# An Adversarial Learning Approach to Medical Image Synthesis for Lesion Detection

Liyan Sun, Jiexiang Wang, Yue Huang, Xinghao Ding, Hayit Greenspan†, and John Paisley‡

*Abstract*—The identification of lesion within medical image data is necessary for diagnosis, treatment and prognosis. Segmentation and classification approaches are mainly based on supervised learning with well-paired image-level or voxel-level labels. However, labeling the lesion in medical images is laborious requiring highly specialized knowledge. We propose a medical image synthesis model named *abnormal-to-normal translation generative adversarial network* (ANT-GAN) to generate a normal-looking medical image based on its abnormal-looking counterpart without the need for paired training data. Unlike typical GANs, whose aim is to generate realistic samples with variations, our more restrictive model aims at producing a normal-looking image corresponding to one containing lesions, and thus requires a special design. Being able to provide a "normal" counterpart to a medical image can provide useful side information for medical imaging tasks like lesion segmentation or classification validated by our experiments. In the other aspect, the ANT-GAN model is also capable of producing highly realistic lesion-containing image corresponding to the healthy one, which shows the potential in data augmentation verified in our experiments.

*Index Terms*—Medical Image Synthesis, Generative Adversarial Network, Unsupervised Learning.



(a) real tumor MRI    (b) generated healthy    (c) difference (a)&(b)

(d) real healthy MRI    (e) generated healthy    (f) difference (d)&(e)

Fig. 1. Results produced by our model. Lesions are isolated, while healthy regions pass through with minimal modification.

## I. INTRODUCTION

Lesions can occur in body tissue as a result of various factors including trauma, infection or cancer. Medical imaging techniques such as magnetic resonance imaging (MRI) and computational tomography (CT) provide detailed information for diagnosing such lesions [1]. With more efficient medical imaging systems being deployed beyond advanced societies, demands on radiologists have also been increasing. Automatic medical analysis systems can help lower the human expert barrier and expedite the diagnosis and treatment process [18].

However, in the current medical image analysis paradigm, machines and human experts differ in their approach. Specifically, radiologists are well-trained using many healthy and unhealthy medical images and transfer their learned internal 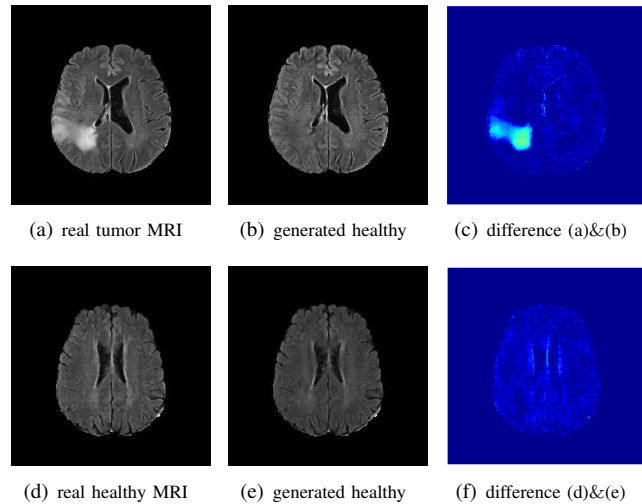representations to new images. Experts search for abnormal regions that differ from their prior knowledge bank of "healthiness" when mentally segmenting the lesions. As for machines, usually a function mapping of the unhealthy medical image to a certain label is learned in a supervised way, either at the image-level or the voxel-level. Since images containing lesions constitute only a small portion of available scans, information in the lesion-free image is often overlooked. Furthermore, the size of medical image datasets is usually limited because labeling requires specialized expertise and is laborious. Such data imbalance and scarcity impacts the performance of medical image analysis models negatively and motivates us to more sufficiently utilize the "healthy" images containing valuable prior information on the appearance of healthy brain structure. We seek to imitate the expert by building a knowledge base of healthy medical images to aid improving diagnostic performance. Predicting a fake healthy version of an image containing lesions can aid in automatic medical image analysis tasks such as segmentation and provide doctors with additional diagnostic information.

Deep neural networks are the state-of-the-art for supervised learning in computer vision and medical imaging tasks such as image classification [10], [25], segmentation [5], [23], [27], object detection [9] and medical image reconstruction [26]. However, in real clinical practice obtaining pairs of normal and abnormal images is unrealistic, since only one can exist at a time, and data augmentation methods are not available here. Thus, instead of formulating a supervised learning framework,

L. Sun, J. Wang, Y. Huang and X. Ding was with the School of Information Science and Engineering, Xiamen University, Xiamen, Fujian, 361005, China, corresponding to: dxh@xmu.edu.cn.

† H. Greenspan was with Department of Biomedical Engineering, Tel Aviv University, Israel

‡ J. Paisley was with Department of Electrical Engineering, Columbia University, New York, NY, USA

we develop an "abnormal-to-normal translation generative adversarial network" (ANT-GAN) model to predict what a lesion-free image should look like that corresponds to an input image; if the model doesn't detect a lesion, the output should be indistinguishable from the input. In the proposed model, an abnormal-to-normal generator (A2N-Generator) converts the input image to its healthy counterpart, and a discriminator is used to decide whether the input is a faked lesion-free image or a real healthy image.

We primarily test our model on the public Multimodal Brain Tumor Segmentation Challenge (BratS) dataset, which is based on magnetic resonance imaging (MRI) of the human brain. We also experiment on the Liver Tumor Segmentation Challenge (LiTS) dataset acquired by computational tomography (CT) on the human liver. Experiments on these two datasets consisting of different imaging modalities and human tissue demonstrate that ANT-GAN can produce highly realistic healthy-looking images corresponding closely to images containing lesions, which can aid the diagnostic work flow. We show some results on the BratS dataset in Figure 1: A real tumor MRI in Figure 1(a) is input into the well-trained A2N-Generator and the corresponding healthy-looking MRI is generated in Figure 1(b). We take the absolute difference between the two images and give the color map in Figure 1(c) where we observe only the tumor regions are highlighted. We also input a real healthy MRI in Figure 1(d) in the A2N-Generator and the corresponding output is shown in Figure 1(e). The colormap of their absolute difference in Figure 1(f) shows little difference, indicating the generator doesn't detect a lesion that isn't there. Since we leverage a cycle consistency strategy in the ANT-GAN model, a normal-to-abnormal generator (N2A-Generator) can also be obtained. Such a regularization can not only stabilize the training and help convergence, but also provide a approach to encode the lesion information in the N2A-Generator, which has the potential in data augmentation where the training data is scarce.

## II. RELATED WORK

The generative adversarial network (GAN) is an emerging deep learning technique for modeling high-dimensional data distributions [8] and has been widely used in computer vision tasks [4], [16], [30], [31]. Image translation is one important application of GAN models [13], [29]. To overcome the need for perfectly aligned input-output pairs, CycleGAN [35] uses a cycle consistency loss; an unsupervised approach is also taken by DualGAN [32] and UNIT [20] for image-to-image translation.

Medical image synthesis is becoming an active research topic in medical imaging. However, most existing work focuses on synthesizing across imaging modalities rather than restoring the image in some way. For example, [12], [22], [34] map from MRI to CT, while [17] map from MRI to PET and [3] map from multiple MRI modalities to other modalities. [24] infer the manifold of normal tissue using a GAN architecture and develop an anomaly scoring scheme to predict abnormal tissue. In [7] the GAN generates synthetic retinal images using segmentation labels, but pathological patterns are not considered.
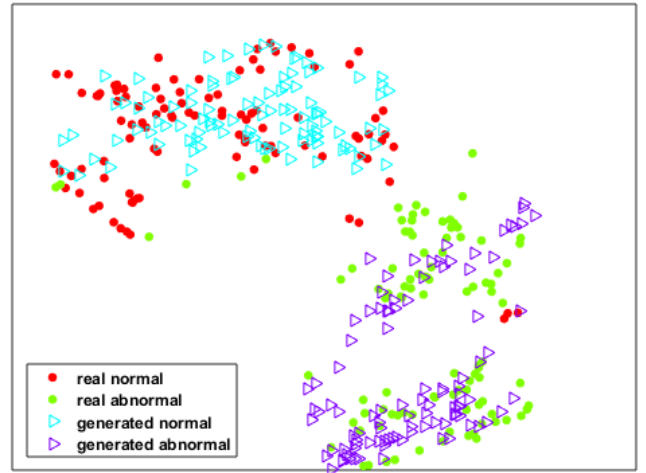


Fig. 2.  100 normal (red) and 100 abnormal (green) medical images from the BratS dataset embedded in $\mathbb{R}^2$ using t-SNE. We also show their embeddings after mapping by their respective generators. Red maps to purple (N2A) and green maps to blue (A2N).

More related to our paper, constrained adversarial auto-encoders are proposed in [6] to detect lesions in brain MRI. The data distribution of brain MRI of healthy subjects are learned using a auto-encoder with unsupervised learning where the constraint that real lesion-containing medical data and its corresponding underlying lesion-free counterpart lie closely in latent space is also imposed. However, this approach is limited by the difficulty in handling high-resolution image synthesis as referred in this work.

## III. METHODS

We denote a normal, healthy medical image as $x^n$ and an abnormal image with lesions as $x^a$. We assume the observed samples are drawn from their corresponding distributions, $x^n \sim p_n(x)$ and $x^a \sim p_a(x)$. In Figure 2, we show a t-SNE embedding of 100 true normal (red dot) and 100 true abnormal images (green dot) in the BratS dataset. To illustrate, we also show the learned A2N-Generator output of each abnormal image (blue triangle) and the learned N2A-Generator output of each normal image (purple triangle). We observe that the distance between the two normal and abnormal manifolds is small and we assume the difference is formed only by the lesions. We next present our ANT-GAN architecture that produced this result.

### A. ANT-GAN architecture overview

The proposed ANT-GAN architecture can be described as an objective function consisting of three different parts, each motivated to capture one aspect desired by the problem. The flow chart of this architecture is shown in Figure 3. The motivation is to take a medical image as input and output a "normal" looking image corresponding to the input. If the input is healthy, we seek an output that is essentially unchanged from the input. The difference between the input and output can then be used to segment abnormal regions and to classify healthy versus unhealthy images.
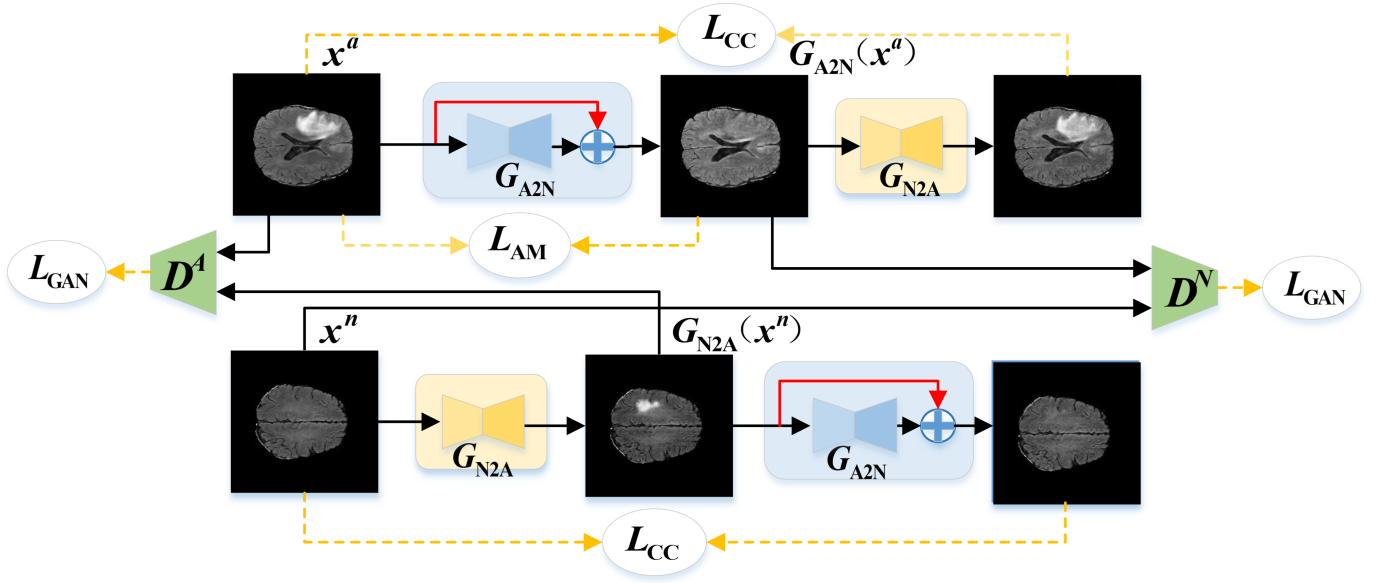
Fig. 3. The flowchart of our proposed ANT-GAN model. The data consists of measured abnormal and normal MRI or CT slices, $x^a$ and $x^n$ respectively. The other images represent intermediate steps within the model and are not measured data.

Our objective function consists of a standard GAN model, plus a cycle consistency loss and a problem-specific loss term to help isolate abnormal regions. The main deep network in our model is the generator $\mathcal{G}_{\text{A2N}}$, which takes in an image $x$, assumed to contain an abnormal region but not necessarily so, and outputs the normal version $\mathcal{G}_{\text{A2N}}(x)$. When $\mathcal{G}_{\text{A2N}}$ is working well, it will produce realistic $\widehat{x}^n = \mathcal{G}_{\text{A2N}}(x^a)$ to fool the discriminator $\mathcal{D}^N$. The generator $\mathcal{G}_{\text{N2A}}$ and discriminator $\mathcal{D}^A$ are used to form the cycle consistency. However, different from many GAN implementations, not any realistic-looking $\widehat{x}^n$ is acceptable, but only one that looks like its corresponding $x^a$ with modifications in the abnormal regions. This motivates the following penalty, which we later validate via an ablation study.

Our full objective function consists of three terms and can be written as

$$\mathcal{L}_{\text{FULL}} = \mathcal{L}_{\text{GAN}} + \lambda_{\text{CC}}\mathcal{L}_{\text{CC}} + \lambda_{\text{AM}}\mathcal{L}_{\text{AM}}, \qquad (1)$$

which we minimize over $\mathcal{G}$ and maximize over $\mathcal{D}$. We break down each of these terms below.

*1) Term 1: Anomaly mask:* During training, we assume that each abnormal image has a corresponding binary mask provided with it that indicates where the abnormal locations are within the image. Let this mask be $\mathbf{M}_x$, which is the same size as the image $x$ being considered in training. We emphasize here that this mask is not available and not needed during testing.

Since we want the generator $\mathcal{G}_{\text{A2N}}$ to automatically isolate and modify the lesions within the image while leaving any healthy region within the image unchanged, we define the penalty

$$\mathcal{L}_{\text{AM}} = \mathbb{E}_{p_a(x)}\left[\|(\mathbf{1} - \mathbf{M}_x) \odot (\mathcal{G}_{\text{A2N}}(x^a) - x^a)\|_2^2\right], \quad (2)$$

where $\odot$ represents element-wise multiplication and $\mathbf{1}$ is an all-ones matrix of the same size of the input image. In

other words, if the generator modifies a pixel or voxel in an abnormal image $x^a$ that does not correspond to the abnormal region, a heavy L2 penalty is paid.

*2) Term 2: GAN:* The main term in our objective is the GAN. Instead of building a unidirectional transform in the abnormal to normal direction, we adopt a bidirectional transform model with two generators $\mathcal{G}_{\text{A2N}}$ and $\mathcal{G}_{\text{N2A}}$ trained simultaneously. This strategy can help stabilize the model training via cycle consistency regularization. The trained $\mathcal{G}_{\text{N2A}}$ can also produce highly realistic lesion-containing medical images, which has the potential for data augmentation not previously available for this problem. We have two generators and two discriminators, depending on whether the input data is normal $x^n$ or abnormal $x^a$.

$$\begin{aligned}\mathcal{L}_{\text{GAN}} = \;& \mathbb{E}_{p_a}\left[\ln \mathcal{D}^A(x^a)\right] + \mathbb{E}_{p_n}\left[\ln \mathcal{D}^N(x^n)\right] \\ & + \mathbb{E}_{p_n}\left[\ln\left(1 - \mathcal{D}^A(\mathcal{G}_{\text{N2A}}(x^n))\right)\right] \\ & + \mathbb{E}_{p_a}\left[\ln\left(1 - \mathcal{D}^N(\mathcal{G}_{\text{A2N}}(x^a))\right)\right].\end{aligned} \qquad (3)$$

We will discuss our selected networks for $\mathcal{D}$ and $\mathcal{G}$ in the following section. While $\mathcal{G}_{\text{A2N}}(x^a)$ is trying to fool the discriminator $\mathcal{D}^N$, the $\mathcal{L}_{\text{AM}}$ term teaches $\mathcal{G}_{\text{A2N}}$ too fool it by only finding and modifying the abnormal regions. The adversarial training strategy is also adopted for $\mathcal{G}_{\text{N2A}}$ and $\mathcal{D}^A$.

*3) Term 3: Cycle consistency:* As motivated by Figure 2, we adopt a cycle consistency term [35] to transform normal and abnormal images into one another, and aid learning of $\mathcal{G}_{\text{N2A}}$ and $\mathcal{G}_{\text{A2N}}$,

$$\begin{aligned}\mathcal{L}_{\text{CC}} = \;& \mathbb{E}_{p_a}\left[\|(\mathcal{G}_{\text{N2A}}(\mathcal{G}_{\text{A2N}}(x^a)) - x^a\|_1\right] \\ & + \mathbb{E}_{p_n}\left[\|(\mathcal{G}_{\text{A2N}}(\mathcal{G}_{\text{N2A}}(x^n)) - x^n\|_1\right].\end{aligned} \qquad (4)$$

This allows for additional information to be shared between normal and abnormal medical images when learning their corresponding generators. As part of this bidirectional regularization, we define the first term to be the abnormality synthesis

(a) Normal Brain    (b) Normal Brain    (c) Brain Tumor    (d) Brain Tumor

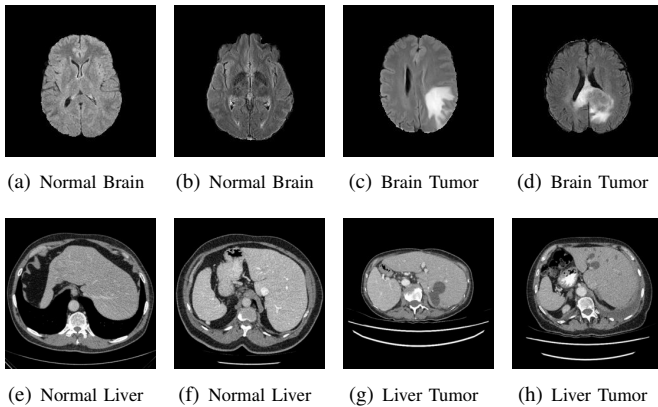(e) Normal Liver    (f) Normal Liver    (g) Liver Tumor    (h) Liver Tumor

Fig. 4. We show some example brain MRI slices from BratS18 (first row) and Liver CT slices from LiTS (second row). Both normal- and abnormal-looking images are provided.

TABLE I
THE NETWORK ARCHITECTURE OF THE GENERATORS. THE "FS" REPRESENTS THE FILTER SIZE, "IN" REPRESENTS THE INSTANCE NORMALIZATION AND "ACT" REPRESENTS THE ACTIVATION FUNCTION.

| Layer | Input | FS | Stride | IN | Act | Output |
|---|---|---|---|---|---|---|
| Conv1 | 240*240*1 | 7*7 | 1 | ✓ | ReLU | 240*240*64 |
| Conv2 | 240*240*64 | 3*3 | 2 | ✓ | ReLU | 120*120*128 |
| Conv3 | 120*120*128 | 3*3 | 2 | ✓ | ReLU | 60*60*256 |
| RB×9 | 60*60*256 | 3*3 | 1 | ✓ | N/A | 60*60*256 |
| Deconv1 | 60*60*256 | 3*3 | 2 | ✓ | ReLU | 120*120*256 |
| Deconv2 | 120*120*256 | 3*3 | 2 | ✓ | ReLU | 240*240*64 |
| Conv4 | 240*240*64 | 7*7 | 1 | | tanh | 240*240*1 |

TABLE II
THE ARCHITECTURE OF THE RESIDUAL BLOCKS (RB).

| Layer | Input | FS | Stride | IN | Act | Output |
|---|---|---|---|---|---|---|
| Conv1 | 60*60*256 | 3*3 | 1 | ✓ | ReLU | 60*60*256 |
| Conv2 | 60*60*256 | 3*3 | 1 | ✓ | N/A | 60*60*256 |

consistency (AC), and the second to be normal synthesis (NC) consistency. In the later ablation study, we show that bidirectional cycle consistency learns a better model than with either unidirectional consistency terms alone.

### B. Implementation

*a) Network Architecture and Training.:* For this medical imaging task in which accuracy is a major requirement of the model, the generator $\mathcal{G}_{A2N}$ needs to detect and modify the lesion region while keeping other parts unchanged. The $\mathcal{L}_{AM}$ penalty is meant to enforce this, but to further help in this task we include a global shortcut (the red arrow in Figure 3) to require the generator to learn a mapping that isolates and removes the lesion.

Following the proposal of [35], we also adopt the network architecture proposed in [14] as the generators $\mathcal{G}_{N2A}$ (except for the added global shortcut) and $\mathcal{G}_{A2N}$. These two generators share the same network architecture but have different parameters. The generators consist of an encoder, residual blocks, and a decoder. We show the architecture of the generators in Table I. In the generators, the 2-stride convolution in shallow layers under-samples the feature into smaller size and the 2-stride deconvolution up-samples the features to the input size. The residual blocks [11] whose architecture is shown in Table II are also adopted to increase model capacity. Similarly, we use PatchGAN [13], [16] as the discriminators $\mathcal{D}^N$ and $\mathcal{D}^A$. In this network, the classification problem is turned into a regression problem. The input image is mapped by the network to a $30 \times 30$ matrix, which is then compared against a matrix of all zeros or ones with an L2 penalty to stabilize training. The architecture of discriminators is shown in Table III. The instance normalization [28] strategy is utilized to help accelerate training.

For training, we update the discriminator using the history of the previous 50 generated images, rather than the output of the most recent generator. For our experiments, we set $\lambda_{CC} = 10$, following [35]. We also set $\lambda_{AM} = 10$. To optimize, we use ADAM with batch size 1. The learning rate is set to $2 \times 10^{-4}$ to train the networks for 20 epochs with 300K iterations.

## IV. RESULTS

### A. Data

We use two popular medical imaging datasets primarily used for the evaluation of lesion segmentation: the Multimodal Brain Tumor Segmentation Challenge 2018 dataset (BratS18) [2], [21] and the Liver Tumor Segmentation Challenge dataset (LiTS).

*a) BratS18.:* The BratS18 dataset provides 210 high grade glioma (HGG) and 75 lower grade glioma (LGG) MRI with binary masks for the tumor (or lack of tumor). Each 3D MRI contains 155 slices of size $240 \times 240$. Not every slice contains a tumor, and therefore healthy MRI are provided by this data as well. We use the FLAIR modality image for all the experiments because the entire tumor is represented well by this modality.However, we also show more experimental results on other modalities, where the ANT-GAN provides impressive visual quality. A more detailed medical description of the data can be found on the challenge website.[1]

*b) LiTS.:* We also experiment with the LiTS data containing a total of 131 contrast enhanced abdominal CT volume images of the liver acquired from 7 different clinical institutions. The in-plane resolution ranges from 0.5mm to 1mm and the slice thickness ranges from 0.7mm to 5.0mm. Each slice is $512 \times 512$ in size and we resize them to $256 \times 256$, and as with the BratS18 MRI not every slice contains a lesion and so these slices are considered to be healthy images. A detailed data description can be found on the challenge website. [2]

Aside from the difference in imaging tissue and modality of these two data sets, the tumor regions on the CT images are of different shape and size, as can be seen in Figure 4. Also, many CT scans are acquired in a way that introduces greater noise-like artifacts than MRI. For each dataset, $80\%$ of randomly selected data are used for training and the resting $20\%$ for testing.

[1] https://www.med.upenn.edu/sbia/brats2018.html
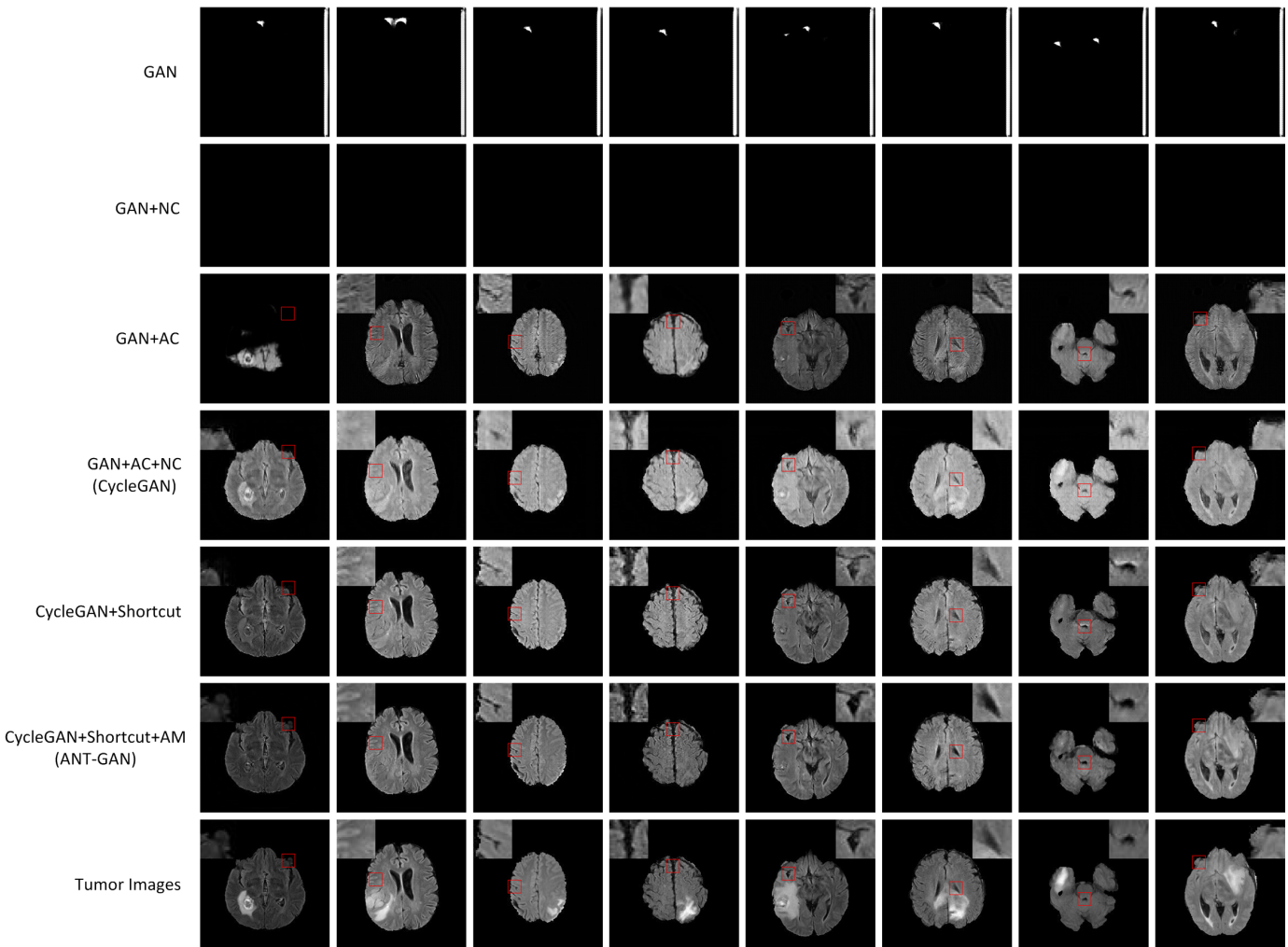[2] https://competitions.codalab.org/competitions/17094

Fig. 5. The results of our ablation study on the BratS18 dataset. Please see text for analysis.

TABLE III
THE ARCHITECTURE OF THE DISCRIMINATORS.

| Layer | Input | FS | Stride | IN | Act | Output |
|---|---|---|---|---|---|---|
| Conv1 | 240*240*1 | 4*4 | 2 | | Leaky ReLU | 120*120*64 |
| Conv2 | 120*120*64 | 4*4 | 2 | ✓ | Leaky ReLU | 60*60*128 |
| Conv3 | 60*60*128 | 4*4 | 2 | ✓ | Leaky ReLU | 30*30*256 |
| Conv4 | 30*30*256 | 4*4 | 1 | ✓ | Leaky ReLU | 30*30*256 |
| Conv5 | 30*30*256 | 4*4 | 1 | | N/A | 30*30*1 |

*B. Qualitative analysis*

Starting from the baseline GAN model, which consists of the $\mathcal{G}_{\text{A2N}}$ generator without shortcut and AM loss and cycle consistency, we conduct the following ablation study to validate the cycle consistency, AM loss and shortcut.

*1) Evaluation of the cycle consistency:* We first conduct experiments to compare the baseline GAN [8] to the same model, but with the abnormality or normality synthesis consistency penalty terms (GAN+AC and GAN+NC). In Figure 5, we show generated (i.e., fake) healthy-looking MRI produced by GAN, GAN+AC and GAN+NC. We observe that model collapse occurs in GAN and GAN+NC where the generator networks have converged to a bad local optimal solution.

GAN+AC produces more meaningful image structures, however it still suffers severe artifacts due to the lack of optimizing constraints. We compare the above models with CycleGAN as a baseline state-of-the-art model for unsupervised image-to-image translation [35]. The cycle consistency term reduces the artifacts by modifying the search space. However, the gray scale shift shows some bias.

*2) Evaluation of the shortcut and anomaly mask:* The global shortcut connection can simplify the function mapping by forcing the generator $\mathcal{G}_{\text{A2N}}$ to focus on the lesion region. We compare with CycleGAN [35] and CycleGAN with a shortcut connection (CycleGAN+shortcut). The difference between CycleGAN+shortcut and our ANT-GAN is the inclusion of the anomaly mask penalty term.

In Figure 5, we observe that the generator can better detect, remove and inpaint the tumor regions without impacting the non-tumor regions by using the proposed global skip connection. However, as shown by the zoomed-in regions, CycleGAN+shortcut still performs less satisfactorily than ANT-GAN in terms of some important details of the healthy regions of the lesion-containing MRI. This is because ANT-GAN contains the anomaly mask term, which forces the generator $\mathcal{G}_{\text{A2N}}$

(a) Conv1 Feature               (b) Conv1 Feature

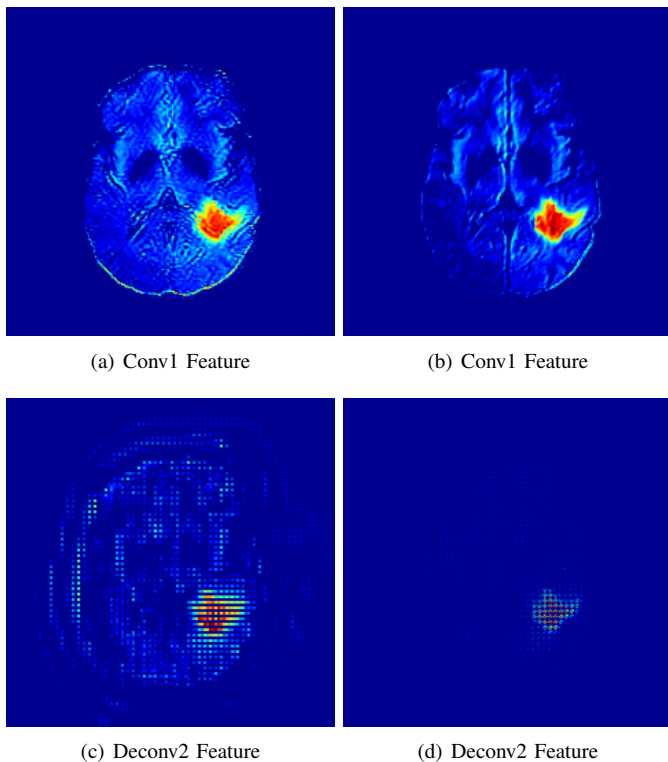(c) Deconv2 Feature           (d) Deconv2 Feature

Fig. 6. A visualization of two features of the generator $\mathcal{G}_{A2N}$ on the BratS MRI data. The tumor is gradually identified.

TABLE IV
THE OBJECTIVE ASSESSMENT ON THE SHORTCUT AND AM LOSS STRATEGY USING PSNR METRIC.

| Model | PSNR dB |
|---|---|
| GAN + AC + NC (CycleGAN) | 21.34 |
| GAN + AC + NC + Shortcut | 27.29 |
| GAN + AC + NC + Shortcut + AM | 28.44 |

to leave healthy portions of an MRI unchanged and only detect and modify lesions. Both CycleGAN and CycleGAN+shortcut do not have this feature since they have different motivation in their design. Though based on the GAN, ANT-GAN is more motivated by image restoration than image generation.

We show feature maps of the generator $\mathcal{G}_{A2N}$ on a BratS MRI data in Figure 6, which clearly shows that the generator $\mathcal{G}_{A2N}$ tries to capture and work on the lesion regions with the global shortcut connection.

We also use PSNR metric to further validate the benefit of shortcut and AM loss strategies in preventing non-lesion regions from distortion. The results evaluated on the testing datasests in BratS18 are shown in Table IV. We observe the proposed ANT-GAN (ANT + AC + NC + Shortcut + AM) achieves the highest PSNR value in the compared models, proving the effectiveness of the shortcut architecture and AM loss in non-lesion region preservation. The lesion regions are not amenable to objective evaluation since their ground truth healthy counterpart are unknown naturally.
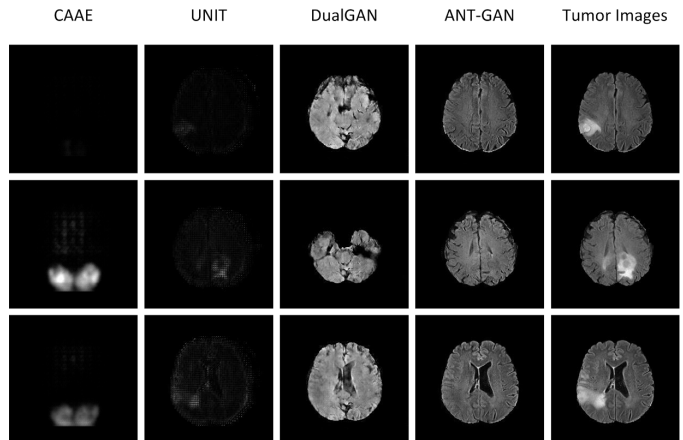


Fig. 7. We compare ANT-GAN with other state-of-the-art GAN methods less tailored to this problem.

*3) Comparison with other GAN formulations:* We compare the proposed ANT-GAN model with the prior work using constrained adversarial auto-encoder model (CAAE) for lesion detection [6] and other two recently proposed state-of-the-art unsupervised GAN models, UNIT [19] and DualGAN [32]. We show these results in Figure 7. We observe that UNIT does not work for this problem. The data sets in this case are too small and the images too large for these models to learn in their less-regularized settings. ANT-GAN also outperforms DualGAN in imaging quality. While DualGAN uses cycle consistency, which makes learning $\mathcal{G}$ easier with less data, no shortcut in $\mathcal{G}$ is used by DualGAN, unlike ANT-GAN. Finally, the stricter regularization of the anomaly mask in ANT-GAN (absent from all other GAN models) not only can enforce greater fidelity to the original image, which is required for this problem, but also aid GAN learning by introducing greater supervision. We observe that CAAE struggles to produce high quality normal-looking medical images in such a high-resolution image synthesis task, which is also the main limitation mentioned in that paper, where evaluations are performed on much smaller images of size $32 \times 32$.

*C. Practical Applications*

In this section, we experiment with using the output of ANT-GAN as input to image segmentation and classification models. We compare with the performance of these same models using the original medical image only without additional information provided by ANT-GAN. The purpose of these experiments is to show how ANT-GAN can supplement existing models to improve their results.

*a) Application to image segmentation.:* The expert radiologist makes diagnoses based on prior knowledge about characteristics of healthy and unhealthy images. To assist in this analysis, automatic segmentation has become an important task in the field of medical image analysis [33]. Since ANT-GAN is trained to isolate abnormal tumor regions and fix those regions only, the difference between an input image $x$ and output image $\mathcal{G}_{A2N}(x)$ can be used to segment MRI for areas of potential concern.
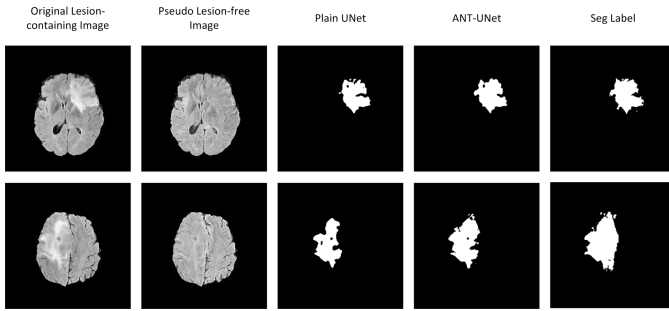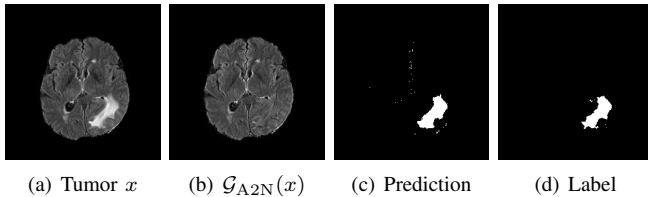
Fig. 8. Example segmentation of two MRI slices.



(a) Tumor $x$    (b) $\mathcal{G}_{\text{A2N}}(x)$    (c) Prediction    (d) Label

Fig. 9. Example segmentation obtained by taking the absolute difference between the real tumor MRI $x$ and the generated normal-looking MRI $\mathcal{G}_{\text{A2N}}(x)$, after binarization at a preset threshold.

For the segmentation model, we use the state-of-art UNet [23] to segment each slice. We input to UNet each MRI slice $x$ and its corresponding generated lesion-free MRI slice $\mathcal{G}_{\text{A2N}}(x)$ from an already-trained ANT-GAN model as multi-channel inputs (referred to as ANT-UNet). We compare with UNet in which only $x$ is input without using information from $\mathcal{G}_{\text{A2N}}(x)$ (referred to as Plain UNet). To ensure that the number of parameters is the same in both models for fair comparison, we use a copy of the MRI slice $x$ to create a multi-channel input for Plain UNet. We train the ANT-GAN model first, and then ANT-UNet and Plain UNet on BratS18 using 5-fold cross validation. In Table V we compare their segmentation performances using the Dice Coefficient (larger is better). The improvement demonstrates that the prior information provided by ANT-GAN on where the lesion may be can significantly aid predictions made by the state-of-the-art segmentation model UNet. We show two visual examples in Figure 8.

TABLE V
THE PERFORMANCE COMPARISON BETWEEN THE PLAIN UNET AND ANT-UNET ON BRATS18.

| DC | Enhancing Tumor | Whole Tumor | Tumor Core |
|---|---|---|---|
| Plain UNet | 75.44% | 87.60% | 77.53% |
| ANT-UNet | 77.77% | 89.10% | 79.77% |

We observe that the generated normal-looking from ANT-GAN can also be used to directly segment the image, since the only difference between a synthesized normal-looking image and its real abnormal counterpart is region with the lesion. To illustrate this, we calculate the absolute difference between $x$ and $\mathcal{G}_{\text{A2N}}(x)$ and show the segmentation after binary thresholding at 0.1 in Figure 9.

*b) Application to image classification.:* We also evaluate the benefits of using the output of ANT-GAN for a lesion classification task. As classifier, we adopt the deep model VGG [25] as the base classifier to predict if the input image contains a lesion or not. Again, to ensure the same number of parameters for comparison, we use a duplicate of $x$ in place of $\mathcal{G}_{\text{A2N}}(x)$ for Plain VGG, while we input both $x$ and $\mathcal{G}_{\text{A2N}}(x)$ as input to the VGG (referred to as ANT-VGG) to see if $\mathcal{G}_{\text{A2N}}(x)$ brings any additional discriminative information. We show the classification results in the Table VI. We use a 5-fold cross validation for evaluation of BratS18. We observe that ANT-VGG outperforms the Plain VGG in all three classification metrics, showing that ANT-GAN can improve the medical image classification task for detecting lesions.

TABLE VI
COMPARISON FOR LESION CLASSIFICATION.

| Methods | Precision | Recall | F1-Measure |
|---|---|---|---|
| Plain-VGG | 89.41% | 89.86% | 0.896 |
| ANT-VGG | 92.35% | 90.96% | 0.917 |

### D. Results on more MRI modalities with brain tumor.

The ANT-GAN model is mainly evaluated and validated on the FLAIR modality. However, we also test the ANT-GAN model on other three MRI setting including T2, T1ce and T1 modalities. The generated pseudo healthy images are shown in Figure 10.

### E. Results on the LiTS Challenge dataset.

We also implement ANT-GAN on the LiTS dataset and show some qualitative results in Figure 11. We observe that the lesions in the liver CT data appears with much lower contrast than in the brain MRI data. While our model can detect and modify the abnormal regions successfully, we note that there are more deformations than with the BratS18 dataset, which is a result of this more difficult task.

## V. DISCUSSIONS

### A. Sensitivity to anomaly mask parameter $\lambda_{\text{AM}}$

We discuss how the regularization parameter $\lambda_{\text{AM}}$ for the anomaly mask term influences the result of our ANT-GAN model and show these results in Figure 12. We observe that setting $\lambda_{\text{AM}} = 10$ leads to a good balance between the preservation of non-lesion regions and the modification of the lesion.

### B. Synthesizing abnormal-looking images

The ability of generating highly realistic pseudo healthy MR image ANT-GAN model using the generator $\mathcal{G}_{\text{A2N}}$ has been verified in our experiments. However, the well-trained ANT-GAN model can produce another generator $\mathcal{G}_{\text{N2A}}$.

We leverage the trained $\mathcal{G}_{\text{N2A}}$ to synthesize abnormal-looking images on BratS18 datasets. Two examples are shown in Figure 13. We observe that the tumor is synthesized near the brain boundary in the second example. In training $\mathcal{G}_{\text{N2A}}$, the discriminator $\mathcal{D}^A$ encodes the core patterns of the real
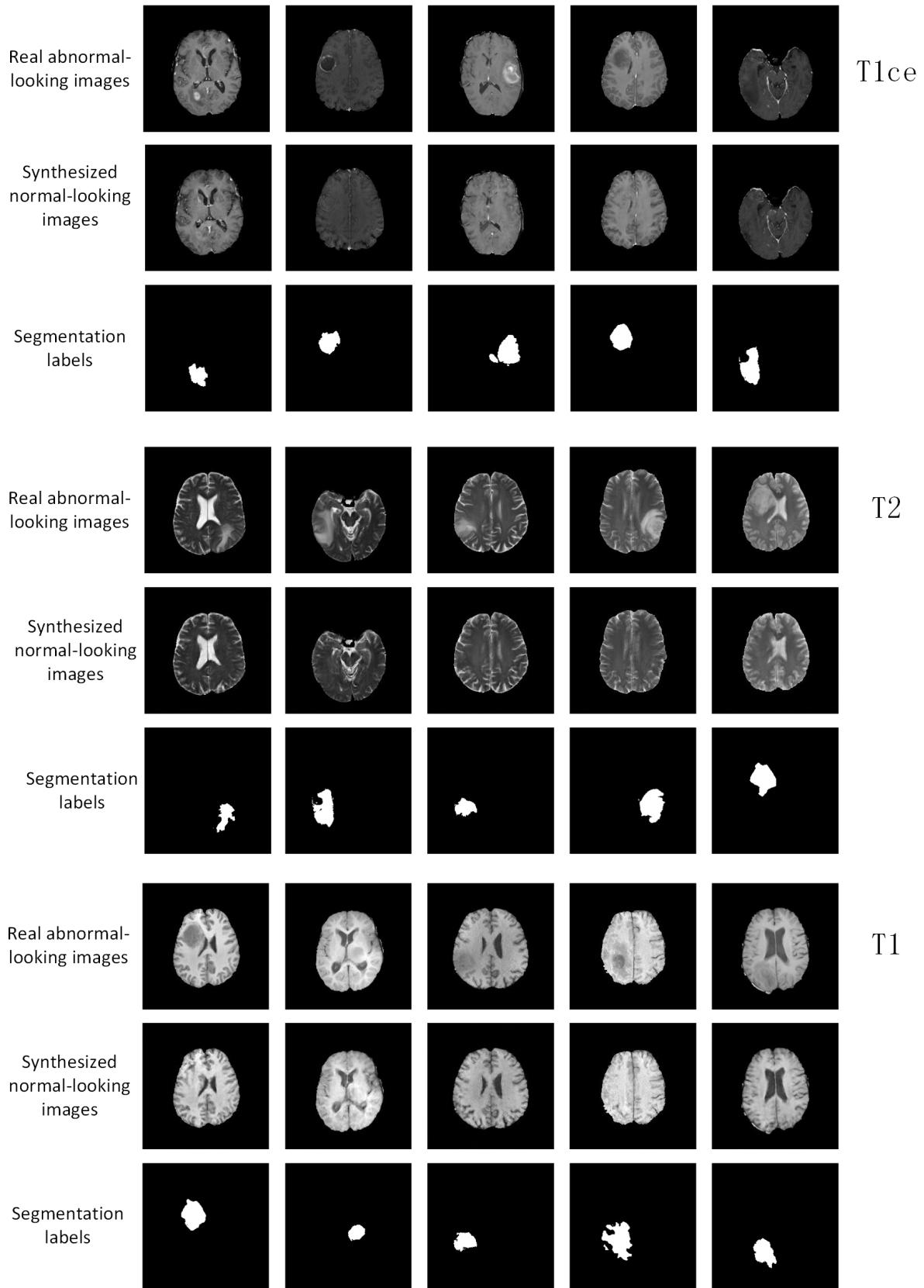
Real abnormal-looking images

Synthesized normal-looking images

Segmentation labels

T1ce

Real abnormal-looking images

Synthesized normal-looking images

Segmentation labels

T2

Real abnormal-looking images

Synthesized normal-looking images

Segmentation labels

T1

Fig. 10. More results on the synthesized normal-looking MRI images on T1ce, T2 and T1 modalities.

Original Lesion-
containing Image     ANT-GAN     Seg Label



Fig. 11. Experimental results on the Liver CT data (LiTS).



(a) Truth       (b) Synthesized lesion

(c) Truth       (d) Synthesized lesion

Fig. 13. Synthesis of a lesion $\mathcal{G}_{\text{N2A}}(x)$ in two healthy MRI $x$.



(a) Tumor      (b) 0.01      (c) 0.1

(d) 1      (e) 10      (f) 100

Fig. 12. The results produced with different $\lambda_{\text{AM}}$.

TABLE VII
THE EVALUATION OF THE ALEXNET MODEL TRAINED ON 2 DATASETS
FOR LESION IDENTIFICATION.

| Training Datasets | Precision | Recall | F1-Measure |
|---|---|---|---|
| Datasets A with 100 samples | 65.56% | 63.36% | 0.644 |
| Datasets B with 200 samples | 77.19% | 77.94% | 0.776 |

images in datasets A to predict their abnormal and normal counterparts using trained $\mathcal{G}_{\text{N2A}}$ and $\mathcal{G}_{\text{A2N}}$, yielding 50 faked healthy images and 50 faked unhealthy ones. We generate a datasets B with the data contained in datasests A plus the pseudo 50 abnormal and 50 normal images.

We train two AlexNet [15] models to identify if the image contains any lesion with training on Datasets A and Datasets B, and the two trained AlexNet models are tested on the resting data in the BratS18 datasets. We show the averaged results on the classification metrics in Table VII.

## VI. CONCLUSION

We proposed an generative adversarial network called ANT-GAN for translating a medical image containing lesions into a corresponding image where the lesion has been "removed" via color correction. We showed how being able to generate these two versions of the same image can help in the medical image segmentation and classification tasks, since the generator can provide additional information about what the image "should" look like if it were healthy. We also showed how our generator was able to not be fooled by healthy MRI, in which case it simply output a near replication of the input image when no lesion is present.

abnormal-looking images data, which can avoid having lesions being generated in physiologically unreasonable locations. We observe the generator $\mathcal{G}_{\text{N2A}}$ produces realistic lesion-containing MR images.

The ANT-GAN model is capable of synthesizing pseudo normal and abnormal images, showing a potential in augmenting data in some cases where the medical images are scarce. To evaluate such potential, we randomly draw 50 real lesion-containing MR images and 50 real lesion-free ones from BratS18 datasets to form a small-scale datasets A. We utilize the chosen 50 real lesion-free and lesion-containing MR
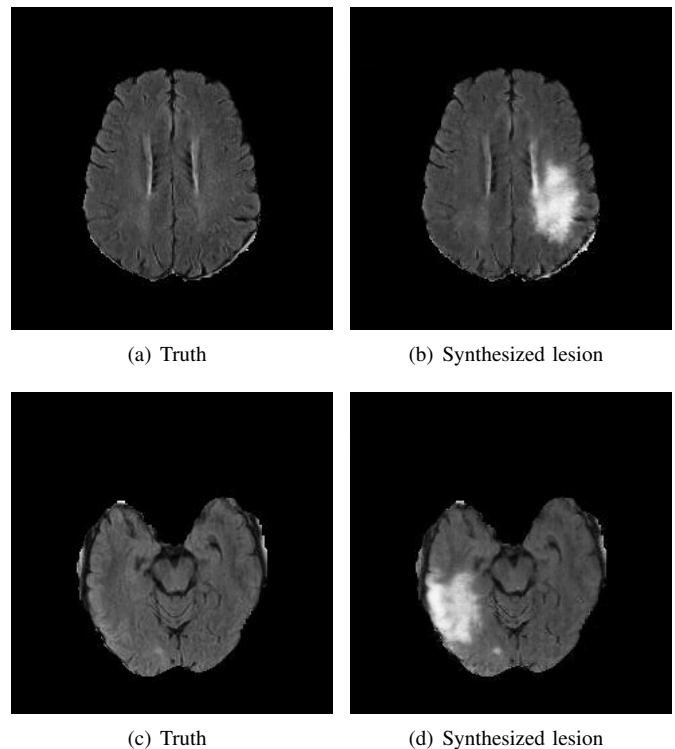
Our objective function was tailored to the problem by introducing an binary anomaly mask term that indicates the lesion location, as well as cycle consistency constraint to regularize the space. (We again note that this mask is not needed for test images.) Comparison with other GAN setups showed how this was a requirement to successfully learn from this small data set. Experiments on the BratS18 and LiTS challenge data sets helped to validate our framework for using computer vision to address fundamental medical image analysis problems of segmentation and classification.

REFERENCES

[1] A. M. Aisen, W. Martel, E. M. Braunstein, K. I. McMillin, W. A. Phillips, and T. Kling. MRI and CT evaluation of primary bone and soft-tissue tumors. *American Journal of Roentgenology*, 146(4):749–756, 1986.

[2] S. Bakas, H. Akbari, A. Sotiras, M. Bilello, M. Rozycki, J. S. Kirby, J. B. Freymann, K. Farahani, and C. Davatzikos. Advancing the cancer genome atlas glioma MRI collections with expert segmentation labels and radiomic features. *Scientific Data*, 4(170117):170117, 2017.

[3] A. Chartsias, T. Joyce, M. V. Giuffrida, and S. A. Tsaftaris. Multimodal MR synthesis via modality-invariant latent representation. *TMI*, 37(3):803–814, 2018.

[4] J. Chen, J. Chen, H. Chao, and M. Yang. Image blind denoising with generative adversarial network based noise modeling. In *CVPR*, June 2018.

[5] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected CRFs. *TPAMI*, 40(4):834–848, 2018.

[6] X. Chen and E. Konukoglu. Unsupervised detection of lesions in brain MRI using constrained adversarial auto-encoders. In *International Conference on Medical Imaging with Deep Learning*, 2018.

[7] P. Costa, A. Galdran, M. I. Meyer, M. Niemeijer, M. Abràmoff, A. M. Mendonça, and A. Campilho. End-to-end adversarial retinal image synthesis. *TMI*, 37(3):781–791, 2018.

[8] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *NIPS*, pages 2672–2680, 2014.

[9] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask R-CNN. In *ICCV*, pages 2980–2988, 2017.

[10] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016.

[11] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016.

[12] T. Huynh, Y. Gao, J. Kang, W. Li, Z. Pei, J. Lian, and D. Shen. Estimating CT image from MRI data using structured random forest and auto-context model. *TMI*, 35(1):174–183, 2015.

[13] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. In *CVPR*, pages 1125–1134. IEEE, 2017.

[14] J. Johnson, A. Alahi, and F. F. Li. Perceptual losses for real-time style transfer and super-resolution. In *ECCV*, pages 694–711, 2016.

[15] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, pages 1097–1105, 2012.

[16] C. Ledig, L. Theis, F. Huszar, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang, and W. Shi. Photo-realistic single image super-resolution using a generative adversarial network. In *CVPR*, July 2017.

[17] R. Li, W. Zhang, H. I. Suk, L. Wang, J. Li, D. Shen, and S. Ji. Deep learning based imaging data completion for improved brain disease diagnosis. In *MICCAI*, pages 305–312. Springer, 2014.

[18] G. Litjens, T. Kooi, B. E. Bejnordi, A. A. A. Setio, F. Ciompi, M. Ghafoorian, J. A. van der Laak, B. Van Ginneken, and C. I. Sánchez. A survey on deep learning in medical image analysis. *MedIA*, 42:60–88, 2017.

[19] M.-Y. Liu, T. Breuel, and J. Kautz. Unsupervised image-to-image translation networks. In *NIPS*, pages 700–708, 2017.

[20] M.-Y. Liu and O. Tuzel. Coupled generative adversarial networks. In *NIPS*, pages 469–477, 2016.

[21] B. H. Menze, A. Jakab, S. Bauer, J. Kalpathy-Cramer, K. Farahani, J. Kirby, Y. Burren, N. Porz, J. Slotboom, R. Wiest, et al. The multimodal brain tumor image segmentation benchmark (BRATS). *TMI*, 34(10):1993–2024, 2015.

[22] D. Nie, R. Trullo, J. Lian, C. Petitjean, R. Su, Q. Wang, and D. Shen. Medical image synthesis with context-aware generative adversarial networks. In *MICCAI*, pages 417–425. Springer, 2017.

[23] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, pages 234–241. Springer, 2015.

[24] T. Schlegl, P. Seebock, S. M. Waldstein, U. Schmidt-Erfurth, and G. Langs. Unsupervised anomaly detection with generative adversarial networks to guide marker discovery. In *IPMI*, pages 146–157, 2017.

[25] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.

[26] L. Sun, Z. Fan, Y. Huang, X. Ding, and J. Paisley. Compressed sensing MRI using a recursive dilated network. In *AAAI*, pages 2444–2451, 2018.

[27] K.-L. Tseng, Y.-L. Lin, W. Hsu, and C.-Y. Huang. Joint sequence learning and cross-modality convolution for 3d biomedical segmentation. In *CVPR*, pages 3739–3746. IEEE, 2017.

[28] D. Ulyanov, A. Vedaldi, and V. S. Lempitsky. Instance normalization: The missing ingredient for fast stylization. *CoRR*, abs/1607.08022, 2016.

[29] T.-C. Wang, M.-Y. Liu, J.-Y. Zhu, A. Tao, J. Kautz, and B. Catanzaro. High-resolution image synthesis and semantic manipulation with conditional GANs. In *CVPR*. IEEE, 2018.

[30] Q. Wu, P. Wang, C. Shen, I. Reid, and A. van den Hengel. Are you talking to me? reasoned visual dialog generation through adversarial learning. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.

[31] W. Yang, W. Ouyang, X. Wang, J. Ren, H. Li, and X. Wang. 3d human pose estimation in the wild by adversarial learning. In *CVPR*, June 2018.

[32] Z. Yi, H. Zhang, P. Tan, and M. Gong. DualGAN: Unsupervised dual learning for image-to-image translation. In *ICCV*, pages 2868–2876. IEEE, 2017.

[33] L. Yu, X. Yang, C. Hao, J. Qin, and P. A. Heng. Volumetric convnets with mixed residual connections for automated prostate segmentation from 3D MR images. In *AAAI*, pages 66–72, 2017.

[34] Z. Zhang, L. Yang, and Y. Zheng. Translating and segmenting multimodal medical volumes with cycle- and shape-consistency generative adversarial network. In *CVPR*, 2018.

[35] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *ICCV*, pages 2242–2251. IEEE, 2017.