

GANai: Standardizing CT Images using Generative Adversarial Network with Alternative Improvement

Gongbo Liang[†], Sajjad Fouladvand^{*,†}, Jie Zhang[‡], Michael A. Brooks[‡], Nathan Jacobs[†], Jin Chen^{*,†,x}

* Institute for Biomedical Informatics, University of Kentucky, USA, Lexington, KY, USA

sajjad.fouladvand@uky.edu, chen.jin@uky.edu

† Department of Computer Science, University of Kentucky, Lexington, KY, USA

liang@cs.uky.edu, jacobs@cs.uky.edu

‡ Department of Radiology, University of Kentucky, Lexington, KY, USA

jie.zhang1@uky.edu, Michael.Brooks@uky.edu

x Department of Internal Medicine, University of Kentucky, Lexington, KY, USA

Abstract—Computed tomography (CT) is a widely-used diagnostic image modality routinely used for assessing anatomical tissue characteristics. However, non-standardized imaging protocols are commonplace, which poses a fundamental challenge in large-scale cross-center CT image analysis. One approach to address the problem is to standardize and normalize CT images using image synthesis algorithms including generative adversarial network (GAN) models. GAN learns the data distribution of training images and generate synthesized images under the same distribution. However, existing GAN models are not directly applicable to this task mainly due to the lack of constraints on the mode of data to generate. Furthermore, they treat every image equally, but in real applications, certain images are more difficult to standardize than the others. All these may lead to the lack-of-detail problem in CT image synthesis. We present a new GAN model called GANai to mitigate the differences in radiomic features across CT images captured using non-standard imaging protocols. Given source images, GANai composes new images by specifying a high-level goal that the image features of the synthesized images should be similar to those of the standard images. GANai introduces a new alternative improvement training strategy to alternatively and gradually improve GAN model performance. The new training strategy enables a series of technical improvements, including phase-specific loss functions, phase-specific training data, and the adoption of ensemble learning, leading to better model performance. The experimental results show that efficiency and stability of GAN models have been much improved in GANai and our model is significantly better than the existing state-of-the-art image synthesis algorithms on CT image standardization.

Index Terms—computed tomography, image synthesis, generative adversarial network, alternative training

I. INTRODUCTION

Computed tomography (CT) is one of the most popular diagnostic image modalities routinely used for assessing anatomical tissue characteristics for disease management [1], [2], [3], [4], [5]. CT scanners provide the flexibility of customizing acquisition and image reconstruction protocols to meet an individual’s clinical needs [6], [7]. However, CT acquisition parameter customization is a double-edged sword [8]. While it enables physicians to capture critical image features towards personalized healthcare, it forms a barrier to analyzing CT images in a large scale, a.k.a. radiomics [9], [10].

Capturing CT images with non-standardized imaging protocols may result in inconsistent radiomic features. As was revealed in a recent study, both intra-CT (by changing CT acquisition parameters) and inter-CT (by comparing different scanners with the same acquisition parameters) tests have demonstrated low reproducibility regarding radiomic features, such as intensity, shape, and texture, for CT imaging [11], [12]. In the example shown in Figure 1, each lung tumor was acquired twice using two different reconstruction kernels (B164 and Br40, Siemens Healthineers, Erlangen, Germany). The figure demonstrates that the appearances (as well as the radiomic features) of the same tumor can be strongly affected by the selection of CT acquisition parameters.

To overcome the barriers that prevent the use of CT images in large-scale radiomic studies, algorithms have been developed aiming to normalize and standardize CT images from multiple sources. Image synthesis is a class of algorithms that generate synthesized images from source images, which satisfy the condition that the feature-based distributions of the synthesized images are similar to that of target images [13]. Mathematically, given a source image x , an image synthesis algorithm composes a synthesized image x' by specifying a *high-level* goal that the image features of x' are significantly more similar to that of the target image y than the source image x . Image synthesis algorithms have been successfully used in image conversion and natural language processing, such as the synthesis of images from text descriptions [14].

Image synthesis algorithms can be roughly classified into two groups based on the techniques they use, i.e., traditional image processing-based algorithms [15], [16] and deep learning-based algorithms [17], [18]. In the first group, the histogram matching-based algorithm has been widely used [19], [20]. In general, it synthesizes images by mapping the histogram of source images to that of target images. However, finding the mapping function requires the presence of the target images, which are often missing or are not well defined in practice. In the second group, generative adversarial network models (GAN), a class of deep learning algorithms, can learn the data distribution of training data and generate

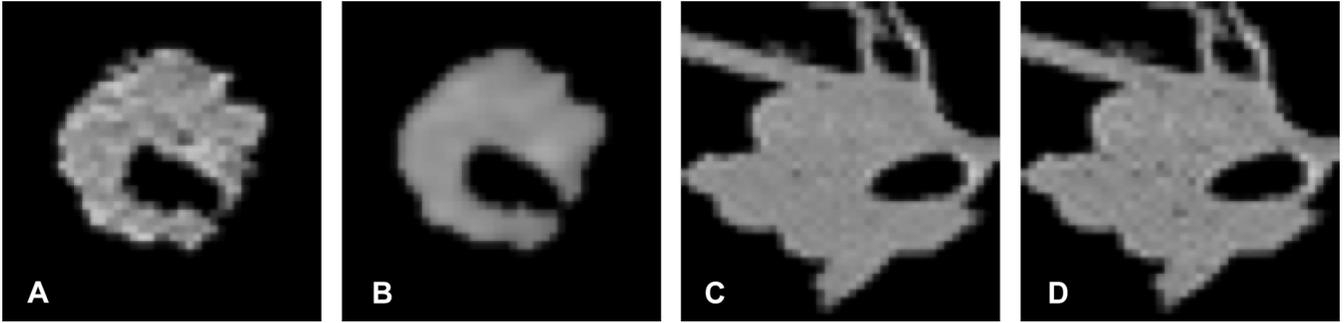


Fig. 1: Lung tumors acquired using two kernels have shown significantly different appearances as well as radiomic features. (A) Lung tumor 1 acquired with kernel BI64. (B) Lung tumor 1 acquired with kernel Br40. (C) Lung tumor 2 acquired with kernel BI64. (D) Lung tumor 2 acquired with kernel Br40.

synthesized examples which fall under the same distribution of the training [21]. In particular, the conditional generative adversarial network (cGAN), a special kind of GANs, learns the conditional distribution of the source image x given the target image y and then performs image transference from one domain to another [22], [23]. GAN and cGAN models have shown promising performance in image to image mapping tasks [23], [24], [25], [26], [27], [28] which motivated us to utilize them in this paper as a baseline. However, GAN models (including cGANs) are not directly applicable to our task mainly due to the limitations in model design and model training. Specifically, GAN models do not contain any constraints to control what modes of data it shall generate. The synthesized images are not guaranteed to be similar to the target images (Figure S1). Also, GAN models treat every image in training equally, but in real applications, some images are more difficult to synthesize than the others (Figure S2). All these limit the functionality of GAN models and may lead to the lack-of-detail problem in image synthesis.

To address the computational challenges in medical image synthesis, where great image details have to be maintained, we propose a novel deep learning framework called “Generative Adversarial Network with Alternative Improvement (GANai)”. GANai has a similar architecture as cGAN, but its training process is significantly different. cGANs characteristics such as training multi-modal models [22], learning conditional generative models, conditioning on an input image and generating its corresponding output image [23] motivated us to utilize cGANs as our base framework. Specifically, GANai introduces an alternative improvement training strategy to alternatively train its deep learning components and steadily improve the whole model performance. The adoption of the new training strategy enables a series of technical improvements, including phase-specific loss functions, component-dedicate training data, adoption of ensemble learning, and so on, leading to a clear improvement on model performance.

While GANai can be deployed in many applications, we adopted and evaluated GANai in standardization of the differences in radiomic features of CT images due to using non-standardized CT imaging protocols. The experimental results show that GANai is significantly better than the state-of-the-

art image synthesis algorithms, such as cGAN and histogram matching, on all the image acquisition parameters that we have tested. In summary, GANai has the following computational advantages:

- 1) GANai introduces an alternative improvement training strategy to alternatively and steadily improve model performance.
- 2) GANai adopts a new phase-specific loss function that allows the discriminator and the generator to collaborate rather than competing with each other.
- 3) GANai improves model training effectiveness by training the discriminator and the generator using specified training images.
- 4) GANai adopts ensemble learning to significantly improve the stability of GAN model training .

II. BACKGROUND

Radiomics is an emerging science to extract and use comprehensive radiomic features from a large volume of medical images for the quantification of overall tumor spatial complexity and the identification of tumor subregions that drive disease transformation, progression, and drug resistance [29], [30], [31], [32]. However, due to the use of non-standardized imaging protocols, variations in acquisition and image reconstruction parameters may cause inconsistency in radiomic features extracted from images, which poses a barrier to the practice of radiomics in large-scale [10], [30], [31].

A. CT Image Acquisition Parameters

In modern CT imaging, there are a large number of imaging protocols, and using non-standardized imaging protocols is common [6]. The CT image acquisition parameters include kV (the x-tube voltage), mAs (the product of x-ray tube current and exposure time), collimation, pitch, reconstruction kernel, field-of-view, and slice thickness [33], [34]. In routine clinical practice, certain parameters are often adjusted to meet the diagnostic needs, i.e., to obtain satisfactory image quality while maintaining low radiation dose to patients. Changing acquisition parameters may significantly affect the resulting images (Figure 1). For example, adjusting kV will change CT numbers (the pixel values of a CT image), changing mAs will

affect image noise rate, and the selection of reconstruction algorithms will result in different image texture features.

B. Histogram Matching

Histogram matching (or called histogram specification) is a widely-used image synthesis tool. It uses the intensity histogram to represent images and then transforms a source image to a target image by matching their intensity histograms [15], [19], [16], [20]. While histograms can represent the density of intensity in the whole image, the major drawback is the loss of location information. A variation of histogram matching is to divide a source image into multiple patches and to apply histogram matching on each patch, expecting that such patch-based representation may lead to location-specific image synthesis. However, patch-based histogram matching may introduce artifacts, esp. on the edges of patches. It is also sensitive to the selection of matching parameters such as the number of bins of a histogram (Figure S3).

C. Generative Adversarial Networks

Recently, deep learning has shown remarkable performance in various medical informatics tasks. For example, it has surpassed the human experts' performance on skin cancer classification by only looking at the dermoscopic images [35].

The generative adversarial network (GAN) is a kind of deep learning models that learns the data distribution of training images and generate synthesized images under the same distribution [21], [36], [37]. A GAN model usually has two components, i.e., the discriminator (D) and the generator (G), where G generates synthesized data from random noise, and D learns a data distribution from the training data and determines whether the synthesized data generated by G fall into the distribution. The goal of G is to generate synthesized data which are good enough to fool D , while D always aims to discriminate the synthesized data and the real data.

The conditional generative adversarial network (cGAN) is a kind of GAN models that learns the *conditional* distribution of the training data and generates synthesized data under the same condition [22], [38], [39]. Among cGAN models, the Image-to-Image model performs the image-to-image transference from one domain to another concerning the given condition, and it has become a widely recognized conditional image synthesis model [23]. Note that the images synthesized by cGANs are not necessarily similar to the target images, although they look "real", meaning having similar semantic meanings as the target images (see Figure S1, S2). However, in medical applications, it is important to maintain authenticity in the synthesized CT images. Specifically, it is expected to generate images with the distribution of radiomic features significantly similar to that of the target images.

While GAN models are advanced in image synthesis [23], [14], image inpainting [40], semantic segmentation [41], etc., GAN models are suboptimal regarding training efficiency and stability. To address the GAN training problem, several ensemble learning-based strategies have been applied to improve model training: 1) to train multiple GANs in parallel using a random initialization of model parameters, and then to

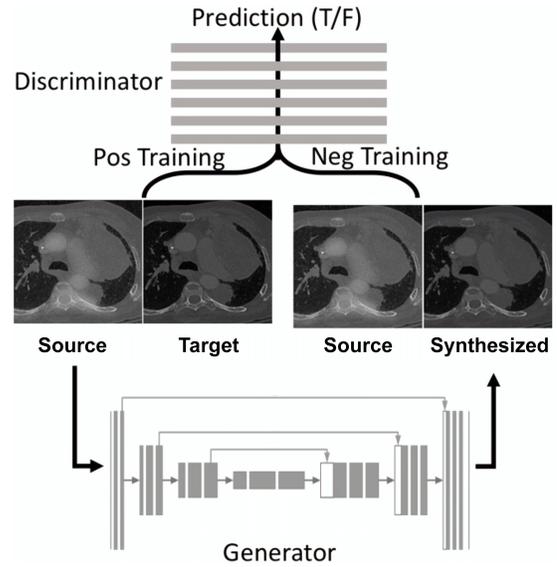


Fig. 2: Architecture of GANai. Given a source image, the generator G synthesizes a new image to fool the discriminator D , while D aims to distinguish the synthesized image and the target image.

randomly choose one of the GANs to generate the synthesized data [42]; 2) to train multiple D s and requires the G to fool a group of D s [43]; and 3) to select training data using boosting and to train a cascade of GANs in sequence. It has been shown that the performance of GANs can be significantly improved by using ensemble learning [43].

III. METHOD

To extend the adversarial learning into the medical image domain and to address the aforementioned challenges, we propose Generative Adversarial Network with Alternative Improvement (GANai).

A. Architecture

GANai consists of two components, i.e., the generator (G) and the discriminator (D), where G is a U-Net with fifteen hidden layers and D is a multilayer perceptron model with six fully connected layers [44]. The architecture of GANai is similar to the cGAN models, shown in Figure 2 [23]. The inputs of D of GANai are image pairs (x, y) and (x, x') , where (x, y) denotes the real pair (positive training), and (x, x') denotes the fake pair (negative training). The goal of D is to distinguish the real pairs from the fake pairs. Given the feedback from D , G learns the mapping from X to Y and generates a synthesized image x' for any given source image x ($x \in X$) in Y 's domain. In contrast to D , G aims to synthesize images that can fool D . If D can distinguish most of the fake pairs from the real pairs, the performance of G needs to be further improved. Otherwise, we conclude that the generative results of G are good enough for the current D .

B. Alternative Improvement

In traditional GAN models, D and G are trained synchronously (D and G trained together) or asynchronously (several batches of D -training followed by several batches of

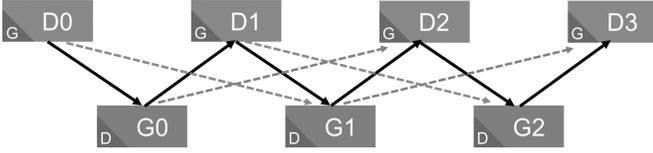


Fig. 3: In each training phase of GANai, D (or G) is trained while the other component is frozen. The name of a block (such as D_0 or G_1) indicates the component in training, and the letter located at the bottom left corner indicates the component that is frozen. The alternative training (solid line) ensures high performance while the ensemble approach (dotted line) improves the training stability.

G -training), based on the assumption that both D and G can be gradually improved together. In practice, however, if D is not well trained to capture the intrinsic features to separate a real and a fake image, G can easily fool D . Similarly, if G is not well “challenged” by D , its model performance is not guaranteed to be improved.

We introduce the alternative training approach for GANs (Figure 3). As the name suggested, GANai has two alternate training phases, i.e., the discriminator training (D -training) and the generator training (G -training). In each training phase, we focus on optimizing one of the components while freezing the other. A training phase will stop if the current component is well trained or the training step exceeds an upper bound. After that, we switch to the other training phase (Figure 3 solid lines). See Section V-B for more details. The alternative training strategy enables a series of technical improvements, including phase-specific loss functions, phase-specific training data, and the adoption of ensemble learning, which will be introduced in the following subsections.

C. Loss Functions

The alternative training of GANai may boost model performance by preventing each component being too strong or too weak. In the literature, strategies have been presented to freeze part of a GAN when the GAN components are imbalanced [45]. However, it is difficult to decide when to freeze/unfreeze a component of GAN. To address this issue, we redesigned the loss functions.

In the D -training phase, G is frozen so that D learns the differences between the synthesized images and the target images and discriminates the synthesized images. Hence, the loss function of D is the same discriminator loss of cGAN [22]:

$$\begin{aligned} \mathbb{L}_{Phase_D}(D) = & \mathbb{E}_{x,y \sim P_{data}(x,y)} [-\log D(x,y)] + \\ & \mathbb{E}_{x \sim p_x, z \sim p_z(z)} [-\log(1 - D(x, G(x,z)))] \end{aligned} \quad (1)$$

where x is the source image; y is the target image; $G(x, z)$ is the synthesized image generated by G , which maps the source image x and a random noise vector z to y ; $D(x, y)$ is the prediction result of the real pair; and $D(x, G(x, z))$ is the prediction result of the fake pair. For $D(x, y)$, the higher the prediction accuracy, the higher the value of $D(x, y)$.

In the G -training phase, D is frozen, and it evaluates the results of G . Since we expect G to fool D , the loss of D in

the G -training phase is defined as:

$$\begin{aligned} \mathbb{L}_{Phase_G}(D) = & \mathbb{E}_{x,y \sim P_{data}(x,y)} [-\log D(x,y)] + \\ & \mathbb{E}_{x \sim p_x, z \sim p_z(z)} [-\log(D(x, G(x,z)))] \end{aligned} \quad (2)$$

Finally, by integrating Eq 1 and Eq 2, the loss function of D in GANai is defined as:

$$\begin{aligned} \mathbb{L}(D) = & \mathbb{E}_{x,y \sim P_{data}(x,y)} [-\log D(x,y)] + \\ & (\mathbb{E}_{x \sim p_x, z \sim p_z(z)} [-\log D(x, G(x,z))])^\alpha + \\ & (\mathbb{E}_{x \sim p_x, z \sim p_z(z)} [-\log(1 - D(x, G(x,z)))]^{1-\alpha} \end{aligned} \quad (3)$$

where parameter $\alpha = 1$ if GANai is in the G -training phase and $\alpha = 0$ in the D -training phase.

The loss function of G is the same as Isola et al. [23]. Also, we adopt the $L1$ loss as the regularization factor.

$$\begin{aligned} \mathbb{L}(G) = & \mathbb{E}_{x,G(x,z)} [-\log D(x, G(x,z))] + \\ & \beta \mathbb{E}_{G(x,z),y} [\|y - G(x,z)\|] \end{aligned} \quad (4)$$

where β is the weight of the regularization term.

To determine when to switch between the D -training phase and the G -training phase, the prediction accuracy on the fake image pairs ($D(x, x')$) is used. The value of $D(x, x')$ is computed at every training step and is compared with two thresholds. More specifically, if $D(x, x') \leq T_l$, GANai will switch from D -training to G -training. If $D(x, x') \geq T_h$, GANai will switch from G -training to D -training. T_l and T_h are the lower and upper thresholds of $D(x, x')$. To improve training stability, the least amount of steps (minibatches) of each training phase is also specified. Note that in GANai, the value of $D(x, x')$ increases and decreases, indicating that the performance of D and G is improved alternatively.

D. Training D and G with Dedicated Training Data

Since the components of GANai are trained separately, one idea motivated by the curriculum learning is to increase model training efficiency by training G and D using different data [46]. More specifically, the images that are potentially synthesizable can be used to accelerate the G -training, while the training of D can benefit from images that are difficult to synthesize.

We develop a procedure to select training data for D and G . First, a cGAN model is trained using all the training data [23]. Second, with the trained cGAN model, we synthesize a new image for every source image and compare every synthesized image with its corresponding target image using Kullback-Leibler divergence [47], normalized mutual information (NMI) [48], and cosine similarity. Finally, the training data is split into two subsets based on z-score, i.e., 1/3 of the source-target image pairs with the highest similarities between synthesized images and target images (called T_{easy}) and 1/3 of the images with the lowest similarities (call T_{hard}). The new procedure allows us to train G using T_{easy} and train D using T_{hard} (see Section V-A for other training set selection strategies).

E. Improving Training Stability using Ensemble Learning

Due to the nature of the generative adversarial concept (i.e., open-ended competition between GAN components), it is not guaranteed that G or D will improve towards the same direction. For example, if the k th state of G fools the $(k-1)$ th state of D , it still may be classified by the older $(k-2)$ th state of D . Therefore, during the two-phase training of GANai, we improve the model stability by adopting the ensemble learning. Simply speaking, a D is required to discriminate multiple G s and a G must fool multiple D s.

Mathematically, the following criteria are specified in GANai: when training the k th G , the G must fool both the $(k-2)$ th state and the $(k-1)$ th state of D , and when training k th D , the D should discriminate both the $(k-2)$ th state and the $(k-1)$ th state of G . For an illustrative example, see the dot lines in Figure 3. These criteria can be further extended to incorporate more historical D s or G s or more sophisticated conditions. In the exception that GANai cannot identify such a D or G that satisfies the criteria after at most T_s steps (the maximum training step in each phase), it will roll back to the previous state, and re-train the current component.

IV. EXPERIMENTAL RESULTS

A. Data

In total 2,448 chest CT image slices of lung cancer patients were collected using Siemens CT Somatom Force at the University of Kentucky Medical Center. For each patient, a CT image was constructed with each of the possible combinations of two image reconstruction parameters, i.e., slice thickness (0.5, 1, 1.5, 3mm) and reconstruction kernels (B157 and B164). With data augmentation, the training data has been extended to 14,958 image patch pairs. Among them, 7,479 assigned as T_{easy} and 7,479 assigned as T_{hard} using the procedure introduced in Section III-D. Each image pair contains a source image x and the target image y . See details of data augmentation in Section S1.A with examples in Figure S4.

The validation data contains 3,554 2.5D images, and multiple radiomic features were extracted for model validation. Specifically, we randomly cropped 2.5D images from the CT images that have not been used as training data, with their dimensions ranging from $5 \times 5 \times 5$ to $60 \times 60 \times 30$ pixels. When cropping the 2.5D validation images, we excluded areas with bone or air, since soft tissues are what physicians are most interested. See Section S1.B for more details.

Given a large number of CT imaging protocols, it is impractical to apply all of them. We selected two image reconstruction parameters (kernel and slice thickness) and used all the combinations for the model performance test. Also, we chose 1mm slice thickness and B164 kernel to be the standard imaging protocol, since it is widely used in the current lung cancer radiomic studies. Note the settings can be easily extended to incorporate more acquisition parameters or to use a different standardized imaging protocol.

B. Implementation Details of GANai

In GANai, G is a fifteen hidden layers U-Net [44], with the size between $128 \times 128 \times 64$ and $1 \times 1 \times 512$ (Figure S5). The input of G are 256×256 images, and the synthesized images have the same image size. D is implemented as a multilayer perceptron model with six fully connected layers with the size between $256 \times 256 \times 3$ and $30 \times 30 \times 1$ (Figure S6).

The training of GANai started with the D -training phase, and all the network weights were randomly initialized. We set the regularization term weight $\beta = 100$ to reduce the visual artifacts [23], and used $T_l = 0.05$ and $T_h = 0.95$ as the training phase switch thresholds, and $T_s = 10$ epochs as the maximum training step. Within each training phase, the model needed to be trained for at least five steps before switching to the other training phase. GANai was trained for 100 epochs with learning rate being 0.0002, momentum being 0.5.

GANai is deployed on Tensorflow [49] on a Linux computer server with eight Nvidia GTX 1080 GPU cards. It took 15 hours to train GANai from scratch using a single GPU card. Using the trained model, it took 0.2 seconds to generate a synthesized image (5 images per second).

Figure 4 shows the discriminator prediction results on all the fake pairs $D(x, x')$ in the first 150 steps of training. With the training of D , $D(x, x')$ decreases. When the value of $D(x, x')$ is below T_l (in our experiment, $T_l = 0.05$), GANai is switched to the G -training phase. In the G -training phase, $D(x, x')$ increases, since D is frozen and the performance of G keeps increasing. When the value of $D(x, x')$ is higher than T_h ($T_h = 0.95$), GANai is switched to the D -training phase.

The training and validation loss of D and G in the first 150 training steps are shown in Figure 5. Both the training and validation loss of D decreased in every training phase, which indicates the model performance of D and G was improved alternatively. In the D -training phase, if the performance of D is increased, the loss of D will reduce, since both $-\log(D(x, y))$ and $-\log(1 - D(x, x'))$ are both reduced (solid lines in Figure 5A). When switching from the D -training to the G -training phase, α in the loss function of D flips from 0 to 1, which immediately turns the loss of D from a small value to a high value (see the jumps located at phase turning points in Figure 5A). In the G -training phase, if the performance of G is increased, the performance of D will decrease, so the loss of D decreases (dotted lines in Figure 5A). Figure 5B shows the loss of G increases in D -training phase (due to the performance improvement of D) and decreases in G -training phase, since the performance of G is improved (See Section V-C).

C. Evaluation Metric

For performance evaluation, we compared GANai with cGAN [23] and the patch-based histogram matching (see details in supplementary section III). Instead of hiring human annotators, we adopt the radiomic features for performance evaluation [50], [51]. Specifically, two classes of radiomic features were used for model performance evaluation, i.e., 2.5D texture features (i.e., gray-level co-occurrence matrix) and 2.5D intensity histogram based features. In total, eight

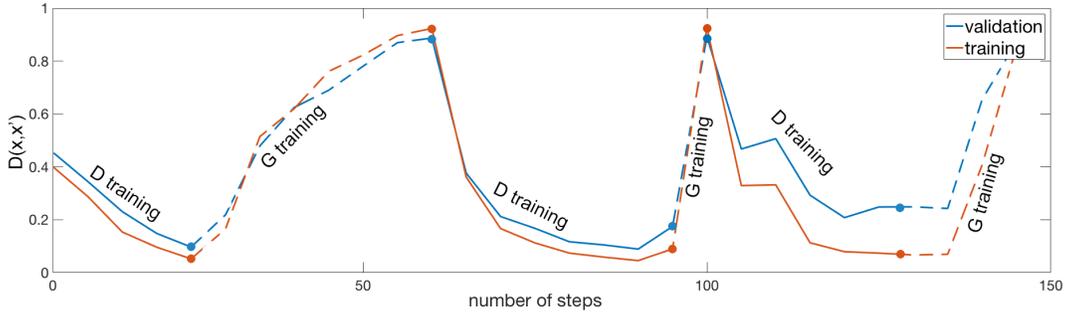


Fig. 4: The prediction results of D on the fake image pairs (x, x') in the first 150 steps of the alternative training. For $D(x, x')$, the higher the prediction accuracy, the lower the value ($D(x, x') \in [0, 1]$).

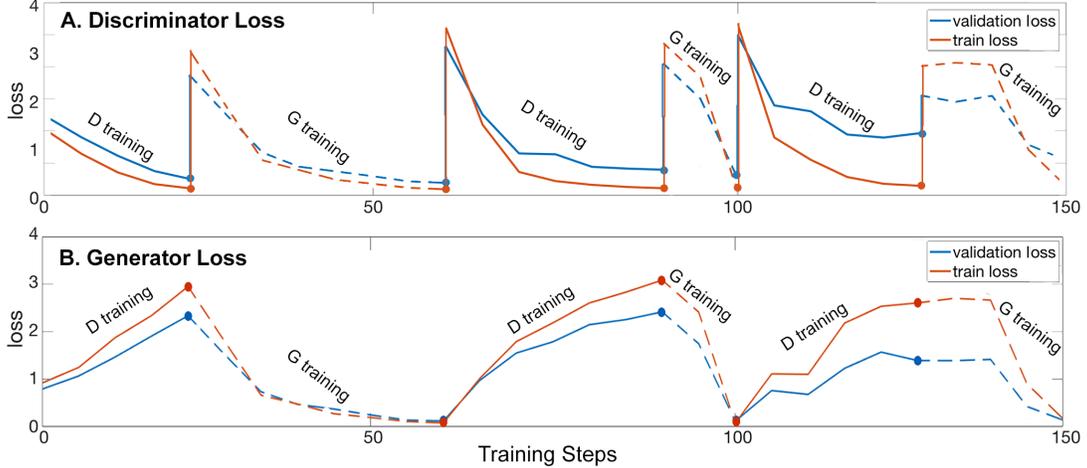


Fig. 5: The training loss and the validation loss of D and G in GANai in first 150 steps of training. The solid lines indicate the loss of D in the D -training phase. The dotted lines indicate the loss of D in the G -training phase. The solid points indicate the time when GANai switches between the D -training phase and the G -training phase.

radiomic features were adopted for performance evaluation (see Section S2 for details).

Per every radiomic feature to test, we compared each synthesized image and its target image, and computed the absolute error and relative error using the following equations:

$$\text{abs_err}(\text{feature}_k, m) = \frac{|\text{feature}(\text{synthesized}, k, m) - \text{feature}(\text{target}, k)|}{\text{feature}(\text{target}, k)} \quad (5)$$

where feature_k is the k th radiomic feature, m is either GANai or a image synthesis model to compare.

$$\text{rel_err}(\text{feature}_k, m_1, m_2) = \frac{\text{abs_err}(\text{feature}_k, m_1) - \text{abs_err}(\text{feature}_k, m_2)}{\text{error}(\text{feature}_k, m_1)} \quad (6)$$

where m_1 and m_2 are two different image synthesis models. For the relative error, a positive value indicates that m_2 has smaller error than m_1 , vice versa.

Model stability is evaluated using the cumulative sum control chart (CUSUM) [52]. CUSUM is a sequential analysis model typically used for monitoring change detection [53]. In CUSUM, the differences between any two adjacent values (in our case, the absolute errors between any two adjacent saved

model states) are measured and are compared with a threshold. CUSUM is computed as the number of the difference values higher than a threshold (called out-of-control points). In our experiment, a series of CUSUM values were generated for each model using multiple thresholds. The normalized sum of the CUSUM values, which is the smaller the better, was used for model stability evaluation.

D. Performance Evaluation Results on Generator

The absolute errors on all the tested radiomic features are shown in Table I. For the detailed feature-based errors, see Figure S7. On the texture features, the mean absolute error of histogram matching over all six features is 0.37. cGAN reduces it to 0.13, and GANai further reduces the absolute error significantly to 0.08 (two sample t-test p -value ≤ 0.01). On the intensity histogram features, GANai decreases the absolute errors by 17.77% from cGAN, and 79.05% from histogram matching. The results indicate that GANai is significantly better than cGAN and patch-based histogram matching.

Table II shows the relative errors of GANai and cGAN on seven sets of the validation data generated using different combinations of CT acquisition parameters. A positive value indicates the error of GANai is lower than cGAN, while a neg-

TABLE I: Averaged absolute errors (SD) of (1) the texture features and (2) the intensity histogram features computed using histogram matching, cGAN, and GANai. In all of them, GANai has the smallest errors (cGAN and GANai two sample t-test p -value ≤ 0.01).

Absolute Error	Hist. Matching	P2P	STAN-CT
Contrast ¹	0.21 \pm 0.15	0.12 \pm 0.08	0.09 \pm 0.06
Correlation ¹	0.18 \pm 0.13	0.18 \pm 0.12	0.09 \pm 0.07
Dissimilarity ¹	0.15 \pm 0.11	0.09 \pm 0.06	0.06 \pm 0.04
Energy ¹	0.47 \pm 0.28	0.19 \pm 0.14	0.14 \pm 0.11
Entropy ¹	0.09 \pm 0.06	0.02 \pm 0.01	0.01 \pm 0.01
Homogeneity ¹	0.28 \pm 0.16	0.10 \pm 0.06	0.07 \pm 0.05
Kurtosis ²	0.54 \pm 0.27	0.18 \pm 0.14	0.15 \pm 0.11
Skewness ²	0.51 \pm 0.27	0.16 \pm 0.12	0.14 \pm 0.11

TABLE II: Averaged relative errors on the texture features¹ and the intensity histogram features² by comparing cGAN and GANai. Positive values mean GANai is better, and negative values mean cGAN is better. Overall, GANai has smaller errors than cGAN.

Relative Error	B17 0.5mm	B17 1mm	B17 1.5mm	B17 3mm	B164 0.5mm	B164 1.5mm	B164 3mm	Overall
Contrast ¹	-0.16	0.00	0.13	0.36	0.42	-0.46	0.34	0.25
Correlation ¹	0.06	0.38	0.37	0.45	0.68	-0.10	0.38	0.50
Dissimilarity ¹	-0.05	0.28	0.32	0.48	0.64	-0.30	0.39	0.33
Energy ¹	-0.19	0.35	0.34	0.55	0.11	-1.61	-0.12	0.26
Entropy ¹	-0.05	0.30	0.30	0.48	0.67	-0.81	0.19	0.50
Homogeneity ¹	-0.34	0.29	0.30	0.48	0.75	-0.85	0.26	0.30
Kurtosis ²	-0.05	0.18	0.26	0.45	-1.66	-1.84	-0.48	0.17
Skewness ²	-0.10	0.17	0.47	0.35	-1.66	-1.84	-0.18	0.13

ative value indicates the error of GANai is higher than cGAN. The results show that GANai outperforms cGAN on five out of seven validation subsets, on which GANai decreased the relative errors by 36.21% on average. For example, on the texture features, GANai reduces the relative error by 54.48% on the B164 kernel with 0.5mm slice thickness images. For the detailed feature-based errors, see Figure S8-S15.

Figure 6 shows an example of the synthesized images using cGAN or GANai generated after 100 training epochs. The GANai synthesized image is more similar to the target image, has sharper edges, and has fewer artifacts than cGAN. Figure 7 shows both cGAN and GANai model reaches their best performance after 20 epochs of training. After that, GANai can still maintain high synthesized image quality, but cGAN started to introduce artifacts.

E. Performance Evaluation Results on Discriminator

To evaluate the performance on the discriminator D , we generated a fake-pair-only dataset and used it to measure the prediction accuracy of all the D s in the model training process. Specifically, given a fixed source image set X_{val} and the correspondent target image set Y_{val} , each having 1,750 images, we generated the synthesized image set X'_{val} using the second last generator of GANai. The accuracy of every discriminator (such as D_0 to D_3 in Figure 3) in the alternative training process of GANai was measured with all image pairs in X_{val} and X'_{val} . Accuracy is defined as the proportion of (x, x') that were correctly classified as the fake image pairs. Figure 8 shows the prediction accuracy of D at every training process. The increasing prediction accuracy shows the performance of D was steadily improving during the training of GANai.

F. Performance Evaluation Results on Training Stability

In GANai, an ensemble learning-based approach is adopted to increase the training stability. To demonstrate the effectiveness of this approach, we designed the following experiment. Three networks (cGAN, $GANai_{singleDG}$, and GANai) were trained for 100 epochs using the same training data, where $GANai_{singleDG}$ is a simplified version of GANai that trains the current component only based on the previous counter component, without using multiple D s or G s. The training state of every 2.5 training epochs was saved. We compared all the three models using the same validation data at every saved model state (Figure 9A). The normalized sum of the CUSUM values of cGAN, $GANai_{singleDG}$, and GANai over all the six texture features are 0.21, 0.15, and 0.13 respectively, indicating GANai is the most stable model among the three. Figure 9 shows the CUSUM on the contrast feature computed using the gray-level co-occurrence matrix.

V. DISCUSSION

A. Training Effectiveness

The training data in GANai are separated into two subsets for the training of G and D . Our assumption is that for certain source images that are difficult to standardize, we should avoid them in the G -training phase. Instead, we use them to train D . To test the assumption, we trained a new GANai model called $GANai_{reverse}$ with the opposite training data assignment (i.e., G trained with T_{hard} and D trained with T_{easy}). Figure 10 shows that the mean absolute errors of $GANai_{reverse}$ are significantly higher than GANai on a majority of the features, indicating that training data assignment is critical for improving GAN performance.

We further tested the effectiveness of the new strategies developed for improving training effect. Two modified cGAN models were trained, one with dedicated training data, i.e., T_{hard} for D and T_{easy} for G , called $cGAN_{SpDa}$, and the other further adopting the alternative training strategy, called $cGAN_{SpDa+AI}$. Experimental results show that 1) $cGAN_{SpDa}$ can effectively reduce the feature-based absolute errors of cGAN on a majority of the texture features, and 2) $cGAN_{SpDa+AI}$ can further reduce the absolute errors on texture features (Figure 11). It indicates that the new training strategies developed in GANai are effective and can be adopted by generic GAN models to further improve their performance.

Though the data separation step may cost as same as the GANai training. The dataset separation step is part of the data collection and preparation task. The cost of this step should not contribute to the cost of GANai model training. In addition, our experiments show that by training D and G with dedicated training data, the model performance improved significantly (e.g., error reduced over 55% on the feature of gray-level co-occurrence matrix correlation). Thus, we believe that the proposal of training D and G with dedicated training data is necessary to GAN model training.

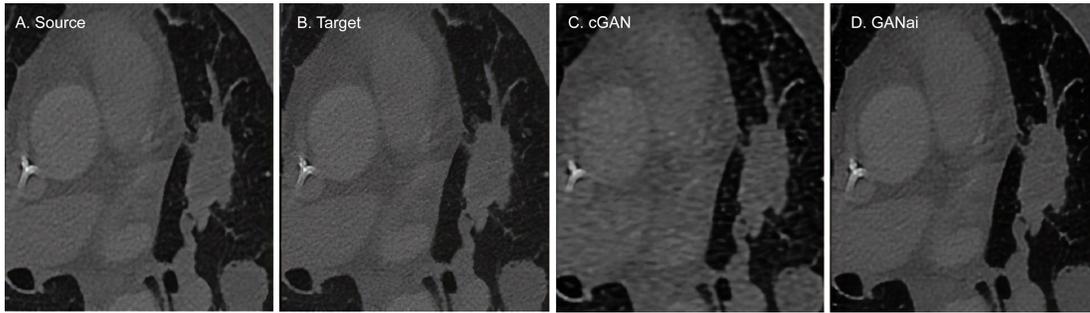


Fig. 6: Examples of the synthesized images generated by cGAN and GANai at 100th epoch. (A) source image. (B) target image. (C) cGAN synthesized image. (D) GANai synthesized image.

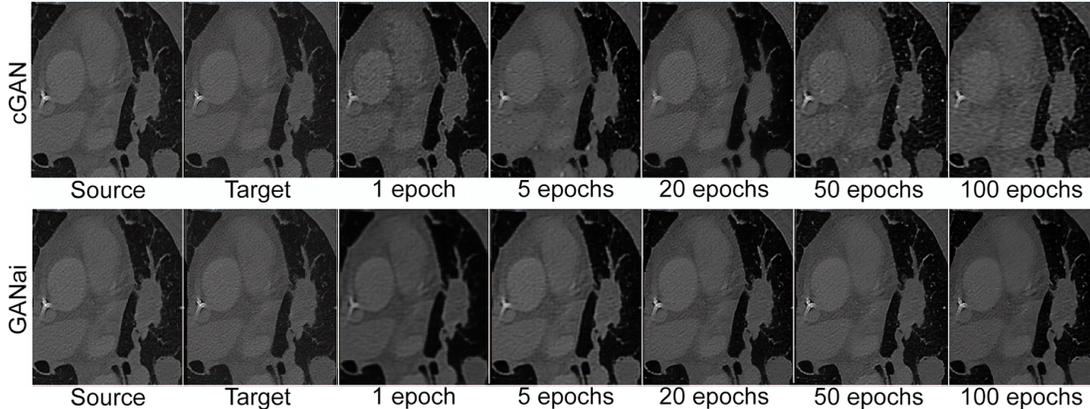


Fig. 7: Examples of the synthesized images generated by cGAN and GANai at multiple training steps. The first two columns are the source and target images. Both cGAN and GANai reached their best performance at about 20 training epochs. The synthesized images generated by cGAN have obvious artifacts and have less sharp edges than that of GANai. Furthermore, GANai maintained a high synthesized image quality in the continuous training after the first 20 epochs, whereas cGAN started to introduce additional artifacts into the synthesized images.

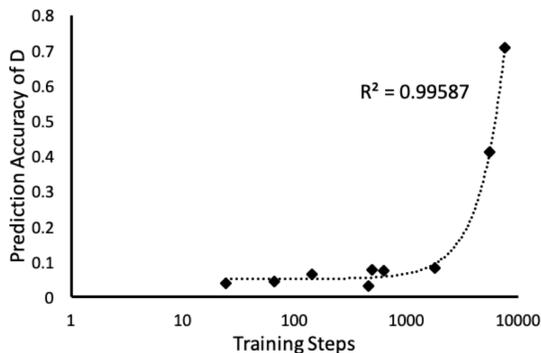


Fig. 8: Prediction accuracy of D_s gradually increases during the alternative training process of GANai.

B. Alternative Improvement and Effectiveness of Ensemble Learning

The alternative improvement strategy of GANai is different from the traditional training procedure of GANs where D and G components are trained synchronously (D and G trained together) or asynchronously (several batches of D -training followed by several batches of G -training). In fact, our freezing strategy and the algorithm we developed to decide on when and how to freeze a component is one of the contributions of this work. Traditional GANs switch to training

the other component once they are done with processing a batch of data (or once they are done with training the current component for a constant number of steps); however, GANai trains components dynamically. GANai switches to training the other component only when the current component which is being trained is efficiently trained based on a predefined performance threshold.

GANai adopts the alternatively improving strategy to train D and G so that both modules can be optimized in each iteration of training. One potential problem of such full optimization is that the model could be trapped at the local minima instead of reaching the global optimization. One such example is shown in Figure S14, where a generator has been trained for more than five epochs, but it still did not result in any significant improvement. It is reasonable to believe that the model was trapped at a local minima. To address this issue, we adopt the ensemble learning approach, i.e., GANai requires a D to discriminate multiple G s and a G to fool multiple D s. Also, we rollback to the previous training phase and then retrain the model, if a satisfactory loss cannot reach in a reasonable amount of time.

C. Validation Loss

The validation loss of G in Figure 5B is constantly lower than the training loss, which is uncommon to machine learning

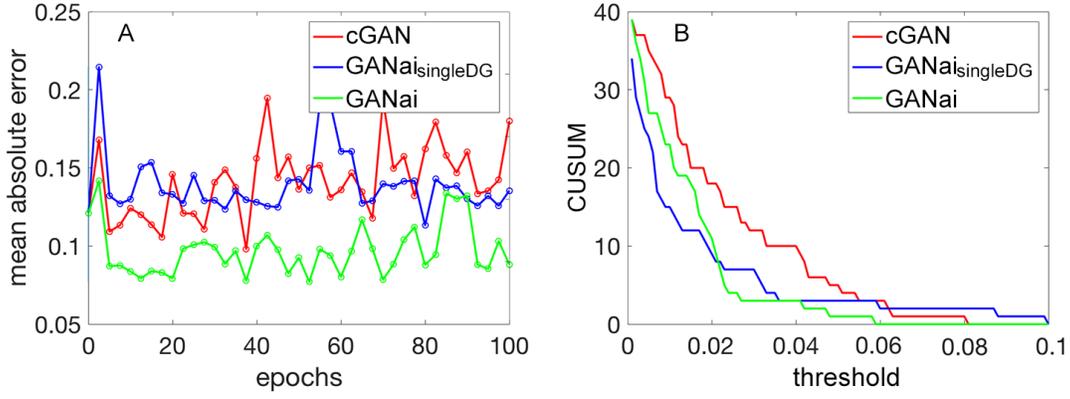


Fig. 9: Performance evaluation on training stability. (A) the mean absolute errors of cGAN, $GANai_{singleDG}$, and GANai on the contrast feature computed using the gray-level co-occurrence matrix. (B) the CUSUM values of cGAN, $GANai_{singleDG}$, and GANai, where the x-axis is the threshold of CUSUM, and the y-axis is the CUSUM value. In general, GANai is the most stable model among the three.

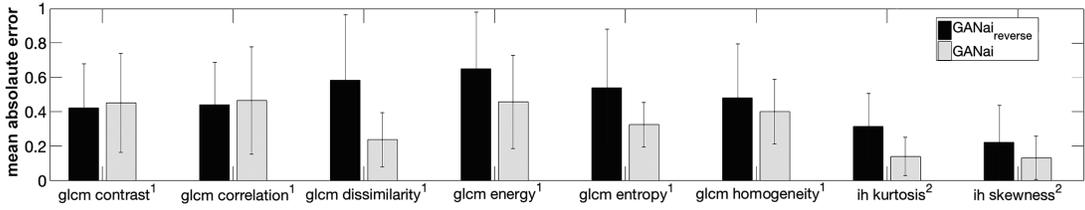


Fig. 10: Averaged feature errors for the data effectiveness test. ¹ Gray-level co-occurrence matrix, ² Intensity Histogram. It shows that the mean absolute errors of $GANai_{reverse}$ are significantly higher than GANai on a majority of the features, indicating that training data assignment is critical for improving GAN performance.

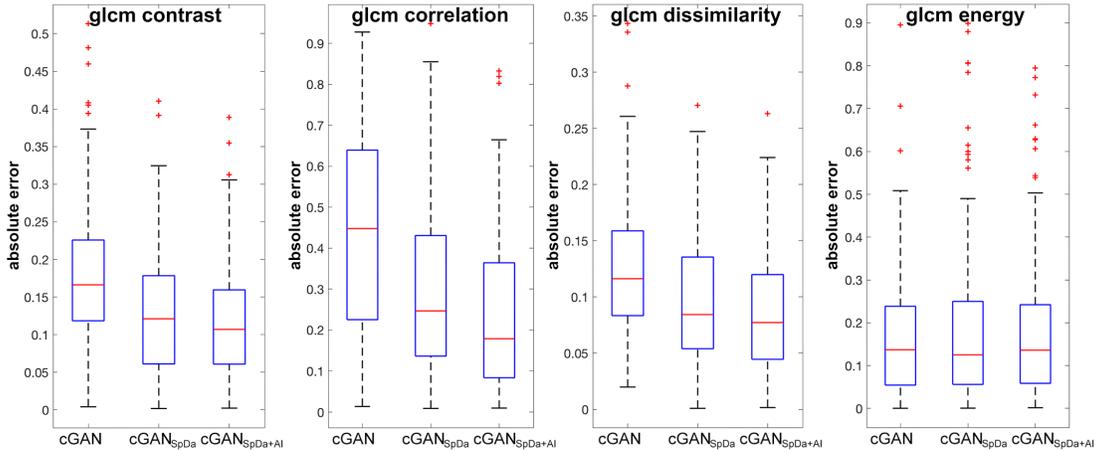


Fig. 11: Gray-level co-occurrence matrix feature errors of different cGAN versions. $cGAN_{SpDa}$ was trained with dedicated training data. $cGAN_{SpDa+AI}$ was further adopting the alternative training strategy.

tasks. This is reasonable because the loss of G is $-\log D(x, x')$ computed using the prediction result on all the fake image pairs. As shown in Figure 4, the value of $D(x, x')$ on the validation dataset is higher than that on the training dataset. After taking the minus log, the validation loss is smaller than the training loss. However, as stated in Gulrajani et al [54], the loss of GANs may not associate with model performance. Thus, the fact that the validation loss of G is smaller than the training loss does not necessarily indicate whether the synthesized images on the validation dataset is better than that on the training dataset. It is also why GANai uses the

prediction of D , rather than using the loss of G , to control the model training phase switch.

D. Limitations

While GANai, in general, performs better than traditional GAN models and histogram matching on texture features, its performance could be suboptimal on shape-based features. Shape-based features, such as volume, are usually determined by the physical setup of CT machines. For instance, a 1.5 mm nodule can be totally omitted in a 3 mm slice thickness scan due to partial volume [55]. **The shape-based features may be captured and addressed by replacing the 2D kernels in GANai**

with 3D kernels. However, the computational cost will be increased dramatically. High computational cost means small batch sizes. This can lead to very stochastic gradients as well as limiting the effectiveness of machine learning techniques such as batch normalization [56], increasing model training time and making the results worse. Finally, even with 3D kernels, whether the model can generate the omitted small nodule back to the synthesized image is still questionable. Thus, in this research, we only focus on texture features. Another inherent limitation of GANai is its tendency to generate over-smoothed images as the number of training epochs increases. This can be something that GANai learns from its loss function which makes GANai more vulnerable to overfitting.

VI. CONCLUSION

As a popular diagnostic image modality, CT is routinely used for assessing anatomical tissue characteristics. However, CT imaging customization poses a fundamental challenge in radiomics, since non-standardized imaging protocols are commonplace. Image synthesis algorithms have been developed to integrate and standardize CT images. Among them, GAN models learn the data distribution of training data and generate synthesized images under the same distribution of the training images. However, GANs are not directly applicable to the CT image standardization task due to the lack-of-detail problem.

We developed a novel GAN model called GANai to mitigate the differences in radiomic features of CT images. Given source images, GANai composes synthesized images by specifying a high-level goal that the image features of the synthesized images should be similar to those of the target images. GANai introduces the alternative training strategy to GAN. In each training phase, the model aims to optimize either G or D while freezing the other component. A training phase will stop if the current component is well trained or the training step exceeds an upper bound. After that, GANai switches to train the counter component. Note that just because of the adoption of the alternative training strategy, new technical improvements become applicable. For example, the inputs of the ensemble learning (multiple states of D s and G s) are the end products of every alternative training phase, and a new loss function and dedicated training data can be specified in different training phases. GANai was compared with the start-of-the-art cGAN model [23] and the patch-based histogram matching method [15]. The experimental results show that GANai is significantly better than cGAN and patch-based histogram matching on the texture and intensity histogram based radiomic features.

In conclusion, GANai is a new GAN model for CT image standardization. Its alternative training strategies are effective, easy to implement, and can be adopted by the other GAN models to further improve their performance. With GANai, CT images from multiple medical centers can be seamlessly integrated and standardized, and large-scale radiomics studies can be conducted to extract comprehensive radiomic features and to identify key tumor characteristics that drive disease transformation, progression, and drug resistance.

REFERENCES

- [1] J. L. Prince and J. M. Links, *Medical imaging signals and systems*. Pearson Prentice Hall Upper Saddle River, 2006.
- [2] J. Beutel, H. L. Kundel, and R. L. Van Metter, *Handbook of medical imaging: Physics and psychophysics*. Spie Press, 2000, vol. 1.
- [3] A. Webb and G. C. Kagadis, "Introduction to biomedical imaging," *Medical Physics*, vol. 30, no. 8, pp. 2267–2267, 2003.
- [4] M. Mahesh, "Fundamentals of medical imaging," *Medical Physics*, vol. 38, no. 3, pp. 1735–1735, 2011.
- [5] J. T. Bushberg and J. M. Boone, *The essential physics of medical imaging*. Lippincott Williams & Wilkins, 2011.
- [6] A. Midya, J. Chakraborty, M. Gönen, R. K. Do, and A. L. Simpson, "Influence of ct acquisition and reconstruction parameters on radiomic feature reproducibility," *Journal of Medical Imaging*, vol. 5, no. 1, p. 011020, 2018.
- [7] S. P. Raman, M. Mahesh, R. V. Blasko, and E. K. Fishman, "Ct scan parameters and radiation dose: practical advice for radiologists," *Journal of the American College of Radiology*, vol. 10, no. 11, pp. 840–846, 2013.
- [8] S. S. Yip and H. J. Aerts, "Applications and limitations of radiomics," *Physics in Medicine & Biology*, vol. 61, no. 13, p. R150, 2016.
- [9] A. J. Buckler, L. Bresolin, N. R. Dunnick, D. C. Sullivan, and Group, "A collaborative enterprise for multi-stakeholder participation in the advancement of quantitative imaging," *Radiology*, vol. 258, no. 3, pp. 906–914, 2011.
- [10] G. Liang, J. Zhang, M. Brooks, J. Howard, and J. Chen, "radiomic features of lung cancer and their dependency on ct image acquisition parameters," *Medical Physics*, vol. 44, no. 6, p. 3024, 2017.
- [11] R. Berenguer, M. d. R. Pastor-Juan, J. Canales-Vázquez, M. Castro-García, M. V. Villas, F. M. Legorburo, and S. Sabater, "Radiomics of ct features may be nonreproducible and redundant: Influence of ct acquisition parameters," *Radiology*, p. 172361, 2018.
- [12] L. A. Hunter, S. Krafft, F. Stingo, H. Choi, M. K. Martel, S. F. Kry, and L. E. Court, "High quality machine-robust image features: Identification in nonsmall cell lung cancer computed tomography images," *Medical physics*, vol. 40, no. 12, 2013.
- [13] M. F. Cohen and J. R. Wallace, *Radiosity and realistic image synthesis*. Elsevier, 2012.
- [14] S. Reed, Z. Akata, X. Yan, L. Logeswaran, B. Schiele, and H. Lee, "Generative adversarial text to image synthesis," in *Proceedings of the 33rd International Conference on Machine Learning*, vol. 48, 2016, pp. 1060–1069.
- [15] R. C. Gonzalez and R. E. Woods, *Digital image processing*. Upper Saddle River, NJ: Prentice Hall, 2012.
- [16] A. Rosenfeld, *Digital picture processing*. Academic press, 1976.
- [17] R. A. Yeh, C. Chen, T. Yian Lim, A. G. Schwing, M. Hasegawa-Johnson, and M. N. Do, "Semantic image inpainting with deep generative models," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 5485–5493.
- [18] R. Mihail, G. Liang, and N. Jacobs, "Automatic hand skeletal shape estimation from radiographs," *IEEE transactions on nanobioscience*, 2019, doi: 10.1109/TNB.2019.2911026.
- [19] A. R. Weeks, L. J. Sartor, and H. R. Myler, "Histogram specification of 24-bit color images in the color difference (cy) color space," *Journal of electronic imaging*, vol. 8, no. 3, pp. 290–301, 1999.
- [20] A. K. Jain, *Fundamentals of digital image processing*. Englewood Cliffs, NJ: Prentice Hall, 1989.
- [21] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in Neural Information Processing Systems 27*, 2014, pp. 2672–2680.
- [22] M. Mirza and S. Osindero, "Conditional generative adversarial nets," *arXiv:1411.1784v1*, 2014.
- [23] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-image translation with conditional adversarial networks," in *Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [24] A. Dosovitskiy and T. Brox, "Generating images with perceptual similarity metrics based on deep networks," in *Advances in Neural Information Processing Systems*, 2016, pp. 658–666.
- [25] C. Li and M. Wands, "Combining markov random fields and convolutional neural networks for image synthesis," in *Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 2479–2486.

- [26] S. P. A. S. P. Donggeun Yoo, Namil Kim and I. S. Kweon, "Pixel-level domain transfer," in *Computer Vision-ECCV 2016*, 2016, pp. 517–532.
- [27] L. Karacan, Z. Akata, A. Erdem, and E. Erdem, "Learning to generate images of outdoor scenes from attributes and semantic layouts," *arXiv preprint arXiv:1612.00215*, 2016.
- [28] X. Wang, K. Yu, S. Wu, J. Gu, Y. Liu, C. Dong, Y. Qiao, and C. C. Loy, "EsrGAN: Enhanced super-resolution generative adversarial networks," in *European Conference on Computer Vision*. Springer, 2018, pp. 63–79.
- [29] B. S. Rosenstein, C. M. West, S. M. Bentzen, J. Alsner, C. N. Andreassen, D. Azria, G. C. Barnett, M. Baumann, N. Burnet, J. Chang-Claude *et al.*, "Radiogenomics: radiobiology enters the era of big data and team science," *International Journal of Radiation Oncology* Biology* Physics*, vol. 89, no. 4, pp. 709–713, 2014.
- [30] Q. Li, J. Kim, Y. Balagurunathan, Y. Liu, K. Latifi, O. Stringfield, A. Garcia, E. G. Moros, T. J. Dilling, M. B. Schabath *et al.*, "Imaging features from pre-treatment ct scans are associated with clinical outcomes in non-small-cell lung cancer patients treated with stereotactic body radiotherapy," *Medical physics*, vol. 44(8), pp. 4341–4349, 2017.
- [31] H. J. Aerts, "The potential of radiomic-based phenotyping in precision medicine: a review," *JAMA oncology*, vol. 2, no. 12, pp. 1636–1642, 2016.
- [32] D. V. Fried, S. L. Tucker, S. Zhou, Z. Liao, O. Mawlawi, G. Ibbott *et al.*, "Prognostic value and reproducibility of pretreatment ct texture features in stage iii non-small cell lung cancer," *International Journal of Radiation Oncology* Biology* Physics*, vol. 90, no. 4, pp. 834–842, 2014.
- [33] K. Thayalan and R. Ravichandran, *The physics of radiology and imaging*. Jaypee Brothers Medical Publishers, 2014.
- [34] A. N. Primak, C. H. McCollough, M. R. Bruesewitz, J. Zhang, and J. G. Fletcher, "Relationship between noise, dose, and pitch in cardiac multi-detector row ct," *Radiographics*, vol. 26, no. 6, pp. 1785–1794, 2006.
- [35] A. Esteva, B. Kuprel, R. A. Novoa, J. Ko, S. M. Swetter, H. M. Blau, and S. Thrun, "Dermatologist-level classification of skin cancer with deep neural networks," *Nature*, no. 542, pp. 115–118, February 2017.
- [36] A. Radford, L. Metz, and S. Chintala, "Unsupervised representation learning with deep convolutional generative adversarial networks," *arXiv preprint arXiv:1511.06434*, 2015.
- [37] J. Zhao, M. Mathieu, and Y. LeCun, "Energy-based generative adversarial network," *arXiv preprint arXiv:1609.03126*, 2016.
- [38] T. Kim, M. Cha, H. Kim, J. K. Lee, and J. Kim, "Learning to discover cross-domain relations with generative adversarial networks," in *Proceedings of the 34th International Conference on Machine Learning*, vol. 70, International Convention Centre, Sydney, Australia, 2017, pp. 1857–1865.
- [39] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *2017 IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 2242–2251.
- [40] U. Demir and G. Unal, "Patch-based image inpainting with generative adversarial networks," *arXiv preprint arXiv:1803.07422*, 2018.
- [41] P. Luc, C. Couprie, S. Chintala, and J. Verbeek, "Semantic segmentation using adversarial networks," *arXiv preprint arXiv:1611.08408*, 2016.
- [42] Y. Wang, L. Zhang, and J. van de Weijer, "Ensembles of generative adversarial networks," *arXiv preprint arXiv:1612.00991*, 2016.
- [43] I. Durugkar, I. Gemp, and S. Mahadevan, "Generative multi-adversarial networks," *arXiv preprint arXiv:1611.01673*, 2016.
- [44] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *Medical Image Computing and Computer-Assisted Intervention*. Springer, 2015, pp. 234–241.
- [45] D. Pfau and O. Vinyals, "Connecting generative adversarial networks and actor-critic methods," *arXiv:1610.01945*, 2017.
- [46] Y. Bengio, J. Louradour, R. Collobert, and J. Weston, "Curriculum learning," in *Proceedings of the 26th annual international conference on machine learning*. ACM, 2009, pp. 41–48.
- [47] S. Kullback and R. A. Leibler, "On information and sufficiency," *The annals of mathematical statistics*, vol. 22, no. 1, pp. 79–86, 1951.
- [48] A. Strehl and J. Ghosh, "Cluster ensembles – a knowledge reuse framework for combining multiple partitions," *The Journal of Machine Learning Research*, vol. 3, pp. 583–617, 2002.
- [49] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard *et al.*, "Tensorflow: A system for large-scale machine learning," in *OSDI*, vol. 16, 2016, pp. 265–283.
- [50] E. Denton, S. Chintala, A. Szlam, and R. Fergus, "Deep generative image models using a laplacian pyramid of adversarial networks," in *NIPS'15 Proceedings of the 28th International Conference on Neural Information Processing Systems*, vol. 1, 2015, pp. 1486–1494.
- [51] I. G. Tim Salimans, W. Zaremba, V. Cheung, A. Radford, and X. Chen, "Improved techniques for training gans," in *NIPS'16 Proceedings of the 30th International Conference on Neural Information Processing Systems*, 2016, pp. 2234–2242.
- [52] E. S. Page, "Continuous inspection schemes," *Biometrika*, vol. 41, no. 1/2, pp. 100–115, 1954.
- [53] O. A. Grigg, V. Farewell, and D. Spiegelhalter, "Use of risk-adjusted csum and rsprcharts for monitoring in medical contexts," *Statistical methods in medical research*, vol. 12, no. 2, pp. 147–170, 2003.
- [54] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. Courville, "Improved training of wasserstein gans," *arXiv:1704.00028*, 2017.
- [55] A. A. Divani, S. Majidi, X. Luo, F. G. Souslian, J. Zhang, A. Abosch, and R. P. Tummala, "The abcs of accurate volumetric measurement of cerebral hematoma," *Stroke*, vol. 42, no. 6, pp. 1569–1574, 2011.
- [56] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," *arXiv preprint arXiv:1502.03167*, 2015.

Supplementary Document of GANai: Standardizing CT Images using Generative Adversarial Network with Alternative Improvement

Gongbo Liang², Sajjad Fouladvand^{1,2}, Jie Zhang³, Michael A. Brooks³, Nathan Jacobs², Jin Chen^{1,2,4}

1 Institute for Biomedical Informatics, University of Kentucky, USA, Lexington, KY, USA

2 Department of Computer Science, University of Kentucky, Lexington, KY, USA

3 Department of Radiology, University of Kentucky, Lexington, KY, USA

4 Department of Internal Medicine, University of Kentucky, Lexington, KY, USA

I. DATASET

This study includes eight datasets of CT images from two lung cancer patients. A CT image was reconstructed with one of the four slice thickness (0.5, 1, 1.5, 3mm) and one of the two reconstruction kernels (B157 and B164) using Siemens CT Samatom Force in the University of Kentucky Medical Center. In total, 2,448 CT slices were obtained.

A. Generating Training Data

Given the CT images, we generated the training data using data for GANai augmentation. First, we randomly cropped fifteen small patches and five large patches from each CT slice. The width (length) of a small patch varied from 10 pixels to 255 pixels. For a large patch, the patch width (length) varied within [255, 500] pixels. None of the patches overlap with each other. Next, for each patch, we generated two rotated images with random rotation degree, two shifted images with random shift, and two rotated and shifted images. Finally, we rescale all the patches to 256×256 . In total, the training data includes 228,480 image patches. An example image patch is shown in Figure S4.

To generate T_{easy} and T_{hard} , we adopted the cross-validation strategy. Nine fold of the data were used to train a cGAN model and the rest were used for model validation. The cGAN model was trained for 20 epochs. Finally, with the strategy in Section 3.4, we generated T_{easy} and T_{hard} , each having 7,479 images.

B. Generating Validation Data

Each validation image is a 2.5d CT image patch. We used the following criteria to generate all the validation images: First, let the image size be between 5 to 60 pixels for the width/height and 5 to 30 pixels for depth, which is similar to the tumor volume. Second, at least 80% of the total pixels in an image has to be soft tissue, soft tissues, fat, or parenchyma. Neither the pixels amount of bones or air can be accessing 10% of the total number of pixels. Finally, the value of the gray-level co-occurrence matrix correlation feature of a validation image needs to be greater than 6 since from our experiment, the gray-level co-occurrence matrix correlation feature value of tumor usually higher than 7.

II. RADIOMIC FEATURES

A. Gray-level Co-occurrence Matrix

Gray-level co-occurrence matrix (GLCM) is one of the most widely used measurements of tumor texture features [1], [6], [3], [2]. Using GLCM, we measured six texture features, i.e. contrast, correlation, energy, homogeneity, dissimilarity, and entropy [4]. The equations are as follows:

$$Contrast_{GLCM} = \sum_{i,j} |i - j|^2 p(i, j) \quad (1)$$

$$Correlation_{GLCM} = \frac{\sum_i \sum_j (ij) p(i, j) - \mu_i \mu_j}{\sigma_i \sigma_j} \quad (2)$$

$$Energy_{GLCM} = \sum_{i,j} p(i, j) p(i, j) \quad (3)$$

$$Entropy_{GLCM} = - \sum_{i,j} p(i, j) \log(p(i, j)) \quad (4)$$

$$Homogeneity_{GLCM} = \sum_{i,j} \frac{p(i, j)}{1 + |i - j|} \quad (5)$$

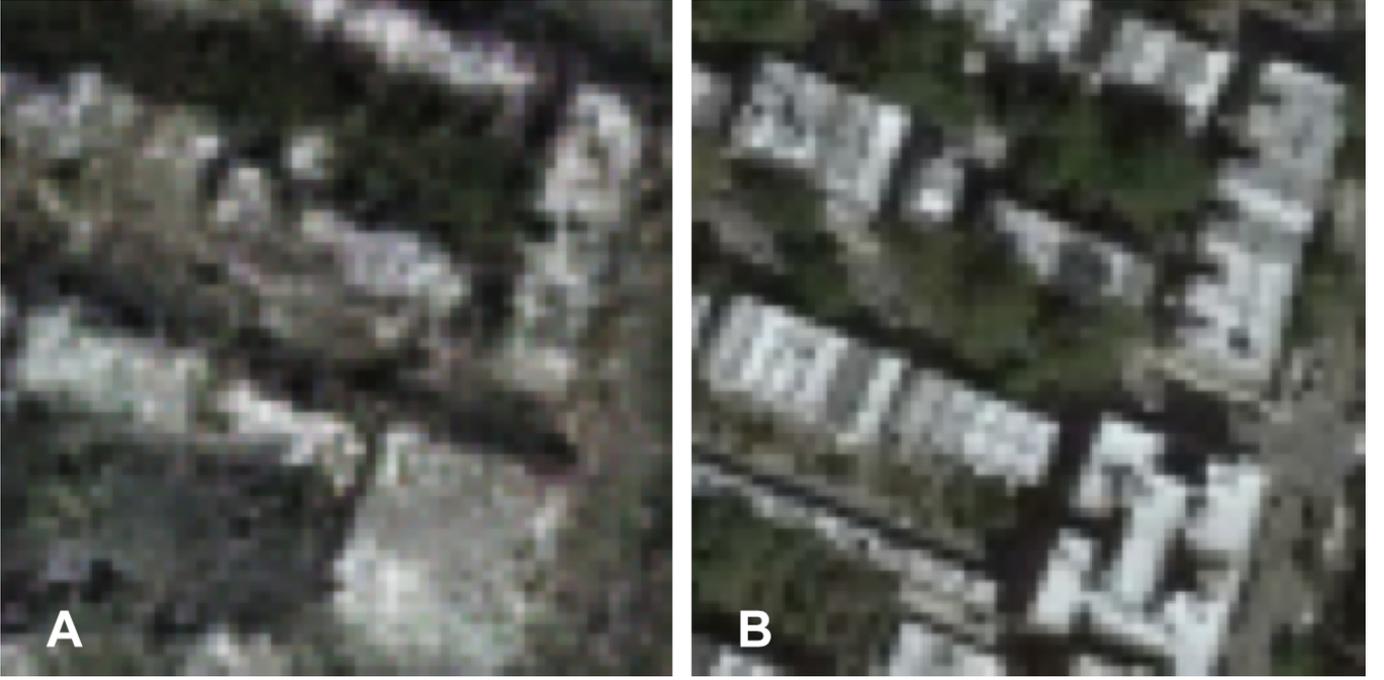


Fig. S1: A cGAN model was trained to convert maps to aerial-view images. (A) the synthesized image generated using the cGAN model. (B) the target image. A side-by-side comparison of A and B shows that many aerial details (such as the boundary of buildings and the shapes of buildings) in B are missing from A.

where $p(i, j)$ is the value of the cell in the co-occurrence matrix, i and j are the row and column indices, μ_i and μ_j are the means of i th row and j th column, σ_i and σ_j are the standard deviation of i th row and j th column.

B. Intensity Histogram

Intensity histogram (IH) is a common method to measure pixel intensity based texture features [2], [3], [6]. Using the intensity histogram, we measured kurtosis and skewness. The equations are as follows:

$$IH_{kurtosis} = \frac{\sum_i^{range} (i - \mu)^4}{(\sum_i^{range} (i - \mu)^2)^2} \quad (6)$$

$$IH_{skewness} = \frac{\sum_i^{range} (i - \mu)^3}{(\sum_i^{range} (i - \mu)^2)^{1.5}} \quad (7)$$

where $range$ indicates the range of intensity (normalized), μ is the mean of weighted sum of IH.

III. MODELS TO COMPARE

For performance evaluation, we compared all the feature-based errors of GANai with cGAN [5] and the patch-based histogram matching. cGAN and patch-based histogram matching were trained using the same training data (the combination of T_{easy} and T_{hard}) and the same validation dataset as GANai. In the cGAN model [5], G is a U-net, and D is a multilayer perceptron model that has the same number of hidden layers and filters as GANai. The learning rate was set to 0.0002, momentum 0.5, and

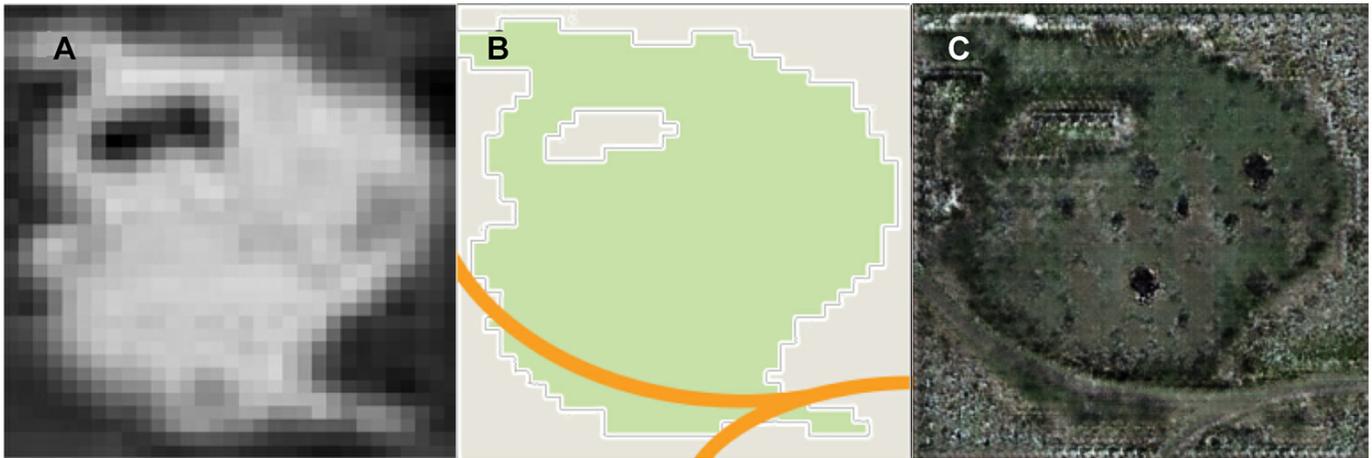


Fig. S2: A more challenging map-to-aerial-view example compared with Fig. S1. (A) a lung tumor image segmented from a CT image. (B) a map-like contour image generated from A. (C) the image synthesized by the cGAN model using B as the source image.

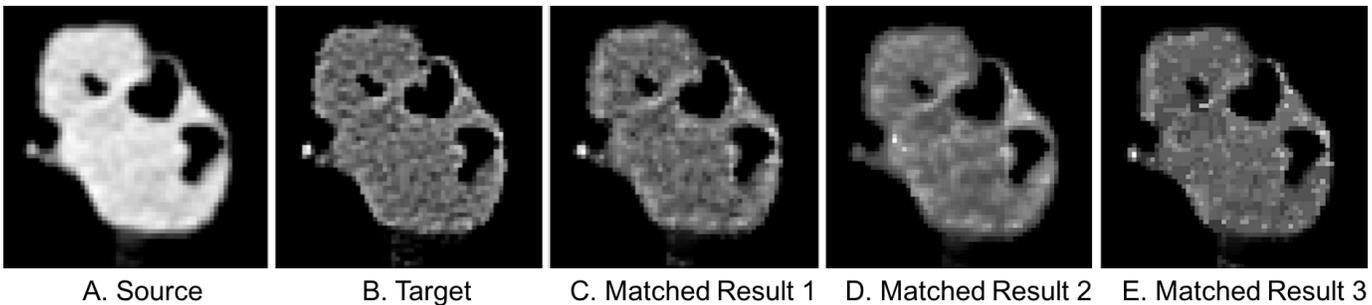


Fig. S3: A patch-based histogram matching example. (A) the source image. (B) the target images. (C)-(E) the synthesized images using different histogram matching parameters (number of bins and patch size). All the synthesized images have artifacts. Specifically, the edges of C are blurred; D has a discontinuous region at the right of the tumor (caused by larger patch sizes); and E has a high noise rate.

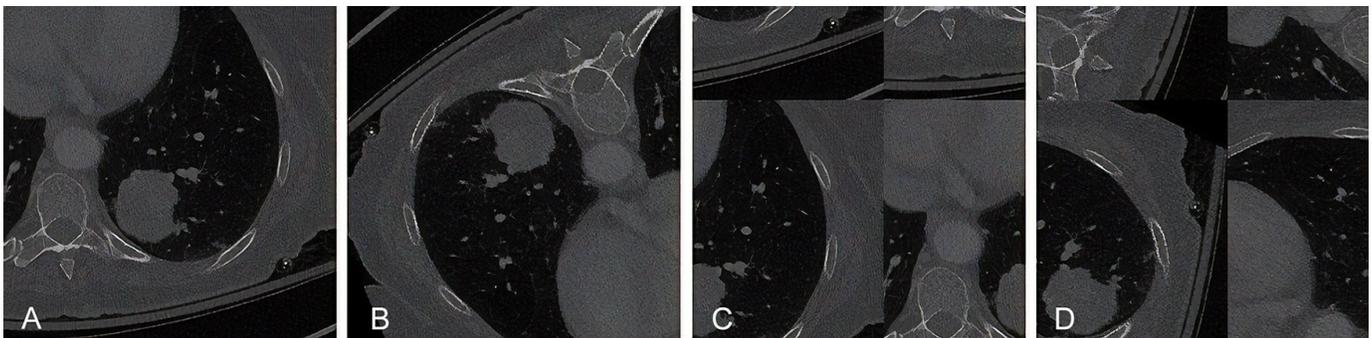


Fig. S4: (A) the cropped patch from CT slice, (B) the rotated patch, (C) the shifted patch, (D) the shifted and rotated patch.

the model was trained for 100 epochs. The patch-based histogram matching model was implemented using the *imhistmatch* function in Matlab¹. The exhaustive search strategy was adopted to find the best mapping between the source and the target images for all the training images. Specifically, for each source image, we applied the patch-based histogram matching model using patch sizes between 32×32 and 128×128 pixels for each patch. We used numbers of bin sizes between 8 and 128. Therefore, for each source image, 15 mapping functions were learned. We then selected the best mapping for each source image (the one generates the synthesized image with the minimum error) and apply it to the rest of the images in the training data

¹<https://www.mathworks.com/help/images/ref/imhistmatch.html>

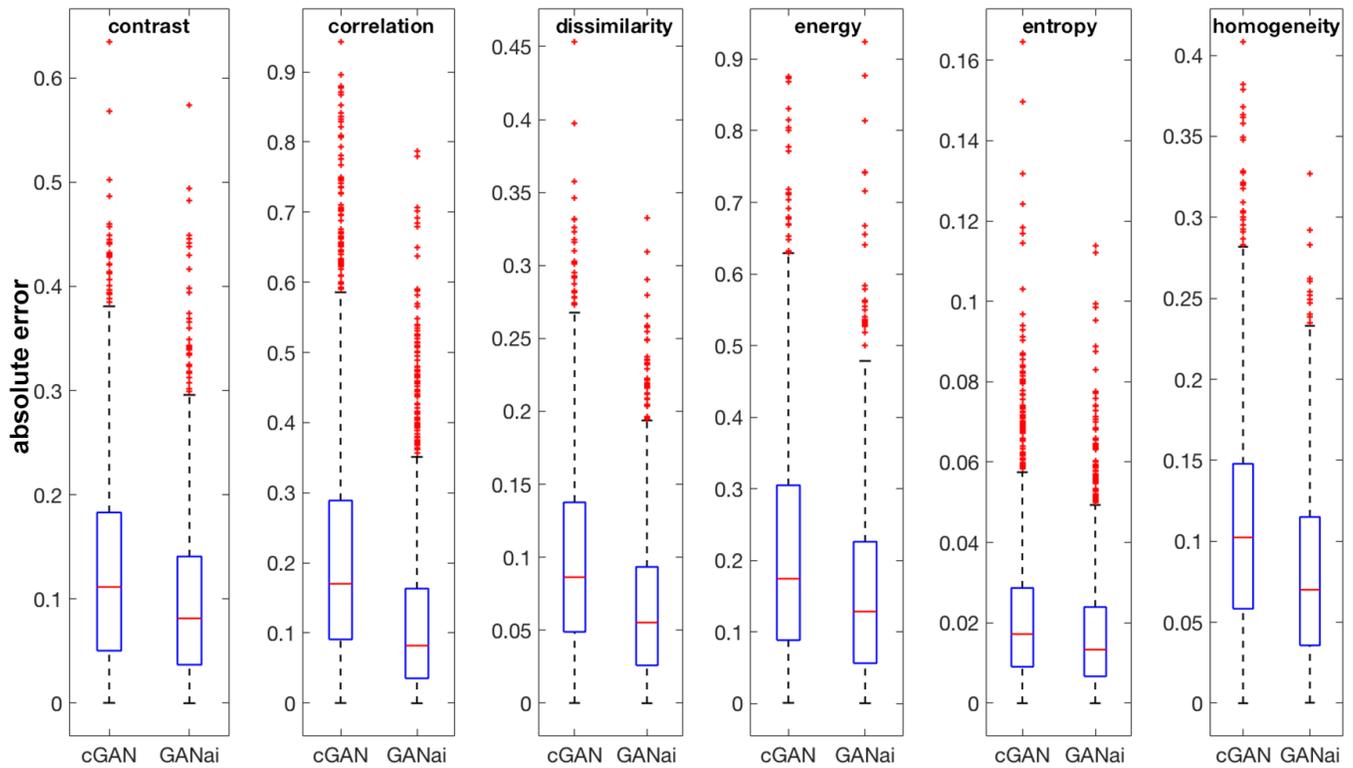


Fig. S5: Comparison of cGAN and GANai on all the features computed using the gray-level co-occurrence matrix. The y-axis is the absolute error. The models were tested on the whole validation data.

and compute the mean error. Finally, we select the mapping function with the minimal mean error in performance evaluation.

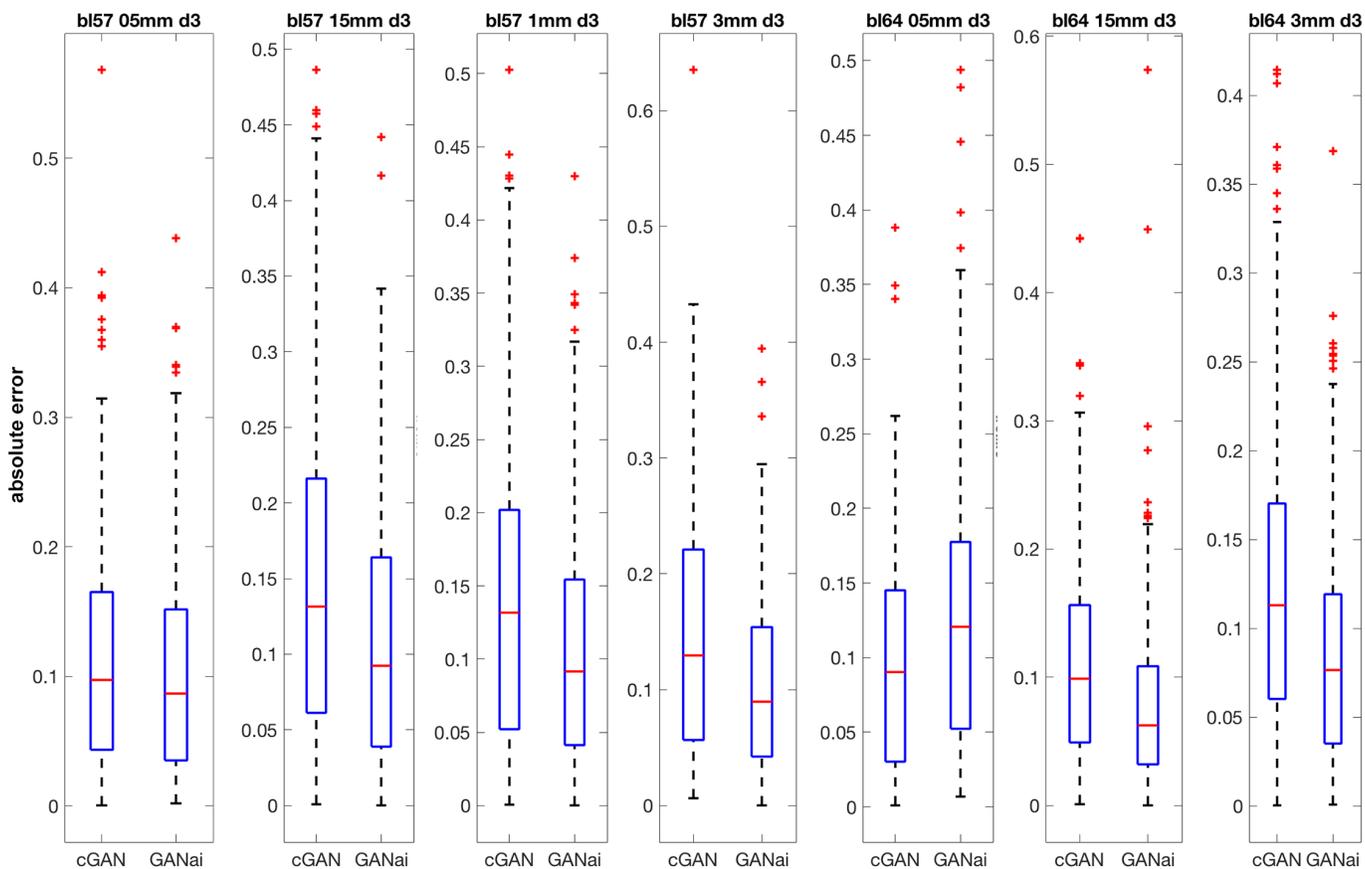


Fig. S6: Comparison of cGAN and GANai on the contrast feature computed using the gray-level co-occurrence matrix. The y-axis is the absolute error. The models were tested on seven subsets of the validation dataset acquired using different acquisition parameters, i.e. B157-0.5mm, B157-1.5mm, B157-1mm, B157-3mm, B164-0.5mm, B164-1.5mm, and B164-3mm.

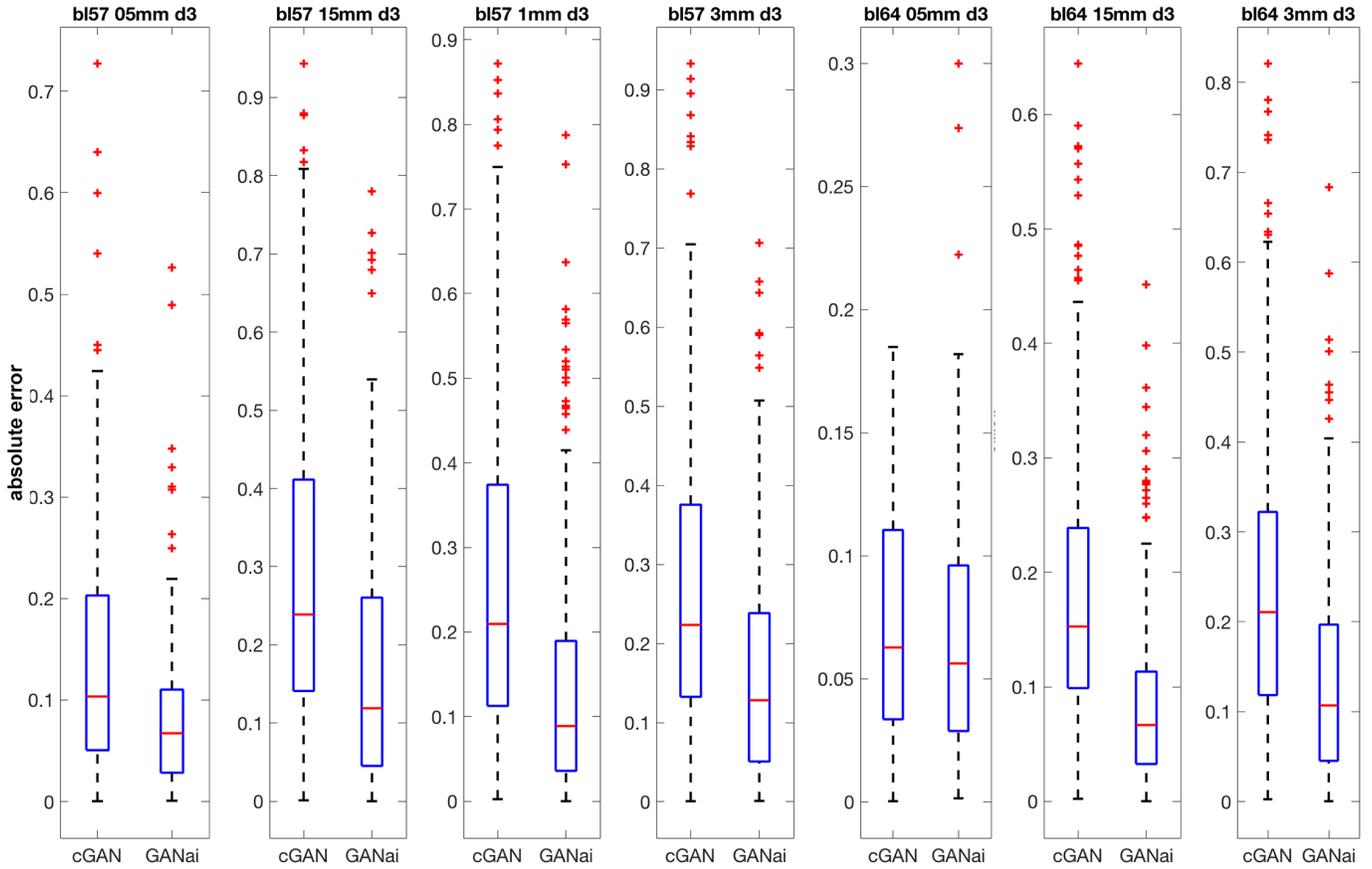


Fig. S7: Comparison of cGAN and GANai on the correlation feature computed using the gray-level co-occurrence matrix. The y-axis is the absolute error. The models were tested on seven subsets of the validation dataset acquired using different acquisition parameters, i.e. B157-0.5mm, B157-1.5mm, B157-1mm, B157-3mm, B164-0.5mm, B164-1.5mm, and B164-3mm.

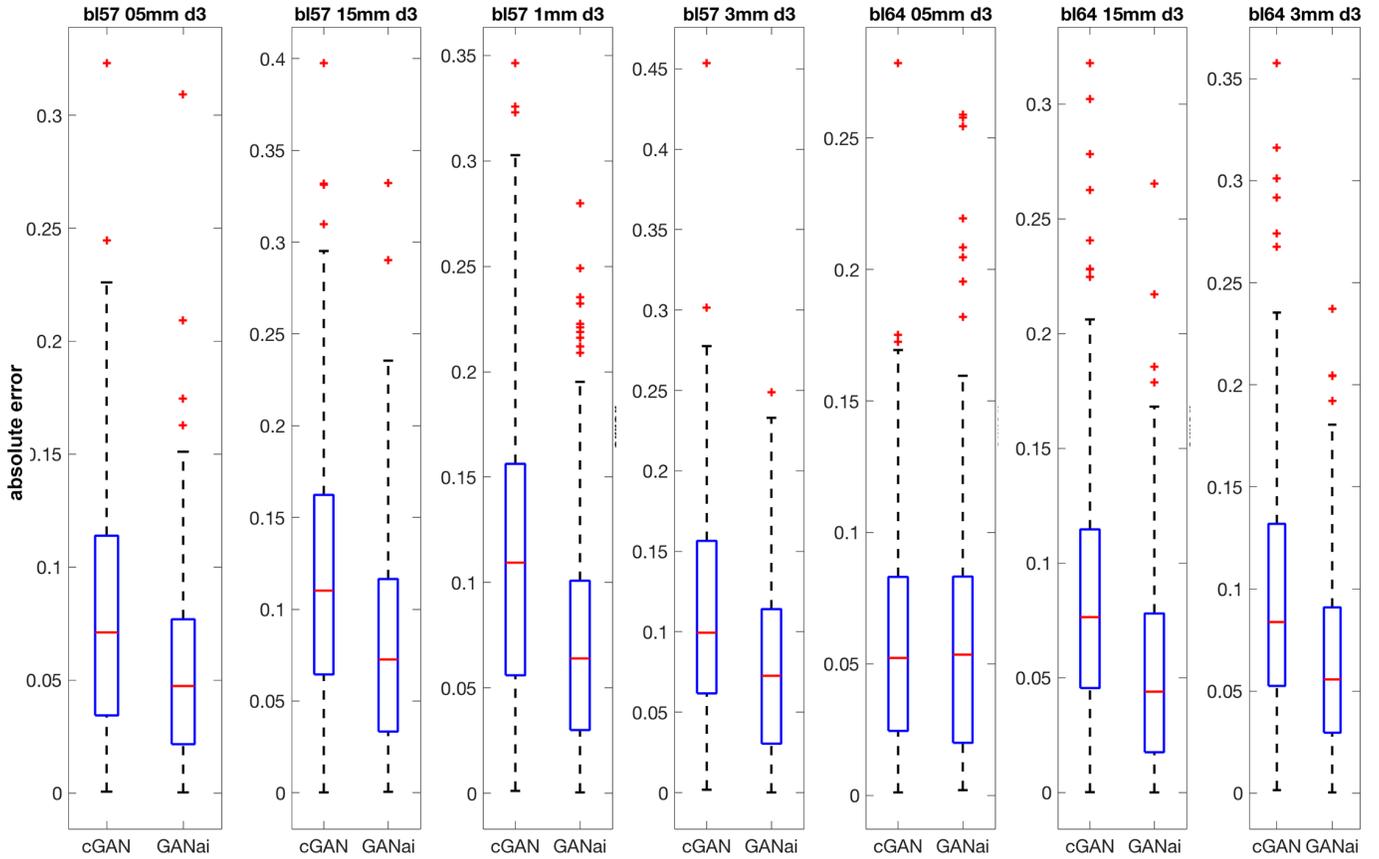


Fig. S8: Comparison of cGAN and GANai on the dissimilarity feature computed using the gray-level co-occurrence matrix. The y-axis is the absolute error. The models were tested on seven subsets of the validation dataset acquired using different acquisition parameters, i.e. B157-0.5mm, B157-1.5mm, B157-1mm, B157-3mm, B164-0.5mm, B164-1.5mm, and B164-3mm.

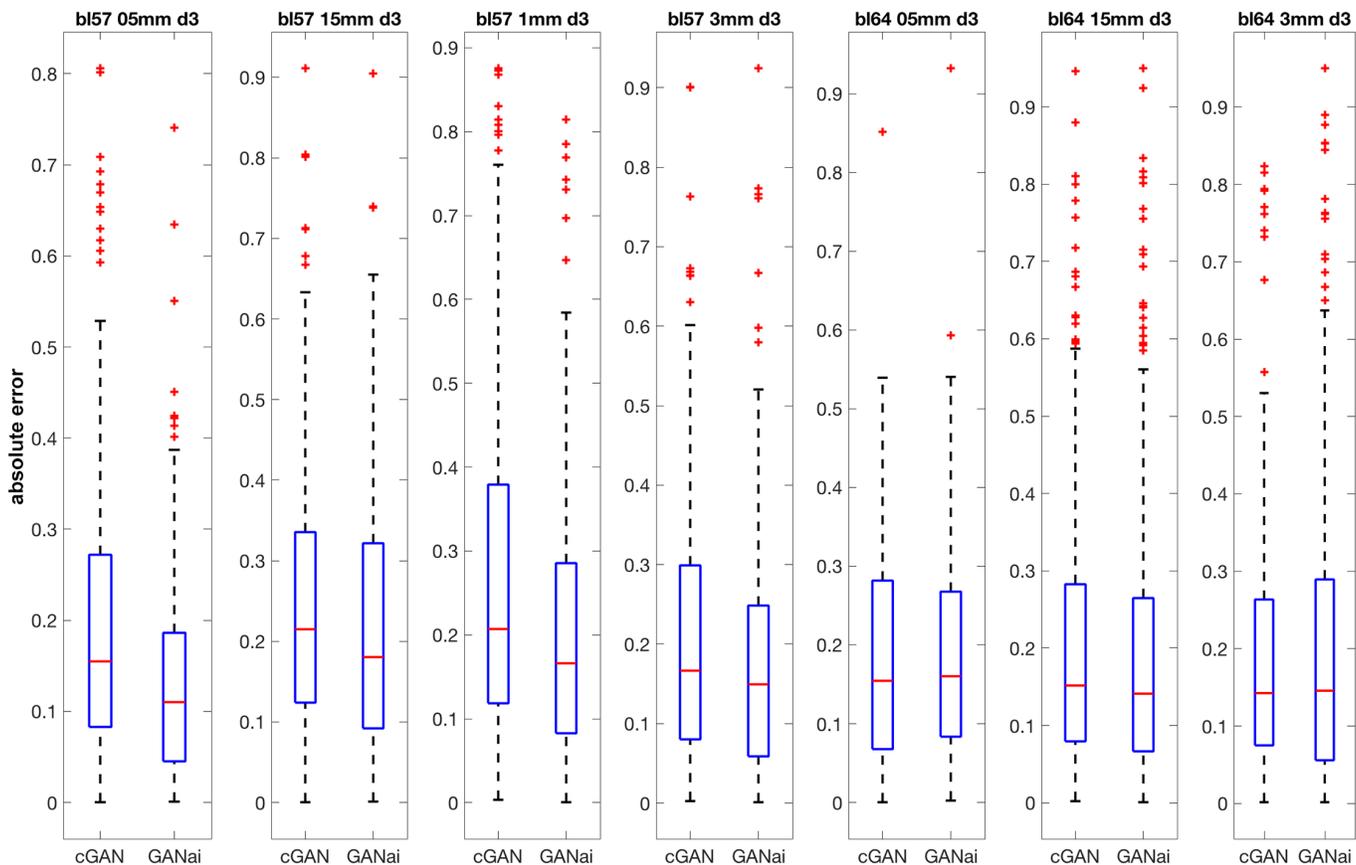


Fig. S9: Comparison of cGAN and GANai on the energy feature computed using the gray-level co-occurrence matrix. The y-axis is the absolute error. The models were tested on seven subsets of the validation dataset acquired using different acquisition parameters, i.e. B157-0.5mm, B157-1.5mm, B157-1mm, B157-3mm, B164-0.5mm, B164-1.5mm, and B164-3mm.

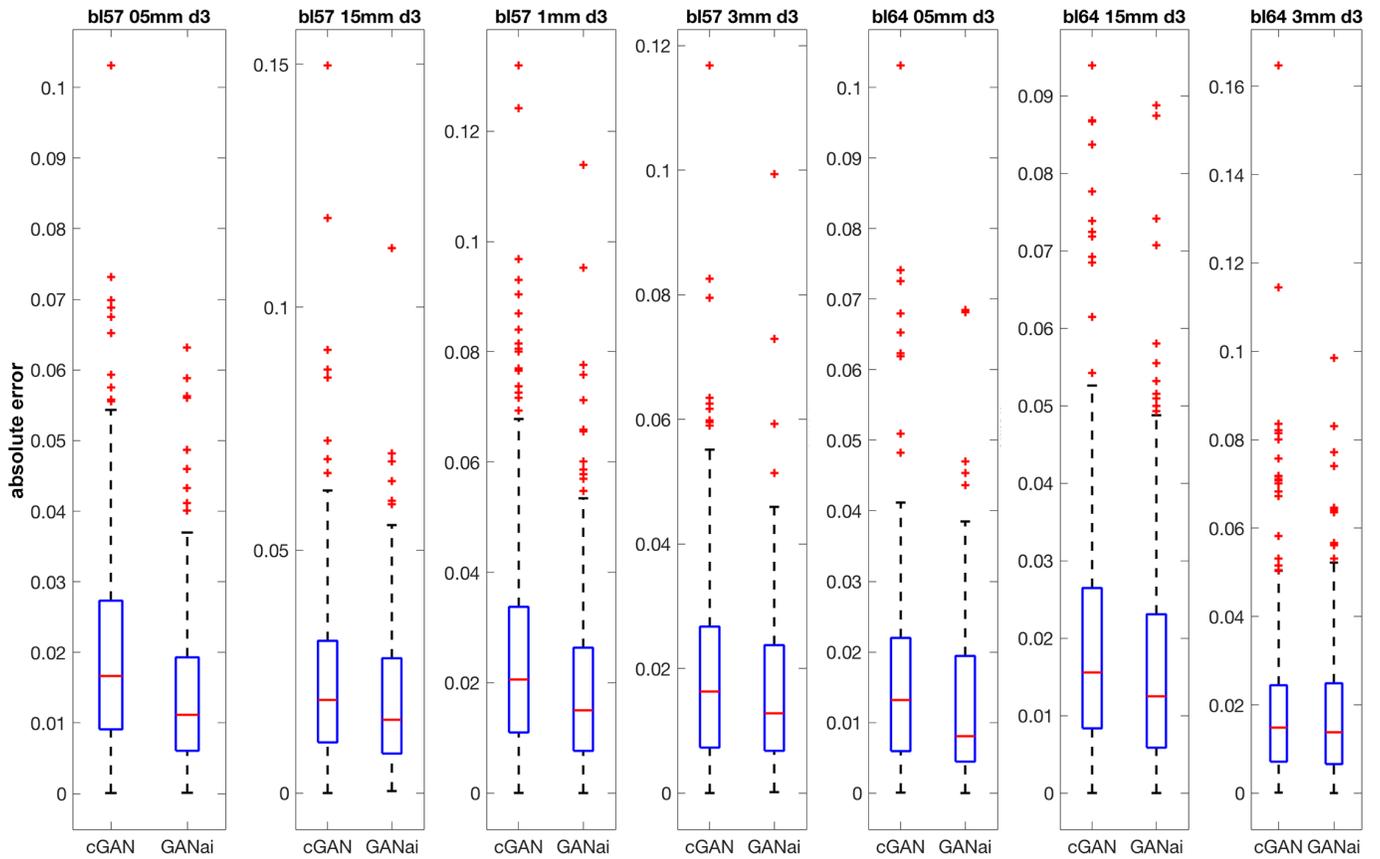


Fig. S10: Comparison of cGAN and GANai on the entropy feature computed using the gray-level co-occurrence matrix. The y-axis is the absolute error. The models were tested on seven subsets of the validation dataset acquired using different acquisition parameters, i.e. B157-0.5mm, B157-1.5mm, B157-1mm, B157-3mm, B164-0.5mm, B164-1.5mm, and B164-3mm.

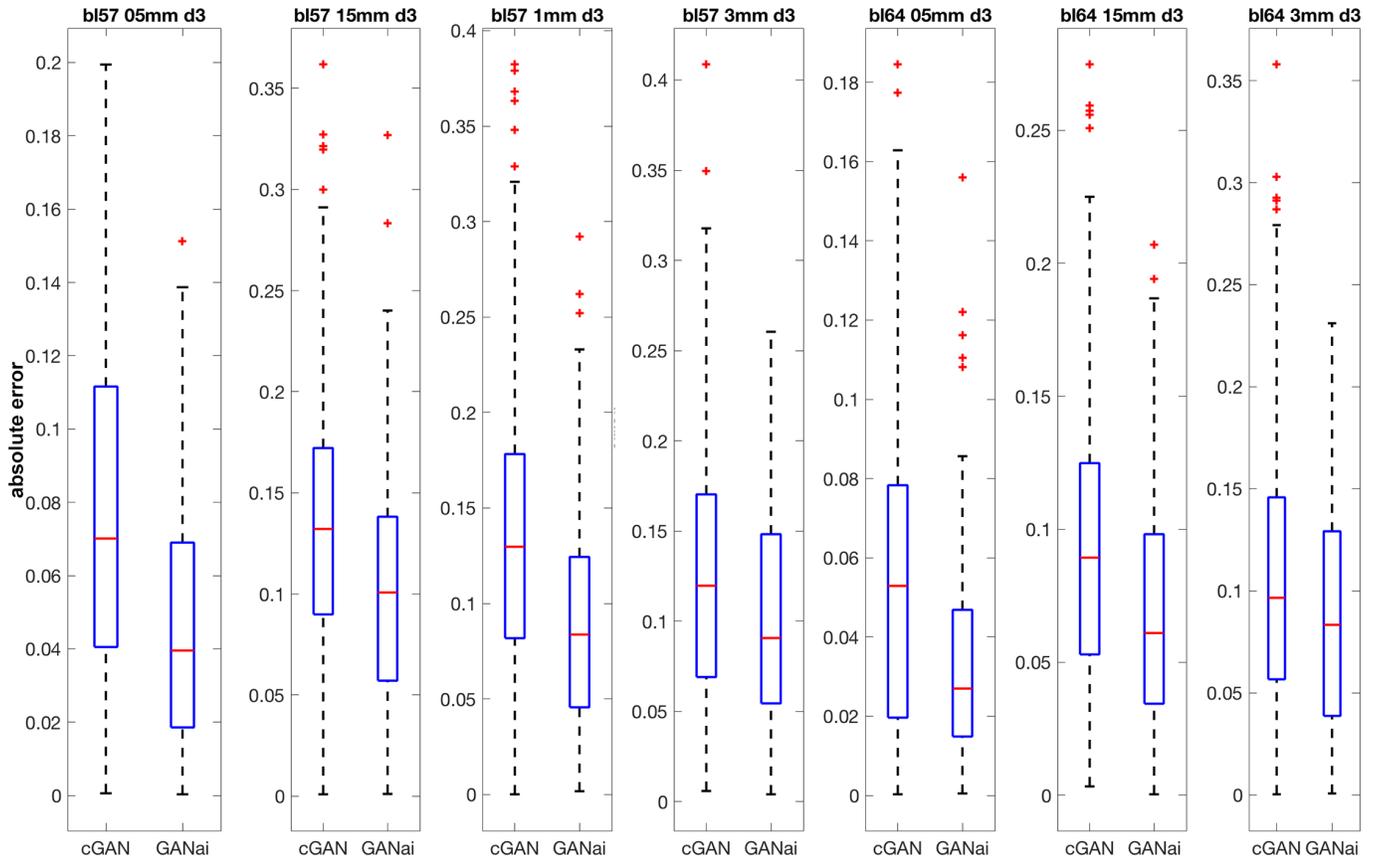


Fig. S11: Comparison of cGAN and GANai on the homogeneity feature computed using the gray-level co-occurrence matrix. The y-axis is the absolute error. The models were tested on seven subsets of the validation dataset acquired using different acquisition parameters, i.e. B157-0.5mm, B157-1.5mm, B157-1mm, B157-3mm, B164-0.5mm, B164-1.5mm, and B164-3mm.

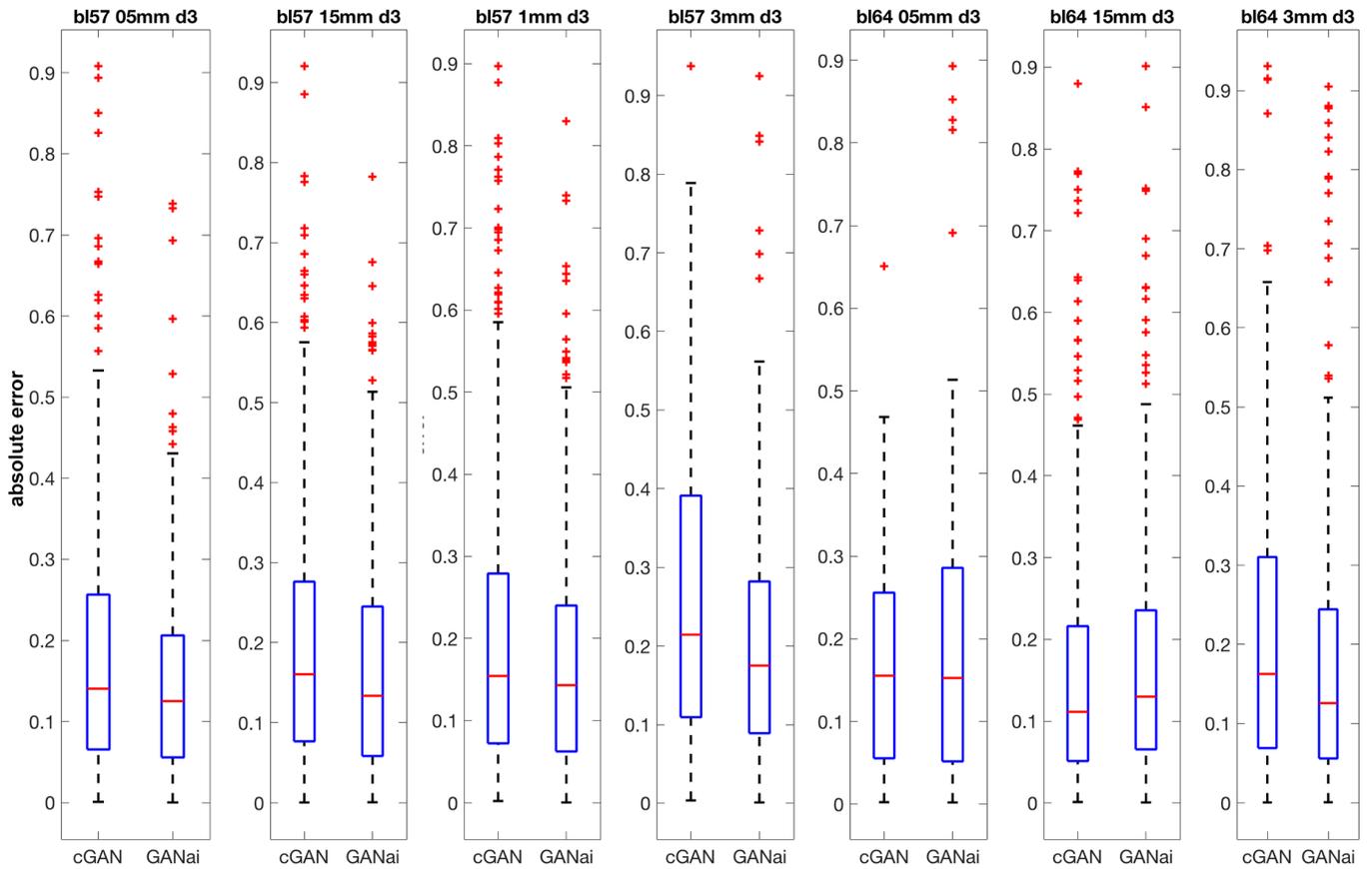


Fig. S12: Comparison of cGAN and GANai on the kurtosis feature computed using the intensity histogram. The y-axis is the absolute error. The models were tested on seven subsets of the validation dataset acquired using different acquisition parameters, i.e. B157-0.5mm, B157-1.5mm, B157-1mm, B157-3mm, B164-0.5mm, B164-1.5mm, and B164-3mm.

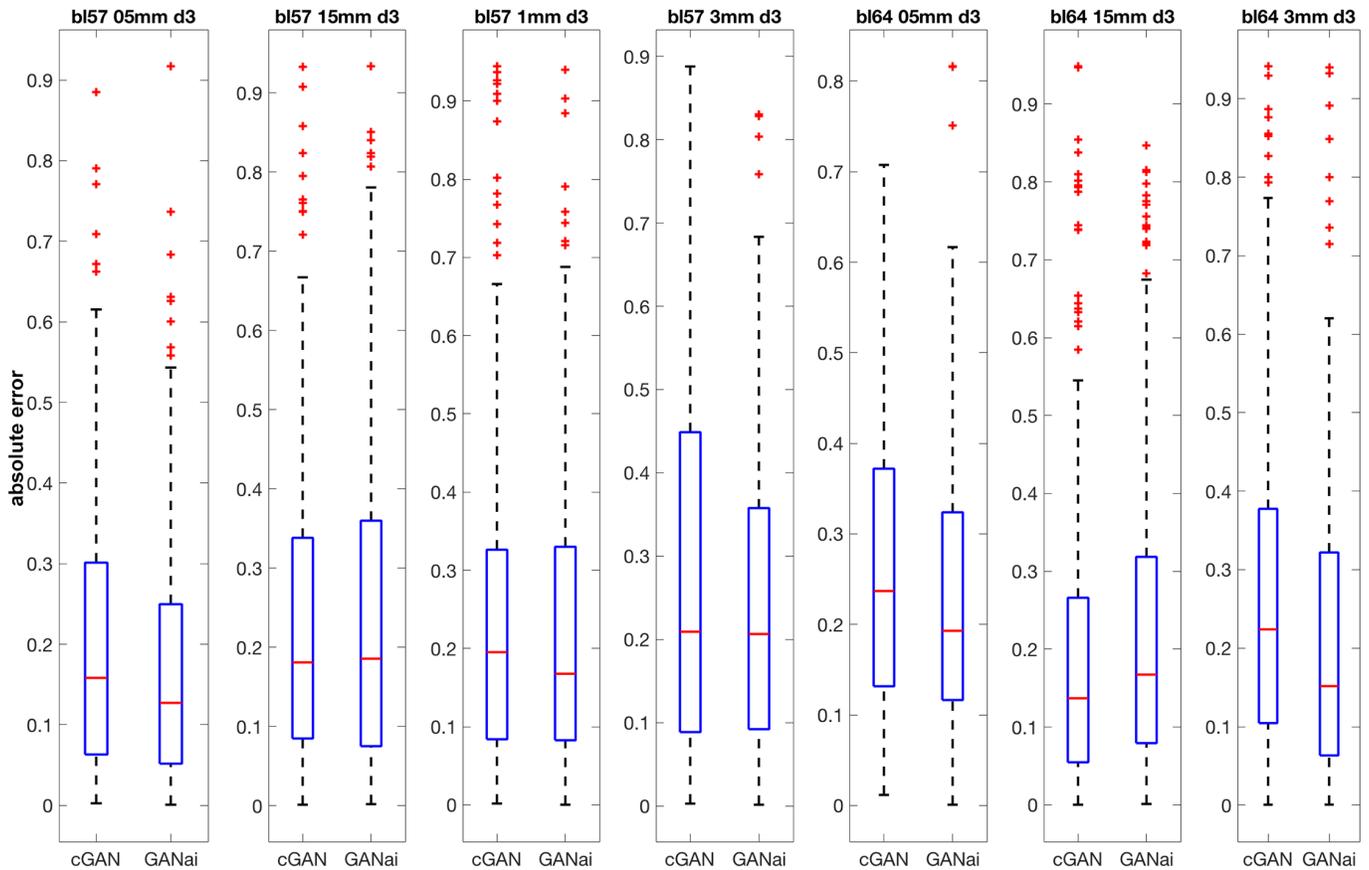


Fig. S13: Comparison of cGAN and GANai on the skewness feature computed using the intensity histogram. The y-axis is the absolute error. The models were tested on seven subsets of the validation dataset acquired using different acquisition parameters, i.e. B157-0.5mm, B157-1.5mm, B157-1mm, B157-3mm, B164-0.5mm, B164-1.5mm, and B164-3mm.

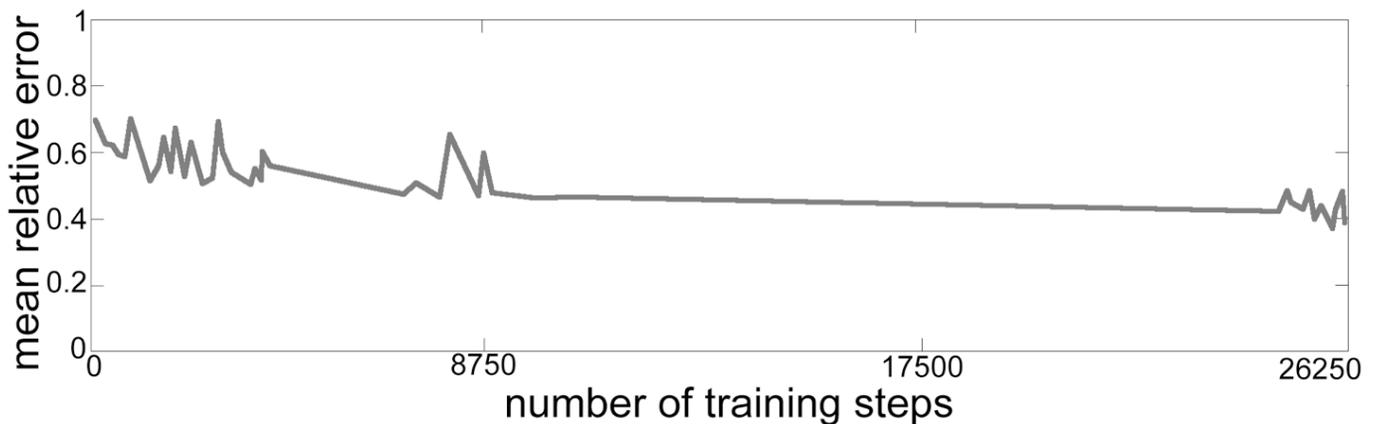


Fig. S14: Results of GAN without ensemble learning. In the example, the generator was trapped at a local minimum during model training. Even after more than five epochs (i.e., about 20000 steps) of training, the generator did not obtain any significant performance improvement. In GANai, by adopting the ensemble learning strategy, the trap-at-local-minimal problem is resolved.

REFERENCES

- [1] Hugo JWL Aerts, Emmanuel Rios Velazquez, Ralph TH Leijenaar, Chintan Parmar, Patrick Grossmann, Sara Carvalho, Johan Bussink, René Monshouwer, Benjamin Haibe-Kains, Derek Rietveld, et al. Decoding tumour phenotype by noninvasive imaging using a quantitative radiomics approach. *Nature communications*, 5:4006, 2014.
- [2] Alexandra R Cunliffe, Samuel G Armato, Xianhan M Fei, Rachel E Tuohy, and Hania A Al-Hallaq. Lung texture in serial thoracic ct scans: Registration-based methods to compare anatomically matched regions. *Medical physics*, 40(6Part1), 2013.
- [3] Balaji Ganeshan, Sandra Abaleke, Rupert CD Young, Christopher R Chatwin, and Kenneth A Miles. Texture analysis of non-small cell lung cancer on unenhanced computed tomography: initial evidence for a relationship with tumour glucose metabolism and stage. *Cancer imaging*, 10(1):137, 2010.
- [4] Robert M Haralick, Karthikeyan Shanmugam, and Its' Hak Dinstein. Textural features for image classification. *IEEE Transactions on systems, man, and cybernetics*, (6):610–621, 1973.
- [5] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A. Efros. Image-to-image translation with conditional adversarial networks. In *Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [6] Marco Ravanelli, Davide Farina, Mauro Morassi, Elisa Roca, Giuseppe Cavalleri, Gianfranco Tassi, and Roberto Maroldi. Texture analysis of advanced non-small cell lung cancer (nsccl) on contrast-enhanced computed tomography: prediction of the response to the first-line chemotherapy. *European radiology*, 23(12):3450–3455, 2013.