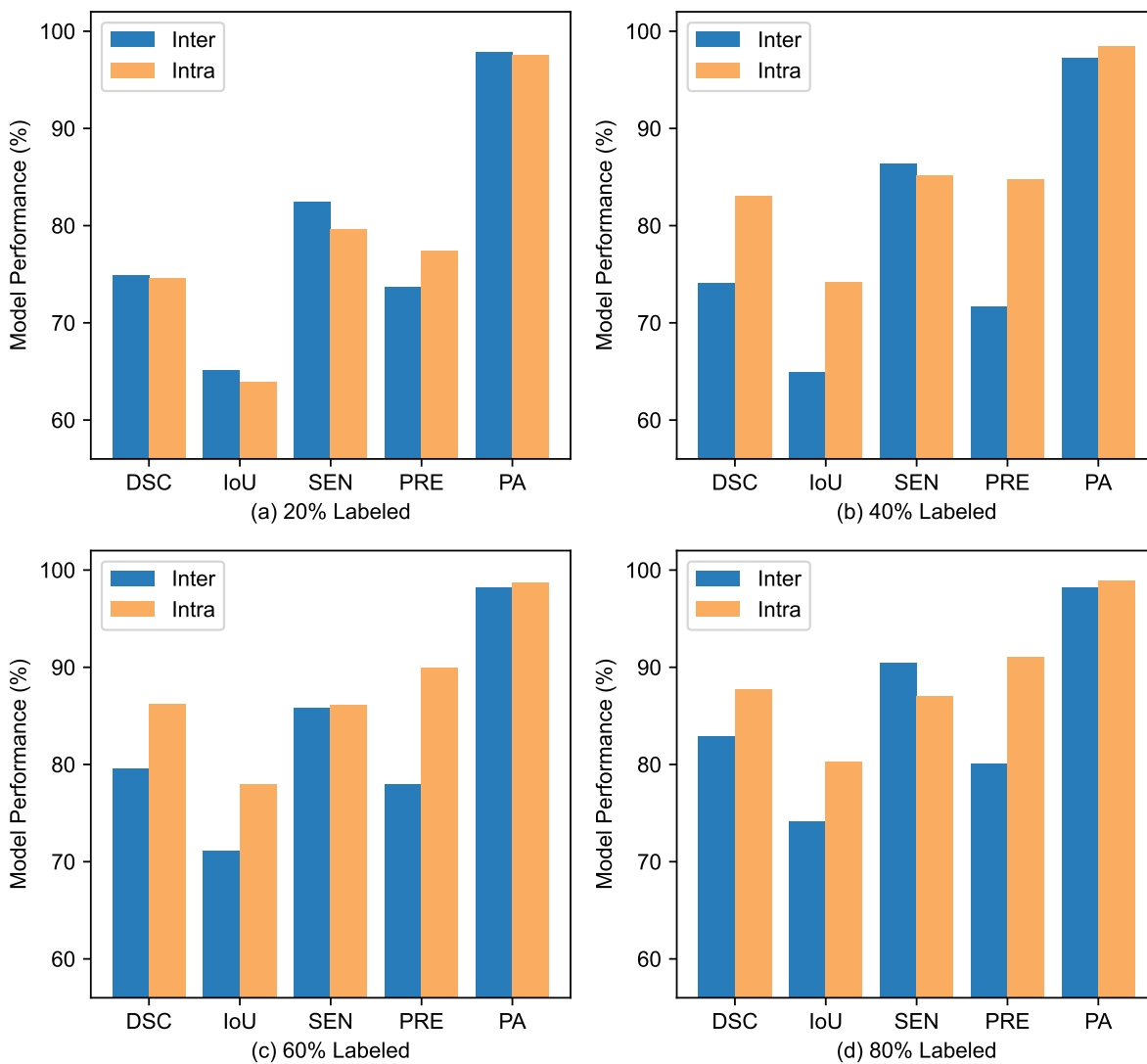


Figure 9: Performance of DEMS on cross-dataset evaluation experiments. The DEMS is trained on the BUS dataset and evaluated on the BUSI and BUS datasets for inter-dataset and intra-dataset configurations, respectively.

increases, the performance disparity between inter-dataset and intra-dataset experiments diminishes. This indicates improved generalization capability as additional visual features are incorporated. When leveraging 80% labeled images for training, the performance gap in metrics reduces to approximately 5%, 6%, 9%, and 1% for DSC, IoU, PRE, and PA, respectively. Similarly, the SEN exhibits an around 3% performance increase. The higher SEN in the inter-dataset experiments compared to the intra-dataset experiments may potentially be attributed to less pronounced features in negative examples within the BUSI dataset, leading to an increased likelihood of negatives being recognized as positives. The relatively minor decline in performance underscores the robust generalization capability of the proposed approach.

To evaluate the computational efficiency across models for practical application, we visualize the relationship between model training time and performance metrics across various methods in Fig. 10. The results show that DEMS demonstrates superior computational efficiency, achieving the highest evaluation metrics with reasonable time consumption. Most methods require approximately 2 hours to complete training, while MGCC constitutes an exception with a requirement of about 10 hours. DEMS achieves the highest performance, followed by MGCC and other methods. Furthermore, it exhibits relatively low standard deviations, indicating superior stability throughout the training process. Considering both training time and model performance, DEMS not only achieves the highest metrics but also maintains reasonable time consumption. Although MGCC demonstrates relatively high performance, it incurs considerably higher computational costs. The remaining methods consume reasonable training time but fall short of achieving ideal performance. The exceptional computational efficiency of DEMS highlights its potential for practical applications.

5.3. Ablation Study

Table 5: Ablation study of the OAA, RRE block, and sensitivity loss on the SNGT dataset with 80% labeled data.

OAA	RRE	L_{sen}	DSC	IoU	SEN	PRE	PA
			56.65±0.71	45.32±0.58	58.70±1.19	62.14±2.07	99.32±0.02
✓			63.50±1.32	52.28±1.19	63.51±1.52	72.32±2.32	99.43±0.04
	✓		62.06±1.41	49.92±0.93	62.93±2.81	70.04±0.95	99.41±0.02
		✓	58.42±0.29	46.94±0.15	60.57±1.81	63.58±2.54	99.36±0.02
✓	✓		66.83±1.27	55.63±1.18	67.86±0.86	71.60±3.93	99.45±0.04
✓		✓	65.12±0.91	54.04±0.98	65.87±3.13	71.82±3.00	99.44±0.01
	✓	✓	63.60±0.81	51.38±1.09	65.57±1.72	68.63±1.95	99.41±0.04
✓	✓	✓	68.50±1.64	56.92±1.22	69.70±2.32	72.30±3.77	99.46±0.04

To investigate the effectiveness of the OAA, RRE block, and L_{sen} within DEMS, we conduct extensive ablation experiments and present the results in Tab. 5. It is worth noting that DA transformations including random rotation and random flip are incorporated when OAA is removed. The reason for this is that exclusively removing

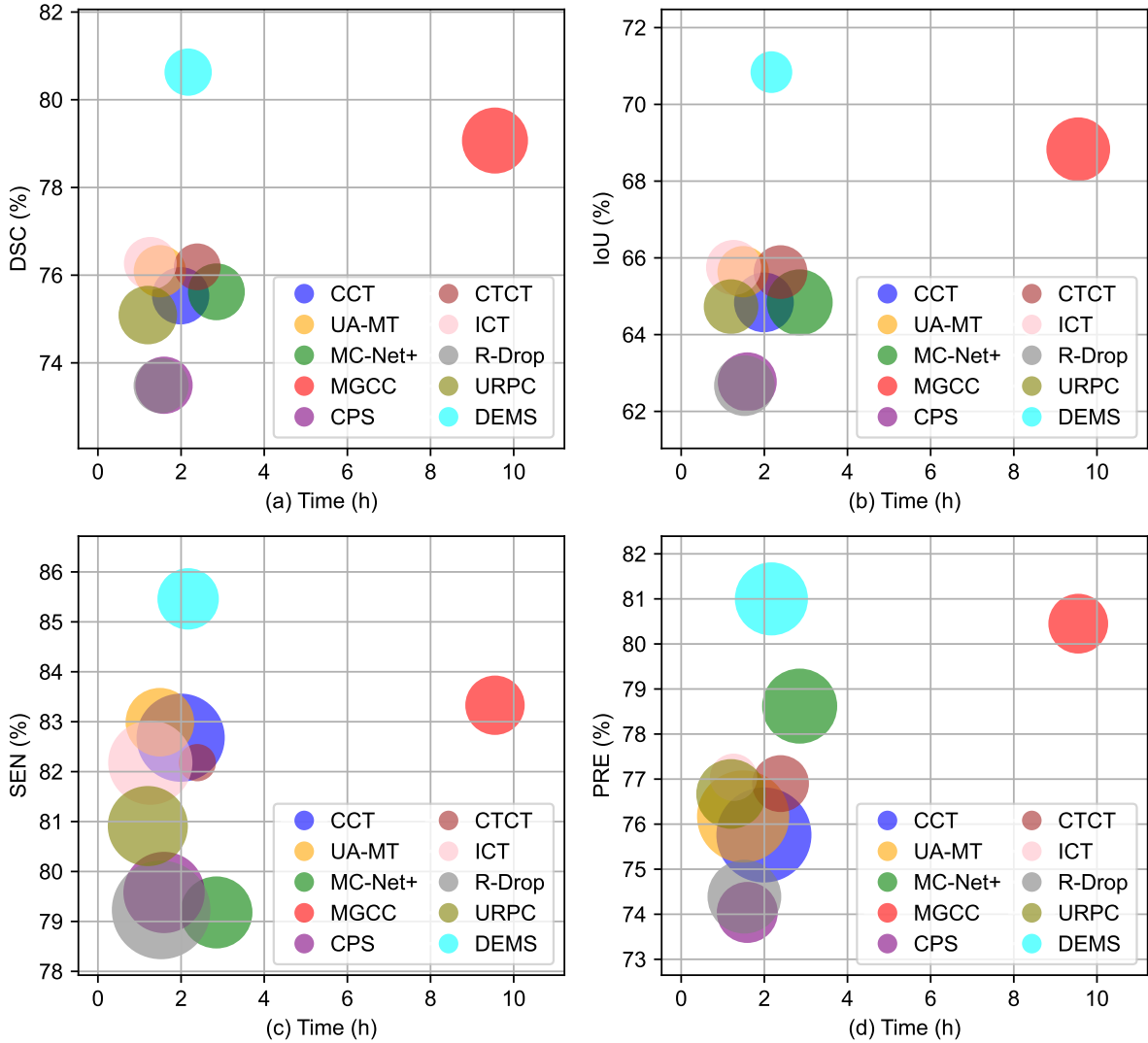


Figure 10: Relationship between time consumption and performance metrics across DEMS and SOTA methods on the DDTI dataset using 80% labeled images. We exclude PA from the visualization for clarity, as its values are consistently similar across methods and can lead to significant overlap of circles. The circle size indicates the standard deviation.

OAA leads to a drastic deterioration in model performance, thereby rendering comparative analysis futile. Detailed observations reveal that the absence of each component results in a significant performance decrement. Specifically, removing OAA, RRE block, and L_{sen} decreases the DSC by 4.90%, 3.38%, and 1.67%, respectively. Additionally, excluding component combinations accelerates performance degradation. When OAA, RRE block, and L_{sen} are excluded, a DSC of 56.65% is observed with a decrease of 11.85%. Although the variant with solely OAA achieves the highest PRE of 72.32%, it merely leads by a marginal 0.02%. However, the DSC, IoU, and SEN significantly fall short by margins of 5.00%, 4.64%, and 6.19%, respectively.

6. Conclusion

6.1. Summary and Discussion

In this manuscript, we introduce a novel semi-supervised segmentation method DEMS to segment medical images with limited data. We devise the OAA to diversify the input data, thereby enhancing the generalization ability. We propose the RRE block to enrich feature diversity and introduce perturbations to produce varying inputs for different decoders, therefore offering greater variability. Furthermore, we propose a novel sensitivity loss to further enhance the consistency across decoders and bolster training stability. Extensive experimental results on both our own and three public datasets demonstrate the superiority of DEMS over SOTA methods. Additionally, DEMS performs exceptionally desirable on severe data shortages, showcasing its remarkable data efficiency and significant advancements in medical segmentation. Despite its impressive performance, adapting DEMS for multi-object segmentation may pose challenges, as its sensitivity loss is formulated specifically for binary segmentation. A potential solution to this limitation involves treating multi-object segmentation as an aggregation of multiple binary segmentation tasks. It should be noted that this kind of solution can come with several challenges such as increased inference times and the complexity of handling overlapping or contiguous objects.

6.2. Future Perspectives

The future perspectives of the DEMS are twofold. Firstly, an enhanced connection structure between the encoder and decoders can be developed. Recently, transformer-based architecture has been proven to be a powerful tool for computer vision tasks. Compared with CNN-based architecture, it prioritizes the assimilation of global features rather than local features. To this end, incorporating transformer and CNN-based architectures is capable of providing richer visual features and therefore further strengthening the model performance. However, integrating the transformer-based architecture may significantly raise computing costs, potentially limiting the application in most of the mobile scenes. One approach to offset the extra computing lies in reducing the number of decoders. Secondly, a novel updating strategy can be formulated to update the coefficients across various loss function terms. In

contrast to fixing the coefficients throughout the training process, adaptively modulating them as training progresses can more effectively maintain the magnitude of each loss term at comparable levels. This can secure the contribution of varying loss terms during the training and is also anticipated to augment the generalization capacity of the model. In addition to the term coefficients, the binarization threshold leveraged to compute the sensitivity loss could also benefit from adaptive adjustments during the training phase.

Acknowledgements

This work is supported by Tan Tock Seng Hospital (A-8001334-00-00).

References

- [1] Gang Hu, Yixuan Zheng, Laith Abualigah, and Abdelazim G Hussien. Detdo: An adaptive hybrid dandelion optimizer for engineering optimization. *Advanced Engineering Informatics*, 57:102004, 2023.
- [2] Mohsen Zare, Mojtaba Ghasemi, Amir Zahedi, Keyvan Golalipour, Soleiman Kadkhoda Mohammadi, Seyedali Mirjalili, and Laith Abualigah. A global best-guided firefly algorithm for engineering problems. *Journal of Bionic Engineering*, pages 1–30, 2023.
- [3] Laith Abualigah, Serdar Ekinci, Davut Izci, and R Abu Zitar. Modified elite opposition-based artificial hummingbird algorithm for designing fopid controlled cruise control system. *Intelligent Automation & Soft Computing*, 2023.
- [4] Jingying Chen, Jinxin Shi, and Ruyi Xu. Dual subspace manifold learning based on gcn for intensity-invariant facial expression recognition. *Pattern Recognition*, 148:110157, 2024.
- [5] Sai Harsha Yelleni, Deepshikha Kumari, PK Srijith, et al. Monte carlo dropout for modeling uncertainty in object detection. *Pattern Recognition*, 146:110003, 2024.
- [6] Zhaoshan Liu, Qiujie Lv, Ziduo Yang, Yifan Li, Chau Hung Lee, and Lei Shen. Recent progress in transformer-based medical image analysis. *Computers in Biology and Medicine*, page 107268, 2023.
- [7] Saidi Guo, Heye Zhang, Yifeng Gao, Hui Wang, Lei Xu, Zhifan Gao, Antonella Guzzo, and Giancarlo Fortino. Survival prediction of heart failure patients using motion-based analysis method. *Computer Methods and Programs in Biomedicine*, 236:107547, 2023.

- [8] Tianfei Zhou, Liulei Li, Gustav Bredell, Jianwu Li, Jan Unkelbach, and Ender Konukoglu. Volumetric memory network for interactive medical image segmentation. *Medical Image Analysis*, 83:102599, 2023.
- [9] Xiwang Xie, Xipeng Pan, Weidong Zhang, and Jubai An. A context hierarchical integrated network for medical image segmentation. *Computers and Electrical Engineering*, 101:108029, 2022.
- [10] Saidi Guo, Xiujian Liu, Heye Zhang, Qixin Lin, Lei Xu, Changzheng Shi, Zhifan Gao, Antonella Guzzo, and Giancarlo Fortino. Causal knowledge fusion for 3d cross-modality cardiac image segmentation. *Information Fusion*, page 101864, 2023.
- [11] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6):84–90, 2017. <https://doi.org/10.48550/10.1145/3065386>.
- [12] Zhaoshan Liu, Qiujie Lv, Chau Hung Lee, and Lei Shen. Gsda: Generative adversarial network-based semi-supervised data augmentation for ultrasound image classification. *Heliyon*, 9(9), 2023.
- [13] Blake Murdoch. Privacy and artificial intelligence: challenges for protecting health information in a new era. *BMC Medical Ethics*, 22(1):1–5, 2021.
- [14] Matthias Eisenmann, Annika Reinke, Vivienn Weru, Minu Dietlinde Tizabi, Fabian Isensee, Tim J Adler, Sharib Ali, Vincent Andrearczyk, Marc Auberville, Ujjwal Baid, et al. Why is the winner the best? *arXiv preprint*, 2023. <https://doi.org/10.48550/arXiv.2303.17719>.
- [15] Xiaokang Li, Menghua Xia, Jing Jiao, Shichong Zhou, Cai Chang, Yuanyuan Wang, and Yi Guo. Hal-ia: A hybrid active learning framework using interactive annotation for medical image segmentation. *Medical Image Analysis*, page 102862, 2023.
- [16] Libo Zhao, Xiaolong Qian, Yinghui Guo, Jiaqi Song, Jinbao Hou, and Jun Gong. Mskd: Structured knowledge distillation for efficient medical image segmentation. *Computers in Biology and Medicine*, page 107284, 2023.
- [17] Fabian Isensee, Paul F Jaeger, Simon AA 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, 2021.
- [18] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014.

- [19] Lu Chai, Zidong Wang, Jianqing Chen, Guokai Zhang, Fawaz E Alsaadi, Fuad E Alsaadi, and Qinyuan Liu. Synthetic augmentation for semantic segmentation of class imbalanced biomedical images: A data pair generative adversarial network approach. *Computers in Biology and Medicine*, 150:105985, 2022.
- [20] Daiqing Li, Junlin Yang, Karsten Kreis, Antonio Torralba, and Sanja Fidler. Semantic segmentation with generative models: Semi-supervised learning and strong out-of-domain generalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8300–8311, 2021.
- [21] Jason Kugelman, David Alonso-Caneiro, Scott A Read, Stephen J Vincent, Fred K Chen, and Michael J Collins. Dual image and mask synthesis with gans for semantic segmentation in optical coherence tomography. In *2020 Digital Image Computing: Techniques and Applications (DICTA)*, pages 1–8. IEEE, 2020.
- [22] Qianli Feng, Chenqi Guo, Fabian Benitez-Quiroz, and Aleix M Martinez. When do gans replicate? on the choice of dataset size. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6701–6710, 2021.
- [23] Wen-Yang Zhou, Guo-Wei Yang, and Shi-Min Hu. Jittor-gan: A fast-training generative adversarial network model zoo based on jittor. *Computational Visual Media*, 7:153–157, 2021.
- [24] Ekin D Cubuk, Barret Zoph, Dandelion Mane, Vijay Vasudevan, and Quoc V Le. Autoaugment: Learning augmentation strategies from data. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 113–123, 2019.
- [25] Dong Yang, Holger Roth, Ziyue Xu, Fausto Milletari, Ling Zhang, and Daguang Xu. Searching learning strategy with reinforcement learning for 3d medical image segmentation. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2019: 22nd International Conference, Shenzhen, China, October 13–17, 2019, Proceedings, Part II 22*, pages 3–11. Springer, 2019.
- [26] Junyan Lyu, Yiqi Zhang, Yijin Huang, Li Lin, Pujin Cheng, and Xiaoying Tang. Aadg: automatic augmentation for domain generalization on retinal image segmentation. *IEEE Transactions on Medical Imaging*, 41(12):3699–3711, 2022.
- [27] Tiexin Qin, Ziyuan Wang, Kelei He, Yinghuan Shi, Yang Gao, and Dinggang Shen. Automatic data augmentation via deep reinforcement learning for effective kidney tumor segmentation. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1419–1423. IEEE, 2020.

- [28] Sangdoon Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6023–6032, 2019.
- [29] Rushi Jiao, Yichi Zhang, Le Ding, Rong Cai, and Jicong Zhang. Learning with limited annotations: a survey on deep semi-supervised learning for medical image segmentation. *arXiv preprint arXiv:2207.14191*, 2022.
- [30] Yassine Ouali, Céline Hudelot, and Myriam Tami. Semi-supervised semantic segmentation with cross-consistency training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12674–12684, 2020.
- [31] Xiaoqiang Li, Yuanchen Wu, and Songmin Dai. Semi-supervised medical imaging segmentation with soft pseudo-label fusion. *Applied Intelligence*, pages 1–13, 2023.
- [32] Han Zheng, Lanfen Lin, Hongjie Hu, Qiaowei Zhang, Qingqing Chen, Yutaro Iwamoto, Xianhua Han, Yen-Wei Chen, Ruofeng Tong, and Jian Wu. Semi-supervised segmentation of liver using adversarial learning with deep atlas prior. In *Medical Image Computing and Computer Assisted Intervention—MICCAI 2019: 22nd International Conference, Shenzhen, China, October 13–17, 2019, Proceedings, Part VI 22*, pages 148–156. Springer, 2019.
- [33] Antti Tarvainen and Harri Valpola. Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results. *Advances in neural information processing systems*, 30, 2017.
- [34] Yichi Zhang, Rushi Jiao, Qingcheng Liao, Dongyang Li, and Jicong Zhang. Uncertainty-guided mutual consistency learning for semi-supervised medical image segmentation. *Artificial Intelligence in Medicine*, 138:102476, 2023.
- [35] Tao Lei, Dong Zhang, Xiaogang Du, Xuan Wang, Yong Wan, and Asoke K Nandi. Semi-supervised medical image segmentation using adversarial consistency learning and dynamic convolution network. *IEEE Transactions on Medical Imaging*, 2022.
- [36] Xuanang Xu, Thomas Sanford, Baris Turkbey, Sheng Xu, Bradford J Wood, and Pingkun Yan. Shadow-consistent semi-supervised learning for prostate ultrasound segmentation. *IEEE Transactions on Medical Imaging*, 41(6):1331–1345, 2021.
- [37] Yongchao Wang, Bin Xiao, Xiuli Bi, Weisheng Li, and Xinbo Gao. Mcf: Mutual correction framework for semi-supervised medical image segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15651–15660, 2023.

- [38] Fenghe Tang, Jianrui Ding, Lingtao Wang, Min Xian, and Chunping Ning. Multi-level global context cross consistency model for semi-supervised ultrasound image segmentation with diffusion model. *arXiv preprint arXiv:2305.09447*, 2023.
- [39] Yicheng Wu, Minfeng Xu, Zongyuan Ge, Jianfei Cai, and Lei Zhang. Semi-supervised left atrium segmentation with mutual consistency training. In *Medical Image Computing and Computer Assisted Intervention—MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part II 24*, pages 297–306. Springer, 2021.
- [40] Yicheng Wu, Zongyuan Ge, Donghao Zhang, Minfeng Xu, Lei Zhang, Yong Xia, and Jianfei Cai. Mutual consistency learning for semi-supervised medical image segmentation. *Medical Image Analysis*, 81:102530, 2022.
- [41] Donghai Zhai, Bijie Hu, Xun Gong, Haipeng Zou, and Jun Luo. Ass-gan: Asymmetric semi-supervised gan for breast ultrasound image segmentation. *Neuro-computing*, 493:204–216, 2022.
- [42] Zihang Xu, Zhenghua Xu, Shuo Zhang, and Thomas Lukasiewicz. Pca: Semi-supervised segmentation with patch confidence adversarial training. *arXiv preprint arXiv:2207.11683*, 2022.
- [43] Xin Yu, Qi Yang, Yinchu Zhou, Leon Y Cai, Riqiang Gao, Ho Hin Lee, Thomas Li, Shunxing Bao, Zhoubing Xu, Thomas A Lasko, et al. Unest: local spatial representation learning with hierarchical transformer for efficient medical segmentation. *Medical Image Analysis*, page 102939, 2023.
- [44] Penghui Li, Rui Zhou, Jin He, Shifeng Zhao, and Yun Tian. A global-frequency-domain network for medical image segmentation. *Computers in Biology and Medicine*, page 107290, 2023.
- [45] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019.
- [46] Siddharth Pandey, Pranshu Ranjan Singh, and Jing Tian. An image augmentation approach using two-stage generative adversarial network for nuclei image segmentation. *Biomedical Signal Processing and Control*, 57:101782, 2020.
- [47] Ahmed Iqbal and Muhammad Sharif. Unet: A semi-supervised method for segmentation of breast tumor images using a u-shaped pyramid-dilated network. *Expert Systems with Applications*, 221:119718, 2023.
- [48] Ziyu Wang, Tom Schaul, Matteo Hessel, Hado Hasselt, Marc Lanctot, and Nando Freitas. Dueling network architectures for deep reinforcement learning. In *International conference on machine learning*, pages 1995–2003. PMLR, 2016.

- [49] Ju Xu, Mengzhang Li, and Zhanxing Zhu. Automatic data augmentation for 3d medical image segmentation. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2020: 23rd International Conference, Lima, Peru, October 4–8, 2020, Proceedings, Part I 23*, pages 378–387. Springer, 2020.
- [50] Zhaoshan Liu, Qiuji Lv, Yifan Li, Ziduo Yang, and Lei Shen. Medaugument: Universal automatic data augmentation plug-in for medical image analysis. *arXiv preprint arXiv:2306.17466*, 2023.
- [51] Amy Zhao, Guha Balakrishnan, Fredo Durand, John V Guttag, and Adrian V Dalca. Data augmentation using learned transformations for one-shot medical image segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8543–8553, 2019.
- [52] Zach Eaton-Rosen, Felix Bragman, Sebastien Ourselin, and M Jorge Cardoso. Improving data augmentation for medical image segmentation. In *Proceedings of the Conference on Medical Imaging with Deep Learning*, 2018.
- [53] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. In *International Conference on Learning Representations*, 2018.
- [54] Jun Lyu, Bin Sui, Chengyan Wang, Qi Dou, and Jing Qin. Adaptive feature aggregation based multi-task learning for uncertainty-guided semi-supervised medical image segmentation. *Expert Systems with Applications*, page 120836, 2023.
- [55] Lequan Yu, Shujun Wang, Xiaomeng Li, Chi-Wing Fu, and Pheng-Ann Heng. Uncertainty-aware self-ensembling model for semi-supervised 3d left atrium segmentation. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2019: 22nd International Conference, Shenzhen, China, October 13–17, 2019, Proceedings, Part II 22*, pages 605–613. Springer, 2019.
- [56] Xiangde Luo, Minhao Hu, Tao Song, Guotai Wang, and Shaoting Zhang. Semi-supervised medical image segmentation via cross teaching between cnn and transformer. In *International Conference on Medical Imaging with Deep Learning*, pages 820–833. PMLR, 2022.
- [57] Hu Cao, Yueyue Wang, Joy Chen, Dongsheng Jiang, Xiaopeng Zhang, Qi Tian, and Manning Wang. Swin-unet: Unet-like pure transformer for medical image segmentation. In *European conference on computer vision*, pages 205–218. Springer, 2022.
- [58] 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: 18th International Conference*,

Munich, Germany, October 5-9, 2015, *Proceedings, Part III 18*, pages 234–241. Springer, 2015.

- [59] Alexander Buslaev, Vladimir I Iglovikov, Eugene Khvedchenya, Alex Parinov, Mikhail Druzhinin, and Alexandr A Kalinin. Alumentations: fast and flexible image augmentations. *Information*, 11(2):125, 2020. <https://doi.org/10.3390/info11020125>.
- [60] 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.
- [61] Jonathan Tompson, Ross Goroshin, Arjun Jain, Yann LeCun, and Christoph Bregler. Efficient object localization using convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 648–656, 2015.
- [62] Huisi Wu, Jiasheng Liu, Fangyan Xiao, Zhenkun Wen, Lan Cheng, and Jing Qin. Semi-supervised segmentation of echocardiography videos via noise-resilient spatiotemporal semantic calibration and fusion. *Medical Image Analysis*, 78:102397, 2022.
- [63] Moi Hoon Yap, Manu Goyal, Fatima Osman, Robert Martí, Erika Denton, Arne Juette, and Reyer Zwiggelaar. Breast ultrasound region of interest detection and lesion localisation. *Artificial Intelligence in Medicine*, 107:101880, 2020.
- [64] Moi Hoon Yap, Gerard Pons, Joan Marti, Sergi Ganau, Melcior Sentis, Reyer Zwiggelaar, Adrian K Davison, and Robert Marti. Automated breast ultrasound lesions detection using convolutional neural networks. *IEEE journal of biomedical and health informatics*, 22(4):1218–1226, 2017.
- [65] Walid Al-Dhabyani, Mohammed Gomaa, Hussien Khaled, and Aly Fahmy. Dataset of breast ultrasound images. *Data in brief*, 28:104863, 2020. <https://doi.org/10.1016/j.dib.2019.104863>.
- [66] 1st place solution in miccai 2020 tn-scui challenge. <https://github.com/WAMAWAMA/TNSCUI2020-Seg-Rank1st>. Accessed 5 Sep 2023.
- [67] Lina Pedraza, Carlos Vargas, Fabián Narváez, Oscar Durán, Emma Muñoz, and Eduardo Romero. An open access thyroid ultrasound image database. In *10th International symposium on medical information processing and analysis*, volume 9287, pages 188–193. SPIE, 2015.
- [68] Bryan C Russell, Antonio Torralba, Kevin P Murphy, and William T Freeman. Labelme: a database and web-based tool for image annotation. *International journal of computer vision*, 77:157–173, 2008.

- [69] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016.
- [70] Xiaokang Chen, Yuhui Yuan, Gang Zeng, and Jingdong Wang. Semi-supervised semantic segmentation with cross pseudo supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2613–2622, 2021.
- [71] Vikas Verma, Kenji Kawaguchi, Alex Lamb, Juho Kannala, Arno Solin, Yoshua Bengio, and David Lopez-Paz. Interpolation consistency training for semi-supervised learning. *Neural Networks*, 145:90–106, 2022.
- [72] Lijun Wu, Juntao Li, Yue Wang, Qi Meng, Tao Qin, Wei Chen, Min Zhang, Tie-Yan Liu, et al. R-drop: Regularized dropout for neural networks. *Advances in Neural Information Processing Systems*, 34:10890–10905, 2021.
- [73] Xiangde Luo, Wenjun Liao, Jieneng Chen, Tao Song, Yinan Chen, Shichuan Zhang, Nianyong Chen, Guotai Wang, and Shaoting Zhang. 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: 24th International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part II 24*, pages 318–329. Springer, 2021.