

Unpaired Brain MR-to-CT Synthesis using a Structure-Constrained CycleGAN

Heran Yang^{1,2}, Jian Sun¹, Aaron Carass², Can Zhao², Junghoon Lee³,
Zongben Xu¹, and Jerry Prince²

¹ School of Mathematics and Statistics, Xi'an Jiaotong University, China
yhr.7017@stu.xjtu.edu.cn

² Department of Electrical and Computer Engineering, Johns Hopkins University,
USA

³ Department of Radiation Oncology and Molecular Radiation Sciences, Johns
Hopkins University, USA

Abstract. The cycleGAN is becoming an influential method in medical image synthesis. However, due to a lack of direct constraints between input and synthetic images, the cycleGAN cannot guarantee structural consistency between these two images, and such consistency is of extreme importance in medical imaging. To overcome this, we propose a structure-constrained cycleGAN for brain MR-to-CT synthesis using unpaired data that defines an extra structure-consistency loss based on the modality independent neighborhood descriptor to constrain structural consistency. Additionally, we use a position-based selection strategy for selecting training images instead of a completely random selection scheme. Experimental results on synthesizing CT images from brain MR images demonstrate that our method is better than the conventional cycleGAN and approximates the cycleGAN trained with paired data.

Keywords: MR-to-CT synthesis, CycleGAN, Deep learning, MIND.

1 Introduction

Magnetic resonance (MR) imaging has been widely utilized to diagnose patients, as it is non-ionizing, non-invasive, and has a range of contrast mechanisms. However, MR images do not directly provide electron density information, which is essential for some applications such as MR-based radiotherapy treatment planning or attenuation correction in hybrid PET/MR scanners. A straightforward solution is to separately scan a computed tomography (CT) image, but this is time-consuming, costly, potentially harmful to patients, and requires accurate MR/CT registrations. Therefore, to avoid the CT scan, a variety of approaches have been proposed to synthesize CT images from available MR images [1, 4–7]. For example, by using paired MR and CT atlases, atlas-based methods [4] first register multiple atlas MR images to a subject MR image, and then the warped atlas CT images are combined to synthesize a subject CT image. Deep learning-based methods [5] have designed different convolutional neural network (CNN) structures to directly learn the MR-to-CT mapping.

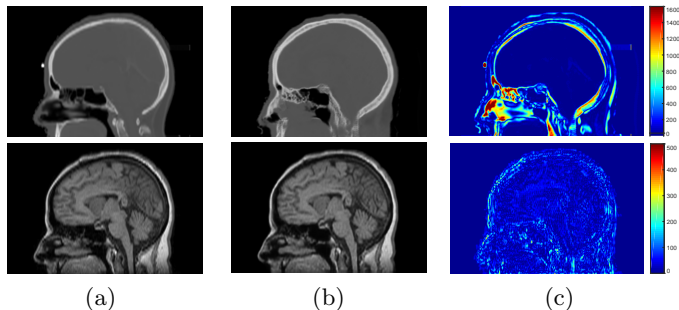


Fig. 1: Visual example of a cycleGAN result. We show (a) ground-truth CT image and input MR image, (b) synthetic CT image and reconstructed MR image, and (c) the relative errors between the ground-truth/synthetic CT images (upper) and the input/reconstructed MR images (lower) .

Although these methods can produce good synthetic images, they rely on a large number of paired CT and MR images, which are hard to obtain in practice, especially for specific MR tissue contrasts. To relax the requirement of paired data, Wolterink et al. [6] and Chartsias et al. [1] used a cycleGAN [8] for MR-to-CT synthesis on unpaired data with promising results. They used a CNN to learn the MR-to-CT mapping with the help of an adversarial loss, which forces synthetic CT images to be indistinguishable from real CT images. To ensure the synthetic CT image correctly corresponds to an input MR image, another CNN is utilized to map synthetic CT back to the MR domain and the reconstructed image should be identical to the input MR image (i.e., cycle-consistency loss).

However, due to a lack of direct constraints between the synthetic and input images, the cycleGAN cannot guarantee structural consistency between these two images. As shown in Fig. 1, the reconstructed MR image is almost identical to the input MR image, indicating the cycle consistency is well kept, but the synthetic CT image is quite different from the ground-truth, especially for the skull region, which illustrates that the structure of the synthetic CT image is not consistent with that of the input MR image. To overcome this, Zhang et al. [7] trained two auxiliary CNNs respectively for segmenting MR and CT images and also defined a loss to force the segmentation of the synthetic image to be the same as the ground-truth segmentation of the input image. This requires a training dataset with ground-truth segmentations of MR and CT images, which further complicates the training data requirements.

In this work, we propose a structure-constrained cycleGAN to constrain structural consistency without requiring ground-truth segmentations. By using the modality independent neighborhood descriptor [3], we define a structure-consistency loss enforcing the extracted features in the synthetic image to be voxel-wise close to the ones extracted in the input image. Additionally, we use a position-based selection strategy for selecting training images instead of a completely random selection scheme. Experimental results on synthesizing CT images from brain MR images show that our method achieves significantly better

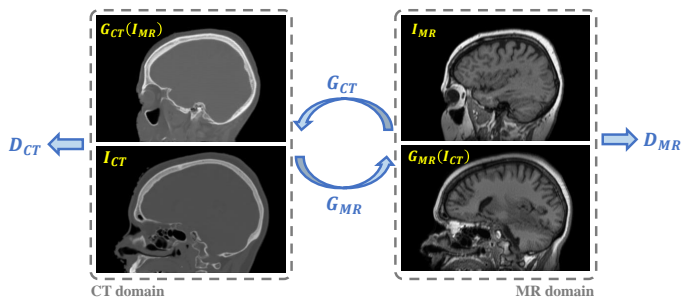


Fig. 2: Illustration of our proposed structure-constrained cycleGAN. Two generators (i.e., G_{CT} and G_{MR}) learn cross-domain mappings between CT and MR domains. The training of these mappings is supervised by adversarial, cycle-consistency, and structure-consistency losses.

results compared to a conventional cycleGAN with various metrics, and approximates the cycleGAN trained with paired data.

2 Method

In this section, we introduce our proposed structure-constrained cycleGAN. As shown in Fig. 2, our method contains two generators G_{CT} and G_{MR} , which provide the MR-to-CT and CT-to-MR mappings, respectively. In addition, discriminator D_{CT} is used to distinguish between real and synthetic CT images, and discriminator D_{MR} is for MR images. Our training loss includes three types of terms: an adversarial loss [2] for matching the distribution of synthetic images to target CT or MR domain; a cycle-consistency loss [8] to prevent generators from producing synthetic images that are irrelevant to the inputs; and a structure-consistency loss to constrain structural consistency between input and synthetic images.

2.1 Adversarial loss

The adversarial loss [2] is applied to both generators. For the generator G_{CT} and its discriminator D_{CT} , the adversarial loss is defined as

$$\mathcal{L}_{GAN}(G_{CT}, D_{CT}) = D_{CT}(G_{CT}(I_{MR}))^2 + (1 - D_{CT}(I_{CT}))^2, \quad (1)$$

where I_{CT} and I_{MR} denote the unpaired input CT and MR images. During the training phase, G_{CT} tries to generate a synthetic CT image $G_{CT}(I_{MR})$ close to a real CT image, i.e., $\max_{G_{CT}} \mathcal{L}_{GAN}(G_{CT}, D_{CT})$, while D_{CT} is to distinguish between a synthetic CT image $G_{CT}(I_{MR})$ and a real image I_{CT} , i.e., $\min_{D_{CT}} \mathcal{L}_{GAN}(G_{CT}, D_{CT})$. Similarly, the adversarial loss for G_{MR} and D_{MR} is defined as

$$\mathcal{L}_{GAN}(G_{MR}, D_{MR}) = D_{MR}(G_{MR}(I_{CT}))^2 + (1 - D_{MR}(I_{MR}))^2. \quad (2)$$

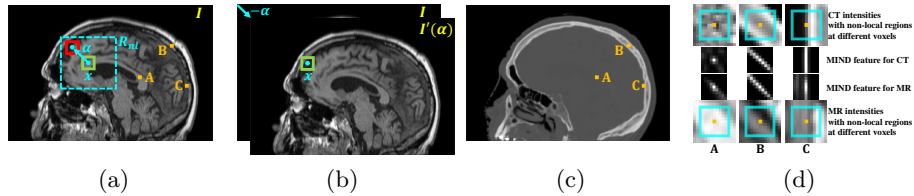


Fig. 3: Illustration of the MIND feature. (a) To extract the MIND feature at x , a patch around $x + \alpha$ is compared with a patch around x for each $x + \alpha \in R_{nl}$; (b) comparison between x and $x + \alpha$ of I in (a) equals a comparison of I and $I'(\alpha)$ at x ; (c) the CT image paired with MR image in (a); (d) visual examples of MIND features extracted at voxels A, B, C within paired MR and CT images in (a) and (c).

2.2 Cycle-consistency loss

To prevent the generators from producing synthetic images that are irrelevant to the inputs, a cycle-consistency loss [8] is utilized for G_{CT} and G_{MR} forcing the reconstructed images $G_{CT}(G_{MR}(I_{CT}))$ and $G_{MR}(G_{CT}(I_{MR}))$ to be identical to their inputs I_{CT} and I_{MR} . This loss is written as

$$\begin{aligned} \mathcal{L}_{\text{cycle}}(G_{CT}, G_{MR}) = & \|G_{CT}(G_{MR}(I_{CT})) - I_{CT}\|_1 \\ & + \|G_{MR}(G_{CT}(I_{MR})) - I_{MR}\|_1. \end{aligned} \quad (3)$$

2.3 Structure-consistency loss

Since the cycle-consistency loss does not necessarily ensure structural consistency (as discussed in Sec. 1), our method uses an extra structure-consistency loss between the synthetic and input images. However, as these two images are respectively in MR and CT domains, we first map these images into a common feature domain by using a modal-independent structural feature, and then the structural consistency between the synthetic and input images is measured in this feature domain. In this work, we use the modality independent neighborhood descriptor (MIND) [3] as the structural feature. MIND is defined using a non-local patch-based self-similarity and depends on image local structure instead of intensity values. It has been previously applied to MR/CT image registration as a similarity metric. Figure 3(d) shows visual examples of MIND features extracted at different voxels in MR and CT images. In the following paragraphs, we introduce the MIND feature and our structure-consistency loss in detail.

The MIND feature extracts distinctive image structure by comparing each patch with all its neighbors in a non-local region. As shown in Fig. 3(a), for voxel x in image I , the MIND feature F_x is an $|R_{nl}|$ -length vector, where R_{nl} denotes a non-local region around voxel x , and each component $F_x^{(\alpha)}$ for a voxel $x + \alpha \in R_{nl}$ is defined as

$$F_x^{(\alpha)}(I) = \frac{1}{Z} \exp\left(-\frac{D_{\mathcal{P}}(I, x, x + \alpha)}{V(I, x)}\right), \quad (4)$$

where Z is a normalization constant so that the maximal component of F_x is 1. $D_{\mathcal{P}}(I, x, x + \alpha)$ denotes the L_2 distance between two image patches \mathcal{P} respectively centered at voxel x and voxel $x + \alpha$ in image I , and $V(I, x)$ is an estimation of local variance at voxel x , which can be written as

$$D_{\mathcal{P}}(I, x, x + \alpha) = \sum_{p \in \mathcal{P}} (I(x + p) - I(x + \alpha + p))^2, \quad (5)$$

$$V(I, x) = \frac{1}{4} \sum_{n \in \mathcal{N}} D_{\mathcal{P}}(I, x, x + n), \quad (6)$$

where \mathcal{N} is the 4-neighborhood of voxel x .

It is difficult to directly compute the operation $D_{\mathcal{P}}$ and its gradient using Eqn. 5 in a deep network. Instead, as shown in Fig. 3(b), $D_{\mathcal{P}}$ can be equivalently computed by using a convolutional operation as

$$D_{\mathcal{P}}(I, x, x + \alpha) = C * (I - I'(\alpha))^2, \quad (7)$$

where C is an all-one kernel of the same size as patch \mathcal{P} , and $I'(\alpha)$ denotes I translated by α . By doing this, the structural feature can be extracted via several simple operations and the gradients of these operations can be easily computed.

Based on the MIND feature introduced above, the structure-consistency loss in our method is defined to enforce the extracted MIND features in the synthetic images $G_{CT}(I_{MR})$ or $G_{MR}(I_{CT})$ to be voxel-wise close to the ones extracted in their inputs I_{MR} or I_{CT} , which can be written as

$$\begin{aligned} \mathcal{L}_{\text{structure}}(G_{CT}, G_{MR}) &= \frac{1}{N_{MR}|R_{nl}|} \sum_x \|F_x(G_{CT}(I_{MR})) - F_x(I_{MR})\|_1 \\ &+ \frac{1}{N_{CT}|R_{nl}|} \sum_x \|F_x(G_{MR}(I_{CT})) - F_x(I_{CT})\|_1, \end{aligned} \quad (8)$$

where N_{MR} and N_{CT} respectively denote the number of voxels in input images I_{MR} and I_{CT} , and $\|\cdot\|_1$ is the L_1 norm. In this work, we use a 9×9 non-local region and a 7×7 patch for computing structure-consistency loss. Furthermore, instead of an all-one kernel C , we utilize a Gaussian kernel C_{σ} with standard deviation $\sigma = 2$ to reweight the importance of voxels within patch \mathcal{P} in Eqn. 7. In preliminary experiments, we tried different non-local regions, patch sizes, and σ values, but did not observe improved performance.

2.4 Training loss

Given the definitions of adversarial, cycle-consistency, and structure-consistency losses above, the training loss of our proposed method is defined as:

$$\begin{aligned} \mathcal{L}(G_{CT}, G_{MR}, D_{CT}, D_{MR}) &= \mathcal{L}_{\text{GAN}}(G_{CT}, D_{CT}) + \mathcal{L}_{\text{GAN}}(G_{MR}, D_{MR}) \\ &+ \lambda_1 \mathcal{L}_{\text{cycle}}(G_{CT}, G_{MR}) + \lambda_2 \mathcal{L}_{\text{structure}}(G_{CT}, G_{MR}), \end{aligned} \quad (9)$$

where λ_1 and λ_2 control the relative importance of the loss terms. During training, λ_1 is set to 10 as per [6, 8] and λ_2 is set to 5. To optimize \mathcal{L} , we alternatively update $D_{MR/CT}$ (with $G_{MR/CT}$ fixed) and $G_{MR/CT}$ (with $D_{MR/CT}$ fixed).

2.5 Network structure

Our method is composed of four trainable neural networks, i.e., two generators, G_{CT} and G_{MR} , and two discriminators, D_{CT} and D_{MR} , and we use the same network structures as [6, 8] in this work. That is, two generators, G_{CT} and G_{MR} , are 2D fully convolutional networks (FCNs) with two stride-2 convolutional layers, nine residual blocks, and two fractionally-strided convolutional layers with stride $\frac{1}{2}$. The two discriminators, D_{CT} and D_{MR} , are 2D FCNs consisting of five convolutional layers to classify whether 70×70 overlapping image patches are real or synthetic. For further details, please refer to [8].

2.6 Position-based selection strategy

Although our input MR and CT slices are unpaired, we can get the positions of their slices within the volumes. Slices in the middle of the volume necessarily have more brain tissue than peripheral slices. Thus, instead of feeding in slices at extremely different positions of the brain, e.g., a peripheral CT slice and a medial MR slice, we input training slices at similar positions; this is referred to as a position-based selection (PBS) strategy. That is, the MR and CT slices are linearly aligned considering their respective numbers of slices within the volumes, and given the i -th MR slice in its volume, the index $T(i)$ of corresponding CT slice selected by our method is determined by

$$T(i) = \begin{cases} \left[i \cdot \frac{K_{CT}-1}{K_{MR}-1} \right] + m, & \text{if } 5 \leq \left[i \cdot \frac{K_{CT}-1}{K_{MR}-1} \right] < K_{CT} - 5, \\ \left[i \cdot \frac{K_{CT}-1}{K_{MR}-1} \right], & \text{otherwise,} \end{cases} \quad (10)$$

where K_{MR} and K_{CT} respectively denote the number of slices in unpaired MR and CT volumes. $[\cdot]$ denotes the rounding function, and m is a random integer within the range of $[-5, 5]$. This strategy forces the discriminators to be stronger at distinguishing synthetic images from real ones, thus avoiding mode collapse. This in turn forces our generators to be better in order to *trick* our discriminators. We evaluate this position-based selection strategy in Sec. 3.

3 Experiments

3.1 Data set

The MR and CT volumes are respectively obtained using a Siemens Magnetom Espree 1.5T scanner (Siemens Medical Solutions, Malvern, PA) and a Philips Brilliance Big Bore scanner (Philips Medical Systems, Netherlands) under a routine clinical protocol for brain cancer patients. Geometric distortions in MR volumes are corrected using a 3D correction algorithm in the Siemens Syngo console workstation. All MR volumes are N4 corrected and normalized by aligning the white matter peak identified by fuzzy C-means.

The data set contains the brain MR and CT volumes of 45 patients, which were divided into a training set containing MR and CT volumes of 27 patients,

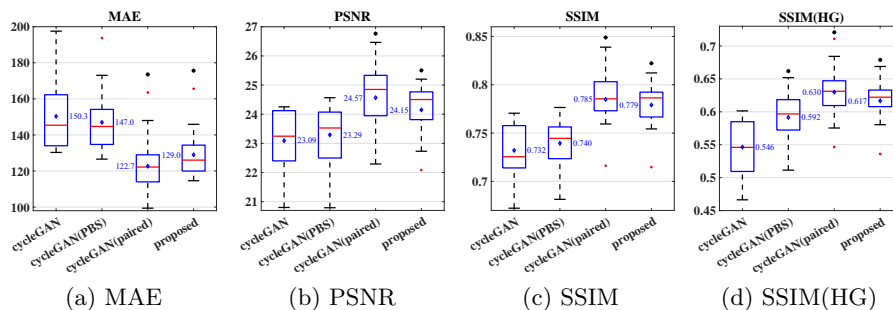


Fig. 4: Comparison of different methods on synthesizing CT images in boxplots, where the diamond and number in blue denote the respective mean and * denotes $p < 0.001$ compared to the conventional cycleGAN using a paired sample t-test.

a validation set of 3 patients for model and epoch selection, and a test set of 15 patients for performance evaluation. As in [6], the experiments were performed on 2D sagittal image slices. Each MR or CT volume contains about 270 sagittal images, which are resized and padded to 384×256 while maintaining the aspect ratio, and the intensity ranges are respectively $[-1000, 3500]$ HU for CT and $[0, 3500]$ for MR. To augment the training set, each image is padded to 400×284 and then randomly cropped to 384×256 as training samples.

3.2 Experimental results

We compare the proposed method to the conventional cycleGAN [6, 8] (denoted as “cycleGAN”) and a cycleGAN trained with paired data (denoted as “cycleGAN (paired)”), which represents the best that a cycleGAN can achieve. To evaluate the position-based selection strategy in Sec. 2.6, a cycleGAN using this strategy during training, denoted as “cycleGAN (PBS)”, is also included in comparison. As in [6, 8], the learning rate is set to 0.0002 for all compared methods.

To quantitatively compare these methods, we use mean absolute error (MAE), peak signal-to-noise ratio (PSNR), and structural similarity (SSIM) between the ground-truth CT volume and the synthetic one, which are computed within the head region mask and averaged over 15 test subjects. Furthermore, SSIM over regions with high gradient magnitudes (denoted as “SSIM(HG)”) is also computed to measure the quality of bone regions in synthetic images. The maximum value in PSNR and the dynamic range in SSIM are set to 4500, as the range of our CT data is $[-1000, 3500]$ HU.

As shown in Fig. 4, our proposed method achieves significantly better performance than conventional cycleGAN in all the metrics ($p < 0.001$) and produces similar results compared to the cycleGAN trained with paired data. Compared to randomly selecting training slices at any position, our proposed position-based selection strategy produces significantly higher SSIM(HG) score ($p < 0.001$) with marginal improvement in the other three metrics. Figure 5 shows visual examples of synthetic CT images by different methods from a test subject.

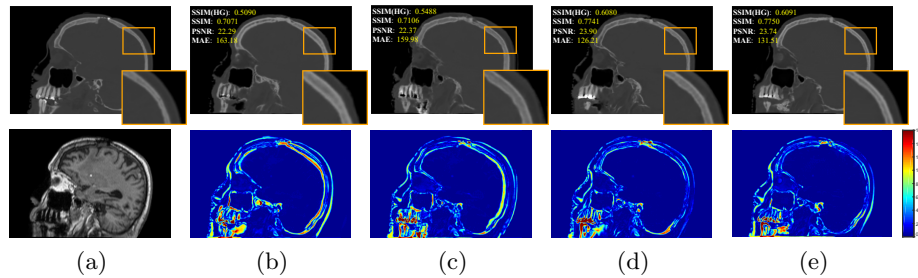


Fig. 5: Visual comparison of synthetic CT images using different methods. For one test subject, we show (a) the ground-truth CT image and input MR image; the synthetic CT image and its difference image (compared to ground-truth CT image) generated by (b) cycleGAN, (c) cycleGAN (PBS), (d) cycleGAN (paired), and (e) proposed method. The small text in each sub-image is the corresponding accuracy on this test subject.

4 Conclusion

We propose a structure-constrained cycleGAN for brain MR-to-CT synthesis using unpaired data. Compared to the conventional cycleGAN [6, 8], we define an extra structure-consistency loss based on the modality independent neighborhood descriptor to constrain structural consistency and also introduce a position-based selection strategy for selecting training images. The experiments show that our method generates better synthetic CT images than the conventional cycleGAN and produces results similar to a cycleGAN trained with paired data.

Acknowledgments. This work is supported by the NSFC (11622106, 11690011, 61721002) and the China Scholarship Council.

References

1. Chatsias, A., Joyce, T., et al.: Adversarial image synthesis for unpaired multi-modal cardiac data. In: SASHIMI. pp. 3–13 (2017)
2. Goodfellow, I., et al.: Generative adversarial nets. In: NIPS. pp. 2672–2680 (2014)
3. Heinrich, M.P., et al.: MIND: Modality independent neighbourhood descriptor for multi-modal deformable registration. *Med. Image Anal.* **16**(7), 1423–1435 (2012)
4. Hofmann, M., Bezrukov, I., et al.: MRI-based attenuation correction for whole-body PET/MRI: quantitative evaluation of segmentation- and atlas-based methods. *J. Nucl. Med.* **52**(9), 1392–1399 (2011)
5. Roy, S., Butman, J.A., Pham, D.L.: Synthesizing CT from ultrashort echo-time MR images via convolutional neural networks. In: SASHIMI. pp. 24–32 (2017)
6. Wolterink, J.M., Dinkla, A.M., et al.: Deep MR to CT synthesis using unpaired data. In: SASHIMI. pp. 14–23 (2017)
7. Zhang, Z., et al.: Translating and segmenting multimodal medical volumes with cycle- and shape-consistency generative adversarial network. In: CVPR (2018)
8. Zhu, J.Y., Park, T., et al.: Unpaired image-to-image translation using cycle-consistent adversarial networks. In: ICCV. pp. 2242–2251 (2017)