

# Distribution Matching Losses Can Hallucinate Features in Medical Image Translation

Joseph Paul Cohen, Margaux Luck, Sina Honari

Montreal Institute for Learning Algorithms, University of Montreal  
{cohenjos, luckmarg, honaris}@iro.umontreal.ca

**Abstract.** This paper discusses how distribution matching losses, such as those used in CycleGAN, when used to synthesize medical images can lead to mis-diagnosis of medical conditions. It seems appealing to use these new image synthesis methods for translating images from a source to a target domain because they can produce high quality images and some even do not require paired data. However, the basis of how these image translation models work is through matching the translation output to the distribution of the target domain. This can cause an issue when the data provided in the target domain has an over or under representation of some classes (e.g. healthy or sick). When the output of an algorithm is a transformed image there are uncertainties whether all known and unknown class labels have been preserved or changed. Therefore, we recommend that these translated images should not be used for direct interpretation (e.g. by doctors) because they may lead to misdiagnosis of patients based on hallucinated image features by an algorithm that matches a distribution. However there are many recent papers that seem as though this is the goal.

**Keywords:** distribution matching, image synthesis, domain adaptation

## 1 Introduction

The introduction of adversarial losses [1] made it possible to train new kinds of models based on implicit distribution matching. Recently, adversarial approaches such as CycleGAN [2], pix2pix [3], UNIT [4], Adversarially Learned Inference (ALI) [5], and GibbsNet [6] have been proposed for un-paired and paired image translation between two domains. These approaches have been used recently in medical imaging research for translating images between domains such as MRI and CT. However, there is a bias when the output of these models are used for interpretation. When translating images from a source domain to a target domain, these models are trained to match the target domain distribution, where they may hallucinate images by adding or removing image features. This can cause a problem when the target distribution during training has over or under representation of known or unknown labels compared to the test time distribution. Due to such a bias, we recommend until better solutions are proposed that maintain the vital information, such translated images should not be used

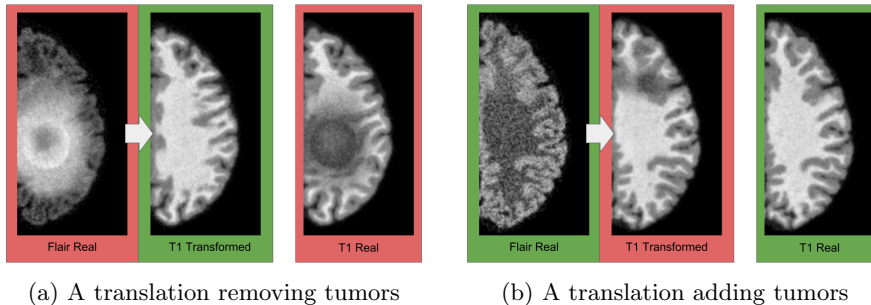


Fig. 1: Examples of two CycleGANs trained to transform MRI images from Flair to T1 types. We show healthy images in green and tumor images in red. In (a) the model was trained with a bias to remove tumors because the target distribution did not have any tumor examples so the transformation was forced to remove tumors in order to match the target distribution. Conversely in (b) the tumors were added to the image to match the distribution which was composed of only tumor examples during training.

for medical diagnosis, since they can lead to mis-diagnosis of medical conditions. This issue should be discussed because recently several papers have been published performing image translation using distribution matching. The main motivation for many of these approaches was to translate images from a source domain to a target domain such that they could be later used for interpretation (e.g. by doctors). Applications include MR to CT [7, 8], CS-MRI [9, 10], CT to PET [11], and automatic H&E staining [12].

We demonstrate the problem with a caricature example in Figure 1 where we *cure cancer* (in images) and *cause cancer* (in images) using a CycleGAN that translates between Flair and T1 MRI samples. In Figure 1(a) the model has been trained only on healthy T1 samples which causes it to remove cancer from the image. This model has learned to match the target distribution regardless of maintaining features that are present in the image. In the following sections, we demonstrate how these methods introduce a bias in image translation due to matching the target distribution.

We draw attention to this issue in the specific use case where the images are presented for interpretation. However, we do not aim to discourage work using these losses for data augmentation to improve the performance of a classification, segmentation, or other model.

## 2 Problem Statement

Our argument is that the composition of the source and target domains can bias the image transformation to cause an unwanted feature hallucination. We systematically review the objective functions used for image translation in Table 1 and discuss how they each exhibit this bias.

Let's first consider a standard GAN model [1] where the generator is a transformation function  $f_{a,b}(a)$  which maps samples from the source domain  $D_a$  to

samples from the target domain  $D_b$ . The discriminator is trained given samples from  $D_b$  through which the transformation function can match the distribution of  $D_b$ .

$$\text{GAN Disc: } \max_{b \sim D_b} \mathbb{E} [\log D(b)] + \mathbb{E}_{a \sim D_a} [\log(1 - D(f_{a,b}(a)))]$$

In order to minimize this objective the transformation function will need to produce images that match real images from the distribution  $D_b$ . Here there are no constraints to force a correct mapping between  $D_a$  and  $D_b$ , so for a non-finite  $D_a$  we can consider it to be equal to a Gaussian noise  $\mathcal{N}$  typically used in a GAN.

In order to better enforce the mapping between the domains CycleGAN [2] extends the generator loss to include cycle consistency terms:

$$\text{Cycle Consistency: } |f_{b,a}(f_{a,b}(a)) - a|$$

Here the function  $f_{a,b}$  is composed of the inverse transformation  $f_{b,a}$  to create a reconstruction loss that will regularize both transformations to not ignore the source image. However, this process does not provide a guarantee that a correct mapping will be made. In order to match the target distribution, image features can be hallucinated and information to reconstruct an image in the other domain can be encoded [13]. Moreover, due to having un-paired source and target data, the target distribution that the generator is trained on may be even distinct from the target distribution that corresponds to the data in the source domain (e.g. having only tumor targets while the source is all healthy). This makes the models such as CycleGAN even more prone to hallucinate features due to the way the data in the target domain is gathered.

Another approach to solve this problem is using a conditional discriminator [3, 14]. The intuition here is that giving the discriminator the source image  $a$  as well as the transformed image  $f_{a,b}(a)$ , we can model the joint distribution. This approach requires paired examples in order to provide real source and target pairs to the discriminator. The dataset  $D_b$  still plays a role in determining what the discriminator learns and therefore how the transformation function operates. The discriminator is trained by:

$$\max_{(a,b) \sim (D_a, D_b)} \mathbb{E} [\log D(b, a)] + \mathbb{E}_{a \sim D_a} [\log(1 - D(f_{a,b}(a), a))]$$

Even in the case of CondGAN that the source and target domain distributions correspond to each other due to having paired data, the discriminator can assign more/less capacity to a feature (e.g. tumors), due to having over/under representation of those features in the target distribution. This can be a source of bias in how those features are translated.

Finally, we look at how to train a transformation using only a L1 loss without any adversarial distribution matching term. With this classic approach we consider transformations based on minimizing the pixel wise error:

$$\mathbb{E}_{(a,b) \sim (D_a, D_b)} \|f_{a,b}(a) - b\|_1$$

Table 1: Loss formulations divided into two phases of training. On the left the discriminator loss is shown (when applicable) and on the right the transformation/generator loss is shown. Note that for GAN losses the generator matches the target distribution indirectly through gradients it receives from the discriminator.

|          | Discriminator Loss (max)   | Domain Transformer/Generator Loss (min)                                      |
|----------|--|--|
| GAN      | $\mathbb{E}_{b \sim D_b} [\log D(b)] + \mathbb{E}_{a \sim D_a} [\log(1 - D(f_{a,b}(a)))]$                  | $\mathbb{E}_{a \sim D_a} [-\log(D(f_{a,b}(a)))]$                             |
| CycleGAN | $\mathbb{E}_{b \sim D_b} [\log D(b)] + \mathbb{E}_{a \sim D_a} [\log(1 - D(f_{a,b}(a)))]$                  | $\mathbb{E}_{a \sim D_a} [-\log(D(f_{a,b}(a))) +  f_{b,a}(f_{a,b}(a)) - a ]$ |
| CondGAN  | $\mathbb{E}_{(a,b) \sim (D_a, D_b)} [\log D(b, a)] + \mathbb{E}_{a \sim D_a} [\log(1 - D(f_{a,b}(a), a))]$ | $\mathbb{E}_{a \sim D_a} [-\log(D(f_{a,b}(a), a))]$                          |
| L1       | -  | $\mathbb{E}_{(a,b) \sim (D_a, D_b)} \ f_{a,b}(a) - b\ _1$                    |

Unlike GAN models that match the target distribution over the entire image, L1 predicts each pixel locally given its receptive field without the need to account for global consistency. As long as some pixels present the category of interest in the image (e.g. tumor), L1 can learn a mapping. However, L1 still can suffer from a bias when the train and test distributions are different, e.g. when no tumor pixels are provided during training, which can be caused by having new known or unknown labels at test time.

With all these approaches to domain translation we find there is the potential for bias in the training data (specifically  $D_b$  for our experiments below).

### 3 Bias Impact

We use the BRATS2013 [15] synthetic MRI dataset because we can visually inspect the presence of a tumor, it is freely available to the public, and we have paired data to inspect results. Our task for analysis is to transform Flair MRI images (source domain) into T1-weighted images (target domain). We start with 1700 image slices where 50% are healthy and 50% have tumors. We use 1400 to construct training sets for the models and 300 as a holdout test set used to test if the transformation added or removed tumors.

In this section, we construct two training scenarios: unpaired and paired. For the CycleGAN we use an unpaired training scenario which keeps the distribution fixed in the source domain (with 50% healthy and 50% tumor samples) and changes the ratio of healthy to cancer samples in the target domain  $D_b$  to simulate how the distribution matching works when the target distribution is irrelevant to the source distribution. For the CondGAN and L1 models we use a paired training scenario where both the source and target domains have the same proportion of healthy to tumor examples because they have to be presented as pairs to the model.

We train 3 models under 11 different percentages of tumor examples in the target distribution, which vary from 0% to 100% with tumors. In place of a doctor to classify the transformed samples we use an impartial CNN classifier (4 convolutional layers with ReLU and Stride-2 convolutions, 1 fully connected layer with no non-linearity, and a two-way softmax output layer) which obtains

80% accuracy on the test set. The results of using this classifier on the generated T1 samples with different target domain composition is shown in Figure 2. As we change the composition of the target domain we can observe the bias impact on the class of the transformed examples from the holdout test set. If there was no bias in matching the target distribution due to the composition of the samples in the target domain, there would be no difference in the percentage of the images diagnosed with a tumor as we change the target domain composition in Figure 2. We also compute the mean absolute pixel reconstruction error between the ground truth image in the target domain and the translated image. If a large feature is added or removed it should produce a large pixel error. If the translation was doing a perfect job, the pixel error should have been 0 for all cases.

We draw the readers attention to CycleGAN which produces the most dramatic change in class labels, since the model learns to map a balanced (tumor to healthy) source domain to an unbalanced composition in the target domain, which encourages the model to add or remove features (see samples in Figure S1). This indicates such models are subject to even more bias due to the composition of the features in the target domain that can be different from the ones in the source domain.

For CondGAN, the pixel error changes across as the composition of tumor/healthy changes, indicating there is a bias due to the training data composition. Perceptually the L1 loss appears the most consistent producing the least bias. However, it has error when it is trained on 0% tumor and the model is asked to translate tumor samples at test time (0% for L1 in Figure 2 bottom row and Figure 3 (a)), which is due to a mis-match between train and test distributions. It indicates that if at test time images with new known or unknown labels (e.g. a new disease) are presented to the model, it cannot transform them properly. In Figure 3 we show examples of the translated images between the models. Note how for GAN based models the cancer tumor gradually appears and gets bigger from left to right. L1 mostly suffers in Figure 3 (a) for 0%. Interestingly, in the case of 100% tumor it can translate healthy images even though it was not trained with healthy images. We believe this is due to having both healthy and tumor regions in each image which allows the network to see healthy sub-regions and learn to translate both categories. Further samples are available in the supplementary information in Figures S2 and S3.

## 4 Conclusion

In this work we discussed concerns about how distribution matching losses, such as those used in CycleGAN, can lead to mis-diagnosis of medical conditions. We have presented experimental evidence that when the output of an algorithm matches a distribution, for unpaired or paired data translation, all known and unknown class labels might not be preserved. Therefore, these translated images should not be used for interpretation (e.g. by doctors) without proper tools to verify the translation process. We illustrate this problem using dramatic examples of tumors being added and removed from MRI images. We hope that future

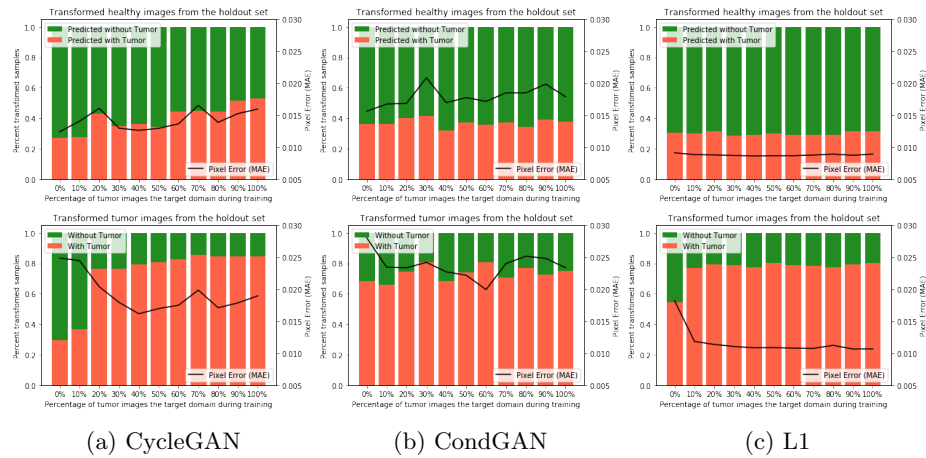


Fig. 2: We plot the classifier’s prediction on 300 (53% tumor) unseen samples (holdout test set) as we vary the distribution of tumor samples in the target domain from 0% to 100% of three models (CycleGAN, CondGAN, L1). This corresponds to 33 trained models. We split the source domain samples of the holdout test set into healthy (top row) and tumor (bottom row) and apply a classifier on the translated images. **Green** represents translated samples predicted by the classifier as healthy and **red** represents samples predicted with tumors. If the translation was without bias the percentage of healthy to tumor images should not change across the 11 models trained for each loss. For CycleGAN, we observe that the percentage of the images diagnosed with tumors increases as the percentage of tumor images in the target distribution increases. The black line represents the mean absolute pixel error between translated and ground truth target samples. While CondGAN seems to have a more stable classification results compared to CycleGAN, the pixel error indicates how much the translated images are away from ground truth samples and subject to change for different percentage of tumor composition in the target domain. L1 loss seem to suffer the least from target distribution matching and produces high error only when the target distribution has 0% of tumors (during training) and is asked to translate tumor samples. This case corresponds to 0% L1 on the bottom row.

methods will take steps to ensure that this bias does not influence the outcome of a medical diagnosis.

## Acknowledgements

We thank Adriana Romero Soriano, Michal Drozdal, and Mohammad Havaei for their valuable input and assistance on the project. This work is partially funded by a grant from the U.S. National Science Foundation Graduate Research Fellowship Program (grant number: DGE-1356104) and the Institut de valorisation des donnees (IVADO). This work utilized the supercomputing facilities managed by the Montreal Institute for Learning Algorithms, NSERC, Compute Canada, and Calcul Quebec.

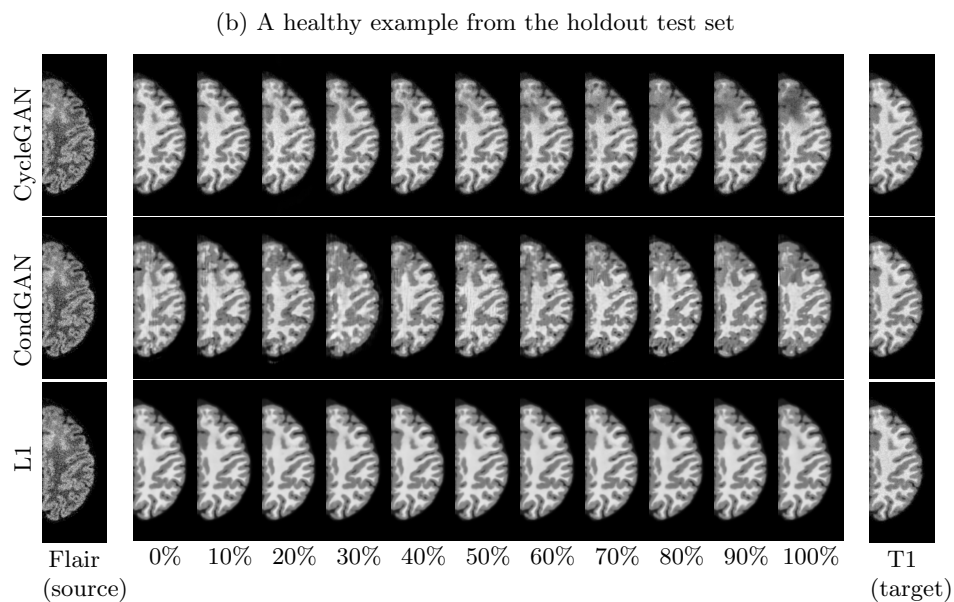
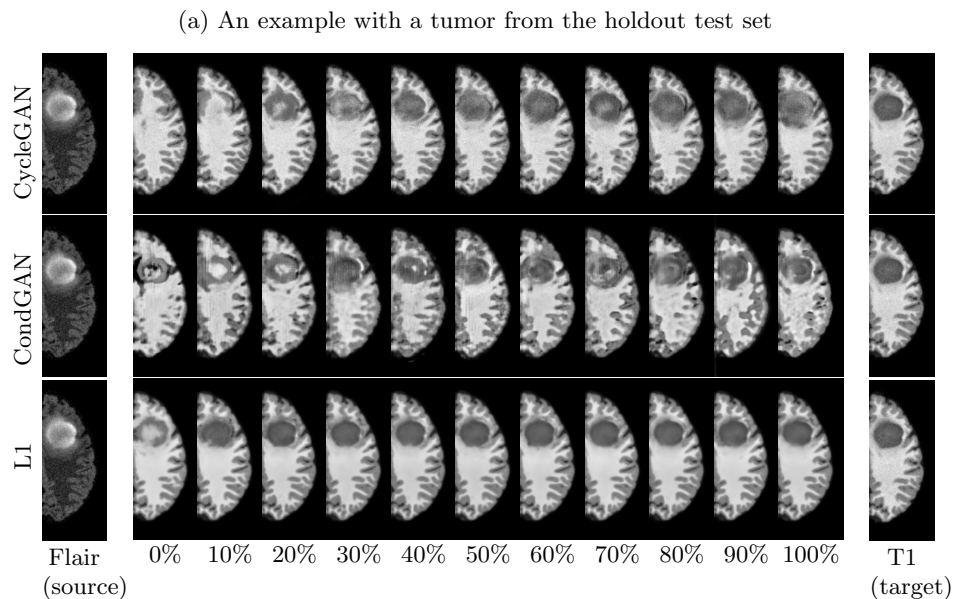


Fig. 3: Illustration of tumor (a) and healthy (b) class change through domain translation while changing the ratio of the healthy to tumor samples in the target domain  $D_b$  for all three models (CycleGAN, CondGAN, L1). We vary the distribution of  $D_b$  from 0% tumor to 100% examples to train 33 different models. We show images of the source domain (Flair) on the left and the corresponding ground truth image in the target domain (T1) on the right. We can observe visually the magnitude of the changes introduced.

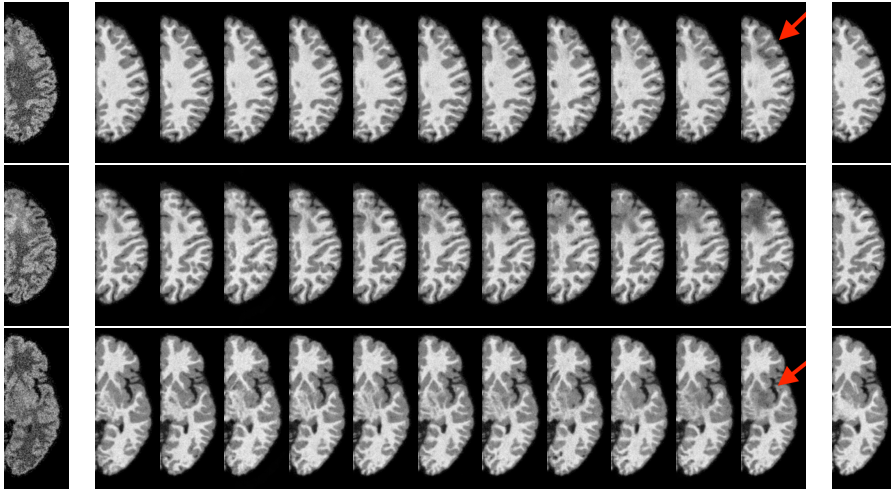
## Bibliography

- [1] Goodfellow, I.J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative Adversarial Networks. In: *Neural Information Processing Systems* (2014)
- [2] Zhu, J.Y., Park, T., Isola, P., Efros, A.A.: Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks. In: *International Conference on Computer Vision* (2017)
- [3] Isola, P., Zhu, J.Y., Zhou, T., Efros, A.A.: Image-to-image translation with conditional adversarial networks. In: *Computer Vision and Pattern Recognition* (2017)
- [4] Liu, M.Y., Breuel, T., Kautz, J.: Unsupervised Image-to-Image Translation Networks. In: *Neural Information Processing Systems* (2017)
- [5] Dumoulin, V., Belghazi, I., Poole, B., Mastropietro, O., Lamb, A., Arjovsky, M., Courville, A.: Adversarially Learned Inference. In: *International Conference on Learning Representations* (2017)
- [6] Lamb, A., Hjelm, D., Ganin, Y., Cohen, J.P., Courville, A., Bengio, Y.: GibbsNet: Iterative Adversarial Inference for Deep Graphical Models. In: *Neural Information Processing Systems* (2017)
- [7] Wolterink, J.M., Dinkla, A.M., Savenije, M.H., Seevinck, P.R., van den Berg, C.A., Išgum, I.: Deep MR to CT synthesis using unpaired data. In: *Workshop on Simulation and Synthesis in Medical Imaging* (2017)
- [8] Nie, D., Trullo, R., Petitjean, C., Ruan, S., Shen, D.: Medical Image Synthesis with Context-Aware Generative Adversarial Networks. In: *Medical Image Computing and Computer-Assisted Intervention* (2016)
- [9] Quan, T.M., Nguyen-Duc, T., Jeong, W.K.: Compressed Sensing MRI Reconstruction using a Generative Adversarial Network with a Cyclic Loss. *IEEE Transactions on Medical Imaging* (2018)
- [10] Yang, G., Yu, S., Dong, H., Slabaugh, G., Dragotti, P.L., Ye, X., Liu, F., Arridge, S., Keegan, J., Guo, Y., Firmin, D.: DAGAN: Deep De-Aliasing Generative Adversarial Networks for Fast Compressed Sensing MRI Reconstruction. *IEEE Transactions on Medical Imaging* (2018)
- [11] Ben-Cohen, A., Klang, E., Raskin, S.P., Amitai, M.M., Greenspan, H.: Virtual PET Images from CT Data Using Deep Convolutional Networks: Initial Results. In: *MICCAI Workshop on Simulation and Synthesis in Medical Imaging* (2017)
- [12] Bayramolu, N., Kaakinen, M., Eklund, L.: Towards Virtual H&E Staining of Hyperspectral Lung Histology Images Using Conditional Generative Adversarial Networks. In: *International Conference on Computer Vision* (2017)
- [13] Chu, C., Zhmoginov, A., Sandler, M.: CycleGAN, a Master of Steganography. In: *Neural Information Processing Systems Workshop on Machine Deception* (2017)
- [14] Mirza, M., Osindero, S.: Conditional Generative Adversarial Nets. *arXiv: 1411.1784* (2014)
- [15] Menze, B.H., Jakab, A., Bauer, S., et al.: The Multimodal Brain Tumor Image Segmentation Benchmark (BRATS). *IEEE Transactions on Medical Imaging* (2015)



## Supplementary Information

(a) Healthy examples from the holdout test set



(b) Examples with a tumor from the holdout test set

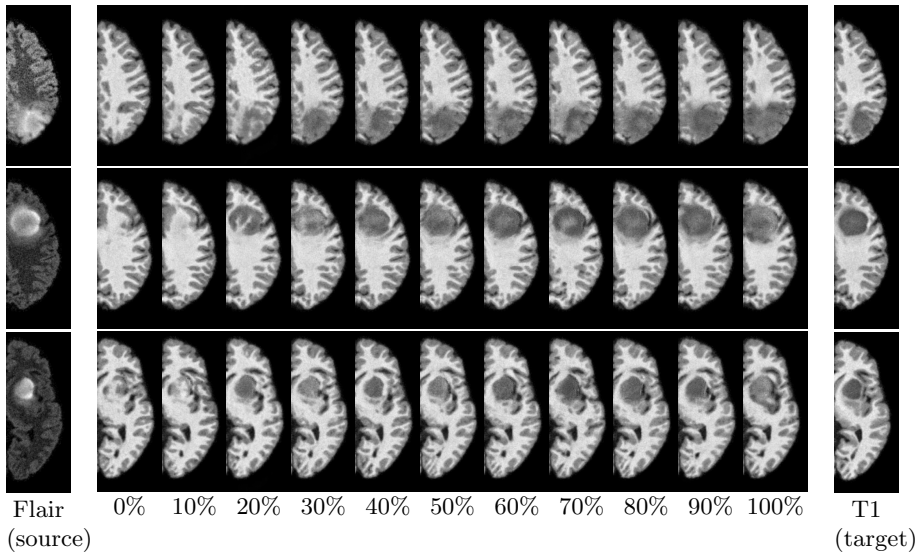
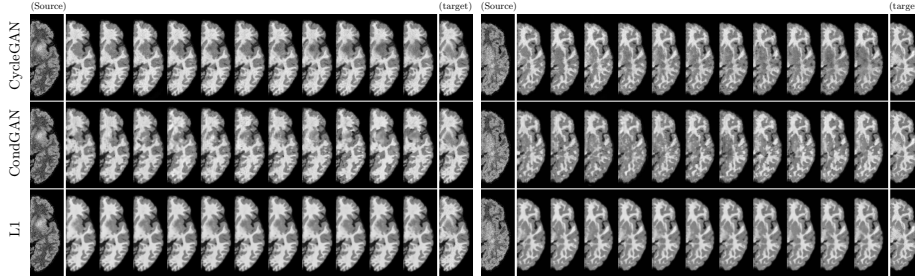


Fig. S1: Illustration of healthy (a) and tumor (b) class change through domain translation while changing the ratio of the healthy to tumor samples in the target domain  $D_b$  for CycleGAN. We vary the distribution of  $D_b$  from 0% tumor examples to 100% tumor examples to train 11 different CycleGAN models. On the left you see the image in the source domain (Flair) and on the right you see the corresponding ground truth image in the target domain (T1). The red arrows point out hard to spot tumors that were introduced.

(a) When source image is healthy and classifier predicts healthy



(b) When source image has tumor and classifier predicts tumor

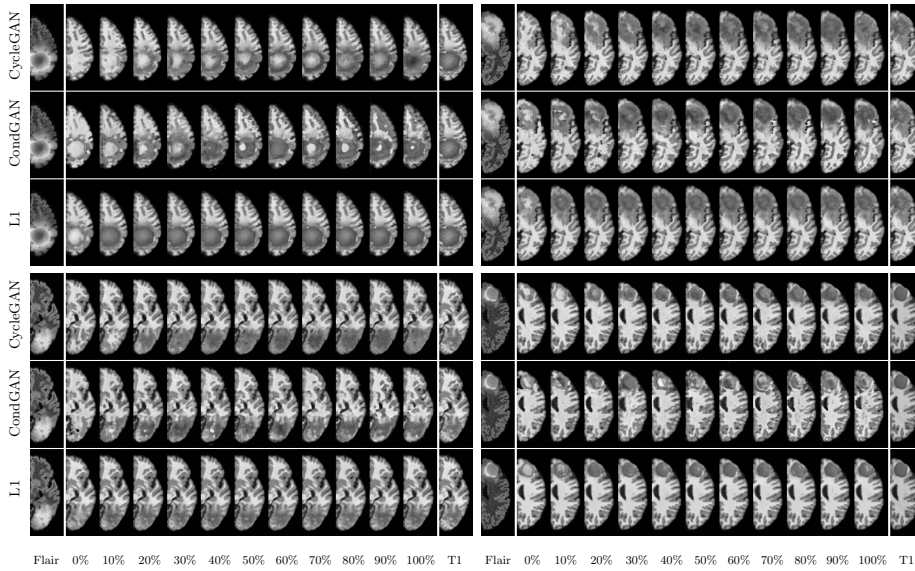


Fig. S2: Showing translation bias when the classifier detects the class properly. On the top row, two examples are shown where the source image is healthy and the classifier predicts healthy for all 33 translated images across all 3 losses (CycleGAN, CondGAN, L1). Note how the translated images are still different from the ground truth T1 image in the target domain. On the bottom two rows, the source image has a tumor and the classifier predicts tumor for all 33 translated images. Note how the tumor seem to disappear in the images on the left and gets bigger in the images on the right. L1 loss suffers less from these artifacts, while adversarial losses produce more erroneous results. These samples show that even when the classifier is doing a perfect job, there is a bias in the translated images, and the bias is different for different losses.

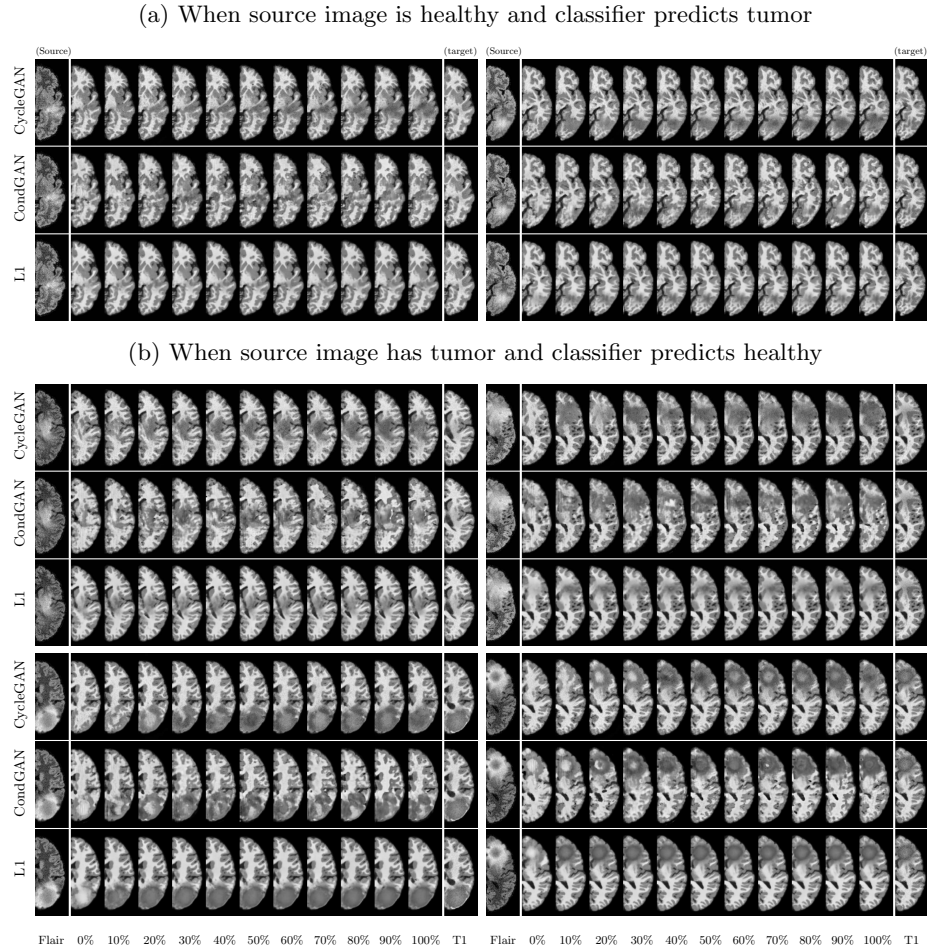


Fig. S3: Showing translation bias when the classifier detects the class wrongly. On the top row, two examples are shown where the source image is healthy and the classifier predicts tumor for all 33 translated images across all 3 losses (CycleGAN, CondGAN, L1). Despite having an imperfect translation, which leads to a tumor-like shape in the translation, see how the degree of bias is different across different loss functions. On the bottom two rows, four samples are shown where the source image has a tumor and the classifier predicts healthy for all 33 translated images. While classifier is doing a poor job here, there is a bias in the translated samples compared to the ground truth T1 samples. L1 suffers less from artifacts. These images show when the classifier detects the class wrongly, which can be due to having a weak classifier or imperfect translation, the bias in translation can still be observed.