# Artifact Disentanglement Network for Unsupervised Metal Artifact Reduction

Haofu Liao[1(✉)], Wei-An Lin[2], Jianbo Yuan[1], S. Kevin Zhou[3], and Jiebo Luo[1]

[1] Department of Computer Science, University of Rochester
`haofu.liao@rochester.edu`
[2] Department of ECE, University of Maryland, College Park
[3] Institute of Computing Technology, Chinese Academy of Sciences

**Abstract.** Current deep neural network based approaches to computed tomography (CT) metal artifact reduction (MAR) are supervised methods which rely heavily on synthesized data for training. However, as synthesized data may not perfectly simulate the underlying physical mechanisms of CT imaging, the supervised methods often generalize poorly to clinical applications. To address this problem, we propose, to the best of our knowledge, the first unsupervised learning approach to MAR. Specifically, we introduce a novel artifact disentanglement network that enables different forms of generations and regularizations between the artifact-affected and artifact-free image domains to support unsupervised learning. Extensive experiments show that our method significantly outperforms the existing unsupervised models for image-to-image translation problems, and achieves comparable performance to existing supervised models on a synthesized dataset. When applied to clinical datasets, our method achieves considerable improvements over the supervised models. The source code of this paper is publicly available at https://github.com/liaohaofu/adn.

## 1 Introduction

Metal artifact is one of the commonly encountered artifacts in computed tomography (CT) images. It is introduced by the metallic implants during the imaging and reconstruction process. The formation of metal artifact involves several mechanisms such as beam hardening, scatter, noise, and the non-linear partial volume effect [1], which make it very challenging to be modeled and removed by traditional methods. Therefore, recent approaches [15,13,4,2,8,7,9] to metal artifact reduction (MAR) propose to use deep neural networks (DNNs) to inherently address the modeling of metal artifacts, and their experimental results show promising MAR performances.

All the existing DNN-based approaches are supervised methods requiring pairs of anatomically identical CT images, one with and the other without metal artifacts, for training. As it is clinically impractical to obtain such pairs of images, most of the supervised methods rely on synthesized images to train their models. However, due to the complexity of metal artifacts and the variations

of CT devices, the synthesized images may not fully simulate the real clinical scenarios, and the performances of these supervised methods may degrade in clinical applications.

In this work, we aim to address the challenging yet more practical unsupervised setting where *no paired CT images are available for training.* To this end, we propose a novel artifact disentanglement network to separate the metal artifacts from clinical CT images in a latent space. The disentanglement enables manipulations between the artifact-affected and artifact-free image domains so that different forms of adversarial- and self-regularizations can be achieved to support unsupervised learning. *To the best of our knowledge, this is the first unsupervised learning approach to MAR.* Extensive experiments show that our method achieves comparable performance to the existing supervised methods on a synthesized dataset. When applied to clinical datasets, all the supervised methods demonstrate certain degrees of degradation, whereas our method outperforms the supervised methods with significantly better clinical MAR results.

## 2    Related work

**Unsupervised image-to-image translation**    Image artifact reduction can be regarded as a form of image-to-image translation. One of the earliest unsupervised works in this category is CycleGAN [16] where a cycle-consistency design is proposed for unsupervised learning. Later works [5,6] improve CycleGAN for diverse and multimodal image generation. However, these unsupervised methods target at image synthesis and do not have suitable components for artifact reduction. Another recent work that is specialized for artifact reduction is deep image prior (DIP) [12], which, however, only works for less structured artifacts such as noise and compression artifacts.

**Deep metal artifact reduction**    A number of studies have recently been proposed to address MAR with DNNs. RL-ARCNN [4] introduces residual learning to a deep convolutional neural network (CNN) and achieves better MAR performance than ordinary CNN. DesteakNet [2] proposes a two-streams approach that can take a pair of NMAR [10] and detail images as the input to jointly reduce metal artifact. CNNMAR [15] uses CNN to generate prior images in the CT image domain to help the correction in the sinogram domain. Both DesteakNet and CNNMAR show significant improvements over the existing non-DNN based methods on synthesized datasets. cGANMAR [13] leverages generative adversarial networks (GANs) [3] to further improve DNN-based MAR performance.

## 3    Methodology

Let $\mathcal{I}$ be the domain of all artifact-free CT images and $\mathcal{I}^a$ be the domain of all artifact-affected CT images, the proposed artifact disentanglement network (ADN) aims to learn a mapping from $\mathcal{I}^a$ to $\mathcal{I}$ without paired data. As illustrated in Figure 1, ADN contains a set of artifact-free image encoder, generator
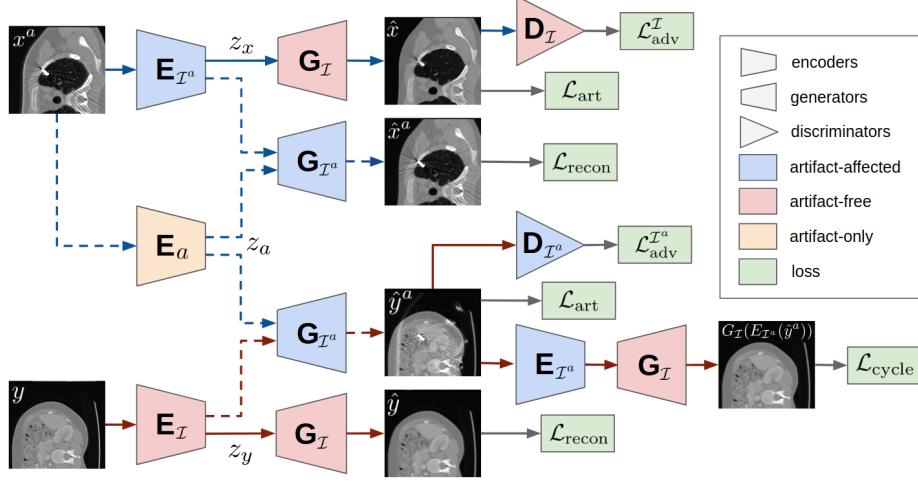
Fig. 1: Overview of the artifact disentanglement network.

and discriminator $\{E_\mathcal{I}, G_\mathcal{I}, D_\mathcal{I}\}$, a set of artifact-affected image encoder, generator and discriminator $\{E_{\mathcal{I}^a}, G_{\mathcal{I}^a}, D_{\mathcal{I}^a}\}$ and an artifact-only encoder $E_a$. The architectures of these building components are inspired from the state-of-the-art studies for image-to-image translation [17,5]. See the supplementary material for their detailed structures.

**Components**    Given two unpaired images $x^a \in \mathcal{I}^a$ and $y \in \mathcal{I}$, the encoders $E_{\mathcal{I}^a}$ and $E_\mathcal{I}$ map the artifact-free content information from $x^a$ and $y$ to a common content space $\mathcal{C}$, respectively. $E_a$ maps the artifact-only information from $x^a$ to an artifact space $\mathcal{A}$,

$$z_x = E_{\mathcal{I}^a}(x^a), z_y = E_\mathcal{I}(y), z_a = E_a(x^a), \quad \{z_x, z_y\} \subset \mathcal{C}, z_a \in \mathcal{A}. \quad (1)$$

The generator $G_{\mathcal{I}^a}$ takes an artifact-free code, $z_x$ or $z_y$, and an artifact-only code $z_a$ as the input and outputs an artifact-affected image. $G_\mathcal{I}$ takes an artifact-free code, $z_x$ or $z_y$, as the input and outputs an artifact-free image,

$$\begin{aligned} \hat{x} &= G_\mathcal{I}(z_x), \quad \hat{x}^a = G_{\mathcal{I}^a}(z_x, z_a), \\ \hat{y} &= G_\mathcal{I}(z_y), \quad \hat{y}^a = G_{\mathcal{I}^a}(z_y, z_a). \end{aligned} \quad (2)$$

During testing, only $E_{\mathcal{I}^a}$ and $G_\mathcal{I}$ are required to obtain an artifact-corrected output, i.e., $\hat{x} = G_\mathcal{I}(E_{\mathcal{I}^a}(x^a))$. The discriminator $D_{\mathcal{I}^a}$ decides whether an input is sampled from $\mathcal{I}^a$ or generated by $G_{\mathcal{I}^a}$. Similarly, $D_\mathcal{I}$ decides whether an input is from $\mathcal{I}$ or $G_\mathcal{I}$.

**Loss functions**    A good MAR model should (i) reduce the artifacts as much as possible and (ii) keep the anatomical content of the input CT images. To remove the artifacts, we train $D_\mathcal{I}$ and $G_\mathcal{I}$ adversarially to encourage the output

$\hat{x}$ to appear similar to an artifact-free image,

$$\mathcal{L}_{\text{adv}}^{\mathcal{I}} = \mathbb{E}_{\mathcal{I}}[\log D_{\mathcal{I}}(y)] + \mathbb{E}_{\mathcal{I}^a}[1 - \log D_{\mathcal{I}}(\hat{x})] \tag{3}$$

To maintain the anatomical content, we apply self-reconstruction to force the encoders and decoders to preserve the content of the inputs,

$$\mathcal{L}_{\text{recon}} = \mathbb{E}_{\mathcal{I},\mathcal{I}^a}[||\hat{x}^a - x^a||_1 + ||\hat{y} - y||_1]. \tag{4}$$

Here, the first term encourages $E_{\mathcal{I}^a}$ encodes all the content information of $x^a$ and the artifact information is not encoded due to the introduction of a separate artifact encoder $E_a$. With the second term, $G_{\mathcal{I}}$ learns how to fully reconstruct the encoded artifact-free content information. Combining these two terms, content persevering for $\hat{x}$ can be achieved.

In addition, we also introduce a *self-reduction design* to further enforce the learning. This idea is carried out in two steps. In the first step, ADN synthesizes "real" metal artifact from $x^a$ and apply it to $y$. Specifically, this is achieved by decoding from $z_y$ and $z_a$, i.e., $\hat{y}^a = G_{\mathcal{I}^a}(z_y, z_a)$, and we use another adversarial loss to guarantee $\hat{y}^a$ looking "real",

$$\mathcal{L}_{\text{adv}}^{\mathcal{I}^a} = \mathbb{E}_{\mathcal{I}^a}[\log D_{\mathcal{I}^a}(x^a)] + \mathbb{E}_{\mathcal{I},\mathcal{I}^a}[1 - \log D_{\mathcal{I}^a}(\hat{y}^a)] \tag{5}$$

In the second step, ADN reduces artifacts from the synthesized data to recover back to $y$. This is regularized by a cycle-consistent loss

$$\mathcal{L}_{\text{cycle}} = \mathbb{E}_{\mathcal{I},\mathcal{I}^a}[||G_{\mathcal{I}}(E_{\mathcal{I}^a}(\hat{y}^a)) - y||_1]. \tag{6}$$

Finally, due to the use of the same metal artifact, the difference map between $x^a$ and $\hat{x}$ and that between $\hat{y}^a$ and $y$ should be close. Thus, we employ an *artifact-consistent* loss to constrain the artifact difference,

$$\mathcal{L}_{\text{art}} = \mathbb{E}_{\mathcal{I},\mathcal{I}^a}[||(x^a - \hat{x}) - (\hat{y}^a - y)||_1]. \tag{7}$$

The full objective function is given by

$$\mathcal{L} = \lambda_{\text{adv}}^{\mathcal{I}} \mathcal{L}_{\text{adv}}^{\mathcal{I}} + \lambda_{\text{adv}}^{\mathcal{I}^a} \mathcal{L}_{\text{adv}}^{\mathcal{I}^a} + \lambda_{\text{recon}} \mathcal{L}_{\text{recon}} + \lambda_{\text{cycle}} \mathcal{L}_{\text{cycle}} + \lambda_{\text{art}} \mathcal{L}_{\text{art}}, \tag{8}$$

where the $\lambda$'s are hyper-parameters that control the importance of each term.

## 4    Experiments

**Datasets.**    We evaluate the proposed method on one synthesized dataset and two clinical datasets. We refer to them as SYN, CL1 and CL2, respectively. For SYN, we randomly select $4,118$ artifact-free CT images from DeepLesion [14] and follow the method from CNNMAR [15] to synthesize metal artifacts. We use $3,918$ of the synthesized pairs for training and validation and the rest 200 pairs for testing.

Table 1: Quantitative evaluation on the SYN dataset.

| | Supervised | | | Unsupervised | | | | |
|---|---|---|---|---|---|---|---|---|
| | CNNMAR[15] | UNet [11] | cGANMAR [13] | Ours | CycleGAN [13] | DIP [12] | MUNIT [5] | DRIT [6] |
| PSNR | 32.5 | **34.8** | **34.1** | <u>33.6</u> | 30.8 | 26.4 | 14.9 | 25.6 |
| SSIM | 91.4 | **93.1** | **93.4** | <u>92.4</u> | 72.9 | 75.9 | 7.5 | 79.7 |



Fig. 2: Qualitative evaluation on the SYN dataset. For better visualization, we obtain the metal region through thresholding and color it with red. See the supplementary material for more qualitative results.

For CL1, we choose the vertebrae localization and identification dataset from Spineweb[1]. We split the CT images from this dataset into two groups, one with artifacts and the other without artifacts. First, we identify regions with HU values greater than $2,500$ as the metal regions. Then, CT images whose largest-connected metal regions have more than 400 pixels are selected as artifact-affected images. CT images with the largest HU values less than $2,000$ are selected as artifact-free images. After this selection, the artifact-affected group contains $6,270$ images and the artifact-free group contains $21,190$ images. We withhold 200 images from the artifact-affected group for testing.

For CL2, we investigate the performance of the proposed method under a more challenging *cross-modality* setting. Specifically, the artifact-affected images of CL2 are from a cone-beam CT (CBCT) dataset collected during spinal interventions. Images from this dataset are very noisy and the majority of them contain metallic implants. There are in total $2,560$ CBCT images from this dataset, among which 200 images are withheld for testing. For the artifact-free images, we reuse the CT images collected from CL1.

**Baselines.** We compare the proposed method with seven state-of-the-art methods that are closely related to our problem. Three of the compared methods are supervised: CNNMAR [15], UNet [11] and cGANMAR [13]. CNNMAR and cGANMAR are two recent approaches that are dedicated to MAR. UNet is a general DNN framework that shows effectiveness in many image-to-image problems. The other four compared methods are unsupervised: CycleGAN [16],

---

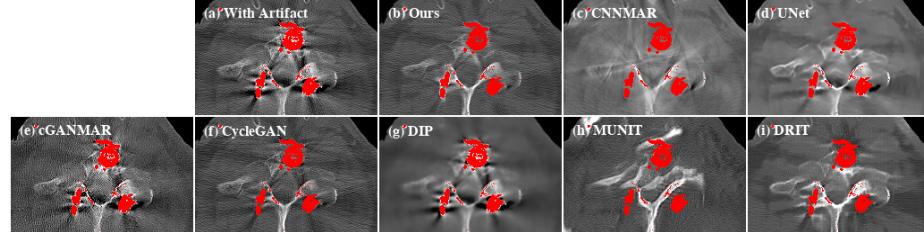[1] spineweb.digitalimaginggroup.ca

Fig. 3: Qualitative evaluation on the CL1 dataset. For better visualization, we obtain the metal region through thresholding and color it with red. See the supplementary material for more qualitative results.

DIP [12], MUNIT [5] and DRIT [6]. These methods are currently state-of-the-art approaches to unsupervised image-to-image translation problems. All the compared methods except UNet are trained with their officially released code. For UNet, a publicly available implementation[2] is used.

**Training and testing.**    We implement our method under the PyTorch deep learning framework[3] and use the Adam optimizer with $1 \times 10^{-4}$ learning rate to minimize the objective function. For the hyper-parameters, we use $\lambda_{\mathrm{adv}}^{\mathcal{I}} = \lambda_{\mathrm{adv}}^{\mathcal{I}^a} = 1.0$, $\lambda_{\mathrm{recon}} = \lambda_{\mathrm{cycle}} = \lambda_{\mathrm{art}} = 20.0$ for SYN and CL1, and use $\lambda_{\mathrm{adv}}^{\mathcal{I}} = \lambda_{\mathrm{adv}}^{\mathcal{I}^a} = 1.0$, $\lambda_{\mathrm{recon}} = \lambda_{\mathrm{cycle}} = \lambda_{\mathrm{art}} = 5.0$ for CL2.

To simulate the unsupervised setting for SYN, we evenly divide the $3,918$ synthesized training pairs into two groups. For one group, only artifact-affected images are used and their corresponding artifact-free images are withheld. For the other group, only artifact-free images are used and their corresponding artifact-affected images are withheld. During training of the unsupervised methods, we randomly select one image from each of the two groups as the input. For the supervised methods, all the $3,918$ synthesized training pairs are used.

To train the supervised methods with CL1, we first synthesize metal artifacts using the images from the artifact-free group of CL1. Then, we train the supervised methods with the synthesized pairs. During testing, the trained models are applied to the testing set containing only clinical metal artifact images. To train the unsupervised methods, we randomly select one image from the artifact-affected group and the other from the artifact-free group as the input.

For CL2, synthesizing metal artifacts is not possible due to the unavailability of artifact-free CBCT images. Therefore, for the supervised methods we directly use the models trained for CL1. In other words, the supervised methods are trained on synthesized CT images (from CL1) and tested on clinical CBCT images (from CL2). For the unsupervised models, each time we randomly select one artifact-affected CBCT image and one artifact-free CT image as the input for training.

---

[2] github.com/milesial/Pytorch-UNet
[3] pytorch.org

**Performance on synthesized data.** SYN contains paired data, allowing for both quantitative and qualitative evaluations. Following the convention in the literature, we use peak signal-to-noise ratio (PSNR) and structural similarity index (SSIM) as the metrics for the quantitative evaluation. For both metrics, the higher the better. Table 1 and Figure 8 show the quantitative and qualitative evaluation results, respectively.

We observe that the proposed method performs significantly better than the other unsupervised methods. MUNIT focuses more on diverse and realistic outputs (Figure 8(i)) with less constraint on structural similarity. CycleGAN and DRIT perform better as both the two models also require the artifact-corrected outputs to be able to transform back to the original artifact-affected images. Although this helps preserve content information, it also encourages the models to keep the artifacts. Therefore, as shown in Figure 8(g) and 2(j), the artifacts cannot be greatly reduced. DIP does not reduce much metal artifact in the input image (Figure 8(h)) as it is not designed to handle the more structured metal artifact.

We also find that the performance of our method is on a par with the supervised methods. The performance of UNet is close to that of cGANMAR which at its backend uses an UNet-like architecture. However, owing to the use of GAN, it produces sharper outputs (Figure 8(e)) than UNet (Figure 8(f)). As for PSNR and SSIM, both methods only slightly outperform our method and, surprisingly, our method performs better than CNNMAR.

**Performance on clinical data.** Next, we investigate the performance of the proposed method on clinical data. Since there are no ground truths available for the clinical images, only qualitative comparisons are performed. The qualitative evaluation results of CL1 are shown in Figure 9. Here, all the supervised methods are trained with paired images that are synthesized from the artifact-free group of CL1. We can see that UNet and cGANMAR generalize poorly when applied to clinical images (Figure 9(d) and 9(e)). CNNMAR is more robust as it corrects the artifacts in the sinogram domain. However, such a sinogram domain correction also introduces secondary artifacts (Figure 9(c)). For the more challenging cross-modality artifact reduction task with CL2 (Figure 10), all the supervised methods fail. This is not totally unexpected as the supervised methods are trained using only CT images because of the lack of artifact-free CBCT images. Similar to the cases with SYN, the other unsupervised methods also show inferior performances when evaluated on both the CL1 and CL2 datasets. By contrast, our method consistently delivers high-quality artifact reduced results on clinical images.

## 5 Conclusion

We presented a novel unsupervised learning approach to MAR. Through the development of an artifact disentanglement network, we showed how to leverage different forms of regularizations to eliminate the requirement of paired images
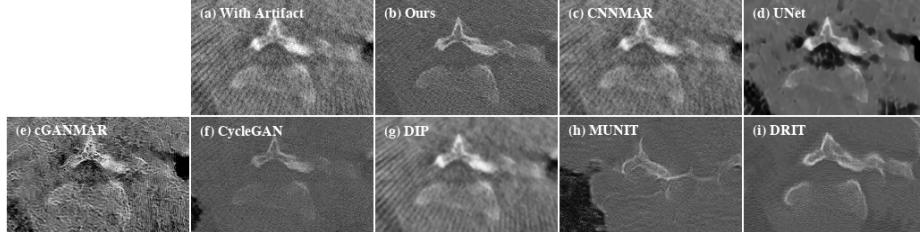
Fig. 4: Qualitative evaluation on the CL2 dataset. See the supplementary material for more qualitative results.

for training. To understand the effectiveness of this approach, we performed extensive evaluations on one synthesized and two clinical datasets. The evaluation results demonstrated the feasibility of using unsupervised learning method to achieve comparable performance to the supervised methods. More importantly, the results also showed that directly learning MAR from clinical CT images under an unsupervised setting was a more feasible and robust approach than transferring the knowledge learned from synthesized data to clinical data. We believe our findings in this work will initiate more applicable research for medical image artifact reduction even under an unsupervised setting.

# References

1. Gjesteby, L., Man, B.D., Jin, Y., Paganetti, H., Verburg, J., Giantsoudi, D., Wang, G.: Metal artifact reduction in CT: where are we after four decades? IEEE Access **4**, 5826–5849 (2016)
2. Gjesteby, L., Shan, H., Yang, Q., Xi, Y., Claus, B., Jin, Y., De Man, B., Wang, G.: Deep neural network for ct metal artifact reduction with a perceptual loss function. In: In Proceedings of The Fifth International Conference on Image Formation in X-ray Computed Tomography (2018)
3. 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 (2014)
4. Huang, X., Wang, J., Tang, F., Zhong, T., Zhang, Y.: Metal artifact reduction on cervical ct images by deep residual learning. Biomedical engineering online **17**(1), 175 (2018)
5. Huang, X., Liu, M., Belongie, S.J., Kautz, J.: Multimodal unsupervised image-to-image translation. In: Computer Vision - ECCV 2018 (2018)
6. Lee, H., Tseng, H., Huang, J., Singh, M., Yang, M.: Diverse image-to-image translation via disentangled representations. In: Computer Vision - ECCV 2018 (2018)
7. Liao, H., Huo, Z., Sehnert, W.J., Zhou, S.K., Luo, J.: Adversarial sparse-view cbct artifact reduction. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 154–162. Springer (2018)

8. Liao, H., Lin, W.A., Huo, Z., Vogelsang, L., Sehnert, W.J., Zhou, S.K., Luo, J.: Generative mask pyramid network for ct/cbct metal artifact reduction with joint projection-sinogram correction. In: International Conference on Medical Image Computing and Computer-Assisted Intervention. pp. 77–85. Springer (2019)
9. Lin, W.A., Liao, H., Peng, C., Sun, X., Zhang, J., Luo, J., Chellappa, R., Zhou, S.K.: Dudonet: Dual domain network for ct metal artifact reduction. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 10512–10521 (2019)
10. Meyer, E., Raupach, R., Lell, M., Schmidt, B., Kachelrieß, M.: Normalized metal artifact reduction (nmar) in computed tomography. Medical physics (2010)
11. Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: International Conference on Medical image computing and computer-assisted intervention (2015)
12. Ulyanov, D., Vedaldi, A., Lempitsky, V.S.: Deep image prior. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2018)
13. Wang, J., Zhao, Y., Noble, J.H., Dawant, B.M.: Conditional generative adversarial networks for metal artifact reduction in ct images of the ear. In: Medical Image Computing and Computer Assisted Intervention – MICCAI 2018 (2018)
14. Yan, K., Wang, X., Lu, L., Summers, R.M.: Deeplesion: automated mining of large-scale lesion annotations and universal lesion detection with deep learning. Journal of Medical Imaging (2018)
15. Zhang, Y., Yu, H.: Convolutional neural network based metal artifact reduction in x-ray computed tomography. IEEE Trans. Med. Imaging **37**(6), 1370–1381 (2018)
16. Zhu, J.Y., Park, T., Isola, P., Efros, A.A.: Unpaired image-to-image translation using cycle-consistent adversarial networks. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2017)
17. Zhu, J., Park, T., Isola, P., Efros, A.A.: Unpaired image-to-image translation using cycle-consistent adversarial networks. CoRR **abs/1703.10593** (2017)

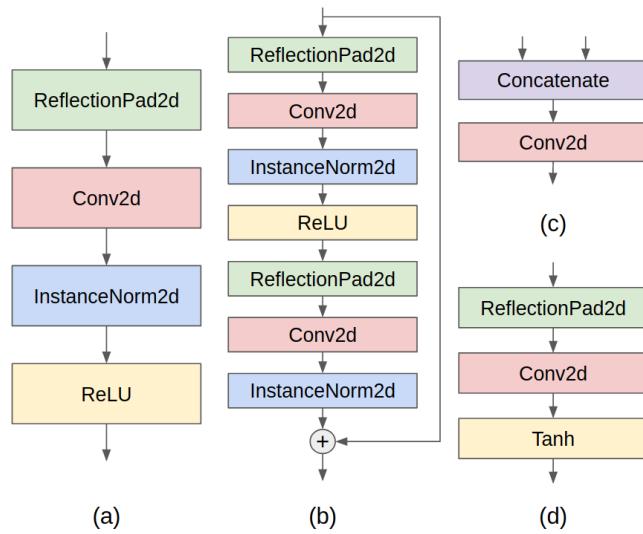# Supplementary Material

## A    Architecture Details



Fig. 5: Basic building blocks of the encoders and generators: (a) convolution block, (b) residual block, (c) merge block, and (d) final block. ReflectionPad2d stands for a reflection padding layer that we use to replace the zero padding of the conventional convolution layer.

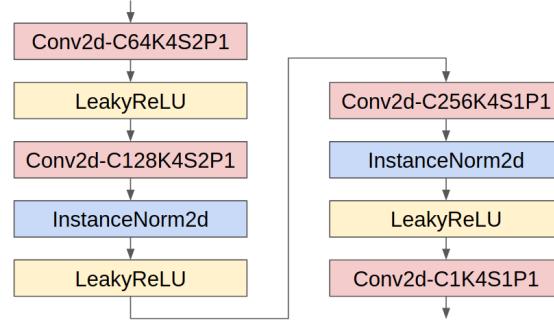Fig. 6: Architecture of the discriminator $D_{\mathcal{I}}$ or $D_{\mathcal{I}^a}$. We use 'C#K#S#P#' to denote the configuration of the convolution layers, where 'K', 'C', 'S' and 'P' stand for the kernel, output channel, stride and padding size, respectively.



Fig. 7: Architecture of the encoders and generators. (a) $E_{\mathcal{I}}$ or $E_{\mathcal{I}^a}$ (b) $G_{\mathcal{I}}$ (c) $E_a$ (d) $G_{\mathcal{I}^a}$. CB, RB, MB and FB are acronyms of the build blocks as illustrated in Fig. 5. The same as in Fig. 6, 'C#K#S#P#' denotes the configurations of the convolution layers in the blocks. For CB, RB, and FB, P is the padding of the reflection padding layer and the padding of the convolutional layer is zero. Note that the artifact code input for $G_{\mathcal{I}^a}$ are the hierarchical features encoded by $E_a$ and are merged with the corresponding outputs from $G_{\mathcal{I}^a}$.
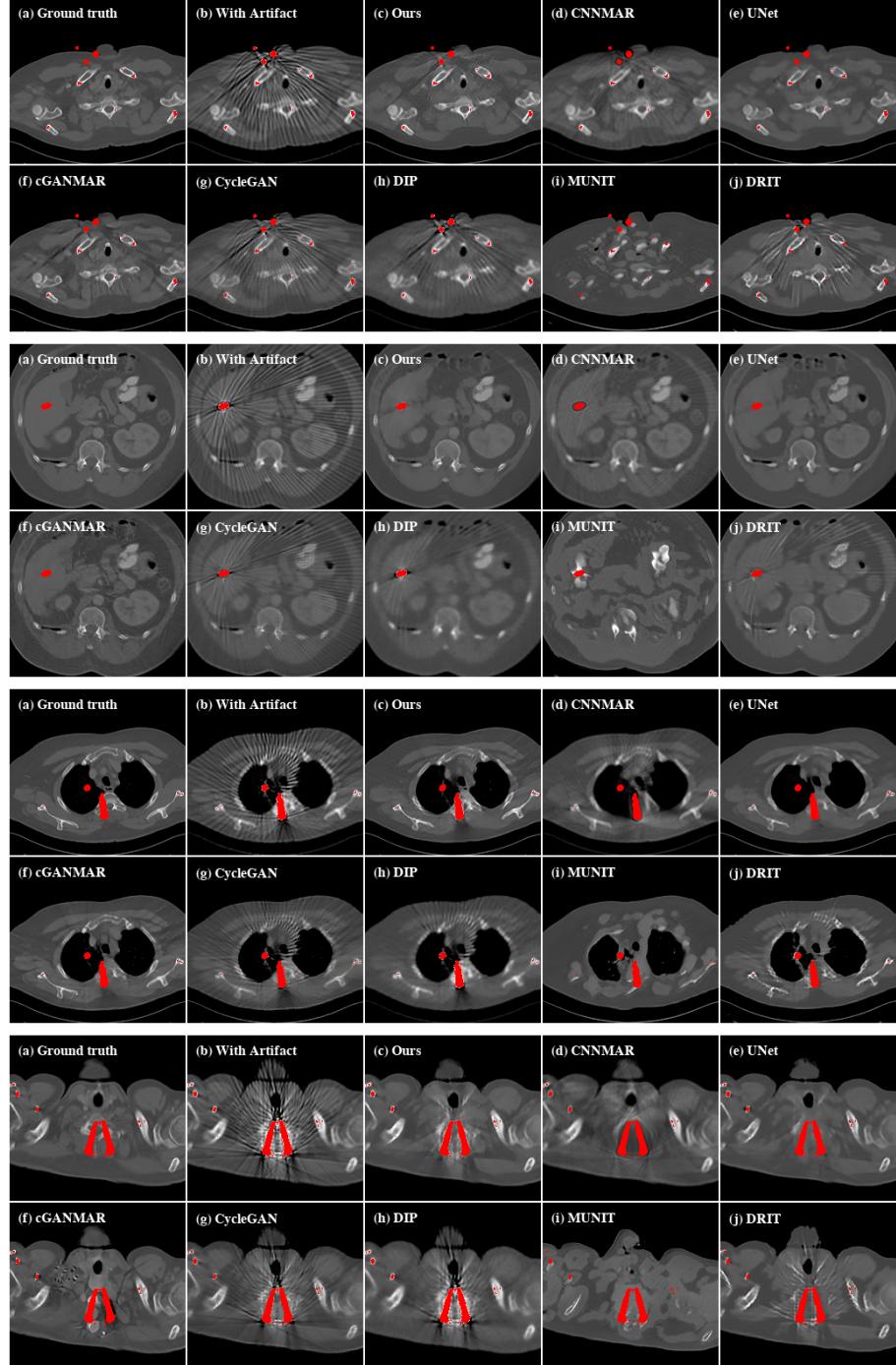
# B    Qualitative Results



Fig. 8: Qualitative evaluation results of SYN. For better visualization, we obtain the metal regions through thresholding and color them with red.
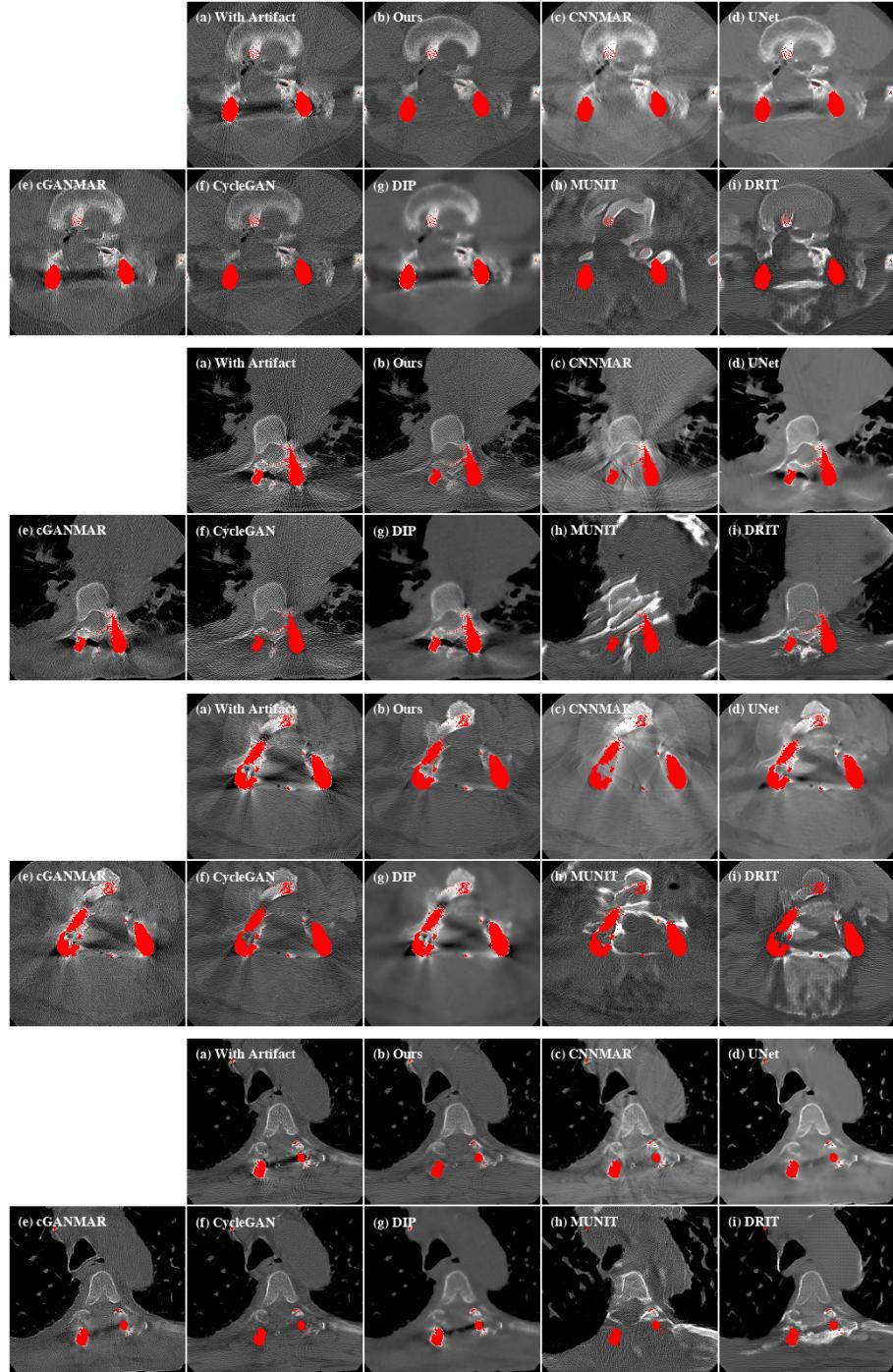
Fig. 9: Qualitative evaluation results of CL1. For better visualization, we obtain the metal regions through thresholding and color them with red.
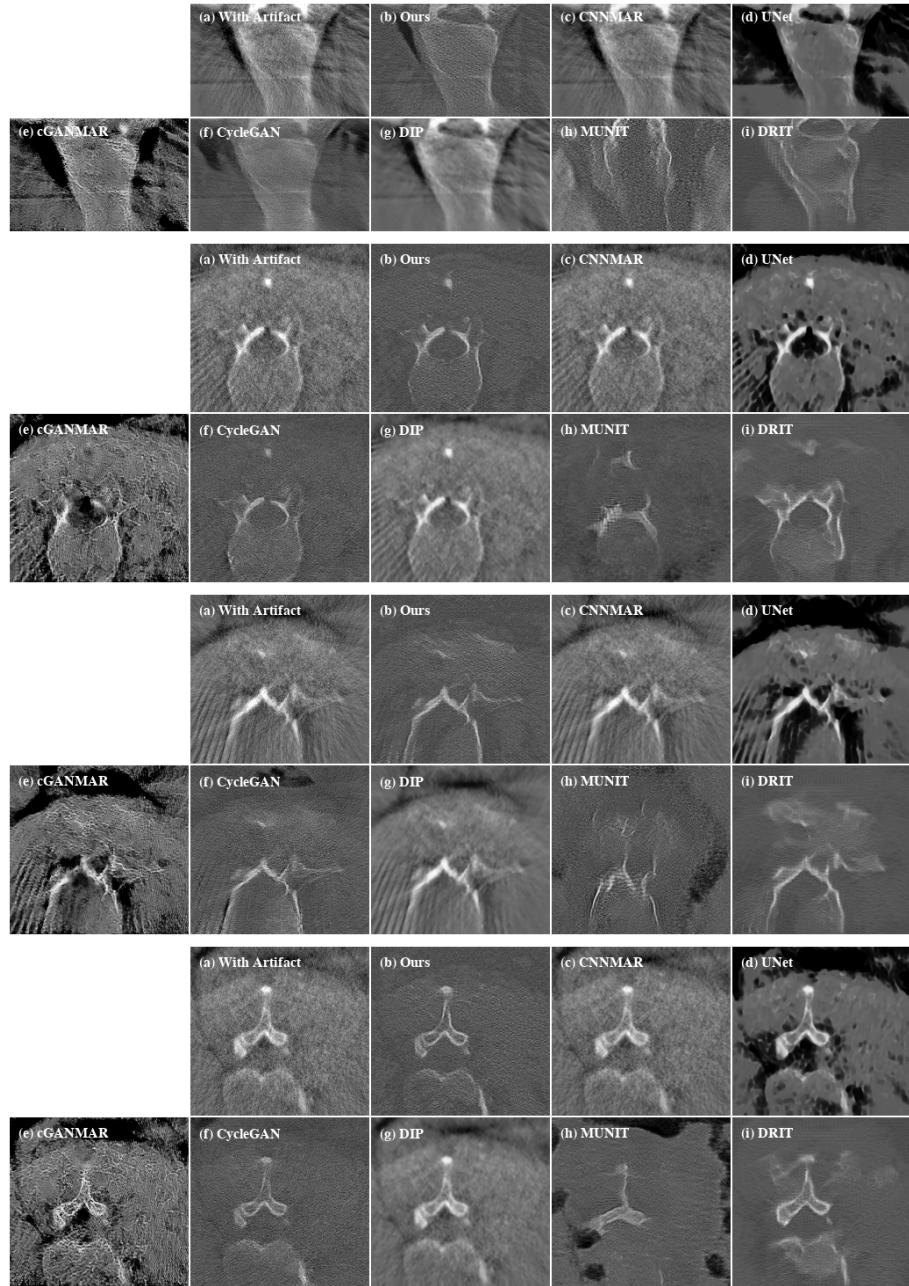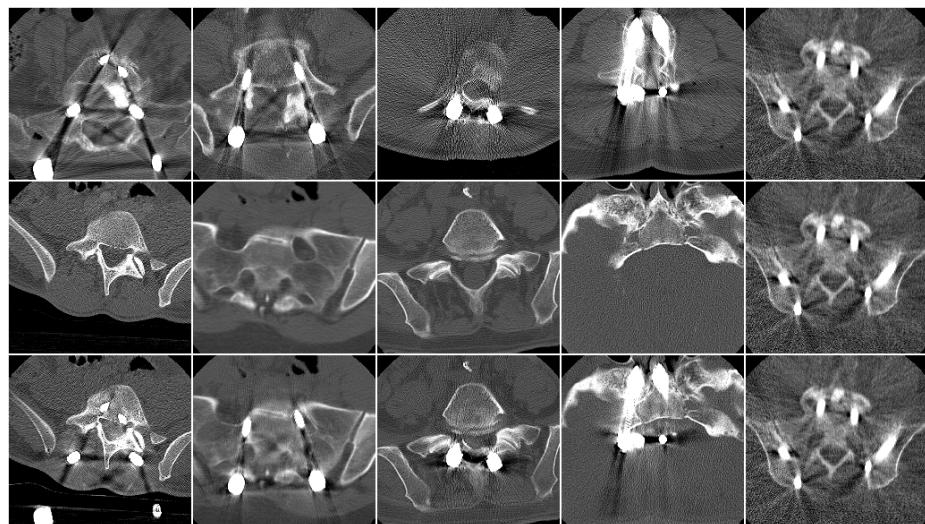
Fig. 10: Qualitative evaluation results of CL2.

Fig. 11: Metal artifact transferring. First row: the clinical images with metal artifacts. Middle row: the clinical images without metal artifacts. Last row: the metal artifacts in the first row transferred to the artifact-free images in the second row.