

An Attentive-based Generative Model for Medical Image Synthesis

Jiayuan Wang¹, Q. M. Jonathan Wu^{1*} and Farhad Pourpanah²

¹ Centre for Computer Vision and Deep Learning, Department of Electrical and Computer Engineering, University of Windsor, Canada.

²Department of Electrical and Computer Engineering, Queens University, Canada.

*Corresponding author(s). E-mail(s): jwu@uwindsor.ca;
Contributing authors: wang621@uwindsor.ca;
farhad.086@gmail.com;

Abstract

Magnetic resonance (MR) and computer tomography (CT) imaging are valuable tools for diagnosing diseases and planning treatment. However, limitations such as radiation exposure and cost can restrict access to certain imaging modalities. To address this issue, medical image synthesis can generate one modality from another, but many existing models struggle with high-quality image synthesis when multiple slices are present in the dataset. This study proposes an attention-based dual contrast generative model, called ADC-cycleGAN, which can synthesize medical images from unpaired data with multiple slices. The model integrates a dual contrast loss term with the CycleGAN loss to ensure that the synthesized images are distinguishable from the source domain. Additionally, an attention mechanism is incorporated into the generators to extract informative features from both channel and spatial domains. To improve performance when dealing with multiple slices, the \mathbf{K} -means algorithm is used to cluster the dataset into \mathbf{K} groups, and each group is used to train a separate ADC-cycleGAN. Experimental results demonstrate that the proposed ADC-cycleGAN model produces comparable samples to other state-of-the-art generative models, achieving the highest PSNR and SSIM values of 19.04385 and 0.68551, respectively. We publish the code at <https://github.com/JiayuanWang-JW/ADC-cycleGAN>.

Keywords: CycleGAN, attention mechanism, deep learning, medical image synthesis, unpaired data

1 Introduction

Medical image analysis plays an important role in many clinical applications. Two types of neuroimaging techniques, including magnetic resonance (MR) and computer tomography (CT), are being widely used to diagnose various diseases and treatment planning. These modalities provide mutually-complementary information. MR images provide excellent soft tissue contrast and have no ionizing radiation, while CT images are suitable for bony structure and chest analysis [1, 2]. Medical image acquisition is a crucial step for medical image analysis. However, obtaining both modalities is a challenging issue due to multiple factors such as the unavailability of certain modalities, radiation dose, or high cost. Therefore, it is important to develop medical image synthesis models to generate one modality from another [3, 4].

Medical image synthesis methods learn a model to transfer knowledge from one modality, i.e., source domain, into another, i.e., target domain, without utilizing extra annotations from the target domain [5]. In other words, it learns a mapping function to map images from the source domain to the target domain. The learning can be done in a supervised manner using paired data [6], i.e., pairs of MR and CT images belong to the same patient and are perfectly registered, or unsupervised manner using unpaired data [7, 8]. Since CT and MR images have different structures, learning a direct mapping function between them is a challenging issue. Thus, the mapping function has to be complex and highly non-linear to bridge the structural differences between the two modalities.

Early medical image synthesis methods are based on segmentation and atlas [9]. Segmentation methods, first, segment an image from the source domain into several tissue classes and then synthesize the corresponding image in the target domain by intensity-filling of each class [10]. In contrast, atlas methods, first, register each image from the source domain into its corresponding atlas via a transformation and then apply the registration to the target domain atlas to synthesize the corresponding sample in the target domain [11]. The quality of the synthesized images by these methods relies on the segmentation and atlas quality.

Recently, convolutional neural networks (CNNs) methods have shown remarkable results in medical areas, such as disease diagnosis [12, 13], especially in synthesizing medical image tasks due to their ability in extracting task-specific features. For example, Li et al. [14] developed a CNNs-based model to use MR images and generate the corresponding positron emission tomography (PET) image for the same subject. Huang et al. [15] proposed a weakly-supervised joint convolutional sparse coding model to simultaneously

conduct super-resolution and medical image synthesis tasks. Zhao et al. [16] designed a CNN architecture to learn a mapping between two modalities utilizing imperfect registered CT-MR pairs.

Generative adversarial networks (GANs) [17] methods have produced promising results in synthesizing medical images. Nie et al. [18] integrated adversarial training strategy into a fully convolutional network (FCN) to model non-linear mapping between two modalities. Dalmaz et al. [19] introduced a novel GAN-based method that leverages the contextual sensitivity of vision transformers, the precision of convolutional operators, and the realism of adversarial learning to improve image generation. However, these generative models need a large number of perfectly registered paired data, which is a challenging issue since paired data are typically scarce. To alleviate the paired data restriction, many studies adopted CycleGAN [20] structure to convert medical image synthesis task into image-to-image translation that can learn from unpaired data. However, CycleGAN cannot perform well in transferring complex texture domain such as CT-MR, and it needs additional term(s) to learn a better mapping function between two modalities and consequently enhance the quality of the synthesized images in the target domain. For example, in [21], an adversarial learning model for synthesizing Ki-67-stained images from H&E-stained images from unpaired data has been proposed. This model attempts to preserve the structural details of the synthesized images by adopting the structural similarity constraint and skip connection. Huo et al. [22] introduced an end-to-end synthetic segmentation model to perform segmentation in the target domain without having access to any manual labels. Chen et al. [3] presented a one-shot generative model for MR image segmentation that utilizes unpaired data in addition to a single paired CT-MR dataset. This model consists of two networks, which are cross-modality image synthesis and MR image segmentation, that are jointly trained.

Moreover, several studies [5, 23–27] integrated attention mechanisms into the model to focus on the most important regions of the image. Studies [23–25] integrated attention mechanism into a conditional GAN to expand the receptive field and extract richer contextual dependencies. In [26], a difficulty-aware attention mechanism that considers the structural information in order to handle hard samples or regions. Tomar et al. [5] introduced a self-attention mechanism for attending various structures of the organ by leveraging an auxiliary semantic segmentation information. Yang et al. [27] incorporated self-attention mechanism into the generators for modelling long-range spatial dependencies in the synthesized images.

In our previous work [28], we proposed DC-cycleGAN, which is a bidirectional generative model, for medical image synthesis. We introduced a new loss term, called dual contrast (DC) loss, to enhance the model performance. DC loss locates the synthesized images far away from the samples of the source domain. To accomplish this, the DC loss uses the samples from the source domain as negative samples. We evaluated the performance of our model using

100 samples selected from a dataset proposed by Han et al. [29]. However, DC-cycleGAN and other methods produce unstable results as this dataset contains various slices. Abu et al. [30] alleviated this problem using auxiliary samples.

In this paper, we propose an *attention-based dual contrast CycleGAN* (ADC-cycleGAN) to further improve the performance of the DC-cycleGAN model and address the above-mentioned limitation of existing methods. The main contributions of our study are as follows:

1. Convolutional block attention module (CBAM) [31] is integrated into the DC-cycleGAN structure, namely ADC-cycleGAN, to extract more informative features from both channel and space dimensions in synthesizing medical images.
2. To generate high-quality images from datasets with multiple structures. To achieve this, the K -means algorithm is employed to cluster the training dataset into K groups and then each group is used to train an ADC-cycleGAN model, i.e., K models are trained. Using the clustering algorithm reduces the complexity of the dataset and alleviates the generator collapse.
3. To evaluate the performance of our proposed bidirectional medical image synthesis method with baseline and other state-of-the-art methods.

This paper contains five sections. Section 2 reviews the existing methods. Section 3 presents the proposed model. Section 4 provides the experimental results and ablation studies. Finally, the concluding remarks and future research directions are presented in Section 5.

2 Related works

This section reviews recent advances in medical image synthesis and attention mechanism.

2.1 Medical image synthesis

Medical image synthesis is an active area of research. On one hand, it aims to reduce time, labor, and cost [32]. On the other hand, some patients have metal devices in their bodies in which they can not scan the MR images. Medical image synthesis techniques can be broadly categorized into traditional and deep learning (DL)-based methods. Traditional methods learn a mapping function between similar patches from the two domains. This category can be grouped into atlas- [33–35] and segmentation- [10, 36–38] methods. Izquierdo et al. [36] first segmented the MR images into 6 tissue classes, and then uses a diffeomorphic approach to non-rigidly co-register. In the same way, the anatomical MR data for new subjects is co-registered with the template. Finally, the inverse transformations were applied to synthesize CT scans. Burgos et al. [39] generated CT scans and attenuation maps to enhance the attenuation correction for PET/MR scanners. As such, CT scans are synthesized using a multi-atlas information propagation scheme, in which a local

image similarity measure is used to locally match the MRI-derived patient's morphology to a database of MRI/CT pairs.

DL methods can be classified into Auto-encoder(AE), GAN, and U-net [32]. DEDIS [40], which is a deep encoder-decoder image synthesizer, performs MR image translation into different modalities, i.e., synthesizes T2 from T1 and DWI from T2. DEDIS is fast, requires a lower computational cost, and can produce comparable results as compared with the traditional methods. Hi-Net [41] solves the missing modality problem in medical imaging by learning a mapping function from the multi-modal domains to the target domain. It consists of three components, including a modality-specific network that learns the features of each modality, a multi-modal fusion network that learns common latent features of multi-modal data, and a multi-modal synthesis network that combines the learned latent features with hierarchical features from each modality to synthesize target images. Auto-GAN [42] is a self-supervised AE structure that obtains target-modality-specific information for the generator to synthesize missing MR image modality from available modalities. SkrGAN [43] integrates a sketch prior constraint into the GAN to synthesize high-quality medical images. It also embeds a sketch-based representation using a color render mapping technique. In another study [44], U-Net is combined with an adversarial training strategy to synthesize 2D PET from MR slices.

A number of models based on CycleGAN have been introduced for medical image synthesis from unpaired data. For example, UC-GAN [45] integrates U-Net into the CycleGAN structure for synthesizing CT scans from MR images. SC-cycleGAN [46] integrates structure-consistency loss with spectral normalization and self-attention mechanism for generating CT from MR images. Nice-GAN [47] utilizes the discriminator's encoding capability to improve the quality of generated images in the target domain. To ensure a stable training process during the adversarial min-max game, a decoupled training strategy is developed. This strategy prevents any training inconsistency and enables the encoder to train effectively by maximizing the loss and keeping it frozen otherwise.

In [48], the CBAM is integrated into β -cycleGAN to focus on the most important channel and spatial features. RegGAN [49] is based on the theory of "loss-correction" that is proposed for translating MR T1 to T2. Although RegGAN can produce promising results, it is a single-direction synthesis model. UGATIT [50] is an unsupervised method that integrates an attention module and a normalization function into its structure.

2.2 Attention

The role of attention in human perception is crucial. The human visual system focuses on the most salient parts of an image to explore the visual structure instead of processing the whole scene at once [51]. Inspired by the human visual system, attention processing has been incorporated into DL to enhance the model's performance [52, 53]. The attention mechanism is a crucial component in computer vision (CV), applied in various application domains such as text

classification [54], machine translation [55] and image classification [56], just to name a few.

The attention mechanism in CV can be broadly divided into four categories: channel attention to focus on “what,” spatial attention to focus on “where,” temporal attention to focus on “when,” and branch attention to focus on “which” [57]. Meantime, several studies have explored the combinations of two or more of them to enhance the model performance, e.g., spatial & temporal, and channel & spatial. The channel & spatial attention simultaneously take advantage of both channel and spatial attention to focus on important objects and regions of the images. The representative models of this category include CBAM [31], dual attention [53], and triplet attention [58]. Among them, CBAM is a simple and effective model that considers both channel and spatial domains. It can be integrated into CNN-based architectures to perform end-to-end training with negligible overhead [59]. Therefore, we integrated CBAM into our generators to extract features from both channel and space domains. We discuss the CBAM structure in Section 3.3.

3 Proposed model

This section introduces ADC-cycleGAN, a method for synthesizing medical images from unpaired data. We begin by formulating the problem and providing an overview of our approach. Next, we present specific instantiations of our method. To aid the reader, we list common symbols along with their corresponding descriptions in Table 1.

Table 1: Table of common symbols.

Symbol	Description
x	real CT image sample
y	real MR image sample
$G(x)$ & \hat{y}	synthesized MR images
$F(y)$ & \hat{x}	synthesized CT images
$D(x)$	discriminator to identify real CT images from the synthetic or randomly selected image from the source domain
$D(y)$	discriminator to identify real MR images from the synthetic or randomly selected image from the source domain
λ	weight of cycle consistency loss in the whole loss functions
β	weight of dual contrast loss in the whole loss functions

3.1 Problem formulation

Assume $X = \{x_i\}_{i=1}^N$ indicates the set of CT images and $Y = \{y_j\}_{j=1}^M$ represents the set of MR images.

The objective is to train a bidirectional projection function between CT and MR images using unpaired data. To accomplish this, ADC-cycleGAN employs two generators: $G : x \rightarrow \hat{y}$ and $F : y \rightarrow \hat{x}$, to learn the CT-to-MR and MR-to-CT mappings, respectively. Here, $\hat{y} = G(x)$ and $\hat{x} = F(y)$ denote the

synthesized MR and CT images, respectively. Additionally, two discriminators D_Y and D_X are used to differentiate between real MR and CT images and synthetic or randomly selected images from the source domain.

3.2 Model overview

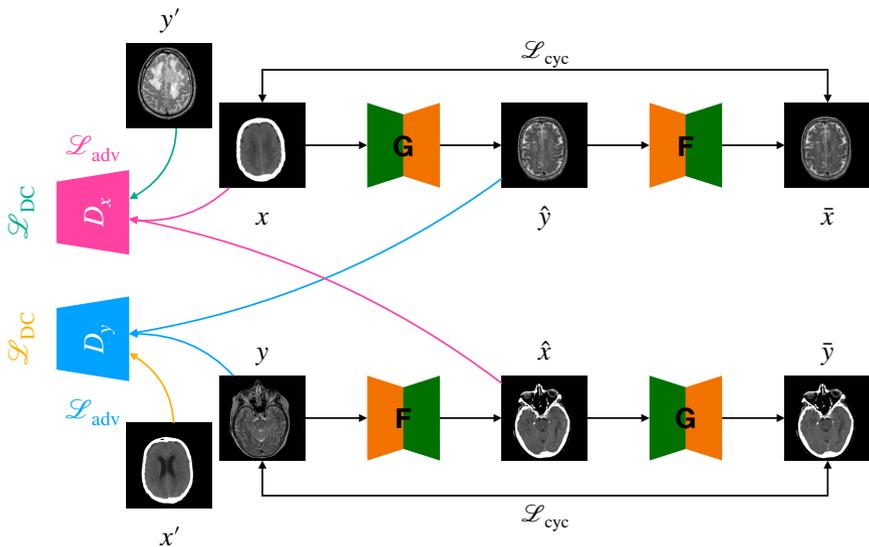


Fig. 1: The structure of ADC-cycleGAN

First, the training set is clustered into a number of groups using K -means algorithm. Then, a model, i.e., ADC-cycleGAN (see Fig. 1), is trained based on samples of each group. Fig. 1 illustrates the overview of ADC-cycleGAN. It is a bidirectional generative model that combines CycleGAN with dual contrast loss [28] to synthesize high-quality images in the target domain. Each component of ADC-cycleGAN is discussed in detail as follows.

3.2.1 CycleGAN

CycleGAN, which is originally proposed by Zhu et al. [20] for unpaired image-to-image translation, consists of a GAN [17] and cycle consistency loss. The GAN is composed of a generator $G : X \rightarrow Y$ to synthesize images in the target domain by transferring knowledge from the source domain, and a discriminator D_Y to identify real images in the target domain from the synthesized ones. These modules play a two-player game in which the generator forces the discriminator to enhance its distinguishing ability, while the discriminator forces the generator to synthesize more realistic images. In another word, the generator learns a mapping between source and target domains $X \rightarrow Y$, and

the discriminator's output is a probability that indicates the input image is real. The generator and discriminator play a min-max game to optimize using adversarial loss function as:

$$\mathcal{L}_{\text{GAN}}(G, D_Y, X, Y) = \mathbb{E}_{y \sim p_{\text{data}}(y)} [\log D_Y(y)] + \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log (1 - D_Y(G(x)))] . \quad (1)$$

where x and y are real images from the source and target domains, respectively, $G(x)$ is the synthesized image in the target domain.

Although adversarial learning can learn a mapping between the source and target domains, it has the limitation of mapping input images to any arbitrary location in the target domain. Consequently, adversarial learning does not guarantee that the learned function can accurately map input x_i to its intended output y_i . Furthermore, acquiring paired image datasets can be challenging in practice. To overcome these issues, it is essential to reconstruct the synthesized images into their target domain via a cycle consistency loss. To accomplish this, an additional generator $F : Y \rightarrow X$ is necessary for reconstructing real source domain images. A discriminator D_X is also required to distinguish real images from the synthesized ones in the target domain

$$\mathcal{L}_{\text{cycle}}(G, F) = \mathbb{E}_{x \sim p_{\text{data}}(x)} [\|F(G(x)) - x\|_1] + \mathbb{E}_{y \sim p_{\text{data}}(y)} [\|G(F(y)) - y\|_1], \quad (2)$$

where $F(G(x))$ and $G(F(y))$ indicate the reconstructed CT and MR images, respectively.

To summarize, the CycleGAN is comprised of four distinct networks: two generators (G and F) and two discriminators (D_X and D_Y). The object function for the CycleGAN is as follows:

$$\mathcal{L}(G, F, D_X, D_Y) = \mathcal{L}_{\text{GAN}}(G, D_Y, X, Y) + \mathcal{L}_{\text{GAN}}(F, D_X, Y, X) + \lambda \mathcal{L}_{\text{cycle}}(G, F), \quad (3)$$

where λ indicates the weight of cycle consistency loss.

The aim to solve is:

$$G^*, F^* = \arg \min_{G, F} \max_{D_X, D_Y} \mathcal{L}(G, F, D_X, D_Y). \quad (4)$$

3.2.2 Dual contrast

Since CT and MR images have different structures and there is no constraints between real source and synthesized images [27], CycleGAN alone can not synthesize high-quality images in the target domain. To mitigate this issue, an additional loss term, known as *dual contrast* (DC) loss, is introduced into discriminators. The DC loss utilizes samples from the source domain as negative

samples (x' and y') to prompt the model to generate images that are distinct from the source domain (refer to Fig. 1). Instead of discriminating between real images (y and x) and synthesized images (\hat{y} and \hat{x}), the discriminators are tasked with distinguishing real images in the target domain from both synthesized images and randomly selected images from the source domain. In other words, the real images are assigned to class 1, whereas synthesized images and randomly selected samples from the source domain are assigned to class 0. The dual contrast loss function, as follows:

$$\mathcal{L}_{\text{DC}}(D_Y, X, Y) = \mathbb{E}_{x' \sim p_{\text{data}}(x')} [\log(1 - D_Y(x'))]. \quad (5)$$

3.2.3 Overall objective function

By adding DC loss into (3), the final objective function can be written as:

$$\begin{aligned} \mathcal{L}(G, F, D_X, D_Y) = & \mathcal{L}_{\text{GAN}}(G, D_Y, X, Y) + \\ & \mathcal{L}_{\text{GAN}}(F, D_X, Y, X) + \beta \mathcal{L}_{\text{DC}}(D_Y, X, Y) + \\ & \beta \mathcal{L}_{\text{DC}}(D_X, Y, X) + \lambda \mathcal{L}_{\text{cycle}}(G, F), \end{aligned} \quad (6)$$

where λ and β indicate the weight of cycle consistency loss and dual contrast loss in the whole loss functions, respectively.

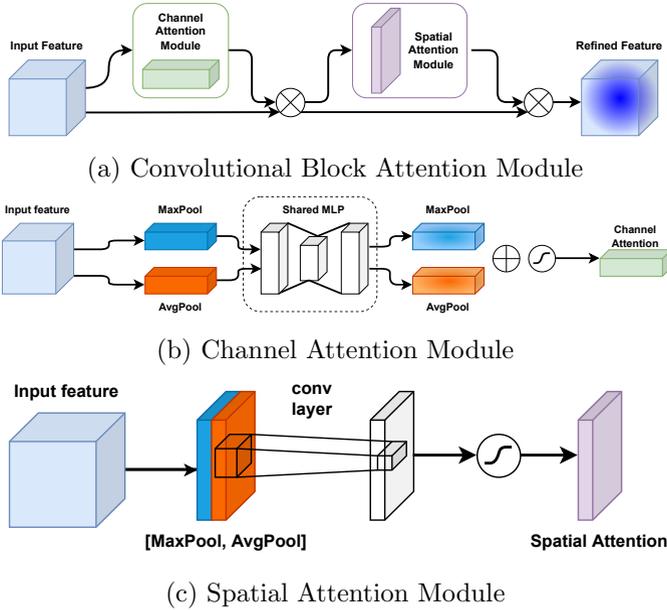
3.3 Attention mechanism

Channel attention, spatial attention, and their combinations, i.e., mixed attention, are often used in computer vision tasks. Channel attention focuses on “*what*” are important channels to pay attention in a feature map. It adds a different weight to each channel, and the weight with a high value is more correlated. In contrast, spatial attention focuses on “*where*” is an informative part to pay attention by learning a weight on a 2D feature map. While, mixed attention combines channel and spatial attentions. CBAM [31] is the most representative model of this category that has been widely integrated into the CNN-based models to improve performance with negligible overheads. It is presented in detail, as follows.

Assume $\mathbb{R}^{C \times H \times W}$ represents an input feature map, where C , H and W indicate channel, high and wide, respectively. As shown in Fig. 2 (a), CBAM consists of two sub-modules, including channel and spatial, that sequentially infers a 1D channel attention map $\mathbf{M}_c \in \mathbb{R}^{C \times 1 \times 1}$ and a 2D spatial attention map $\mathbf{M}_s \in \mathbb{R}^{1 \times H \times W}$. The CBAM whole process can be expressed as:

$$\begin{aligned} \mathbf{F}' &= \mathbf{M}_c(\mathbf{F}) \otimes \mathbf{F}, \\ \mathbf{F}'' &= \mathbf{M}_s(\mathbf{F}') \otimes \mathbf{F}', \end{aligned} \quad (7)$$

where \otimes is the element-wise multiplication, \mathbf{F}' indicates a feature map after combining the input feature \mathbf{F} and channel attention map $\mathbf{M}_c(\mathbf{F})$, and \mathbf{F}'' is the final refined feature.

**Fig. 2:** Structure of CBAM and sub modules

It first uses max-pooling and average-pooling operations to produce two spatial context descriptors $\mathbf{F}_{\text{avg}}^c$ and $\mathbf{F}_{\text{max}}^c$, and then, the channel attention map $\mathbf{M}_c \in \mathbb{R}^{C \times 1 \times 1}$ is produced by feeding the descriptors into a shared network, i.e., an MLP with one hidden layer, as:

$$\begin{aligned} \mathbf{M}_c(\mathbf{F}) &= \sigma(MLP(\text{AvgPool}(\mathbf{F})) + MLP(\text{MaxPool}(\mathbf{F}))) \\ &= \sigma(\mathbf{W}_1(\mathbf{W}_0(\mathbf{F}_{\text{avg}}^c)) + \mathbf{W}_1(\mathbf{W}_0(\mathbf{F}_{\text{max}}^c))), \end{aligned} \quad (8)$$

where σ is sigmoid activation function, and $\mathbf{W}_0 \in \mathbb{R}^{C/r \times C}$ and $\mathbf{W}_1 \in \mathbb{R}^{C \times C/r}$ are the weights of the shared network, where r indicates the reduction ratio which is set to 8 in this study.

Spatial attention module: Fig. 2 (c) depicts the spatial attention module, which utilizes the inter-spatial relationship of features to produce the spatial attention map. Unlike the channel attention module, this module applies max-pooling and average-pooling operations to the channel axis, resulting in two 2D maps $\mathbf{F}_{\text{avg}}^s \in \mathbb{R}^{1 \times H \times W}$ and $\mathbf{F}_{\text{max}}^s \in \mathbb{R}^{1 \times H \times W}$. These maps are concatenated to generate a feature descriptor. Then, a convolution layer is applied to produce a 2D spatial attention map $\mathbf{M}_s(\mathbf{F})$, as follows:

$$\begin{aligned} \mathbf{M}_s(\mathbf{F}) &= \sigma(f^{7 \times 7}([\text{AvgPool}(\mathbf{F}); \text{MaxPool}(\mathbf{F})])) \\ &= \sigma(f^{7 \times 7}([\mathbf{F}_{\text{avg}}^s; \mathbf{F}_{\text{max}}^s])) \end{aligned} \quad (9)$$

where σ indicates the sigmoid activation functions, $f^{7 \times 7}$ indicate a convolution layer with a 7 x 7 filter.

3.4 Model structure

This subsection discusses the structures of K -means algorithm, generators and discriminators in detail.

3.4.1 k-means algorithm

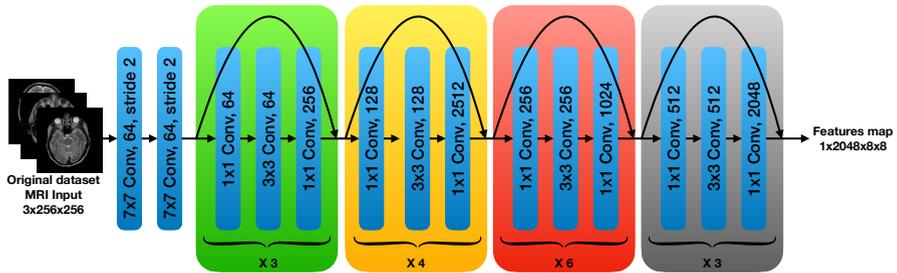


Fig. 3: The preprocessing of cluster.

To cluster the training set, we first use the Resnet50 [60] to extract features (see Fig. 3), and then, the ‘KMeans’ function from the sklearn package is used for clustering. We initialized the cluster centers randomly using the ‘KMeans’ function in the sklearn package. Additionally, we used the K -means algorithm with the same clustering results for all methods by clustering the training dataset into K groups. We repeated this procedure five times to ensure the consistency of the results and avoid dependency on the initial random seed. The KMeans function uses the Euclidean distance to continuously compute the distance between samples and the clusters’ center, and update their centroids by optimizing:

$$E = \sum_{i=1}^K \sum_{\mathbf{x} \in C_i} \|\mathbf{x} - \boldsymbol{\mu}_i\|_2^2, \quad (10)$$

where K is the number of clusters, $\boldsymbol{\mu}_i$ represents the center of the i -th cluster, and $\mathbf{x} \in C_i$ means x belong to the i -th cluster.

3.4.2 Generator

In this study, CBAM [31] is integrated into the generators to extract more informative features from channel and space domains (see Fig. 4). Each generator has one instance normalization (IN) convolutional layer with ReLU activation function, two stride-2 IN-convolutional layers with ReLU activation function, nine residual blocks, in which CBAM is added into the first eight

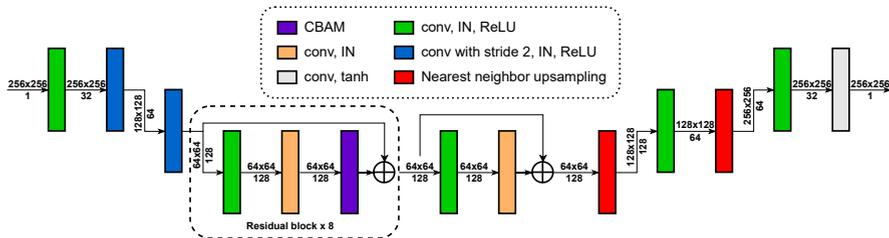


Fig. 4: The structure of generator integrated into the ADC-cycleGAN model.

blocks, two blocks of nearest neighbor up-sampling and IN-convolutional layers with ReLU activation function, followed by a convolutional layer with tanh activation function as output layer.

The vanilla CycleGAN utilizes mean absolute error (MAE) loss in generators. However, using MAE may result in synthesizing low-quality images that can not be recognized by human [28, 61]. To alleviate this issue, Snell et al. [61] showed that using SSIM instead of MAE improves the quality of the reconstructed images. This is mainly due to the considering luminance l , contrast c , and structure s , as follows:

$$SSIM(x_1, x_2) = \frac{(2\mu_{x_1}\mu_{x_2} + c_1)(2\sigma_{x_1x_2} + c_2)}{(\mu_{x_1}^2 + \mu_{x_2}^2 + c_1)(\sigma_{x_1}^2 + \sigma_{x_2}^2 + c_2)}, \quad (11)$$

where c_1 and c_2 and $c_3 = \frac{c_2}{2}$ are constant values, μ_{x_i} is the mean of the i -th image ($i = 1, 2$), σ_{x_i} is the standard deviation of the i -th image ($i = 1, 2$), and $\sigma_{x_1x_2}$ represents the covariance of x_1 and x_2 . while:

$$l(x_1, x_2) = \frac{2\mu_{x_1}\mu_{x_2} + c_1}{\mu_{x_1}^2 + \mu_{x_2}^2 + c_1}, \quad (12)$$

$$c(x_1, x_2) = \frac{2\sigma_{x_1}\sigma_{x_2} + c_2}{\sigma_{x_1}^2 + \sigma_{x_2}^2 + c_2}, \quad (13)$$

$$s(x_1, x_2) = \frac{\sigma_{x_1x_2} + c_3}{\sigma_{x_1}\sigma_{x_2} + c_3}. \quad (14)$$

Therefore, in this study, we use SSIM instead of MAE in the cycle consistency loss. The Eq. (2) can be written as:

$$\mathcal{L}_{\text{cycle}}(G, F) = (1 - SSIM(F(G(x)), x)) + (1 - SSIM(G(F(y)), y)). \quad (15)$$

3.4.3 Discriminator

Fig. 5 shows the structure of discriminators integrated into the ADC-cycleGAN structure. Each discriminator has a stride-2 convolutional layer and three stride-2 IN-convolutional layers with LReLU activation function, and the output layer is 94×94 overlapping patches for identifying whether the input image belongs

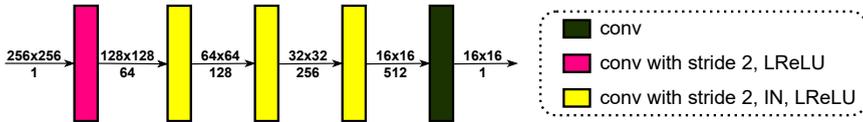


Fig. 5: The structure of discriminator integrated into the ADC-cycleGAN model.

to class 1 (real) or class 0 (synthesized or a random image from the source domain).

Unlike the conventional CycleGAN that uses mean squared error (MSE) loss in the discriminators, in this study, we use cross entropy (CE) because CE converges fast with a small back-propagation error as compared with MSE [62]. On one hand, MSE suffers from the problem of gradient vanishing problem in the output layer of networks [63]. In this study, we employed binary CE (BCE) loss at the output layer of the discriminators to differentiate between input images belonging to class 1 (real) and class 0 (synthesized or from the source domain). The BCE is defined as follows:

$$CE(t, y) = -\frac{1}{N} \sum_{i=1}^N t_i \cdot \log(y_i) + (1 - t_i) \cdot \log(1 - y_i), \quad (16)$$

where $N = 256$ due to our discriminator output has 16×16 patches, t_i is a label for i -th patch and y_i is the predicted probability for i -th patch.

3.5 Summary of proposed framework

During the learning phase, we first use the K -means algorithm to cluster the training set into K groups, and then, an ADC-cycleGAN model is trained using the samples of each group, i.e., totally K ADC-cycleGAN models is trained. During the test phase, each test sample is clustered based on the centroids of the k -means algorithm obtained by the training set, and then, the corresponding ADC-cycleGAN with respect to the cluster number is used to synthesize image.

4 Experimental results

This section conducts a number of experiments to evaluate the effectiveness of ADC-cycleGAN in synthesizing MR images from CT scans and vice versa from unpaired data and compare it with other medical image synthesis methods

such as CycleGAN [20], NiceGAN [47], UGATIT [50], RegGAN [49] and DC-cycleGAN [28]. Their codes are obtained from their official GitHub ^{1 2 3 4 5}. In this study, CycleGAN is used as the baseline model, because it has been widely used in generating medical images from unpaired data.

4.1 Evaluation indexes

We used three evaluation metrics, namely SSIM [64], MAE, and PSNR for performance evaluation and compression.

SSIM is a common evaluation metric that measures the similarity between two images (x_1 and x_2) by considering luminance l , contrast c , and structure s . Eq. (11) can be used to compute SSIM.

MAE measures the average absolute between two images. In other words, MAE computes the distance between real and synthesized images. It can be written as:

$$MAE = \frac{1}{HW} \sum_{i=1}^H \sum_{j=1}^W |x_1(i, j) - x_2(i, j)|, \quad (17)$$

where H and W are the high and wide of the images, respectively.

PSNR is another indicator that can be used to assess the quality of the synthesized image, which can be computed as:

$$PSNR = 10 \log_{10}(L/MSE), \quad (18)$$

where L is the dynamic range of the pixel values, and:

$$MSE = \frac{1}{HW} \sum_{i=1}^H \sum_{j=1}^W (x_1(i, j) - x_2(i, j))^2. \quad (19)$$

A large PSNR indicates that the two images are closer.

4.2 Dataset

Although the aim is medical image synthesis from unpaired data, it is required to use paired data for computing the evaluation metrics, i.e., SSIM, MAE and PSNR. To achieve this, the dataset introduced by Han et al. [29] is obtained from their project website⁶. This dataset includes 367 paired CT and MR images from multiple slices, and the size of each image is 512×256 . CT scans have some head frames due to Gamma Knife treatment. As the publicly available dataset we utilized only included images in png format, we were

¹<https://github.com/simontomaskarlsson/CycleGAN-Keras>

²<https://github.com/alpc91/NICE-GAN-pytorch>

³<https://github.com/taki0112/UGATIT>

⁴<https://github.com/Kid-Liet/Reg-GAN>

⁵<https://github.com/JiayuanWang-JW/DC-cycleGAN>

⁶<https://github.com/ChengBinJin/MRI-to-CT-DCNN-TensorFlow>

unable to follow the approach used by Han et al. [29] for calculating masks based on Hounsfield units and removing head frames. Instead, we manually removed the head frames by cropping the CT scans.

4.3 Parameters setting

ADC-cycleGAN has a similar structure to CycleGAN [20]. It consists of two generators G and F and two discriminators D_x and D_y . We followed the CycleGAN and set λ , batch size and the number of epochs to 10, 1 and 200, respectively. The number of clusters of K -means algorithm based on the ablation study is set to 4. To have a fair comparison, we follow our previous work for $\beta=0.5$. This value is obtained based on a sensitivity analysis. During the training phase, generators are updated five times followed by updating discriminators once.

We normalized all images in the range of -1 to 1 and resized them to 256×256 . In addition, 90% and 10% of the dataset are used for training and testing, respectively. In order to have a fair comparison, the K -means algorithm is used for all methods, i.e., for each method, we first used K -means algorithm to cluster training set into K groups, and then, each group is trained by a model. This procedure is repeated five times for each comparison method and used the mean value along with the standard deviation (SD) as the final result. The experiments were carried out on a server equipped with an Intel(R) Xeon(R) E5-2650 CPU and an Nvidia GTX 1080TI GPU.

4.4 Ablation study

In this section, we conduct ablation studies to show the effects of using different components of the proposed ADC-cycleGAN model. Two experiments are conducted. The first experiment studies the effect of a different number of clusters K and finds the best value, As such, K is varied from 2 to 5. Tables 2 and 3 show the results of ADC-cycleGAN with different number of clusters for MR and CT synthesis, respectively. ADC-cycleGAN with $K=4$ outperforms other cluster numbers. Therefore, for the rest of experiments K is set to 4.

Table 2: The ablation results for MR synthesis. “Mean (standard deviation)” for the different number of clusters.

# of clusters	MAE ↓	PSNR ↑	SSIM ↑
2	0.13360 (0.00783)	18.32465 (0.53110)	0.59565 (0.01087)
3	0.12458 (0.00995)	18.81969 (0.72182)	0.62591 (0.01589)
4	0.11566 (0.00746)	19.69240 (0.75123)	0.64330 (0.01975)
5	0.11897 (0.01354)	19.74868 (1.37066)	0.63929 (0.03150)

Table 3: The ablation results for CT synthesis. “Mean (standard deviation)” for the different number of clusters.

# of clusters	MAE ↓	PSNR ↑	SSIM ↑
2	0.12722 (0.00632)	16.74714 (0.81303)	0.69372 (0.01155)
3	0.11166 (0.00571)	17.58067 (0.60216)	0.71480 (0.01184)
4	0.10809 (0.00292)	17.95930 (0.52605)	0.71589 (0.00655)
5	0.12328 (0.01482)	17.52259 (1.75595)	0.69916 (0.02603)

The second experiment shows the impact of using DC loss, attention mechanism and clustering method in the ADC-cycleGAN structure. Four different combinations are used, as follows:

- CycleGAN (wo): the baseline model.
- CycleGAN (w): CycleGAN with clustering technique.
- A-cycleGAN (wo): CycleGAN with CBAM.
- A-cycleGAN (w): CycleGAN with CBAM and clustering technique.
- DC-cycleGAN (wo): CycleGAN with DC loss.
- DC-cycleGAN (w): CycleGAN with DC loss and clustering technique.
- ADC-cycleGAN (wo): CycleGAN with DC loss and CBAM.
- ADC-cycleGAN (w): CycleGAN with DC loss, CBAM and clustering technique.

Table 4: The ablation study for CT and MR synthesis. The mean (standard deviation) for different conditions. The notions “(w)” and “(wo)” indicate with or without K -means algorithm.

Method	MAE ↓	PSNR ↑	SSIM ↑
CycleGAN (wo)	0.13789 (0.01073)	16.47120 (0.57922)	0.64637 (0.00892)
CycleGAN (w)	0.12245 (0.01008)	17.70577 (0.79366)	0.66171 (0.01509)
A-cycleGAN (wo)	0.14080 (0.00989)	16.28680 (0.55052)	0.64795 (0.00870)
A-cycleGAN (w)	0.12537 (0.00847)	17.39279 (0.52072)	0.66055 (0.01228)
DC-cycleGAN (wo)	0.14157 (0.00646)	16.30498 (0.36758)	0.63576 (0.00617)
DC-cycleGAN (w)	0.11069 (0.00401)	18.80157 (0.35545)	0.67907 (0.00947)
ADC-cycleGAN (wo)	0.13622 (0.00235)	16.57089 (0.11158)	0.64359 (0.00412)
ADC-cycleGAN (w)	0.11005 (0.00450)	19.04385 (0.48771)	0.68551 (0.00849)

Table 4 presents the results of the second ablation study. As can be seen, ADC-cycleGAN outperforms all other combinations. This finding confirms that every component plays a crucial role in synthesizing both MR and CT images. Moreover, our proposed methods with attention mechanisms, i.e., ADC-cycleGAN (wo) and ADC-cycleGAN (w), performed better and more stable results as compared with those without attention mechanisms, i.e., DC-cycleGAN (wo) and DC-cycleGAN (w). However, the baseline models, i.e., CycleGAN (wo) and CycleGAN (w), perform slightly better than those with attention mechanisms, i.e., A-CycleGAN (wo) and A-CycleGAN (w). In addition, all methods that use the clustering mechanism, i.e., CycleGAN(w), A-cycleGAN(w), DC-cycleGAN(w), and ADC-cycleGAN(w) outperform those methods without the clustering algorithm, i.e., CycleGAN(wo), A-cycleGAN(wo), DC-cycleGAN(wo), and ADC-cycleGAN(wo). This proves the capability of the K -means algorithm in synthesizing high-quality images from datasets that contain images with various structures. However, the K -means algorithm reduces the stability of the model.

4.5 Comparison with other methods

This section compares the performance of ADC-cycleGAN with CycleGAN [20], NiceGAN [47], UGATIT [50], RegGAN [49] and DC-cycleGAN [28]. For all methods, the K -means clustering algorithm (with $K = 4$) is used to group the training samples into K clusters, and then, each group is trained by a model, i.e., $K = 4$ models are trained for each method. Note that, during the test phase, the centroids of the K -means algorithm obtained during training are used to cluster the test set. Then, for each test sample, the trained model with respect to its cluster is used to generate the image.

Table 5: The mean (standard deviation) values of the MR synthesis quality evaluation metrics for different methods.

Method	MAE ↓	PSNR ↑	SSIM ↑
CycleGAN	0.12636 (0.01343)	18.52153 (1.06237)	0.62686 (0.02044)
NiceGAN	0.12562 (0.00145)	18.34712 (0.12351)	0.62077 (0.00280)
UGATIT	0.10647 (0.00320)	19.92874 (0.20880)	0.64863 (0.00463)
RegGAN	0.09704 (0.00162)	20.45722 (0.18531)	0.66536 (0.00191)
DC-cycleGAN	0.11395 (0.00406)	19.85754 (0.38057)	0.64502 (0.01302)
ADC-cycleGAN	0.11080 (0.00571)	20.12068 (0.55074)	0.65568 (0.01215)

The quantitative results, i.e., MAE, PSNR, and SSIM, for MR and CT synthesis are shown in Tables 5 and 6, respectively. RegGAN produces the

Table 6: The mean (standard deviation) values of the CT synthesis quality evaluation metrics for different methods.

Method	MAE ↓	PSNR ↑	SSIM ↑
CycleGAN	0.11853 (0.00673)	16.89000 (0.52494)	0.69656 (0.00973)
NiceGAN	0.12742 (0.00120)	15.99856 (0.15009)	0.68270 (0.00175)
UGATIT	0.10825 (0.00328)	17.04600 (0.18058)	0.70392 (0.00612)
RegGAN	0.13940 (0.04498)	16.21307 (1.03060)	0.67735 (0.02595)
DC-cycleGAN	0.10742 (0.00395)	17.74560 (0.33032)	0.71312 (0.00592)
ADC-cycleGAN	0.10931 (0.00329)	17.96701 (0.42467)	0.71534 (0.00482)

Table 7: The results (Mean (standard deviation)) of bidirectional MR and CT synthesis.

Method	MAE ↓	PSNR ↑	SSIM ↑
CycleGAN	0.12244 (0.01008)	17.70576 (0.79366)	0.66171 (0.01508)
NiceGAN	0.12652 (0.00132)	17.17284 (0.13680)	0.65174 (0.00228)
UGATIT	0.10736 (0.00324)	18.48737 (0.19469)	0.67628 (0.00537)
RegGAN	0.11822 (0.02330)	18.33515 (0.60796)	0.67136 (0.01393)
DC-cycleGAN	0.11069 (0.00401)	18.80157 (0.35545)	0.67907 (0.00947)
ADC-cycleGAN	0.11005 (0.00450)	19.04385 (0.48771)	0.68551 (0.00849)

best results in synthesizing MR images from CT scans, and ADC-cycleGAN is ranked as the second-best method in terms of PSNR and SSIM, and the third method in term of MAE. In contrast, ADC-cycleGAN outperforms all methods in terms of PSNR and SSIM in generating CT scans from MR images, while DC-cycleGAN is ranked as the first method in term of MAE. However, RegGAN produces inferior results in synthesizing CT scans from MR images and it is not able to perform bidirectional learning, i.e., it is required to train two times, one for CT-to-MR and one for MR-to-CT. Overall, our proposed ADC-cycleGAN model is able to produce the best results in generating both MR and CT images, i.e., bidirectional learning, in terms of PSNR and SSIM, while UGATIT and ADC-cycleGAN perform similarly in terms of MAE (see Table 7). While our method can produce comparable results as compared with other models, there are several limitations that must be addressed. First, The effectiveness of ADC-cycleGAN and other methods is largely influenced by the quality of the training dataset. When there exist artifacts on the head frames

of CT scans, all methods produce inferior results. Second, our ADC-cycleGAN model still suffers from the model collapse issue that causes by the presence of various slices in the dataset. Although we mitigated this issue by employing the K -mean clustering algorithm, further research is required to find more effective solutions.

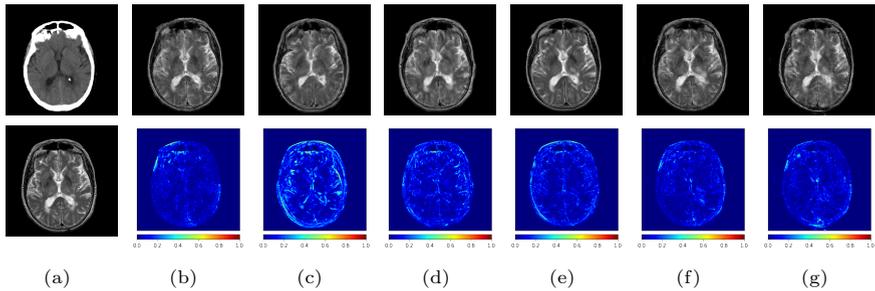


Fig. 6: Synthesized MR images along with absolute error maps between groundtruth and synthesized images by different methods. (a) Real image, (b) CycleGAN, (c) NiceGAN, (d) UGATIT, (e) RegGAN, (f) DC-cycleGAN, (g) ADC-cycleGAN

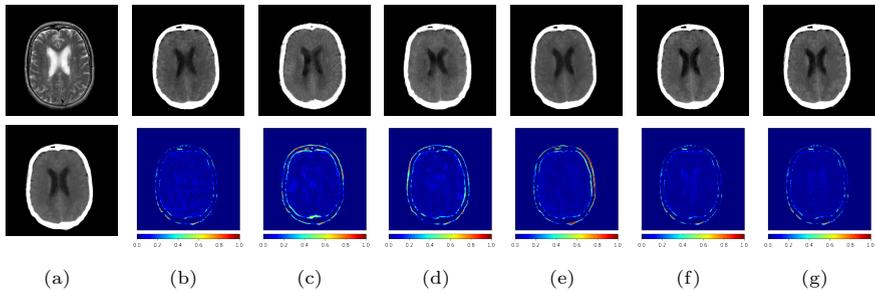


Fig. 7: Synthesized CT images and their corresponding absolute error maps for different methods. (a) Real image, (b) CycleGAN, (c) NiceGAN, (d) UGATIT, (e) RegGAN, (f) DC-cycleGAN, (g) ADC-cycleGAN

In addition, the synthesized CT images and their corresponding absolute error maps for different methods are shown in Figs. 6 and 7, respectively. As compared with other methods, the synthesized MR and CT images by ADC-cycleGAN demonstrate a higher level of fidelity to the real images. In particular, it can effectively capture the intricate details in soft tissue within the synthesized MR scans, and accurately replicate the edge structures in the synthesized CT scans. These results indicate the effectiveness of the proposed ADC-cycleGAN in synthesizing MR and CT images.

5 Conclusion

In this study, we proposed a bidirectional generative model based on CycleGAN with an attention mechanism for synthesizing medical images from unpaired data. We introduced a dual contrast loss that leverages samples from the source domain as negative samples, pushing the synthesized images further away from the source domain. To capture important features in both channel and space domains, we integrated CBAM into the generators. Additionally, to enable the model to synthesize high-quality images from datasets with varying slice numbers, we employed the K -means algorithm to cluster the training set into groups and trained a model for each group. The experimental results, along with the ablation study, demonstrate the effectiveness of the proposed method in synthesizing MR from CT scans and vice versa.

In our future research, we plan to improve the quality of synthesized images by developing novel structures for generators and discriminators, aiming to enhance model performance and stability. Additionally, we aim to utilize the synthesized images for solving other tasks such as segmentation. Furthermore, we plan to develop a new dataset for performance evaluation.

Declarations

The dataset analyzed during the current study are available in the Github repository, <https://github.com/ChengBinJin/MRI-to-CT-DCNN-TensorFlow>.

References

- [1] Xu, L., Zeng, X., Zhang, H., Li, W., Lei, J., Huang, Z.: Bpgan: Bidirectional ct-to-mri prediction using multi-generative multi-adversarial nets with spectral normalization and localization. *Neural Networks* **128**, 82–96 (2020)
- [2] Yang, H., Lu, X., Wang, S.-H., Lu, Z., Yao, J., Jiang, Y., Qian, P.: Synthesizing multi-contrast mr images via novel 3d conditional variational auto-encoding gan. *Mobile Networks and Applications* **26**(1), 415–424 (2021)
- [3] Chen, X., Lian, C., Wang, L., Deng, H., Fung, S.H., Nie, D., Thung, K.-H., Yap, P.-T., Gateno, J., Xia, J.J., *et al.*: One-shot generative adversarial learning for mri segmentation of craniomaxillofacial bony structures. *IEEE transactions on medical imaging* **39**(3), 787–796 (2019)
- [4] Lee, J.H., Han, I.H., Kim, D.H., Yu, S., Lee, I.S., Song, Y.S., Joo, S., Jin, C.-B., Kim, H.: Spine computed tomography to magnetic resonance image synthesis using generative adversarial networks: a preliminary study. *Journal of Korean Neurosurgical Society* **63**(3), 386–396 (2020)

- [5] Tomar, D., Lortkipanidze, M., Vray, G., Bozorgtabar, B., Thiran, J.-P.: Self-attentive spatial adaptive normalization for cross-modality domain adaptation. *IEEE Transactions on Medical Imaging* **40**(10), 2926–2938 (2021)
- [6] Mérida, I., Costes, N., Heckemann, R.A., Drzezga, A., Förster, S., Hammers, A.: Evaluation of several multi-atlas methods for pseudo-ct generation in brain mri-pet attenuation correction. In: 2015 IEEE 12th International Symposium on Biomedical Imaging (ISBI), pp. 1431–1434 (2015). IEEE
- [7] Lian, C., Li, X., Kong, L., Wang, J., Zhang, W., Huang, X., Wang, L.: Cocyclereg: Collaborative cycle-consistency method for multi-modal medical image registration. *Neurocomputing* (2022)
- [8] Li, X., Jia, M., Islam, M.T., Yu, L., Xing, L.: Self-supervised feature learning via exploiting multi-modal data for retinal disease diagnosis. *IEEE Transactions on Medical Imaging* **39**(12), 4023–4033 (2020)
- [9] Jiao, J., Namburete, A.I., Papageorghiou, A.T., Noble, J.A.: Self-supervised ultrasound to mri fetal brain image synthesis. *IEEE Transactions on Medical Imaging* **39**(12), 4413–4424 (2020)
- [10] Berker, Y., Franke, J., Salomon, A., Palmowski, M., Donker, H.C., Temur, Y., Mottaghy, F.M., Kuhl, C., Izquierdo-Garcia, D., Fayad, Z.A., *et al.*: Mri-based attenuation correction for hybrid pet/mri systems: a 4-class tissue segmentation technique using a combined ultrashort-echo-time/dixon mri sequence. *Journal of nuclear medicine* **53**(5), 796–804 (2012)
- [11] Sjölund, J., Forsberg, D., Andersson, M., Knutsson, H.: Generating patient specific pseudo-ct of the head from mr using atlas-based regression. *Physics in Medicine & Biology* **60**(2), 825 (2015)
- [12] Bhosale, Y.H., Patnaik, K.S.: Application of deep learning techniques in diagnosis of covid-19 (coronavirus): a systematic review. *Neural Processing Letters*, 1–53 (2022)
- [13] Bhosale, Y.H., Patnaik, K.S.: Puldi-covid: Chronic obstructive pulmonary (lung) diseases with covid-19 classification using ensemble deep convolutional neural network from chest x-ray images to minimize severity and mortality rates. *Biomedical Signal Processing and Control* **81**, 104445 (2023)
- [14] Li, R., Zhang, W., Suk, H.-I., Wang, L., Li, J., Shen, D., Ji, S.: Deep learning based imaging data completion for improved brain disease diagnosis. In: *International Conference on Medical Image Computing and Computer-assisted Intervention*, pp. 305–312 (2014). Springer

- [15] Huang, Y., Shao, L., Frangi, A.F.: Simultaneous super-resolution and cross-modality synthesis of 3d medical images using weakly-supervised joint convolutional sparse coding. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 6070–6079 (2017)
- [16] Zhao, Y., Liao, S., Guo, Y., Zhao, L., Yan, Z., Hong, S., Hermosillo, G., Liu, T., Zhou, X.S., Zhan, Y.: Towards mr-only radiotherapy treatment planning: synthetic ct generation using multi-view deep convolutional neural networks. In: International Conference on Medical Image Computing and Computer-Assisted Intervention, pp. 286–294 (2018)
- [17] Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. In: Advances in Neural Information Processing Systems, pp. 2672–2680 (2014)
- [18] Nie, D., Trullo, R., Lian, J., Wang, L., Petitjean, C., Ruan, S., Wang, Q., Shen, D.: Medical image synthesis with deep convolutional adversarial networks. *IEEE Transactions on Biomedical Engineering* **65**(12), 2720–2730 (2018)
- [19] Dalmaz, O., Yurt, M., Çukur, T.: Resvit: Residual vision transformers for multi-modal medical image synthesis. *IEEE Transactions on Medical Imaging* (2022)
- [20] Zhu, J.-Y., Park, T., Isola, P., Efros, A.A.: Unpaired image-to-image translation using cycle-consistent adversarial networks. In: Proceedings of the IEEE International Conference on Computer Vision, pp. 2223–2232 (2017)
- [21] Liu, S., Zhang, B., Liu, Y., Han, A., Shi, H., Guan, T., He, Y.: Unpaired stain transfer using pathology-consistent constrained generative adversarial networks. *IEEE Transactions on Medical Imaging* **40**(8), 1977–1989 (2021)
- [22] Huo, Y., Xu, Z., Moon, H., Bao, S., Assad, A., Moyo, T.K., Savona, M.R., Abramson, R.G., Landman, B.A.: Synseg-net: Synthetic segmentation without target modality ground truth. *IEEE transactions on medical imaging* **38**(4), 1016–1025 (2018)
- [23] Liu, Y., Lei, Y., Wang, T., Fu, Y., Tang, X., Curran, W.J., Liu, T., Patel, P., Yang, X.: Cbct-based synthetic ct generation using deep-attention cyclegan for pancreatic adaptive radiotherapy. *Medical physics* **47**(6), 2472–2483 (2020)
- [24] Huang, Z., Chen, Z., Zhang, Q., Quan, G., Ji, M., Zhang, C., Yang, Y.,

- Liu, X., Liang, D., Zheng, H., *et al.*: Cagan: A cycle-consistent generative adversarial network with attention for low-dose ct imaging. *IEEE Transactions on Computational Imaging* **6**, 1203–1218 (2020)
- [25] Xu, Z., Qi, C., Xu, G.: Semi-supervised attention-guided cyclegan for data augmentation on medical images. In: *IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*, pp. 563–568 (2019)
- [26] Nie, D., Shen, D.: Adversarial confidence learning for medical image segmentation and synthesis. *International journal of computer vision* **128**(10), 2494–2513 (2020)
- [27] Yang, H., Sun, J., Carass, A., Zhao, C., Lee, J., Prince, J.L., Xu, Z.: Unsupervised mr-to-ct synthesis using structure-constrained cyclegan. *IEEE transactions on medical imaging* **39**(12), 4249–4261 (2020)
- [28] Wang, J., Wu, Q., Pourpanah, F.: Dc-cyclegan: Bidirectional ct-to-mr synthesis from unpaired data. *arXiv preprint arXiv:2211.01293* (2022)
- [29] Han, X.: Mr-based synthetic ct generation using a deep convolutional neural network method. *Medical physics* **44**(4), 1408–1419 (2017)
- [30] Abu-Srhan, A., Almallahi, I., Abushariah, M.A., Mahafza, W., Al-Kadi, O.S.: Paired-unpaired unsupervised attention guided gan with transfer learning for bidirectional brain mr-ct synthesis. *Computers in Biology and Medicine* **136**, 104763 (2021)
- [31] Woo, S., Park, J., Lee, J.-Y., Kweon, I.S.: CBAM: Convolutional block attention module. In: *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 3–19 (2018)
- [32] Wang, T., Lei, Y., Fu, Y., Wynne, J.F., Curran, W.J., Liu, T., Yang, X.: A review on medical imaging synthesis using deep learning and its clinical applications. *Journal of applied clinical medical physics* **22**(1), 11–36 (2021)
- [33] Hofmann, M., Steinke, F., Scheel, V., Charpiat, G., Farquhar, J., Aschoff, P., Brady, M., Schölkopf, B., Pichler, B.J.: Mri-based attenuation correction for pet/mri: a novel approach combining pattern recognition and atlas registration. *Journal of nuclear medicine* **49**(11), 1875–1883 (2008)
- [34] Chen, M., Jog, A., Carass, A., Prince, J.L.: Using image synthesis for multi-channel registration of different image modalities. In: *Medical Imaging 2015: Image Processing*, vol. 9413, pp. 462–468 (2015). SPIE
- [35] Dowling, J.A., Lambert, J., Parker, J., Salvado, O., Fripp, J., Capp, A., Wratten, C., Denham, J.W., Greer, P.B.: An atlas-based electron density

- mapping method for magnetic resonance imaging (mri)-alone treatment planning and adaptive mri-based prostate radiation therapy. *International Journal of Radiation Oncology* Biology* Physics* **83**(1), 5–11 (2012)
- [36] Izquierdo-Garcia, D., Hansen, A.E., Förster, S., Benoit, D., Schachoff, S., Fürst, S., Chen, K.T., Chonde, D.B., Catana, C.: An spm8-based approach for attenuation correction combining segmentation and nonrigid template formation: application to simultaneous pet/mr brain imaging. *Journal of Nuclear Medicine* **55**(11), 1825–1830 (2014)
- [37] Delpon, G., Escande, A., Ruef, T., Darréon, J., Fontaine, J., Noblet, C., Supiot, S., Lacornerie, T., Pasquier, D.: Comparison of automated atlas-based segmentation software for postoperative prostate cancer radiotherapy. *Frontiers in oncology* **6**, 178 (2016)
- [38] Hsu, S.-H., Cao, Y., Huang, K., Feng, M., Balter, J.M.: Investigation of a method for generating synthetic ct models from mri scans of the head and neck for radiation therapy. *Physics in Medicine & Biology* **58**(23), 8419 (2013)
- [39] Burgos, N., Cardoso, M.J., Thielemans, K., Modat, M., Pedemonte, S., Dickson, J., Barnes, A., Ahmed, R., Mahoney, C.J., Schott, J.M., *et al.*: Attenuation correction synthesis for hybrid pet-mr scanners: application to brain studies. *IEEE transactions on medical imaging* **33**(12), 2332–2341 (2014)
- [40] Sevetlidis, V., Giuffrida, M.V., Tsaftaris, S.A.: Whole image synthesis using a deep encoder-decoder network. In: *International Workshop on Simulation and Synthesis in Medical Imaging*, pp. 127–137 (2016). Springer
- [41] Zhou, T., Fu, H., Chen, G., Shen, J., Shao, L.: Hi-net: hybrid-fusion network for multi-modal mr image synthesis. *IEEE transactions on medical imaging* **39**(9), 2772–2781 (2020)
- [42] Cao, B., Zhang, H., Wang, N., Gao, X., Shen, D.: Auto-gan: self-supervised collaborative learning for medical image synthesis. In: *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, pp. 10486–10493 (2020)
- [43] Zhang, T., Fu, H., Zhao, Y., Cheng, J., Guo, M., Gu, Z., Yang, B., Xiao, Y., Gao, S., Liu, J.: Skrgan: Sketching-rendering unconditional generative adversarial networks for medical image synthesis. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 777–785 (2019). Springer
- [44] Hu, S., Yuan, J., Wang, S.: Cross-modality synthesis from mri to pet

- using adversarial u-net with different normalization. In: 2019 International Conference on Medical Imaging Physics and Engineering (ICMIPE), pp. 1–5 (2019). IEEE
- [45] Wu, H., Jiang, X., Jia, F.: Uc-gan for mr to ct image synthesis. In: Workshop on Artificial Intelligence in Radiation Therapy, pp. 146–153 (2019). Springer
- [46] Yang, H., Sun, J., Carass, A., Zhao, C., Lee, J., Prince, J.L., Xu, Z.: Unsupervised mr-to-ct synthesis using structure-constrained cyclegan. *IEEE transactions on medical imaging* **39**(12), 4249–4261 (2020)
- [47] Chen, R., Huang, W., Huang, B., Sun, F., Fang, B.: Reusing discriminators for encoding: Towards unsupervised image-to-image translation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 8168–8177 (2020)
- [48] Lee, J., Gu, J., Ye, J.C.: Unsupervised ct metal artifact learning using attention-guided β -cyclegan. *IEEE Transactions on Medical Imaging* **40**(12), 3932–3944 (2021)
- [49] Kong, L., Lian, C., Huang, D., Hu, Y., Zhou, Q., et al.: Breaking the dilemma of medical image-to-image translation. *Advances in Neural Information Processing Systems* **34** (2021)
- [50] Kim, J., Kim, M., Kang, H., Lee, K.H.: U-gat-it: Unsupervised generative attentional networks with adaptive layer-instance normalization for image-to-image translation. In: International Conference on Learning Representations
- [51] Larochelle, H., Hinton, G.E.: Learning to combine foveal glimpses with a third-order boltzmann machine. *Advances in neural information processing systems* **23** (2010)
- [52] Fukui, H., Hirakawa, T., Yamashita, T., Fujiyoshi, H.: Attention branch network: Learning of attention mechanism for visual explanation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 10705–10714 (2019)
- [53] Fu, J., Liu, J., Tian, H., Li, Y., Bao, Y., Fang, Z., Lu, H.: Dual attention network for scene segmentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 3146–3154 (2019)
- [54] Liu, G., Guo, J.: Bidirectional lstm with attention mechanism and convolutional layer for text classification. *Neurocomputing* **337**, 325–338 (2019)

- [55] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. *Advances in neural information processing systems* **30** (2017)
- [56] Chen, C.-F.R., Fan, Q., Panda, R.: Crossvit: Cross-attention multi-scale vision transformer for image classification. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 357–366 (2021)
- [57] Guo, M.-H., Xu, T.-X., Liu, J.-J., Liu, Z.-N., Jiang, P.-T., Mu, T.-J., Zhang, S.-H., Martin, R.R., Cheng, M.-M., Hu, S.-M.: Attention mechanisms in computer vision: A survey. *Computational Visual Media*, 1–38 (2022)
- [58] Misra, D., Nalamada, T., Arasanipalai, A.U., Hou, Q.: Rotate to attend: Convolutional triplet attention module. In: *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 3139–3148 (2021)
- [59] Wang, S.-H., Fernandes, S.L., Zhu, Z., Zhang, Y.-D.: Avnc: attention-based vgg-style network for covid-19 diagnosis by cbam. *IEEE Sensors Journal* **22**(18), 17431–17438 (2021)
- [60] He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778 (2016)
- [61] Snell, J., Ridgeway, K., Liao, R., Roads, B.D., Mozer, M.C., Zemel, R.S.: Learning to generate images with perceptual similarity metrics. In: *2017 IEEE International Conference on Image Processing (ICIP)*, pp. 4277–4281 (2017). IEEE
- [62] Zhou, Y., Wang, X., Zhang, M., Zhu, J., Zheng, R., Wu, Q.: Mpce: a maximum probability based cross entropy loss function for neural network classification. *IEEE Access* **7**, 146331–146341 (2019)
- [63] Zhong, Y., Liu, L., Zhao, D., Li, H.: A generative adversarial network for image denoising. *Multimedia Tools and Applications* **79**(23), 16517–16529 (2020)
- [64] Wang, Z., Bovik, A.C., Sheikh, H.R., Simoncelli, E.P.: Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing* **13**(4), 600–612 (2004)