

CDDSA: Contrastive Domain Disentanglement and Style Augmentation for Generalizable Medical Image Segmentation

Ran Gu^a, Guotai Wang^{a,i,*}, Jiangshan Lu^a, Jingyang Zhang^{b,c}, Wenhui Lei^{d,i}, Yanan Chen^{e,g}, Wenjun Liao^f, Shichuan Zhang^f, Kang Li^g, Dimitris N. Metaxas^h, Shaoting Zhang^{a,e,i,*}

^a*School of Mechanical and Electrical Engineering, University of Electronic Science and Technology of China, Chengdu, China*

^b*School of Biomedical Engineering, Shanghai Jiao Tong University, Shanghai, China*

^c*School of Biomedical Engineering, ShanghaiTech University, Shanghai, China.*

^d*School of Electronic Information and Electrical Engineering, Shanghai Jiao Tong University, Shanghai, China*

^e*Sensetime Research, Shanghai, China*

^f*Department of Radiation Oncology, Sichuan Cancer Hospital and Institute, University of Electronic Science and Technology of China, Chengdu, China*

^g*West China Hospital-SenseTime Joint Lab, West China Biomedical Big Data Center, Sichuan University, Chengdu, China*

^h*Department of Computer Science, Rutgers University, Piscataway NJ 08854, USA*

ⁱ*Shanghai AI Lab, Shanghai, China*

Abstract

Generalization to previously unseen images with potential domain shifts and different styles is essential for clinically applicable medical image segmentation, and the ability to disentangle domain-specific and domain-invariant features is key for achieving Domain Generalization (DG). However, existing DG methods can hardly achieve effective disentanglement to get high generalizability. To deal with this problem, we propose an efficient Contrastive Domain Disentanglement and Style Augmentation (CDDSA) framework for generalizable medical image segmentation. First, a disentangle network is proposed to decompose an image into a domain-invariant anatomical representation and a domain-specific style code, where the former is sent to a segmentation model that is not affected by the domain shift, and the disentangle network is regularized by a decoder that combines the anatomical and style codes to reconstruct the input image. Second, to achieve better disentanglement, a contrastive loss is proposed to encourage the style codes from the same domain and different domains to be compact and divergent, respectively. Thirdly, to further improve generalizability, we propose a style augmentation method based on the disentanglement representation to synthesize images in various unseen styles with shared anatomical structures. Our method was validated on a public multi-site fundus image dataset for optic cup and disc segmentation and an in-house multi-site Nasopharyngeal Carcinoma Magnetic Resonance Image (NPC-MRI) dataset for nasopharynx Gross Tumor Volume (GTV_{nx}) segmentation. Experimental results showed that the proposed CDDSA achieved remarkable generalizability across different domains, and it outperformed several state-of-the-art methods in domain-generalizable segmentation. Code is available at <https://github.com/HiLab-git/DAG4MIA>

Keywords:

Disentanglement, Domain Generalization, Contrastive Learning, Medical Image Segmentation

1. Introduction

Deep learning with Convolutional Neural Networks (CNNs) has achieved remarkable performance in medical image segmentation [1, 2, 3], and most existing models are built on the assumption that training and testing images are from the same domain and have very similar, if not the same, distributions. However, in clinical practice, this assumption does often not hold due to several factors such as the differences in scanning devices, imaging protocols, patient groups and image quality between training and testing images, where the testing images are usually acquired from a different medical center than the training set. Such differences (a.k.a domain shift [4]) can substantially degrade the model's performance at test time [5, 6].

To address this problem, many Domain Adaptation (DA) methods have been explored to transfer knowledge from a set of labeled images in a source domain to images in a target domain [7, 8, 9]. However, the DA methods need to tune the model's parameters based on a set of images in the target domain, which is not only time-consuming but also impractical if the target domain is not known in advance [10]. What's more, the model needs to be adapted to each target domain respectively, and is faced with the problem of catastrophic forgetting on previous domains, which is not scalable when applied to a range of new unknown domains.

In contrast to DA, Domain Generalization (DG) that encourages a model to be generalizable to unseen domains is more appealing and efficient as it does not need to tune the model after training. Recently, domain generalization has attracted increasing attentions in the field of both computer vision [11] and medical image analysis [12, 13]. Existing DG methods mainly

*Corresponding author.

Email addresses: Guotai.Wang@uestc.edu.cn (Guotai Wang), Zhangshaoting@uestc.edu.cn (Shaoting Zhang)

include image- and feature-based approaches. For image-based approaches, data augmentation has been widely used for improving the generalizability of a model, and [14] proposed BigAug that uses a series of stacked transformations to augment the training images, with the assumption that the shift between source and target domains could be simulated through extensive data augmentation. However, the configuration of data augmentation requires empirical settings and could be data-specific. In contrast, feature-based methods mainly focus on representation learning to extract most representative features for better generalization across domains [15, 13]. [15] introduced a domain-oriented feature embedding method that dynamically enriches image features with domain prior knowledge learned from multi-site domains to make the semantic features more discriminative. [13] developed a Domain Composition and Attention Network (DCA-Net) that represents features in a certain domain as a linear combination of a set of basis representations in a representation bank, where the combination coefficients are obtained by an attention module. Both methods rely on a domain knowledge pool or a representation bank to infer the domain-specific knowledge to make the network aware of the domain of an input image, which helps to improve the generalizability. However, they are limited by the capacity and representation power of the knowledge pool/bank, and have a limited ability to recognize the invariant features across different domains.

Recently, disentanglement has been introduced to computer vision that aims to explicitly decompose features into domain-invariant contents and domain-specific styles [16]. It has also been employed to learn domain-invariant features for domain adaptation on multi-modality medical image segmentation datasets. [17] applied disentangled representations to unsupervised domain adaptation for liver segmentation. They decomposed the images from two domains into a shared domain-invariant content space and a domain-specific style space, and used representations in the content space for segmentation. [18] used disentangled domain-invariant and domain-specific features for cardiac image segmentation across two modalities, and introduced a zero-loss to enhance the disentanglement. However, most existing disentangling methods are based on Generative Adversarial Networks (GAN), where a content encoder and a style encoder need to be trained for each known modality/domain, and multiple discriminators are involved, leading to a complex training process. Despite their suitability for domain adaptation, the GAN-based disentanglement methods are not scalable, as the number of required content/style encoders and discriminators will increase with the grow of domain number. What's more, such a paradigm cannot be applied to unseen domains as it requires the encoders for each domain to be trained in advance. Therefore, they are not applicable to DG problems.

In this work, we propose a novel GAN-free disentanglement framework named as Contrastive Domain Disentanglement and Style Augmentation (CDDSA) for domain-generalizable medical image segmentation. As shown in Fig. 1, it decomposes medical images in different domains into domain-invariant anatomical representations and domain-specific style codes with only one pair of anatomy Encoder and style encoder,

which is regularized by a decoder that accepts an anatomical representation and a style code to reconstruct an image. The encoders and decoder are shared across different domains, without adversarial learning and domain-specific training, which is efficient and scalable to multiple domains. Our method was inspired by Spatial Decomposition Network (SDNet) [19] that implements feature disentanglement without GAN. Note that SDNet [19] was proposed for semi-supervised learning, modality transformation and multi-modal image segmentation, and it can only perform disentanglement and image reconstruction on seen domains with poor generalizability in unseen domains. The main reason is that SDNet lacks effective constraints on the style codes to encourage them to be domain-specific, which limits the ability to extract domain-invariant feature representations. In addition, it restricts the anatomical representations as binary codes, leading to a limited representation ability for effective image reconstruction.

Differently from SDNet [19], our CDDSA is proposed for domain-generalizable segmentation of medical images. To improve the disentanglement performance, we relax the anatomical representation to soft values and propose a domain style contrastive learning loss to encourage the style codes in different domains to be discriminative from each other, which improves the model's ability to recognize domain-invariant anatomical representations that is sent to a segmentor to obtain segmentation results. As the segmentor is not affected by domain-specific features, it has a high generalizability across different domains. In addition, based on the extracted style codes in training domains, we can generate a new random style code and combine it with an existing anatomical representation to simulate images in an unseen domain with a new styles using the decoder, i.e., style augmentation, which further improves the generalizability of our framework.

To the best of our knowledge, this is the first work in the literature to propose feature disentanglement learning for domain-generalizable medical image segmentation. The contributions of our method are summarised as follows:

- 1) We introduce a novel framework CDDSA using GAN-free disentanglement for domain generalization in medical image segmentation. It achieves generalizability by segmentation from decomposed domain-invariant representations extracted by a single anatomy Encoder that is shared across domains and more efficient and scalable than GAN-based disentanglement.
- 2) To make the disentangled domain-specific style codes more representative and distinguishable, we propose domain style contrastive learning, which forces the style codes from the same domain and different domains to be similar and dissimilar, respectively.
- 3) We propose style augmentation based on the disentangled anatomical representations and style codes to simulate images from unseen domains with different styles, which further improves the generalizability of the disentanglement and segmentation models.

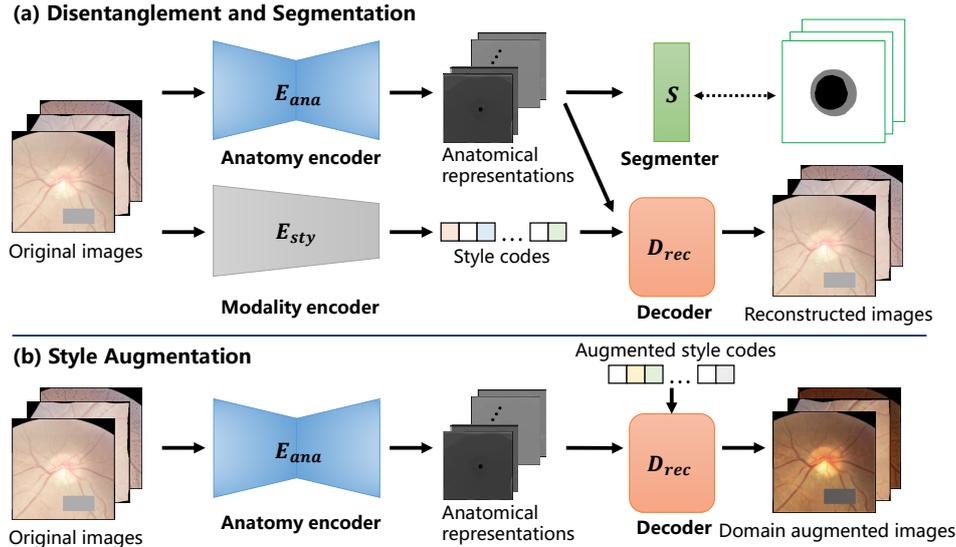


Figure 1: Workflow of our proposed Contrastive Domain Disentanglement and Style Augmentation (CDDSA) method. (a) shows the disentanglement and segmentation networks, where an anatomy encoder and a style encoder obtains anatomical representations and style codes respectively, and they are regularized by a decoder to reconstruct the input image. The segmentor takes domain-invariant anatomical representations as input to obtain the segmentation results. (b) represents the style augmentation strategy, where we combine anatomical representations from a given image with augmented style codes to generate images in a new domain.

- 4) Comprehensive experimental results on multi-domain fundus images and multi-domain nasopharyngeal carcinoma magnetic resonance images (NPC-MRI) showed that our proposed CDDSA achieved high generalization on unseen domains, and it outperformed several state-of-the-art domain generalization methods.

2. Related Works

2.1. Domain Generalization for Medical Image Analysis

Recently, domain generalization has attracted increasing attentions to avoid dramatic performance degradation when inferring with images from unseen domains [20]. It aims to learn a model from a single or multiple source domains to make it directly applicable for unseen target domains without extra training [21, 12, 11]. Existing DG methods mainly include meta-learning methods, data-based methods and feature-based methods. Meta-learning [22, 12] splits a set of source domains into meta-train and meta-test subsets, and adopts meta-optimization that iteratively updates model parameters to improve performance on the meta-test subset to simulate the situation when inferring on unseen domains. [23] combined meta-learning with federated learning to achieve privacy-preserving generalizable segmentation through continuous frequency space interpolation across clients. However, meta-optimization process is highly time-consuming since all potential splitting results of meta-train and meta-test should be considered during training [12].

Data-based approaches usually use different data augmentation strategies for improving the model’s generalizability. [14] a deep stacked transformation assuming that the shift between different domains can be simulated by extensive data augmentation on a single domain. [24] utilized Cycle-GAN [25] to transform images from one certain domain to other domains for

augmentation. [26] proposed Mixed Task Sampling (MTS) to enhance the variety of task-level training samples. Mixup in frequency domains [23, 27] has also been used to synthesize new images for model generalization. However, the efficiency of data augmentation largely depends on the ability to cover the data distribution in unseen domains, hence requiring empirical settings and even data-specific modifications.

Feature-based approaches use domain-adaptive feature calibration or learn domain-invariant features to deal with domain generalization [15, 28, 29]. [15] introduced a domain-oriented feature embedding framework that dynamically updates the domain-specific prior knowledge to make the semantic features more discriminative. [30] proposed a dynamic convolutional head to make the model’s convolutional parameters adaptive to unseen target domains. [13] proposed a domain composition and attention method that calibrates the input feature based on attention coefficients represented by a representation bank. However, these methods did not explicitly obtain domain-invariant features for domain generalization, and they did not separate features into purely domain-specific and domain-invariant representations well, leading to limited performance on domain generalization.

2.2. Disentanglement Representation Learning

Disentanglement explicitly decomposes features into domain-invariant contents and domain-specific styles [31, 32]. In addition to applications such as image synthesis [19], artifact removal and multi-task learning [33] for medical image analysis, it is widely adopted for domain adaptation [17, 18]. [17] used disentanglement to obtain domain-invariant content features for liver segmentation with domain adaptation. [34] used disentanglement to improve the performance of image translation for domain adaptation, and they disentangled the content features from domain information for both the source

and translated images. [18] applied disentanglement-based domain adaptation for cardiac image segmentation, and introduced a zero loss to enhance disentanglement. [35] proposed a bidirectional unsupervised DA framework based on disentangled representation learning for equally competent two-way DA performances on cardiac image segmentation. Despite their good performance on DA, these works achieve disentanglement based on GAN, where multiple discriminators are needed in the adversarial training process that is complex and tricky to optimize. What’s more, they need to have access to images for target domains during training, and are not applicable to DG tasks that involves unseen domains. [19] proposed a GAN-free Spatial Decomposition Network (SDNet) that decomposes an input image into a spatial factor (anatomy) and a non-spatial factor (style), and applied it to semi-supervised segmentation and image synthesis. However, it performs disentanglement and reconstruction well only on seen domains can hardly deal with unseen domains that are not involved in training.

2.3. Contrastive Learning

Contrastive learning is a self-supervised learning method to learn feature representations by enforcing positive pairs to be close and negative pairs to be distant [36]. Previous contrastive learning methods were mainly proposed to pre-train a powerful and representational feature extractor that can distinguish similar and dissimilar samples [37, 38]. For computer vision and medical image analysis, contrastive learning has been mainly used for annotation-efficient learning. For example, [39] proposed a contrastive adaptation network that minimizes the intra-class domain discrepancy and maximizes the inter-class domain discrepancy. [40] used contrastive learning of global and local features sequentially for 3D medical image segmentation with limited annotations. [41] proposed contrastive learning of relative position regression for one-shot object localization in 3D medical images. [42] proposed a contrastive voxel-wise representation learning to effectively learn low-level and high-level features for semi-supervised medical image segmentation. Unlike these works, we design a contrastive learning strategy to enhance disentanglement between domain-invariant and domain-specific features to deal with domain generalization problems.

3. Methods

For the domain generalization problem, the training set consists of images from D domains and can be denoted as $\mathcal{D} = \{(x_i^d, y_i^d)\}_{i=1}^{N_d}$ ($d = 1, 2, \dots, D$), where x_i^d depicts the i -th training sample from the d -th source domain with its corresponding ground-truth annotation y_i^d . N_d denotes the number of training samples in domain d .

Our proposed Contrastive Domain Disentanglement and Style Augmentation (CDDSA) framework is illustrated in Fig. 2. Firstly, we employ a disentangle network containing an anatomy encoder E_{ana} and a style encoder E_{sty} to decompose an image into a domain-invariant anatomical representation and a domain-specific modality representation (i.e., style

code), and they can be used to reconstruct the input image based on a decoder D_{rec} . We further send the disentangled anatomical representation into a segmenter S to predict the segmentation mask. Secondly, to boost the disentanglement performance with more discriminative style codes across different domains, we introduce domain style contrastive learning that forces the decomposed modality representations to have low intra-domain discrepancy and high inter-domain discrepancy. Thirdly, to further enhance model generalization, we proposed a style augmentation strategy to randomly generate style codes and combine them with given anatomical representations to reconstruct images with new styles that are not present in the training set.

3.1. Domain Disentangle Network.

As shown in Fig. 2, for an input image x_i^d , we send it to an anatomy encoder E_{ana} and a style encoder E_{sty} to obtain an anatomical representation $f_{i,a}^d$ and a modality representation (style code) $f_{i,s}^d$, respectively. Then $f_{i,a}^d$ and $f_{i,s}^d$ are sent to a decoder D_{rec} to reconstruct an input-like images \hat{x}_i^d , and a reconstruction loss \mathcal{L}_{rec} is used to encourage the consistency between x_i^d and \hat{x}_i^d . A segmentor S takes $f_{i,a}^d$ as input to obtain the segmentation result.

3.1.1. Anatomy Encoder and Segmenter

To decompose domain-invariant anatomical representations, we employ U-Net [1] as the backbone to implement E_{ana} . We modify U-Net by setting the output channel of the last layer as T and use tanh as the activation function in that layer. Let H and W represent the height and width of the input image x_i^d respectively, the output of E_{ana} is denoted as $f_{i,a}^d \in [-1, 1]^{H \times W \times T}$, and we assume that each channel of $f_{i,a}^d$ emphasizes some anatomical information. Differently from SDNet [19] that constrains $f_{i,a}^d$ to take binary values that may lose many details of object structures, we aim to reserve enough structural information for accurate image reconstruction and further style augmentation, and therefore soften the anatomical representation with a $\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$ activation function in the last layer. The anatomical representation extraction procedure is formulated as:

$$f_{i,a}^d = E_{ana}(x_i^d) \quad (1)$$

Then, the decomposed anatomical representation $f_{i,a}^d$ is fed into a segmentation network S to obtain a segmentation probability map $p_i^d = S(f_{i,a}^d)$. Let y_i^d denote the ground truth, and the supervised segmentation loss for domain d is:

$$\mathcal{L}_{seg} = \frac{1}{2N_d} \sum_{i=1}^{N_d} (\mathcal{L}_{Dice}(p_i^d, y_i^d) + \mathcal{L}_{ce}(p_i^d, y_i^d)) \quad (2)$$

where we use a hybrid segmentation loss that consists of a Dice loss \mathcal{L}_{Dice} and a cross-entropy loss \mathcal{L}_{ce} .

3.1.2. Style Encoder

The domain-specific modality representations are obtained by a style encoder E_{sty} that is implemented by a Variational Autoencoder (VAE) [43]. The VAE learns a low dimensional

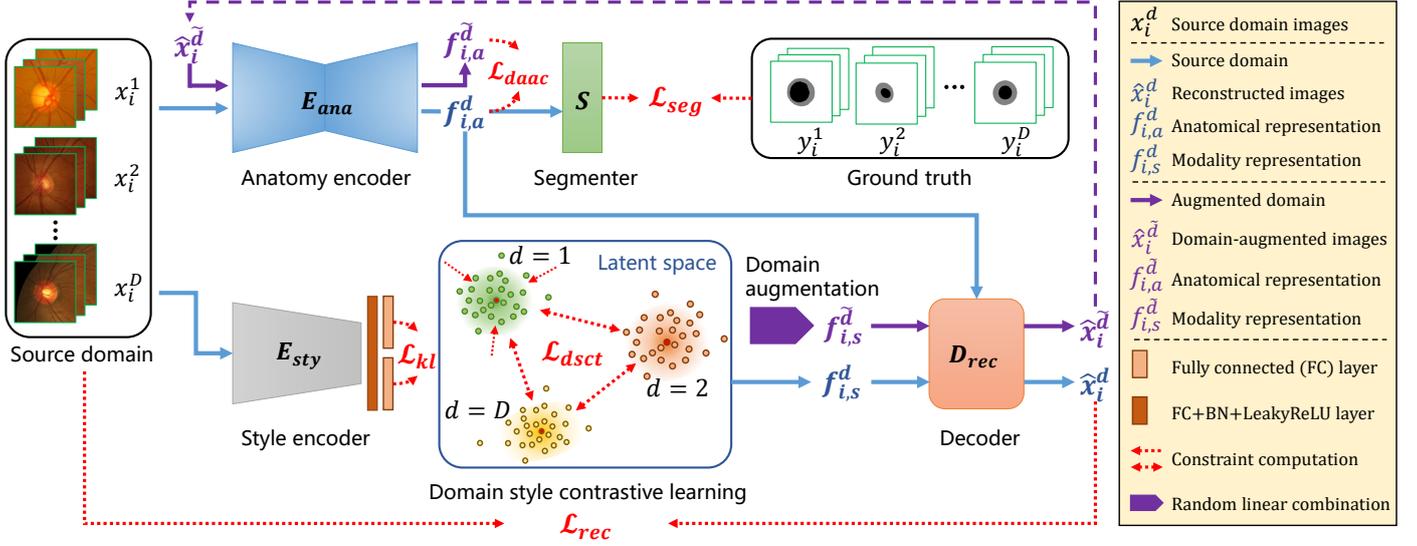


Figure 2: Overview of the proposed Contrastive Domain Disentanglement and Style Augmentation (CDDSA) network for multi-domain generalizable segmentation. We use an anatomy encoder E_{ana} and a modality encoder E_{sty} to extract anatomical representations f_a^d and style codes f_s^d , respectively. A reconstruction decoder D_{rec} takes f_a^d and f_s^d as input and obtains a reconstructed image \hat{x}^d . The decomposed anatomical representations f_a^d is further used for segmentation. \hat{x}^d is our simulated image using domain-augmentation strategy. \hat{x}^d will further input into E_{ana} to get its anatomical representation $f_a^{\bar{d}}$. \mathcal{L}_{saac} is used to encourage the consistency between $f_a^{\bar{d}}$ and f_a^d .

latent space so that the learned latent representations match a prior distribution of an isotropic multivariate Gaussian $p(z) = \mathcal{N}(0, 1)$. Given the input x_i^d , E_{sty} predicts the mean u_i^d and variance v_i^d of the distribution of a latent code $z \in \mathbb{R}^{1 \times Z}$, where Z is the length of the latent code. The style code $f_{i,s}^d$ of an input image x_i^d is sampled from the distribution characterized by mean u_i^d and variance v_i^d . VAE is trained to minimize a reparameterization error, and a KL divergence loss is computed between the estimated Gaussian distribution $q(z|u_i^d, v_i^d)$ and the unit Gaussian $p(z)$:

$$\mathcal{L}_{kl} = D_{kl}(q(z|u_i^d, v_i^d) \| p(z)) \quad (3)$$

where $D_{kl}(p \| q) = \sum p(x) \log \frac{p(x)}{q(x)}$. When training is finished, sampling a vector from the unit Gaussian over a latent space can obtain a new style code, and we send it together with an anatomical representation to the decoder to obtain a reconstructed image, where the decoder is used as a generative model, as detailed in the following.

3.1.3. Reconstruction Decoder

Fig. 3 shows the structure of our reconstruction decoder D_{rec} to generate an image \hat{x}_i^d given two decomposed representations $f_{i,a}^d$ and $f_{i,s}^d$. The collaboration of the two representations acts as a repainting mechanism where the anatomical representation $f_{i,a}^d$ is used to derive the anatomical content, and the modality representation $f_{i,s}^d$ is used to color the style distribution on the whole image [44].

Specifically, the decoder uses four convolutional blocks to map $f_{i,a}^d$ to a reconstructed image conditioned on three Style Reconstruction Modules (SRM), as shown in Fig. 3. For the intermediate feature map obtained by each convolutional

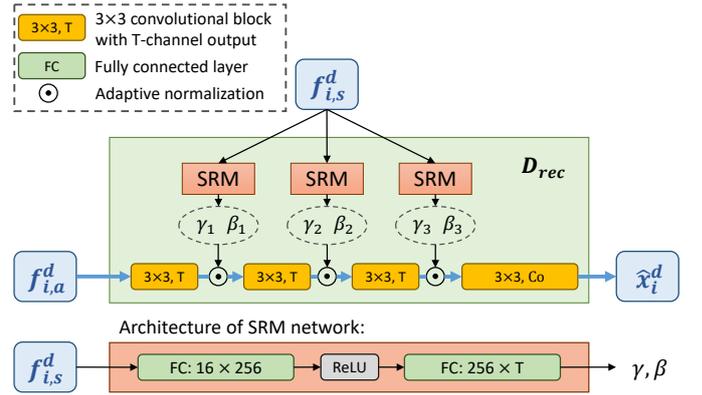


Figure 3: Framework of the reconstruction decoder D_{rec} . T is the channel of feature maps and Co represents the channel of output reconstructed image.

block in the decoder, we apply Adaptive Instance Normalization (AdaIN) to control the output style, where the affine transformation parameters (scale and bias) are predicted by an SRM that takes $f_{i,s}^d$ as input. Let $F_{i,c}$ represent the c -th channel of the intermediate feature map, we use two Fully Connected (FC) layers with a ReLU activation to implement the SRM that maps the modality representation $f_{i,s}^d$ to the scale $\gamma_{i,c}$ and bias $\beta_{i,c}$ that are used by affine transformation of AddIN:

$$AdaIN(F_{i,c} | \gamma_{i,c}, \beta_{i,c}) = \gamma_{i,c} \frac{F_{i,c} - \mu(F_{i,c})}{\sigma(F_{i,c})} + \beta_{i,c} \quad (4)$$

where each channel of the intermediate feature map is normalized separately, and we apply AddIN with SRM after each of three convolutional blocks in the decoder respectively, as shown in Fig. 3. By mapping $f_{i,s}^d$ to the scale and bias values for each intermediate feature map, the reconstruction decoder D_{rec} adap-

tively repaints style distribution on the anatomical representation $f_{i,a}^d$ in a coarse-to-fine manner. We use \hat{x}_i^d to denote the reconstructed image based on $f_{i,s}^d$ and $f_{i,a}^d$, and it is obtained by:

$$\hat{x}_i^d = D_{rec}(f_{i,s}^d, f_{i,a}^d) \quad (5)$$

As $f_{i,s}^d$ and $f_{i,a}^d$ are obtained from x_i^d , the reconstructed image \hat{x}_i^d should be as close as possible to x_i^d . Therefore, a reconstruction loss is employed to train the anatomy encoder E_{ana} , style encoder E_{sty} and reconstruction decoder D_{rec} :

$$\mathcal{L}_{rec} = \frac{1}{N_d} \sum_{i=1}^{N_d} |x_i^d - \hat{x}_i^d| \quad (6)$$

where we simply define the reconstruction loss as the Mean Absolute Error (MAE) loss due to its robustness to outliers.

3.2. Domain Style Contrastive Learning

An effective disentanglement expects that the style code $f_{i,s}^d$ to be domain-specific, but the reconstruction loss \mathcal{L}_{rec} does not provide sufficient supervision for achieving domain-specific style codes. To address the problem and make the model decompose more discriminative modality representations for different domains, we propose a domain style contrastive learning strategy to explicitly constrain the disentangled style code $f_{i,s}^d$.

Let x_i^d and x_j^d represent two different samples from the same domain d in the training set, and their style codes obtained by E_{sty} are denoted as $f_{i,s}^d$ and $f_{j,s}^d$, respectively. We define $(f_{i,s}^d, f_{j,s}^d)$ as a positive pair for maximizing their similarity. At the same time, for N samples each from a different domain d' ($d' \in [0, 1, \dots, D]$ and $d' \neq d$), their corresponding style codes compose a negative set \mathcal{N}_i^d for $f_{i,s}^d$, and each element in \mathcal{N}_i^d should have a minimized similarity compared with $f_{i,s}^d$. Following the standard formula of self-supervised contrastive loss InfoNCE [45, 46], we define our domain style contrastive loss as:

$$\mathcal{L}_{dsct} = -\log \frac{e^{\text{sim}(f_{i,s}^d, f_{j,s}^d)/\tau}}{e^{\text{sim}(f_{i,s}^d, f_{j,s}^d)/\tau} + \sum_{f \in \mathcal{N}_i^d} e^{\text{sim}(f_{i,s}^d, f)/\tau}} \quad (7)$$

where $\text{sim}(\cdot, \cdot)$ is the cosine similarity, and $\tau = 0.1$ is the temperature scaling parameter. In practice, to save the computational cost during training, we fetch b samples for each domain in a mini-batch, and their style codes are saved in a list Q . Let \tilde{Q} denote a permuted version of Q , the corresponding elements with the same index in the two lists are used as a positive pair, i.e., b positive pairs are considered for each domain in a mini-batch. The style codes of the $b(D-1)$ samples from domains other than domain d are used as the negative set for x_i^d .

3.3. Style Augmentation with Anatomical Consistency

Based on the disentangled anatomical representations, style codes and the decoder, we can augment the style of an image by replacing its style code during image reconstruction, and therefore propose a style augmentation strategy to automatically generate images in new domains with different styles. At each iteration of training, we denote the style codes of a batch

Table 1: Statistics of retinal fundus images in four domains used in our experiment following [15].

Domain No.	Dataset	Cases (train / test)	Scanner
Domain 1	Drishiti-GS	50 / 51	Aravind eye hospital
Domain 2	RIM-ONE-r3	99 / 60	Nidek AFC-210
Domain 3	REFUGE-train	320 / 80	Zeiss Visucam 500
Domain 4	REFUGE-val	320 / 80	Canon CR-2

as a style code bank $\mathcal{F} = \{f_{i,s}^d | i = 1, 2, \dots, B; d = 1, 2, \dots, D\}$, where the batch has B samples for each domain. Based on the style codes in \mathcal{F} , we obtain a new style code using a linear combination of them with random weights:

$$f_s^{\tilde{d}} = \sum_{i=1}^{|\mathcal{F}|} \alpha_i \mathcal{F}_i \quad (8)$$

where $f_s^{\tilde{d}}$ is a generated style code that is assumed to be from an unseen domain \tilde{d} . \mathcal{F}_i is the i -th element in the style code bank, and the weight $\alpha_i \in [-1, 1]$ is randomly sampled from a uniform distribution.

Given an anatomical representation $f_{i,a}^d$ from an image in the source domain, we repaint it with the new style code $f_s^{\tilde{d}}$ to generate a new image $\hat{x}_i^{\tilde{d}}$:

$$\hat{x}_i^{\tilde{d}} = D_{rec}(f_s^{\tilde{d}}, f_{i,a}^d) \quad (9)$$

Since the generated image $\hat{x}_i^{\tilde{d}}$ and the real image x_i^d share the same anatomical representation $f_{i,a}^d$, we introduce an anatomical consistency loss \mathcal{L}_{saac} that forces the anatomy encoder E_{ana} to obtain domain-invariant anatomical representations in spite of the different styles between $\hat{x}_i^{\tilde{d}}$ and x_i^d :

$$\mathcal{L}_{saac} = \frac{1}{N_d} \sum_{i=1}^{N_d} |f_{i,a}^d - E_{ana}(\hat{x}_i^{\tilde{d}})| \quad (10)$$

where the MAE loss is used for the anatomical consistency.

3.4. Overall Loss

As a summary of the proposed CDDSA framework, the overall loss function for training is formulated as:

$$\mathcal{L} = \mathcal{L}_{seg} + \lambda_1 \mathcal{L}_{kl} + \lambda_2 \mathcal{L}_{rec} + \lambda_3 \mathcal{L}_{dsct} + \lambda_4 \mathcal{L}_{saac} \quad (11)$$

where \mathcal{L}_{seg} is the supervised segmentation loss (Eq. 2), \mathcal{L}_{kl} is the KL divergence loss for style encoder (Eq. 3), \mathcal{L}_{rec} is the image reconstruction loss (Eq. 6), \mathcal{L}_{dsct} and \mathcal{L}_{saac} are the domain style contrastive loss (Eq. 7) and style augmentation-based anatomical consistency loss (Eq. 10), respectively. $\lambda_1, \lambda_2, \lambda_3$ and λ_4 act as trade-off parameters for different loss terms.

4. Experiments and Results

4.1. Datasets and Implementation Details

In this study, we evaluated our proposed CDDSA and compared it with several state-of-the-art DG methods on a public multi-domain fundus image dataset and an in-house multi-domain nasopharyngeal carcinoma MRI dataset.

Table 2: Statistics of the in-house multi-domain nasopharyngeal carcinoma MRI dataset [47]

Domain No.	Sequence	Slice thickness (mm)	Volumes (train / test)	Slices (train / test)	Total Volumes (train / test)	Total slices (train / test)	Scanner
Domain 1	T1	6 - 7.75	39 / 26	305 / 201			SPPH - Siemens
Domain 2	CE-T1	3	27 / 18	359 / 234	114 / 75	1427 / 994	WCH - Siemens
Domain 3	T1-water	3	24 / 15	402 / 302			WCH - Siemens
Domain 4	T2-water	3	24 / 16	361 / 257			WCH - Siemens

Multi-domain Fundus Image Dataset: For a fair comparison with state-of-the-art DG methods, we evaluated our approach for Optic Cup (OC) and Disc (OD) segmentation on a public multi-domain retinal fundus image dataset¹ [15]. The dataset was collected from four public fundus image datasets obtained by different scanners at different sites that have distinct domain discrepancies in visual appearance and image quality: Domain 1 is from the Drishti-GS [48] dataset containing 50 and 51 images for training and testing, respectively; Domain 2 is from the RIM-ONE [49] dataset containing 99 and 60 images for training and testing, respectively; and the Domain 3 and 4 are from REFUGE [50] challenge’s training and validation datasets, respectively, and both of them contain 320 and 80 images for training and testing.

To evaluate generalizability of OC/OD segmentation models in unseen domains, we followed the leave-one-domain-out cross validation strategy in DoFE [15], where each time three domains were used for training and the other domain was used as the unseen testing domain. In total, there are 789 and 271 images for training and testing, respectively. The statistics of these multi-domain retinal fundus images are summarized in Table 1. For preprocessing, we adopted a series of basic data augmentations to enhance the diversity of training samples as conducted by DoFE [15], and the images were randomly cropped with a size of 256×256 during training.

Multi-domain Nasopharyngeal Carcinoma MRI Dataset: We collected an in-house multi-domain Nasopharyngeal Carcinoma (NPC) MRI dataset for nasopharynx Gross Tumor Volume (GTVnx) segmentation. It was collected from two hospitals with four different imaging protocols [47] (i.e., four domains): T1-weighted imaging, gadolinium contrast-enhanced T1-weighted (CE-T1) imaging, T1 water imaging and T2 water imaging, respectively. Images in Domain 1 were collected from Sichuan Provincial People’s Hospital (SPPH) with slice thickness of 6 - 7.75 mm, and images in Domain 2-4 were collected from West China Hospital (WCH) with slice thickness of 3 mm. In total, there were 189 volumes each from a specific patient, and they were split to 114 for training and 75 for testing. The corresponding slice numbers for training and testing were 1427 and 994, respectively. The volume and slice numbers for each domain are detailed in Table 2.

For preprocessing, we unified the orientation of different volumes into the standard RAI (right to left, anterior to posterior, inferior to superior in the x-, y-, and z-axes, respectively). The voxel intensity was clipped by the 0.1 and 99.9 percentiles of each volume and then normalized to [0, 255]. Each volume

was firstly cropped along z axis based on the slices containing GTVnx delineation, and then center-cropped with a 256×256 window in x-y plane. We used 2D networks for the GTVnx segmentation in each slice and stacked the results into a 3D volume for evaluation, and a leave-one-domain-out cross validation strategy was also employed during the experiment.

Implementation Details: Training and inference were implemented on one NVIDIA GeForce GTX 1080 Ti GPU. The anatomical representation E_{ana} was implemented by U-Net [1] as the backbone, with channel numbers of 16, 32, 64, 128 and 256 at five resolution scales, respectively. We set the channel number of anatomical representations as $T = 8$. The segmenter S consists of two convolutional blocks. The first block has a convolution layer with a kernel size of 3×3 followed by BN and LeakyReLU (sloop = 0.2), and the second block has a 1×1 convolution layer followed by Softmax to obtain a segmentation probability map. The style encoder E_{sty} has convolutional blocks each with a down-sampling layer to reduce the resolution, and the output of the last convolutional block is sent to two fully connected layers to obtain the mean and variance of a Gaussian distribution for the latent style code, and the size of the latent style code was set as $Z = 16$.

The weights in the total loss function were: $\lambda_1 = 1.0$, $\lambda_2 = 0.001$, $\lambda_3 = 0.01$ and $\lambda_4 = 1.0$, respectively. The networks were trained with the Adam optimizer, and the learning rate was initialized to 10^{-3} and decayed to 95% when the performance did not improve in 8 epochs. In a mini-batch, the image/slice number for each domain was 8 and 6 for the fundus image and NPC-MRI datasets, respectively. The epoch number was 200 and 400 for the fundus image and NPC-MRI datasets, respectively. To measure the segmentation performance quantitatively, we adopt the Dice score (Dice) and Average Symmetric Surface Distance (ASSD) for evaluation.

4.2. Fundus Image Segmentation

4.2.1. Comparison with State-of-the-art DG Methods

For domain generalization study, we conducted leave-one-domain-out cross validation on the multi-domain fundus image dataset. We first considered all the available training domains as a single dataset (i.e., ignoring the domain shift in training set) and trained a U-Net [1] using a standard Dice loss, and directly applied it to the unseen domain, which is referred to as ‘**Inter-domain**’ and serves as a lower bound of the experiment. Then, for each domain, we trained and tested the U-Net [1] with the training and testing sets respectively, i.e., no unseen domain involved, which serves as the upper bound for DG and is referred as ‘**Intra-domain**’. For DG methods, we compared our proposed CDDSA with four representative state-

¹<https://github.com/emmasjwang/Dofe>

Table 3: Comparison of Dice (%) by different DG methods on the multi-site fundus image dataset. CDDSA \diamond means that the new style code for style augmentation was randomly sampled from a Gaussian distribution rather than obtained by a random linear combination of style codes in the source domains.

Methods	Domain 1		Domain 2		Domain 3		Domain 4		Avg	
	cup	disc	cup	disc	cup	disc	cup	disc	cup	disc
Lower bound (Inter-domain)	74.38±12.96	96.67±2.04	77.71±20.84	85.05±14.67	79.72±9.51	90.01±5.81	86.63±8.52	89.55±3.26	79.61	90.32
Upper bound (Intra-domain)	83.35±13.99	96.10±1.88	81.53±9.42	94.62±3.01	87.57±7.59	95.91±1.85	88.88±7.10	95.58±1.98	85.33	95.55
BigAug [14]	82.36±11.74	93.73±9.29	75.45±15.01	87.83±11.17	84.32±9.45	91.99±10.72	85.32±7.50	92.97±6.58	81.86	91.63
DoFE [15]	80.25±10.84	95.61±1.45	78.97±14.80	88.74±4.58	84.81±7.71	92.81±2.63	86.65±6.39	93.46±2.43	82.67	92.66
FedDG [23]	79.84±13.55	93.50±4.11	76.57±13.95	88.74±4.91	84.23±6.80	93.73±3.22	85.33±10.19	94.03±4.14	81.49	92.50
DCA-Net [13]	82.16±12.23	94.39±2.94	80.63±15.58	91.50±2.78	84.48±7.77	91.63±4.38	87.11±12.67	93.05±4.98	83.60	92.64
Baseline	80.63±11.55	95.02±2.65	79.35±13.66	89.76±3.20	83.29±8.04	93.67±3.36	84.12±11.33	93.51±4.03	81.85	92.99
+ \mathcal{L}_{dset}	80.64±11.94	96.11±2.95	80.13±17.39	88.27±1.23	85.20±8.06	92.36±2.59	86.33±9.12	93.36±2.60	83.08	92.53
+ \mathcal{L}_{saac}	84.28±11.56	96.13±1.35	81.95±12.60	88.18±3.50	84.59±8.11	92.95±3.30	86.49±9.69	93.27±3.40	84.33	92.63
+ \mathcal{L}_{dset} + \mathcal{L}_{saac} (CDDSA \diamond)	85.53±11.37	96.74±1.70	76.39±17.45	88.25±6.91	84.60±7.76	92.34±4.08	86.56±9.75	92.63±3.39	83.27	92.49
+ \mathcal{L}_{dset} + \mathcal{L}_{saac} (CDDSA)	85.75±12.31	96.79±1.53	81.04±13.63	89.71±3.60	86.94±7.94	93.25±3.55	86.86±8.97	94.44±3.96	85.15	93.55

Table 4: Comparison of ASSD (pixel) by different DG methods on the multi-site fundus image dataset.

Methods	Domain 1		Domain 2		Domain 3		Domain 4		Avg	
	cup	disc	cup	disc	cup	disc	cup	disc	cup	disc
Lower bound (Inter-domain)	22.35±9.74	6.47±3.80	15.77±20.21	18.25±19.60	12.30±5.82	12.33±5.03	7.45±4.60	9.27±2.62	14.47	11.58
Upper bound (Intra-domain)	16.04±6.65	7.84±3.87	13.10±7.68	8.55±5.80	8.41±5.02	6.32±4.02	6.07±3.41	5.46±2.48	10.91	7.04
BigAug [14]	17.91±10.11	8.67±4.08	22.33±15.26	19.77±6.69	13.51±7.67	14.46±4.96	8.90±5.02	8.77±6.63	15.66	12.92
DoFE [15]	17.16±9.40	7.62±2.38	15.28±12.94	14.52±5.36	10.73±6.22	10.11±5.11	7.18±3.23	7.60±3.64	12.59	9.96
FedDG [23]	18.97±12.82	7.83±3.11	15.34±9.33	13.74±6.79	12.21±5.57	9.71±5.63	9.21±6.62	8.15±5.89	13.93	9.86
DCA-Net [13]	17.19±7.64	9.32±4.70	12.39±12.32	10.46±3.23	11.28±5.21	11.32±5.54	7.37±6.51	7.22±4.75	12.06	9.58
Baseline	17.33±8.58	8.21±4.61	13.10±6.66	11.79±3.90	11.03±5.22	9.31±4.23	8.04±6.42	7.43±4.89	12.38	9.18
+ \mathcal{L}_{dset}	18.21±8.05	7.52±5.85	13.33±13.43	14.10±13.50	10.09±5.42	9.98±3.11	7.30±3.88	7.03±2.57	12.23	9.66
+ \mathcal{L}_{saac}	15.77±6.93	7.14±2.59	11.04±5.91	12.97±3.94	10.58±5.21	9.60±3.71	7.22±5.31	7.51±4.10	11.15	9.31
+ \mathcal{L}_{dset} + \mathcal{L}_{saac} (CDDSA \diamond)	14.85±6.87	6.78±3.47	15.35±12.98	14.61±12.51	10.72±5.25	9.92±4.07	7.35±4.44	7.75±3.77	12.07	9.77
+ \mathcal{L}_{dset} + \mathcal{L}_{saac} (CDDSA)	14.65±8.39	6.54±3.74	12.91±10.79	13.06±8.60	9.38±5.40	9.32±4.11	7.28±5.85	6.87±5.03	11.06	8.95

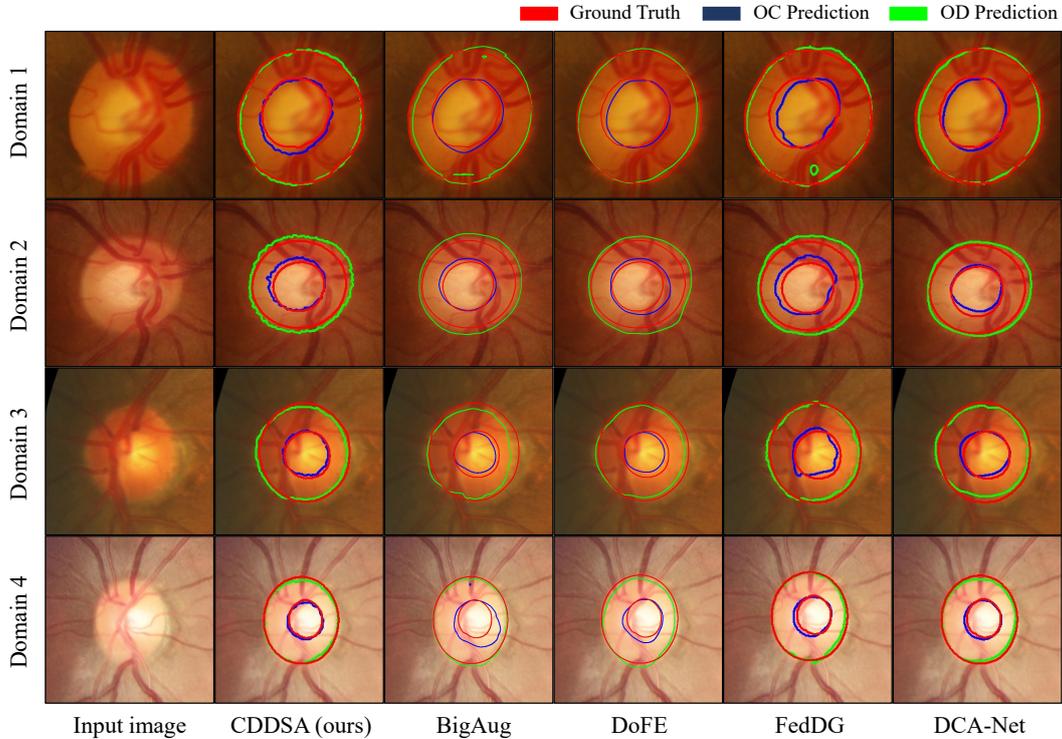


Figure 4: Visual comparison between our proposed CDDSA and BigAug [14], DoFE [15], FedDG [23] and DCA-Net [13] on multi-domain fundus image segmentation.

Table 5: Comparison between different activation functions used by the output of E_{ana} for multi-domain OC/OD segmentation. Gumbel-H and Gumbel-S are two variants of gumbel softmax that return discrete one-hot values and soft continuous values, respectively.

Metric	Activation	Domain 1		Domain 2		Domain 3		Domain 4		Avg	
		cup	disc	cup	disc	cup	disc	cup	disc	cup	disc
Dice	Gumbel-H	82.29±11.97	96.46±2.11	78.67±18.79	86.21±5.06	85.40±8.14	93.36±3.19	87.89±7.67	93.22±3.16	83.56	92.31
	Gumbel-S	82.88±11.28	96.71±1.65	80.22±13.76	89.11±3.52	83.20±8.90	92.59±4.91	86.41±9.88	94.39±3.53	83.18	93.20
	softmax	85.22±10.46	96.94±1.26	80.72±16.09	88.42±12.16	85.22±7.43	93.13±3.75	85.34±10.07	93.23±3.78	84.13	92.93
	tanh	85.75±12.31	96.79±1.53	81.04±13.63	89.71±3.60	86.94±7.94	93.25±3.55	86.86±8.97	94.44±3.96	85.15	93.55
ASSD	Gumbel-H	18.31±8.49	6.65±4.11	14.60±17.21	18.13±16.60	9.92±5.45	9.10±3.66	6.67±3.96	6.97±3.02	12.38	10.21
	Gumbel-S	17.22±6.91	6.33±2.83	12.30±6.02	12.98±5.03	10.94±5.57	9.78±5.85	7.52±5.42	6.89±4.77	12.00	9.00
	softmax	14.40±6.36	6.67±2.53	12.61±12.71	13.28±12.93	10.01±5.14	9.90±5.09	8.28±5.90	7.67±4.73	11.33	9.38
	tanh	14.65±8.39	6.54±3.74	12.91±10.79	13.06±8.60	9.38±5.40	9.32±4.11	7.28±5.85	6.87±5.03	11.06	8.95

of-the-art approaches: BigAug [20] based on data augmentation, DoFE [15] based on domain-oriented feature embedding, DCA-Net [13] based on domain composition and attention, and FedDG [23] that is a federated learning-based domain generalization method.

Table 3 and Table 4 show the quantitative evaluation results of OC/OD segmentation in terms of Dice and ASSD, respectively. Intra-domain achieved the highest performance among the compared methods, with an average Dice of 85.33% and 95.55% for the OC and OD across the four domains. In contrast, the average Dice achieved by Inter-domain was only 79.61% and 90.32% in OC and OD segmentation, respectively, showing the performance gap caused by domain shift. BigAug [20] obtained a slight improvement from Inter-domain, suggesting that aimlessly conducting data augmentation in the image domain has a limited performance. Among the compared existing methods, DCA-Net [13] achieved the highest performance, with an average Dice of 83.60% and 92.64% for OC and OD, respectively. In contrast, our proposed CDDSA outperformed the existing methods, with an average Dice of 85.15% and 93.55% for OC and OD, respectively. The average ASSD obtained by our method was 11.06 and 8.95 pixels for OC and OD, respectively, which also outperformed the compared methods, as shown in Table 4. Fig. 4 shows a visual comparison between our proposed CDDSA and BigAug, DoFE, FedDG and DCA-Net for images from the four testing domains, respectively. It shows that the segmentation results obtained by our proposed CDDSA had boundaries that are closer to the ground truth, while the other DG methods have more over- and under-segmented regions than ours.

4.2.2. Ablation Studies

Effectiveness of Domain Style Contrastive Learning and Style Augmentation: We conducted ablation studies to evaluate the effectiveness of the components of our CDDSA framework, where the baseline was only training E_{ana} , E_{sty} , D_{rec} and S with basic loss functions of \mathcal{L}_{seg} , \mathcal{L}_{kl} and \mathcal{L}_{rec} , following SDNet [19]. We use $+\mathcal{L}_{dsc}$ and $+\mathcal{L}_{saac}$ to denote adding the domain style contrastive learning and domain style augmentation with anatomical consistency to the baseline, respectively. $+\mathcal{L}_{dsc} + \mathcal{L}_{saac}$ means our proposed CDDSA.

Quantitative evaluation results in terms of Dice and ASSD of these variants are shown in the last section of Table 3 and Table 4, respectively. The baseline’s average Dice across OC

and OD was 87.42%, and combining the baseline with \mathcal{L}_{dsc} improved it to 88.81%, indicating that encouraging the network to obtain more discriminative style codes leads to better generalization performance. Baseline + \mathcal{L}_{saac} also improved the two classes’ average Dice to 88.48%, and our method using \mathcal{L}_{dsc} and \mathcal{L}_{saac} further improved the average Dice to 89.35% (85.15% for OC and 93.55% for OD), showing the extra improvement brought by the proposed style augmentation.

To additionally evaluate the effectiveness of our proposed random linear combination for generating new style code during style augmentation, we compared it with an alternative method that randomly samples the style code from a Gaussian distribution, which is denoted as CDDSA^o in Table 3 and Table 4. The results showed CDDSA^o performed slightly worse than CDDSA, but outperformed most existing DG methods, proving that our proposed random linear combination was better for style augmentation than randomly sampling the style code from a Gaussian distribution.

Reconstruction and Style Augmentation Quality: Since the decomposed anatomical representation f_a serves as the input of the segmentor S and the reconstruction decoder D_{rec} , the quality of the anatomical representation has an impact on the performance of S and D_{rec} . To explore the influence of different formats of the anatomical representation on reconstruction and segmentation quality, we compared four activation functions at the end of E_{ana} : 1) the gumbel softmax [51] returning discrete one-hot values, which is referred to as Gumbel-H; 2) the gumbel softmax returning continuous soft values, which is referred to as Gumbel-S; 3) Softmax and 4) Tanh. Quantitative comparison of these activation functions in the fundus image segmentation task is shown in Table 5. We found that Gumbel-S obtained a better segmentation performance than Gumbel-H (83.56% and 92.31% vs 83.18% and 93.20% of average Dice score). Softmax and tanh further improved the model’s performance. Notably, Tanh achieved the highest average Dice score of 85.15% in OC and 93.55% in OD and the lowest average ASSD (11.06 pixels in OC and 8.95 pixels in OD) among the compared activation functions. The results show that using continuous soft values for anatomical representations led to better segmentation performance, as soft representations are more informative compared with binary representations [52].

To further investigate how the activation function used by E_{ana} affects the reconstructed images and style augmentation, we compared the original training image with the im-

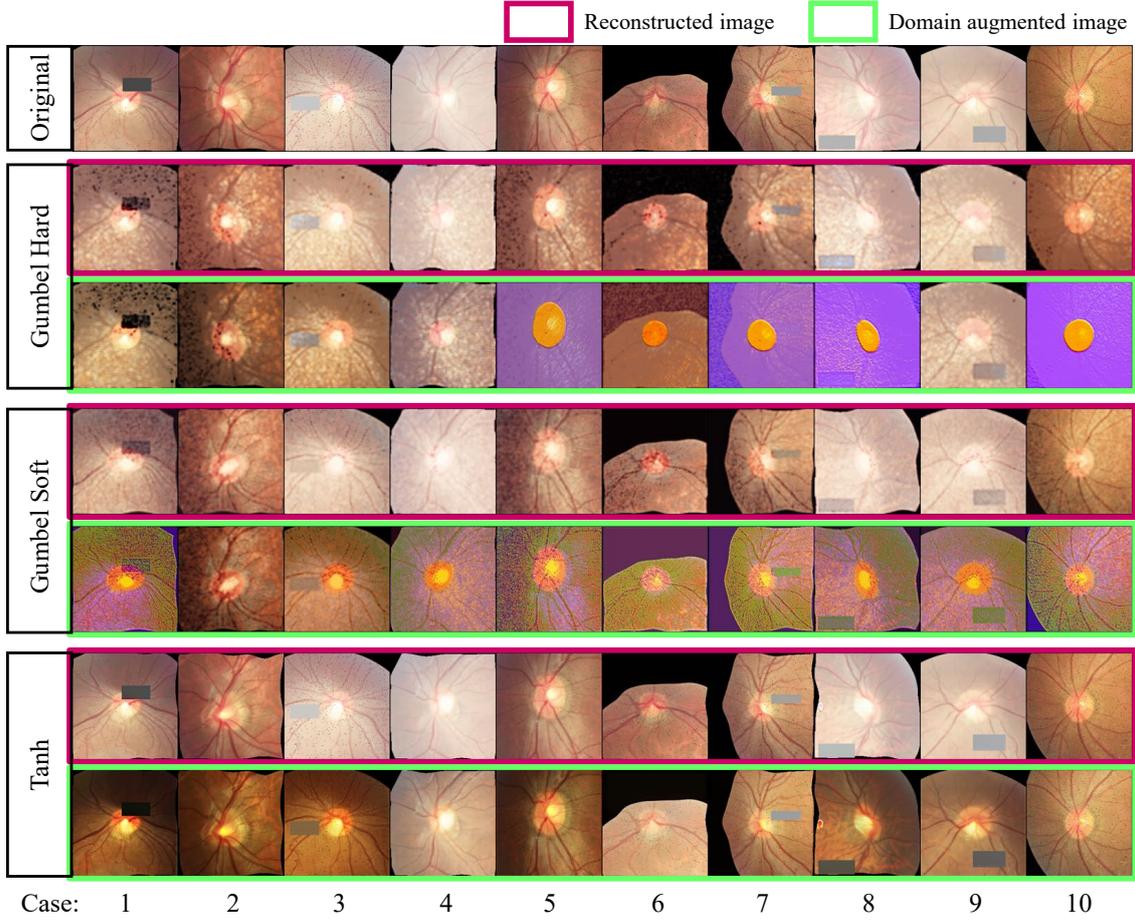


Figure 5: Visual comparison of reconstructed and augmented fundus images with different activation functions at the end of anatomy Encoder. The original images are from domain 4. For each method, the first row (red rectangles) shows images reconstructed from the disentangled anatomical representation and style code, and the second row (green rectangles) shows style-augmented images that are generated based on the anatomical representation from the original images and changed style codes.

Table 6: Quantitative comparison of different DG methods on the multi-domain NPC-MRI image dataset for GTVnx segmentation.

Methods	Domain 1		Domain 2		Domain 3		Domain 4		Avg	
	Dice (%)	ASSD (pix)	Dice (%)	ASSD (pix)	Dice (%)	ASSD (pix)	Dice (%)	ASSD (pix)	Dice	ASSD
Lower bound (Inter-domain)	65.44±13.35	2.37±1.52	76.65±7.75	1.70±1.06	81.17±7.01	1.61±0.77	61.62±13.66	3.38±1.47	71.22	2.27
Upper bound (Intra-domain)	79.19±6.35	1.13±0.54	82.89±7.54	1.58±1.75	86.30±3.25	1.26±0.48	79.46±7.40	2.09±1.48	81.96	1.52
BigAug [14]	75.63±5.97	1.67±1.01	78.51±7.89	1.65±0.99	82.30±5.05	1.82±0.66	63.88±12.32	4.05±1.71	75.08	2.30
DoFE [15]	78.44±6.97	1.27±0.99	75.00±4.95	1.80±0.85	79.66±5.77	1.66±0.68	64.71±15.06	2.57±1.91	74.45	1.83
FedDG [23]	65.07±11.59	2.58±1.63	78.90±6.11	1.67±1.09	81.57±6.05	1.79±0.79	72.38±12.13	3.91±2.64	74.48	2.49
DCA-Net [13]	77.27±6.66	1.27±0.99	77.14±7.53	1.80±0.85	81.63±6.20	1.66±0.68	69.32±10.08	2.57±1.91	76.34	1.83
Baseline	76.63±5.74	1.53±1.23	77.93±5.99	1.41±0.72	82.77±4.65	1.55±0.55	62.71±11.63	3.07±1.91	75.01	1.89
+ \mathcal{L}_{dset}	77.02±6.09	1.56±1.09	77.11±7.15	1.71±0.89	83.02±3.62	1.45±0.41	63.88±12.74	2.81±1.66	75.26	1.88
+ \mathcal{L}_{saac}	77.74±5.67	1.35±0.82	76.33±8.37	1.66±0.83	83.23±5.38	1.44±0.56	68.43±12.92	3.45±1.99	76.43	1.98
+ \mathcal{L}_{dset} + \mathcal{L}_{saac} (CDDSA \diamond)	77.77±6.65	1.36±0.97	78.00±7.74	1.52±0.86	83.10±5.18	1.60±0.63	66.39±10.88	3.50±1.62	76.32	2.00
+ \mathcal{L}_{dset} + \mathcal{L}_{saac} (CDDSA)	78.34±5.14	1.37±0.82	79.16±6.68	1.61±1.19	83.53±4.55	1.48±0.54	69.53±10.28	2.46±1.50	77.64	1.73

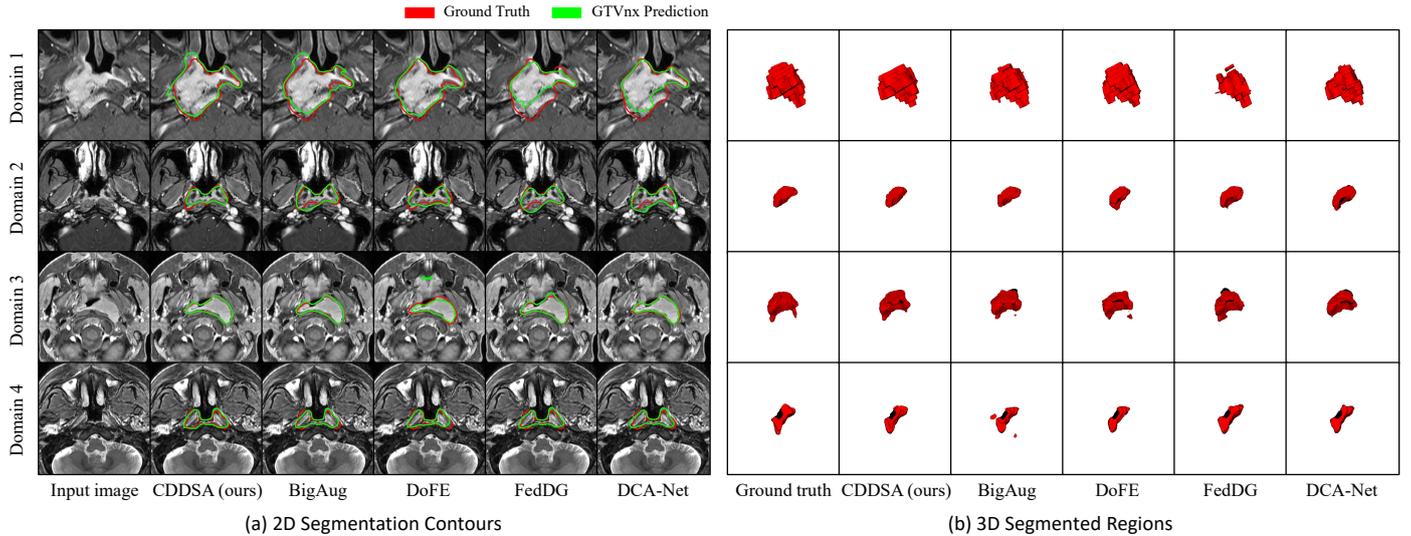


Figure 6: Visual comparison between different DG methods for multi-domain NPC GTVnx segmentation.

age reconstructed from disentangled f_a and f_s and the style-augmented image in Fig 5, where we show the differences between Gumbel-H, Gumbel-S and Tanh. First, for image reconstruction (even rows), it can be observed that using Gumbel-H only reconstructed coarse-grained images and only roughly retained the overall content without detailed structures. Gumbel-S achieved a better quality with more details than Gumbel-H. However, the reconstructed images have some noticeable artefacts compared with the original images. In contrast, Tanh obtained a much higher quality than Gumbel-H, and the reconstructed images were closer to the original inputs in terms of both anatomical structures and styles. Second, for style augmentation (odd rows), Gumbel-H can not keep the same anatomical structure after changing the style, and led to unrealistic images in the augmented domain, especially in cases 5, 7, 8 and 10 in the third row of Fig. 5. Gumbel-S has a better ability to remain the anatomical structures, but the augmented images have a lot of artefacts with unrealistic appearance. In contrast, Tanh achieved very high quality in the style-augmented images with realistic appearances. They have quite different styles with shared anatomical structures compared with the original images, as shown in the last row of Fig. 5. The results show the advantage of our proposed style augmentation strategy, which can successfully generate new samples in an unknown domain with anatomical structures unchanged, which is beneficial for enhancing the model’s generalizability.

4.3. NPC GTVnx Image segmentation

4.3.1. Comparison with State-of-the-art DG Methods

For the multi-domain GTVnx segmentation task, we employed the same set of methods as in Section 4.2.1 for comparison, and the quantitative evaluation results are shown in Table 6. First, Intra-domain (upper bound) achieved the highest performance with average Dice of 81.96% and average ASSD of 1.52 mm across the four domains. In contrast, Inter-domain (lower bound) only obtained an average Dice of 71.22% and

ASSD of 2.27 mm. The performance gap between them was over 10% in average Dice, indicating the large shift among the different domains. All the four existing DG methods achieved great improvements compared with the Inter-domain that does not consider the differences across domains. DoFE [15] and FedDG [23] had a similar segmentation performance, with average Dice score of 74.45% and 74.48%, respectively. BigAug [14] achieved an average Dice of 75.08%, indicating that some augmentation strategies are beneficial for GTVnx segmentation in cross-modality MRI images. DCA-Net [13] obtained an average Dice of 76.34%, which outperformed the other three existing DG methods. In contrast, our proposed CDDSA based on domain-invariant feature learning obtained higher generalizability, achieving an average Dice of 77.64% and ASSD of 1.73 mm, which outperformed the state-of-the-art DG methods.

Fig. 6 provides a visual comparison between our proposed CDDSA and the four state-of-the-art DG methods on the multi-domain NPC GTVnx segmentation dataset. Fig. 6(a) shows that the 2D segmentation boundaries of our CDDSA are closer to the ground truth than those of the other methods. The 3D visualization in Fig. 6 (b) shows that our CDDSA achieved high-quality segmentation results, while the other DG methods have more noises in the results.

4.3.2. Ablation Studies

Effectiveness of Domain Style Contrastive Learning and Style Augmentation: Similar with multi-site fundus image segmentation. We also proved the effectiveness of our proposed domain style contrastive learning and style augmentation strategy in multi-site NPC GTVnx segmentation. Quantitative results are shown in Table 6. The baseline, i.e., re-implementation of SDNet [19] based on our network structures, obtained an average Dice of 75.01%, and combining it with our domain style contrastive learning \mathcal{L}_{dsc} improved it to 75.26%. Combining it with our domain augmentation method \mathcal{L}_{saac} achieved an average Dice of 76.43%. In contrast, our proposed method

Table 7: Comparison between different activation functions used by the output of E_{ana} for multi-domain NPC GTVnx segmentation.

Activation	Domain 1		Domain 2		Domain 3		Domain 4		Avg	
	Dice (%)	ASSD (pix)	Dice (%)	ASSD (pix)	Dice (%)	ASSD (pix)	Dice (%)	ASSD (pix)	Dice (%)	ASSD (pix)
Gumbel-H	75.60±6.20	1.61±1.15	72.69±9.42	1.60±0.76	82.81±7.67	1.62±0.58	64.73±13.53	2.90±1.88	73.96	1.93
Gumbel-S	76.54±6.39	1.48±1.10	78.40±6.48	1.57±0.88	83.23±4.54	1.47±0.52	65.24±11.78	2.73±1.43	75.88	1.81
softmax	77.46±5.40	1.57±1.01	79.87±4.70	1.46±0.74	82.85±4.89	1.89±0.91	62.59±12.48	3.08±1.59	75.69	2.00
tanh	78.34±5.14	1.37±0.82	79.16±6.68	1.61±1.19	83.53±4.55	1.48±0.54	69.53±10.28	2.46±1.50	77.64	1.73

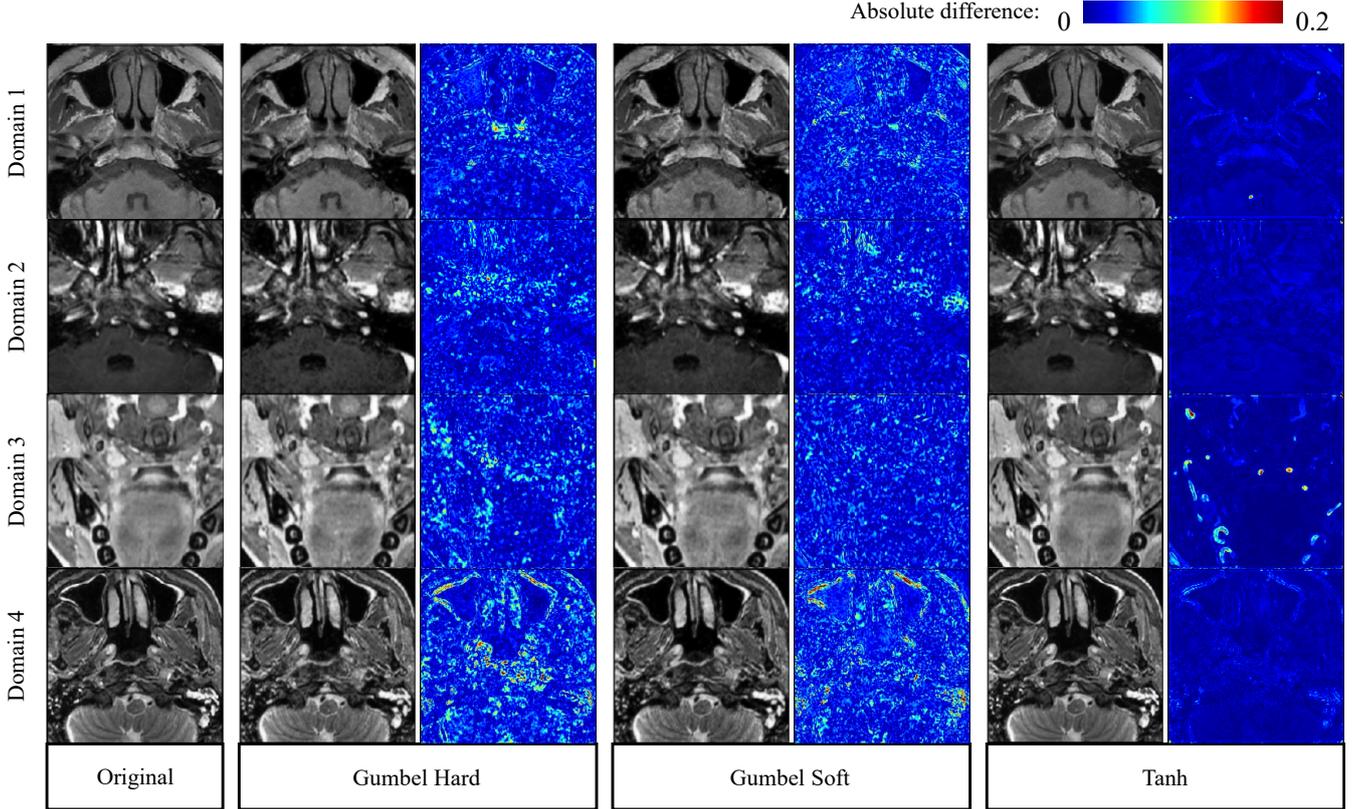


Figure 7: Comparison of reconstructed images with different activation functions at the end of E_{ana} . For each method, the first column shows the reconstructed images based on the disentangled anatomical representation and style code, and the second column shows the absolute difference between the reconstructed and original images.

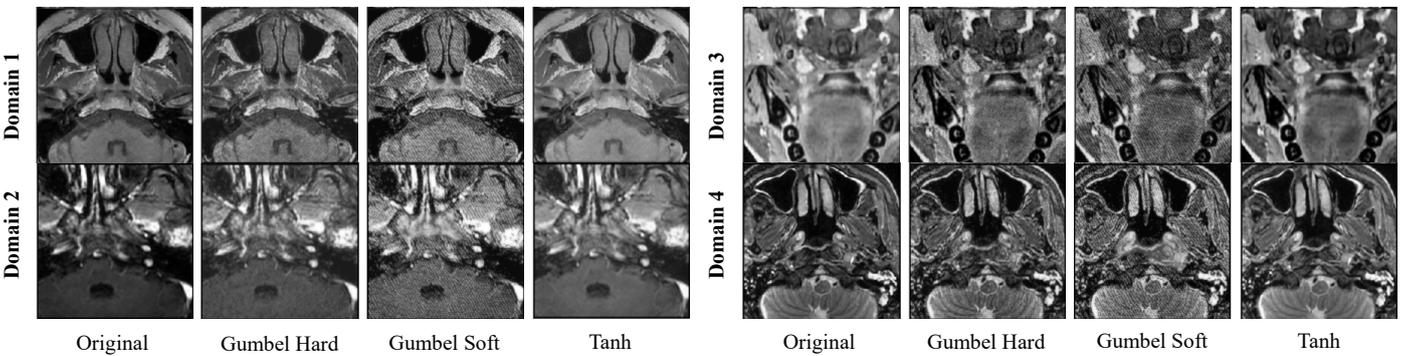


Figure 8: Visual comparison of style-augmented images with different activation functions at the end of E_{ana} .

that uses \mathcal{L}_{dset} and \mathcal{L}_{saac} simultaneously improved the average Dice to 77.64%, which is the highest among the compared variants and significantly better than the baseline (p -value < 0.05). Table 6 also shows that CDDSA performed better than CDDSA \diamond (77.64% vs 76.32%) in terms of Dice, indicating that our style augmentation based on random linear combination of the style codes was better than directly sampling style codes from a Gaussian distribution for style augmentation.

Reconstruction and Style Augmentation Qualities: Similar to Section 4.2.2, we compared four different activation functions at the end of E_{ana} to represent f_a on the NPC-MRI dataset. The corresponding NPC-MRI GTVnx segmentation results are shown in Table 7. It can be observed that Gumbel-S had a higher performance than Gumbel-H (75.88% vs 73.96% in terms of average Dice). Using Tanh further improved the average Dice to 77.64%, which was significantly better than the other activations.

Fig. 7 shows a visual comparison of these activation functions in reconstructing the original image after disentanglement. It can be observed that when Gumbel-H is used, the reconstructed images have a large difference from the original images. Gumbel-S has a lower reconstruction error than Gumbel-H. However, it is inferior to our method using Tanh, showing that Tanh is more suitable to obtaining anatomical representation in disentanglement for high-fidelity reconstruction.

In addition, Fig. 8 shows a visual comparison of style-augmented images when different activation functions are used at the end of E_{ana} . We found that all the methods can generate new-style images based on the augmented domain style code \tilde{f}_s and the anatomical representation f_a of the input. However, Gumbel-H and Gumbel-S led to obvious artifacts in the augmented images. In contrast, our method can change the style of an input image while better retaining the anatomical structures.

5. Conclusion

In this paper, we present a Contrastive Domain Disentanglement and Style Augmentation (CDDSA) framework to tackle the domain generalization problem in medical image segmentation. We introduce a GAN-free efficient disentangle method to decompose medical images from multiple domains into a domain-invariant anatomical representation and a domain-specific style code, where a segmentor works on the anatomical representation to achieve generalizability. To improve the disentanglement and segmentation performance, we use a soft representation for the anatomical representation based on Tanh, and propose domain style contrastive learning to minimize the similarity of style codes in different domains. Based on the disentanglement, we propose a style augmentation strategy that changes the style of an image with remained structure information for augmentation, which can further improve the model’s generalizability. Quantitative experimental results on a multi-site fundus image dataset and a multi-domain NPC MRI dataset showed that our CDDSA outperformed several state-of-the-art multi-domain generalization methods. In the future, it is of interest to apply our CDDSA framework to other multi-domain medical image analysis tasks.

References

- [1] O. Ronneberger, P. Fischer, T. Brox, U-net: Convolutional networks for biomedical image segmentation, in: International Conference on Medical image computing and computer-assisted intervention, Springer, pp. 234–241.
- [2] D. Shen, G. Wu, H.-I. Suk, Deep learning in medical image analysis, Annual review of biomedical engineering 19 (2017) 221–248.
- [3] R. Gu, G. Wang, T. Song, R. Huang, M. Aertsen, J. Depreest, S. Ourselin, T. Vercauteren, S. Zhang, CA-Net: Comprehensive attention convolutional neural networks for explainable medical image segmentation, IEEE Transactions on Medical Imaging 40 (2021) 699–711.
- [4] H. Guan, M. Liu, Domain adaptation for medical image analysis: a survey, IEEE Transactions on Biomedical Engineering 69 (2021) 1173–1185.
- [5] Y. Ganin, E. Ustinova, H. Ajakan, P. Germain, H. Larochelle, F. Laviolette, M. Marchand, V. Lempitsky, Domain-adversarial training of neural networks, The journal of machine learning research 17 (2016) 2096–2030.
- [6] K. Kamnitsas, C. Baumgartner, C. Ledig, V. Newcombe, J. Simpson, A. Kane, D. Menon, A. Nori, A. Criminisi, D. Rueckert, et al., Unsupervised domain adaptation in brain lesion segmentation with adversarial networks, in: International conference on information processing in medical imaging, Springer, pp. 597–609.
- [7] E. Tzeng, J. Hoffman, K. Saenko, T. Darrell, Adversarial discriminative domain adaptation, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 7167–7176.
- [8] J. Wu, R. Gu, G. Dong, G. Wang, S. Zhang, FPL-UDA: Filtered pseudo label-based unsupervised cross-modality adaptation for vestibular schwannoma segmentation, in: 2022 IEEE 19th International Symposium on Biomedical Imaging (ISBI), IEEE, pp. 1–5.
- [9] R. Gu, J. Zhang, G. Wang, W. Lei, T. Song, X. Zhang, K. Li, S. Zhang, Contrastive semi-supervised learning for domain adaptive segmentation across similar anatomical structures, IEEE Transactions on Medical Imaging (2022) 1–1.
- [10] C. Chen, Q. Dou, H. Chen, P.-A. Heng, Semantic-aware generative adversarial nets for unsupervised domain adaptation in chest x-ray segmentation, in: Y. Shi, H.-I. Suk, M. Liu (Eds.), Machine Learning in Medical Imaging, Springer International Publishing, Cham, 2018, pp. 143–151.
- [11] J. Wang, C. Lan, C. Liu, Y. Ouyang, T. Qin, W. Lu, Y. Chen, W. Zeng, P. Yu, Generalizing to unseen domains: A survey on domain generalization, IEEE Transactions on Knowledge and Data Engineering (2022).
- [12] Q. Dou, D. Coelho de Castro, K. Kamnitsas, B. Glocker, Domain generalization via model-agnostic learning of semantic features, Advances in Neural Information Processing Systems 32 (2019).
- [13] R. Gu, J. Zhang, R. Huang, W. Lei, G. Wang, S. Zhang, Domain composition and attention for unseen-domain generalizable medical image segmentation, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer, pp. 241–250.
- [14] L. Zhang, X. Wang, D. Yang, T. Sanford, S. Harmon, B. Turkbey, B. J. Wood, H. Roth, A. Myronenko, D. Xu, et al., Generalizing deep learning for medical image segmentation to unseen domains via deep stacked transformation, IEEE Transactions on Medical Imaging 39 (2020) 2531–2540.
- [15] S. Wang, L. Yu, K. Li, X. Yang, C.-W. Fu, P.-A. Heng, Dofe: Domain-oriented feature embedding for generalizable fundus image segmentation on unseen datasets, IEEE Transactions on Medical Imaging (2020).
- [16] L. Tran, X. Yin, X. Liu, Disentangled representation learning gan for pose-invariant face recognition, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 1415–1424.
- [17] J. Yang, N. C. Dvornek, F. Zhang, J. Chapiro, M. Lin, J. S. Duncan, Unsupervised domain adaptation via disentangled representations: Application to cross-modality liver segmentation, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer, pp. 255–263.
- [18] C. Pei, F. Wu, L. Huang, X. Zhuang, Disentangle domain features for cross-modality cardiac image segmentation, Medical Image Analysis 71 (2021) 102078.
- [19] A. Chartsias, T. Joyce, G. Papanastasiou, S. Semple, M. Williams, D. E. Newby, R. Dharmakumar, S. A. Tsaftaris, Disentangled representation learning in cardiac image analysis, Medical Image Analysis 58 (2019) 101535.

- [20] H. Li, S. J. Pan, S. Wang, A. C. Kot, Domain generalization with adversarial feature learning, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 5400–5409.
- [21] D. Li, Y. Yang, Y.-Z. Song, T. Hospedales, Learning to generalize: Meta-learning for domain generalization, in: Proceedings of the AAAI Conference on Artificial Intelligence, volume 32.
- [22] Q. Liu, Q. Dou, P.-A. Heng, Shape-aware meta-learning for generalizing prostate mri segmentation to unseen domains, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer, pp. 475–485.
- [23] Q. Liu, C. Chen, J. Qin, Q. Dou, P.-A. Heng, FeddG: Federated domain generalization on medical image segmentation via episodic learning in continuous frequency space, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 1013–1023.
- [24] R. H. Fick, A. Moshayedi, G. Roy, J. Dedieu, S. Petit, S. B. Hadj, Domain-specific cycle-gan augmentation improves domain generalizability for mitosis detection, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer, pp. 40–47.
- [25] J.-Y. Zhu, T. Park, P. Isola, A. A. Efros, Unpaired image-to-image translation using cycle-consistent adversarial networks, in: Proceedings of the IEEE international conference on computer vision, pp. 2223–2232.
- [26] C. Li, X. Lin, Y. Mao, W. Lin, Q. Qi, X. Ding, Y. Huang, D. Liang, Y. Yu, Domain generalization on medical imaging classification using episodic training with task augmentation, *Computers in Biology and Medicine* 141 (2022) 105144.
- [27] Z. Zhou, L. Qi, Y. Shi, Generalizable medical image segmentation via random amplitude mixup and domain specific image restoration, in: Proceedings of the European conference on computer vision (ECCV).
- [28] K. Muandet, D. Balduzzi, B. Schölkopf, Domain generalization via invariant feature representation, in: International Conference on Machine Learning, PMLR, pp. 10–18.
- [29] Y. Li, M. Gong, X. Tian, T. Liu, D. Tao, Domain generalization via conditional invariant representations, in: Proceedings of the AAAI conference on artificial intelligence, volume 32.
- [30] S. Hu, Z. Liao, J. Zhang, Y. Xia, Domain and content adaptive convolution for domain generalization in medical image segmentation, *IEEE Transactions on Medical Imaging* (2022) 1–1.
- [31] Y. Bengio, A. Courville, P. Vincent, Representation learning: A review and new perspectives, *IEEE transactions on pattern analysis and machine intelligence* 35 (2013) 1798–1828.
- [32] L. A. Gatys, A. S. Ecker, M. Bethge, Image style transfer using convolutional neural networks, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 2414–2423.
- [33] Q. Meng, N. Pawlowski, D. Rueckert, B. Kainz, Representation disentanglement for multi-task learning with application to fetal ultrasound, in: *Smart Ultrasound Imaging and Perinatal, Preterm and Paediatric Image Analysis*, Springer, 2019, pp. 47–55.
- [34] X. Xie, J. Chen, Y. Li, L. Shen, K. Ma, Y. Zheng, MI²GAN: Generative adversarial network for medical image domain adaptation using mutual information constraint, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer, pp. 516–525.
- [35] M. Ning, C. Bian, D. Wei, S. Yu, C. Yuan, Y. Wang, Y. Guo, K. Ma, Y. Zheng, A new bidirectional unsupervised domain adaptation segmentation framework, in: International Conference on Information Processing in Medical Imaging, Springer, pp. 492–503.
- [36] R. Hadsell, S. Chopra, Y. LeCun, Dimensionality reduction by learning an invariant mapping, in: 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, volume 2, IEEE, pp. 1735–1742.
- [37] K. He, H. Fan, Y. Wu, S. Xie, R. Girshick, Momentum contrast for unsupervised visual representation learning, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pp. 9729–9738.
- [38] T. Chen, S. Kornblith, M. Norouzi, G. Hinton, A simple framework for contrastive learning of visual representations, in: International conference on machine learning, PMLR, pp. 1597–1607.
- [39] G. Kang, L. Jiang, Y. Yang, A. G. Hauptmann, Contrastive adaptation network for unsupervised domain adaptation, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 4893–4902.
- [40] K. Chaitanya, E. Erdil, N. Karani, E. Konukoglu, Contrastive learning of global and local features for medical image segmentation with limited annotations, *Advances in Neural Information Processing Systems* 33 (2020).
- [41] W. Lei, W. Xu, R. Gu, H. Fu, S. Zhang, S. Zhang, G. Wang, Contrastive learning of relative position regression for one-shot object localization in 3d medical images, in: International Conference on Medical Image Computing and Computer-Assisted Intervention, Springer, pp. 155–165.
- [42] C. You, R. Zhao, L. H. Staib, J. S. Duncan, Momentum contrastive voxel-wise representation learning for semi-supervised volumetric medical image segmentation, in: *Medical Image Computing and Computer Assisted Intervention – MICCAI 2022*, Springer Nature Switzerland, Cham, 2022, pp. 639–652.
- [43] D. P. Kingma, M. Welling, Auto-encoding variational bayes, in: International Conference on Learning Representations.
- [44] X. Huang, M.-Y. Liu, S. Belongie, J. Kautz, Multimodal unsupervised image-to-image translation, in: Proceedings of the European conference on computer vision (ECCV), pp. 172–189.
- [45] A. v. d. Oord, Y. Li, O. Vinyals, Representation learning with contrastive predictive coding, *arXiv preprint arXiv:1807.03748* (2018).
- [46] X. Wang, R. Zhang, C. Shen, T. Kong, L. Li, Dense contrastive learning for self-supervised visual pre-training, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 3024–3033.
- [47] W. Liao, J. He, X. Luo, M. Wu, Y. Shen, C. Li, J. Xiao, G. Wang, N. Chen, Automatic delineation of gross tumor volume based on magnetic resonance imaging by performing a novel semi-supervised learning framework in nasopharyngeal carcinoma, *International Journal of Radiation Oncology* Biology* Physics* (2022).
- [48] J. Sivaswamy, S. Krishnadas, A. Chakravarty, G. Joshi, A. S. Tabish, et al., A comprehensive retinal image dataset for the assessment of glaucoma from the optic nerve head analysis, *JSM Biomedical Imaging Data Papers* 2 (2015) 1004.
- [49] F. Fumero, S. Alayón, J. L. Sanchez, J. Sigut, M. Gonzalez-Hernandez, Rim-one: An open retinal image database for optic nerve evaluation, in: 2011 24th international symposium on computer-based medical systems (CBMS), IEEE, pp. 1–6.
- [50] J. I. Orlando, H. Fu, J. B. Breda, K. van Keer, D. R. Bathula, A. Diaz-Pinto, R. Fang, P.-A. Heng, J. Kim, J. Lee, et al., Refuge challenge: A unified framework for evaluating automated methods for glaucoma assessment from fundus photographs, *Medical Image Analysis* 59 (2020) 101570.
- [51] E. Jang, S. Gu, B. Poole, Categorical reparameterization with gumbel-softmax, in: 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24–26, 2017, Conference Track Proceedings.
- [52] A. Chartsias, G. Papanastasiou, C. Wang, S. Semple, D. E. Newby, R. Dharmakumar, S. A. Tsaftaris, Disentangle, align and fuse for multimodal and semi-supervised image segmentation, *IEEE transactions on medical imaging* 40 (2020) 781–792.