

MouseGAN++: Unsupervised Disentanglement and Contrastive Representation for Multiple MRI Modalities Synthesis and Structural Segmentation of Mouse Brain

Ziqi Yu, Xiaoyang Han, Shengjie Zhang, Jianfeng Feng, Tingying Peng* and Xiao-Yong Zhang*

Abstract—Segmenting the fine structure of the mouse brain on magnetic resonance (MR) images is critical for delineating morphological regions, analyzing brain function, and understanding their relationships. Compared to a single MRI modality, multimodal MRI data provide complementary tissue features that can be exploited by deep learning models, resulting in better segmentation results. However, multimodal mouse brain MRI data is often lacking, making automatic segmentation of mouse brain fine structure a very challenging task. To address this issue, it is necessary to fuse multimodal MRI data to produce distinguished contrasts in different brain structures. Hence, we propose a novel disentangled and contrastive GAN-based framework, named MouseGAN++, to synthesize multiple MR modalities from single ones in a structure-preserving manner, thus improving the segmentation performance by imputing missing modalities and multi-modality fusion. Our results demonstrate that the translation performance of our method outperforms the state-of-the-art methods. Using the subsequently learned modality-invariant information as well as the modality-translated images, MouseGAN++ can segment fine brain structures with averaged dice coefficients of 90.0% (T2w) and 87.9% (T1w), respectively, achieving around +10% performance improvement compared to the state-of-the-art algorithms. Our results demonstrate that MouseGAN++, as a simultaneous image synthesis and segmentation method, can be used to fuse cross-modality information in an unpaired manner and yield more robust performance in the absence of multimodal data. We release our method as a mouse brain structural segmentation tool for free academic usage at <https://github.com/you2019>.

Index Terms—Segmentation, Mouse brain, MRI, Generative adversarial network, Disentangled representations.

I. INTRODUCTION

ACCURATE segmentation of brain structures on magnetic resonance (MR) images is crucial for delineating morphological structures, analyzing brain functions, and understanding their relationships. As one of the most important model organisms, the mouse plays an important role in neuroscience, drug discovery, and translational medicine. Since the mouse and human brains are evolutionarily conserved, the mouse brain has proven to be a powerful model for understanding human brain.

Currently, automatic segmentation methods of brain structures have been developed to segment human brains. However, when they are applied to MRI data of mouse brains, their performance remarkably suffers from the differences in image contrast, image size, and anatomy. Segmenting fine mouse brain structures on MRI data using automatic methods has remained a challenging task until now. The main reason is that accurate segmentation of fine brain structures usually requires multi-modality MRI with high-resolution, which provides much more complementary feature information than a single modality [1]. For example, T1w images provide contrast differences between grey and white matter; T2w images are more sensitive to water-rich tissues; and quantitative susceptibility mapping (QSM) is suitable for differentiating deep brain tissues. These modalities described above provide complementary information for accurate segmentation of brain anatomical structures. However, such multi-modal mouse brain MRI data are often lacking because collecting such data takes too much scan time, which is impractical in most preclinical MR facilities. Therefore, mouse brain fine structure segmentation suffers from missing some MRI modalities in practice, but multi-modality fusion is expected to alleviate this dilemma.

Due to the restriction on acquiring a series of multi-modality images followed by segmentation, imputing the missing modality and decoupling semantic features has emerged as a crucial research area. In the computer vision field, image-to-image translation methods based on a generative adversarial network (GAN) have been demonstrated to be successful

This work was supported in part by National Natural Science Foundation of China under Grants 81873893, 82171903, 92043301; in part by Shanghai Municipal Science and Technology Major Project under Grant 2018SHZDZX01 and ZJLab. *Corresponding authors: Xiao-Yong Zhang: xiaoyong.zhang@fudan.edu.cn and Tingying Peng: tingying.peng@helmholtz-muenchen.de

Z. Yu, X. Han, S. Zhang, J. Feng and X.-Y. Zhang are with Institute of Science and Technology for Brain-Inspired Intelligence, Fudan University, Shanghai, China. They are also with MOE Key Laboratory of Computational Neuroscience and Brain-Inspired Intelligence, and MOE Frontiers Center for Brain Science, Fudan University, Shanghai, China (e-mail: xiaoyong_zhang@fudan.edu.cn)

T. Peng is with Helmholtz AI, Helmholtz Zentrum Muenchen, Munich, Germany (e-mail: tingying.peng@helmholtz-muenchen.de).

© 20xx IEEE. Personal use of this material is permitted. Permission from IEEE must be obtained for all other uses, in any current or future media, including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works

in image synthesis [2]–[4]. Particularly, [5] attempted to disentangle the modality-specific information and modality-invariant representations. In the medical image field, several existing cross-modality works have been applied to segment brain lesions [6], [7]. Whilst performing satisfactorily for human brain lesion segmentation tasks, these methods are not readily applicable to our mouse brain segmentation problem as they have several unsolved issues, such as patch-level context might not be appropriate for mouse brain regions due to their finer structures. Moreover, unlike lesions with distinguishable contrast on one specific MRI modality, there are often no significant contrast differences between different brain structures on single-modality MRI images. As a result, it is difficult to accurately delineate their anatomical boundaries using single-modal MRI data. To solve this problem, we aim to propose a synthesis-and-segmentation deep learning framework to first synthesize multi-modality MRI data, hence achieving more accurate segmentation of fine structures in the mouse brain via multimodal image fusion.

Specifically, we propose a disentangled and contrastive GAN-based framework, termed as MouseGAN++, which is a unified model that combines modality synthesis and structure segmentation. MouseGAN++ contains a modality translation module with two novel contrastive losses to project multi-modality image features into a shared latent content space that encodes modality-invariant brain structures and modality-specific attributes. Subsequently, the latent contents are combined with other modality-specific attributes to impute images of other modalities. Concretely, in MouseGAN++, the content contrastive loss is proposed to compel the network to avoid confusing structure-wise information during translation. We reuse attribute and content encoders during adversarial training to produce contrastive learning concurrently. The shared content space also facilitates decoder training in the segmentation module. Moreover, imputing modality with this model can enlarge the dataset, enabling the network to jointly learn modality-invariant semantic representations, thereby enhancing attention to multi-modality fusion.

This work significantly improves our previous model, MouseGAN, published in MICCAI conference [8]. Compared with MouseGAN and other related state-of-the-art (SOTA) work, MouseGAN++ has the following contributions:

- MouseGAN++ uses an unsupervised disentanglement and contrastive framework that concurrently optimizes adversarial loss and contrastive loss without designing and pre-training on additional pretext tasks, thus eliminating the gap between pretext and translation tasks.
- MouseGAN++ disentangles MR images into attribute and content spaces, where the modality-agnostic content features can be flexibly recovered by various attributes or segmented by the decoder. Notably, two novel attribute and content contrastive losses introduced in MouseGAN++ further improve the disentanglement and the downstream segmentation performance.
- We performed extensive experiments to evaluate both translation and segmentation tasks for the mouse brain, and provided the packaged MouseGAN++ as a pipeline that could further assist the mouse community.

II. RELATED WORK

A. Mouse Brain Segmentation

Brain structure segmentation in mouse MR images is a fundamental step in neuroimaging and preclinical studies. With the development of the active mouse research community, there have been many efforts to address this segmentation problem. To date, they can be mainly grouped into two categories based on their methodologies: atlas-based methods and deep learning methods. The atlas-based methods, which use image registration to propagate structural label information from a specific atlas volume to a native space corresponding to each subject volume, are primarily used in mouse brain segmentation pipelines [9], [10]. So far, only a few registration-based toolboxes, e.g., [11], [12], can be used for brain segmentation. However, small and elongated regions are more susceptible to biases introduced by the registration procedure. Besides, the accessibility of suitable atlases and the complexity of making them further impede the spread of these atlas-based methods.

Recently, despite the successful application of deep learning in the field of segmentation, only a few attempts have been made to segment mouse brain structures, e.g., MU-Net [13]. However, due to the lack of prior knowledge of multi-modality contrasts, there is still a considerable performance gap when dealing with fine structure segmentation. In addition, how to disentangle these semantic representations in multimodal MRI data also presents challenges for the network design.

B. Image Synthesis and Disentangled Representation

Image synthesis techniques have made significant strides in recent years. Since modality synthesis can be deemed as an image-to-image translation task, GAN is considered the ideal model as it can be constrained on the conditional image, thus it has been widely used in medical images including MRI, CT, and PET [14]–[17]. Additionally, GAN brings strong disentanglement prior because of its hierarchical structure [18], which naturally serves our purpose.

CycleGAN is a seminal paper dealing with this task [2]. Following the success of CycleGAN, recent efforts to exploit GAN mainly focus on the different construction strategies for the encoder and decoder, such as DiscoGAN [4] and DualGAN [3]. UNIT [19] proposes a shared latent space assumption for a better perceptual appearance. As most previous models consider a mapping between two domains, the scalability is limited when dealing with multiple-domain translation. Then MUNIT [20], DRIT++ [5] and StarGAN-v2 [21] are proposed to address this limitation. Nevertheless, [5], [20], [22] assume the latent attribute spaces are under the constraint of the Gaussian prior, and [21] uses a Gaussian noise-into-style mapping network, which might lead to inadequate disentanglement. Thus, unlike the existing works, we explicitly impose a disentanglement prior to our MouseGAN++.

The core idea of disentangled representation learning approaches follows a prior that one model can learn to embed images into two spaces: a modality-specific attribute space and a modality-independent content space. By assuming that a shared modality-independent content space that maintains

structural information can be exploited for both modalities, previous work mainly concentrated on unsupervised domain adaptation. For instance, Chen et al. [23] proposes a method that conducts synergistic alignment of modalities from both image and feature perspectives. For fusing the complementary information from multimodal data, Chen et al. [7] proposes a learning framework with feature disentanglement and gated feature fusion for robust segmentation and Yang et al. [6] adopts contrastive loss for brain tumor segmentation, however, the paired training data they used is prohibitively expensive in clinics for the research. Recently, extensive efforts on GAN-based disentanglement methods have been explored. In [14], [24]–[30], they share the key idea of disentangling input images into modality-specific and modality-independent features. Since a majority of them introduce Gaussian priors into their latent attribute spaces, it might cause confusion in multi-modality feature distribution. In contrast to previous works, our method proposes contrastive learning prior into latent spaces by combining attribute features and shared content features with contrastive loss to facilitate learning various cross-modality mappings and improve segmentation performance.

C. Contrastive Learning

As a popular branch of self-supervised learning, the network guided by contrastive learning can learn hierarchical features from vast unlabeled datasets, which benefits the training of downstream tasks. Specifically, a contrastive loss is used to push positive pair representations to be similar and negative pair representations to be differentiated. SimCLR [31] and MoCo [32] are two representative methods that provide two training strategies with SOTA performance. Note that most contrastive learning methods focus on image classification, assuming that the instances in two individual images have distinct characteristics. However, when applied to brain MR images, the structural similarity patterns may confuse the network. For example, a recent work [33] utilizes the patch-wise contrastive losses for training. However, patch-level information fails to consider global context, as the same anatomical structure belonging to symmetrical left and right brain hemispheres should have similar features, so if we expect the network to distinguish them, it will lead to false negative results. Besides, the same structure in different modalities may present different features, which possibly results in a large number of false negative pairs. Several other works have also deployed contrastive learning in medical images successfully [30], [34], however, the design of pretext tasks or the network architectures are decoupled from downstream applications. Inspired by SimCLR [31], we reuse two encoders in our model to seamlessly provide embedding features from each batch to conduct contrastive learning to guide both disentangled representation and the quality of generated images.

III. METHODS AND CONCEPT FORMULATION

A. Overview

Fig. 1 depicts our proposed method for multi-modality image translation via unsupervised disentangled and contrastive

representation, which improves the downstream segmentation model as it can leverage well-learned anatomical knowledge combined with imputed modalities which were originally missing.

The details of disentanglement learning are given in Section 3.B. To advance the disentangling operation and alleviate mode collapse during the training stage, we further propose two novel unsupervised contrastive learning strategies: one for attribute learning and the other for content learning as add-on inductive priors detailed in Section 3.C and 3.D, respectively. Together, they complement disentanglement learning as the attribute contrastive learning forces the modality-specific features to be separable in latent attribute space, whilst the content contrastive learning ensures subject-specific features are aligned in content space. The mathematical formulation of each loss function in the modality transfer module is given in Section 3.E. After successful feature disentanglement by the above-mentioned modality transfer module, we could utilize the learned content encoder as the encoder of our segmentation model, as stated in Section 3.F. As the building blocks are adopted from our previous MouseGAN [8], which uses a modified version of DRIT++ [5] as the translation module, the content and attribute encoders' architectures remain the same as [8] and then be reused seamlessly for contrastive losses.

B. Basic Disentanglement Building Blocks

Based on the consensus that unsupervised disentanglement is impossible unless appropriate inductive biases are imposed [18], we introduce our basic assumption here that anatomical semantic information should be included in content factors, whereas style factors only differ in image appearance. We believe this assumption should fit well with the fundamental nature of multi-modality MR images. Thus, in this section, we describe how we decouple multi-modality images into modality-specific attribute space and subject-specific (and modality-independent) content simultaneously, and then further introduce the derivation of contrastive learning priors in the subsequent sections.

Denoting k as the number of modalities in the dataset, $\{M_i\}_{i=1-k}$ be the descriptor for each modality, $(m_i, m_j) \in M$ be the two randomly-sampled images from two exemplary modalities, respectively, and their corresponding modality codes in one-hot format $(z_{m_i}^d, z_{m_j}^d)$, where $z^d \subset R^k$. As shown in Fig. 3, the model consists of modality-specific content encoders $\{E_{m_i}^c, E_{m_j}^c\}$, attribute encoders $\{E_{m_i}^a, E_{m_j}^a\}$, generators $\{G_{m_i}^c, G_{m_j}^c\}$, and modality discriminators $\{D_{m_i}, D_{m_j}\}$ as well as a content discriminator D^c . In the encoding path, for given input images (m_i, m_j) , we obtain the disentangled modality-independent content features $z_{m_i}^c = E_{m_i}^c(m_i)$, $z_{m_j}^c = E_{m_j}^c(m_j)$ in the shared content space $\{\Omega_i^c\}_{i=1-k}$ via the content encoders and simultaneously, the modality-specific attribute codes $z_{m_i}^a = E_{m_i}^a(m_i)$ and $z_{m_j}^a = E_{m_j}^a(m_j)$ in the latent attribute space $\{\Omega_i^a\}_{i=1-k}$ via the attribute encoders. The first forward image style translation cycle is achieved by a generator by combining the content feature from m_i as well as the swapped attribute feature from m_j , i.e., $m_{i \rightarrow j} = G_{m_j}(z_{m_i}^c, z_{m_j}^a, z_{m_j}^d)$. The disentanglement of

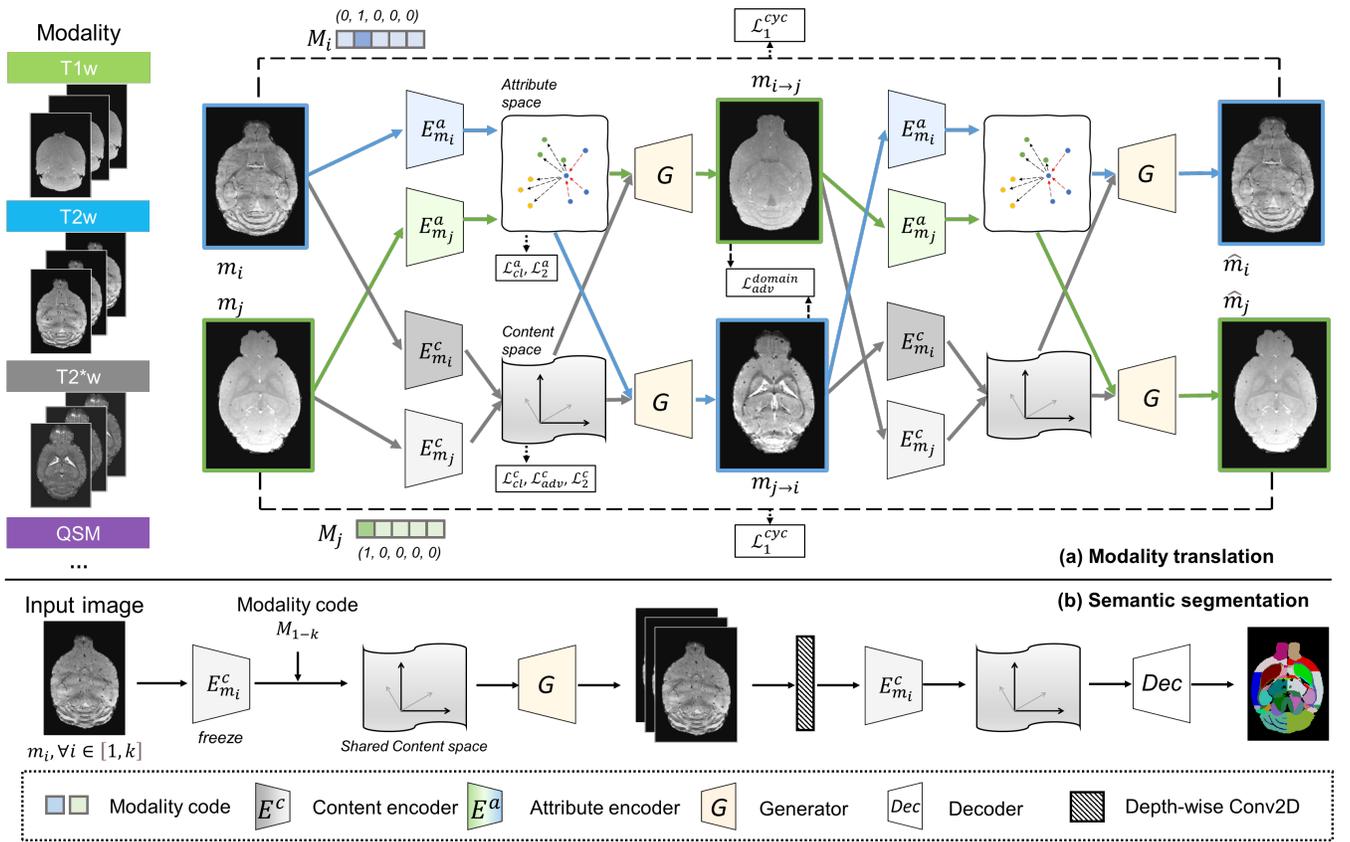


Fig. 1. Schematic of MouseGAN++. Given an input which can be any of the five modalities (T1w, T2w, T2*w, QSM, and Mag), (a) we first train a modality translation module based on feature disentanglement and contrastive learning to synthesize all modalities. (b) Then, we employ this modality translation module as an auxiliary network in the subsequent segmentation pipeline by reusing the content encoder parameters and imputed missing modalities. We demonstrate that the learning of structural features shared between different modalities in a self-supervised manner can better characterize brain structures, thereby leading to better segmentation.

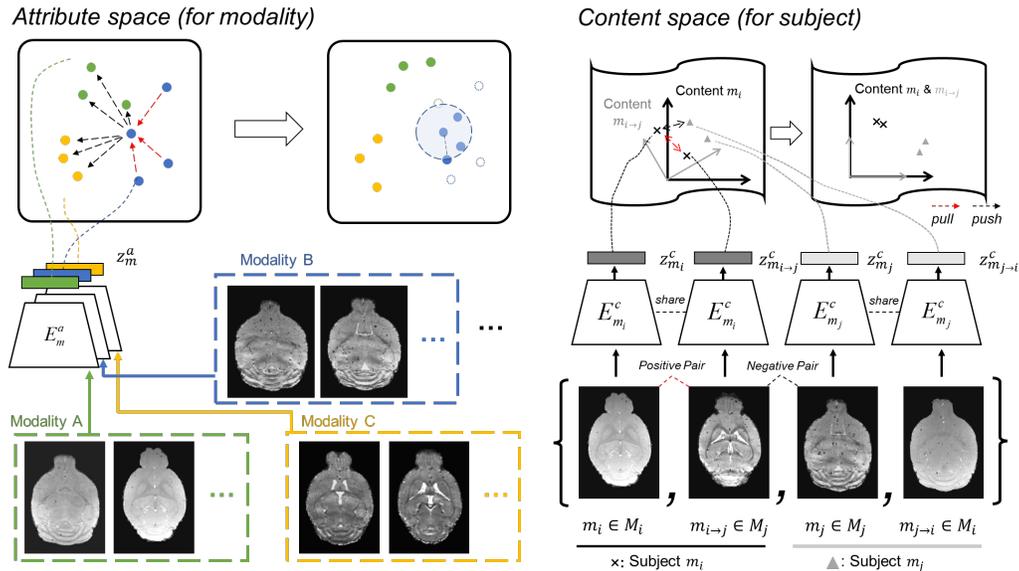


Fig. 2. Illustration of the attribute discrimination (left) and content discrimination (right) achieved by contrastive learning. In the attribute spaces, samples from the same modalities are defined as positive pairs and are driven close to each other, whilst samples from different modalities are negative pairs and are pushed away from each other. In the content space, samples from the same subject are defined as positive pairs and should have similar features in the embedding space, whilst samples from different subjects form negative pairs with distinct features.

modality-independent content features $z_{m_i}^c$ is desirable for us as it well conserves the structural properties which could be

crucial for medical image synthesis and downstream applications, albeit the generated image $m_{i \rightarrow j}$ have a similar high-

level style representation like m_j . Similarly, we can generate $m_{j \rightarrow i} = G_{m_i}(z_{m_j}^c, z_{m_i}^a, z_{m_i}^d)$ by exchanging its attribute and modality codes $z_{m_j}^a, z_{m_j}^d$ for $z_{m_i}^a, z_{m_i}^d$ while preserving its own content feature $z_{m_j}^c$.

C. Modality-Specific Attribute Discrimination

To promote learning a better disentangled feature representation, we introduce a task-specific prior here, as our desirable attribute embedding space should satisfy the following requirements:

- The attribute features of the same modality, even those from different subjects, are close to each other in latent attribute space.
- The attribute features of different modalities, even those from the same subject, are pushed away.

These requirements are in accordance with the key idea from contrastive learning theory [35], [36], which learns feature representation for positive pairs to be similar while pushing features from negative pairs apart. To this end, inspired by the work of [31], [37]–[39], we introduce contrast learning to guide the attribute encoder to achieve a more effective and discriminative modality-specific attribute embedding. Specifically, considering we have two input images m_i and m_j in one batch. Note that i could be equal to j here, when the two randomly sampled images are coincidentally from the same modality. After feeding them into the attribute encoder $E_{m_i}^a$ and $E_{m_j}^a$, we have embedding attribute features $z_{m_i}^a = E_{m_i}^a(m_i)$ and $z_{m_j}^a = E_{m_j}^a(m_j)$, as shown in Fig. 2. Unlike the most commonly used contrastive learning technology that employs instance discrimination, e.g., regarding each image as a separate individual, we extend the contrastive learning concept to attribute (modality) discrimination here. More specifically, the positive pairs contributed to $\Omega_a^+ = \{m_i, m_{j \rightarrow i} | \forall i, j\}$ are defined by the images that share the same modality code $z_{m_i}^d$, which could be either originated from or translated to the same modality, whereas the negative ones $\{m_j, m_{i \rightarrow j} | \forall i, j \text{ and } i \neq j\}$ from different modalities. We optimize the attribute space by minimizing the contrastive loss, which could be formulated as follows:

$$\mathcal{L}_{cl}^a = \sum_{i=1}^N \frac{-1}{|\Omega_a^+|} \sum_{m^+ \in \Omega_a^+} \log \frac{\exp(CL^+/\tau_a)}{\sum_{m_a \in \Omega_a(m_a)} \exp(CL/\tau_a)}, \quad (1)$$

where $CL^+ = \text{sim}(z_{m_i}^a, z_{m^+}^a)$, $CL = \text{sim}(z_{m_i}^a, z_{m_a}^a)$ and $\Omega_a(m_a) = \Omega_a \setminus \{m_i\}$. τ_a is the temperature scaling parameter. $\text{sim}(\cdot, \cdot)$ is pairwise similarity function which calculates the similarity of two vectors in the attribute space and is defined by cosine distance:

$$\text{sim}(x, y) = \frac{x \cdot y^\top}{\|x\| + \|y\|}. \quad (2)$$

D. Instance-Specific Content Discrimination

In medical image synthesis, one of the grimmest issues to contend with is how to retain the anatomical structure information fidelity during translation. Many recent image synthesis methods mainly focus on global visual similarity

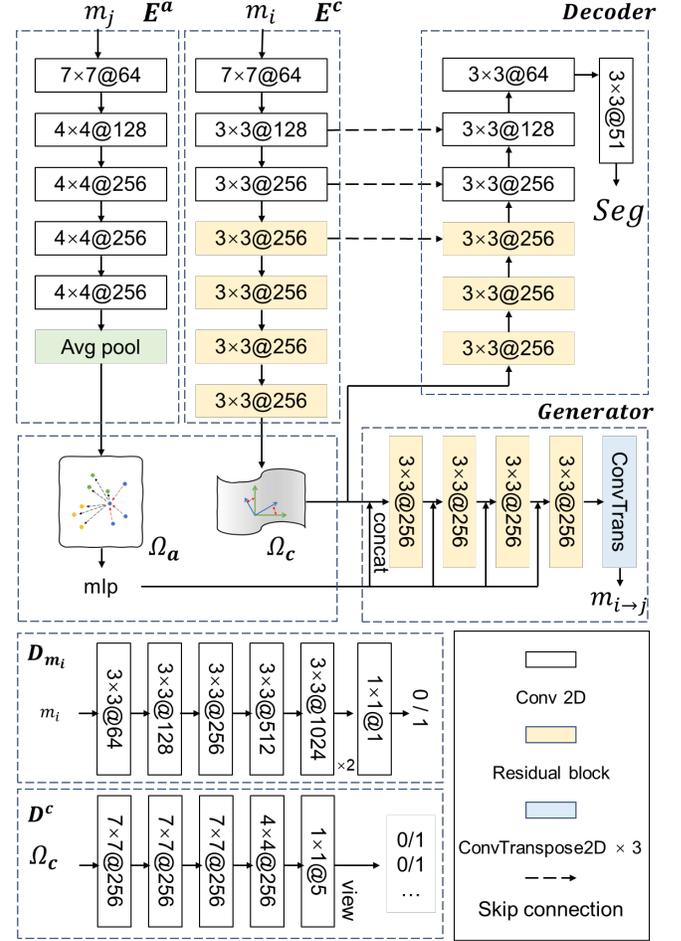


Fig. 3. The network architecture of attribute (E^a), content (E^c) encoder, decoder, generator and discriminator (D_{m_i}/D_{m_j} and D^c) in our framework.

or the rationality of perception. These may work well for natural images but are difficult to transfer to medical images in practice. Here, we propose an inductive prior that the content features $z_{m_i}^c$ are modality-agnostic and preserve only structural information, which should be subject-specific. Hence, we formulate two criteria to constrain the content space:

- For a specific subject, the content features of input image m_i and translated images $m_{i \rightarrow j}$ are close to each other in latent content space since they share the same intrinsic structures.
- For different subjects, the content features are pushed away, regardless of their image modalities.

Concretely, as shown in Fig. 2, given the images m_i^1 and m_i^2 from two subjects, the embedding content features $z_{m_i^1}^c = E^c(m_i^1)$ and $z_{m_i^2}^c = E^c(m_i^2)$ are generated via content encoders, whereas their translated content features are denoted by $z_{m_{i \rightarrow j}^1}^c$ and $z_{m_{i \rightarrow j}^2}^c$, respectively. Similarly as attribute discrimination in above section, we design the content discrimination pretext task here, which denotes m_i^1 and $m_{i \rightarrow j}^1$ as a positive pair $\in \Omega_c^+$ (same subject, different modalities) but m_i^1 and $m_{j \rightarrow i}^2$ as a negative pair (different subject, same

modalities). Mathematically, we have:

$$\mathcal{L}_{cl}^c = \sum_{i=1}^N \frac{-1}{|\Omega_c^+|} \sum_{m_1^+ \in \Omega_c^+} \log \frac{\exp(CL^+/\tau_c)}{\sum_{m_c^1 \in \Omega_c(m_c^1)} \exp(CL/\tau_c)}. \quad (3)$$

where $CL^+ = \text{sim}(z_{m_1^+}^c, z_{m_1^+}^c)$ and $CL = \text{sim}(z_{m_1^+}^c, z_{m_c^1}^c)$ and $\Omega_c(m_c^1) = \Omega_c \setminus \{m_1^1\}$. τ_c is the temperature scaling parameter. In our experiments, we set temperature of $\tau_c, \tau_a = 0.1$.

E. The Overall Loss Function of Modality Translation

As bidirectional reconstruction loss further encourages disentanglement by synergistically updating the gradient propagation among blocks, we reconstruct synthetic images $(m_{i \rightarrow j}, m_{j \rightarrow i})$ back via a second modality transfer step and enforce (\hat{m}_i, \hat{m}_j) to match the real input image (m_i, m_j) with cross-cycle consistency loss:

$$\mathcal{L}_1^{cyc}(E_m^c, E_m^a, G_m) = \|\hat{m}_i - m_i\|_1 + \|\hat{m}_j - m_j\|_1, \quad (4)$$

In addition, to compel the task of unsupervised translation during training, we use a self-reconstruction loss $\mathcal{L}_1^{self-recon}$ to constrain encoders and decoders for the generation of high quality translations, and can be formulated as follows:

$$\mathcal{L}_1^{self-recon}(E_m^c, E_m^a, G_m) = \|G_{m_i}(z_{m_i}^c, z_{m_i}^a) - m_i\|_1 + \|G_{m_j}(z_{m_j}^c, z_{m_j}^a) - m_j\|_1, \quad (5)$$

The latent space reconstruction loss \mathcal{L}_1^{latent} is adopted from [40]. The $\mathcal{L}_{cls}^{domain}$ is used to classify the modality of images [22]. To enforce the translated images to be as real as possible, following the successful practice of [2], [5], [8], the discriminators $\{D^c, D_{m_i}, D_{m_j}\}$ and corresponding losses $\{\mathcal{L}_{adv}^{cont}, \mathcal{L}_{adv}^{domain}\}$ are implemented for adversarial learning: the former discriminator enforces that one cannot infer image modality from content codes only, whilst the later one guarantees the synthesized images has the modality-specific features as the real ones. We formulate these losses as follows:

$$\begin{aligned} \min_{(E^c, G)} \max_{D^c} \mathcal{L}_{adv}^{cont} = & \mathbb{E}_{m_i} \left[\frac{1}{2} \log D^c(z_{m_i}^c) + \frac{1}{2} \log(1 - D^c(z_{m_i}^c)) \right] + \\ & \mathbb{E}_{m_j} \left[\frac{1}{2} \log D^c(z_{m_j}^c) + \frac{1}{2} \log(1 - D^c(z_{m_j}^c)) \right], \end{aligned} \quad (6)$$

$$\begin{aligned} \min_{(E^c, E^a, G)} \max_D \mathcal{L}_{adv}^{domain} = & \mathbb{E}_{m_i} [\log D_{m_i}(m_i)] + \mathbb{E}_{m_j \rightarrow i} [\log(1 - D_{m_i}(m_{j \rightarrow i}))] + \\ & \mathbb{E}_{m_j} [\log D_{m_j}(m_j)] + \mathbb{E}_{m_i \rightarrow j} [\log(1 - D_{m_j}(m_{i \rightarrow j}))], \end{aligned} \quad (7)$$

Besides, we impose L2 regularization \mathcal{L}_2^a to make latent space compact. The overall objective loss function of the translation module is:

$$\begin{aligned} \mathcal{L}_{trans} = & \mathcal{L}_{adv}^{domain} + \mathcal{L}_{adv}^{cont} + \lambda_1 \mathcal{L}_1^{self-recon} + \lambda_2 \mathcal{L}_1^{cyc} + \\ & \lambda_3 \mathcal{L}_1^{latent} + \lambda_4 \mathcal{L}_2^a + \lambda_5 \mathcal{L}_{cls}^{domain} + \lambda_6 \mathcal{L}_{cl}^a + \lambda_7 \mathcal{L}_{cl}^c. \end{aligned} \quad (8)$$

F. Semantic Segmentation Model

The pretrained modality translation model above serves as an auxiliary network for image segmentation by imputing missing modality in the input, since it is too expensive and time-consuming to collect all modality data for every mouse. After modality translation, an original image and its synthesized ones are fed to a depth-wise convolution before an encoder to preprocess the context modality-by-modality (Fig. 1b). We reuse the architecture and the parameters of the content encoder obtained in the modality translation training stage as the encoder of the segmentation module and update the decoder first. This is motivated by the fact that the content encoder in the modality translation model has already learned modality-independent anatomical features in a self-supervised manner and could extract and distill these representative features in a shared latent content space, leading to a better segmentation of anatomical structures for different modalities. Moreover, in contrast to a typical segmentation task, which is only trained by labelled data in a supervised fashion, our segmentation model also leverages information from unlabeled data via the unsupervised modality translation model. Our segmentation loss is defined as:

$$\mathcal{L}_{seg} = - \sum_{k=1}^K (y_i)_j^k \log(\text{Decoder}(E^c(m_i)))_j^k. \quad (9)$$

where y_i is the segmentation ground truth of m_i , $\text{Decoder}(\cdot)_j^k$ denote the probability prediction of voxel j for class k .

IV. EXPERIMENTAL SETUP

A. Dataset Acquisition

1) **Multi-Modality Dataset:** The multi-modality structure images of the mouse brain were generated by a 3D-mGRE sequence acquired with a 11.7T MR scanner, including T2w imaging, T1w imaging, T2*w imaging, quantitative susceptibility mapping (QSM), and Magnitude MR images (Mag). 3D-mGRE was acquired with the following parameters: $TR/TE/\Delta TE = 100/2/2$ ms, echo number = 12, flip angle = 15° , and resolution = $0.07 \times 0.07 \times 0.07$ mm³. The multi-modality structure images covered 75 mouse brains with 56,970 MR slices. All images were preprocessed by skull stripping and bias field correction. The ground truth of 50 brain structures was generated via an atlas-based method and then manually corrected by a brain anatomy expert.

2) **MRM NeAt Dataset:** To further verify our proposed method, we also tested it on external data from the MRM NeAt dataset [41], which includes 10 T2w MR images acquired with a 16.7T MR scanner. Each scan was manually annotated into 37 structures. All images had the resolution = $0.10 \times 0.10 \times 0.10$ mm³, covered 960 slices, and were preprocessed by skull stripping, denoising, and bias field correction. Since this dataset included mono-modality, we only executed segmentation tasks on this dataset for comparison.

B. Implementation Details

1) **Multi-Modality Dataset:** The architecture of MouseGAN++ is presented in Fig. 3. To conduct 5-fold cross-validation, we

randomly split 80% scans (60 subjects) as training sets and 20% scans (15 subjects) as testing sets at subject-level in each modality and further utilize unpaired data in experiments for both translation module and segmentation module training. In the preprocessing process, each slice was resized to 256×256 matrix size and the intensity distribution was normalized into zero mean and unit variance. Data augmentation was also introduced to alleviate overfitting, including random cropping to 216×216 matrix size and random flip. In addition, to train the translation module, we used one-hot modality codes to indicate each modality of input images in the training stage. In the test stage, the input single-modality images were encoded into shared content space and generated to all modalities with different one-hot modality codes. Empirically, we set $\lambda_1, \lambda_2, \lambda_3 = 10, \lambda_4 = 0.01, \lambda_5, \lambda_6, \lambda_7 = 1$.

For segmentation model training, we first froze the parameters of the shared content encoder and then trained only the decoder and depth-wise convolution filter. The images are resized to 224×160 . For all training procedures, we used the Adam optimizer with a learning rate of 0.0001, and set the size of the batch to 16. Finally, we finetuned the segmentation model with a learning rate of 0.00001 to update and refine all parameters, including the content encoder. The proposed framework was deployed in the Pytorch library and trained on a NVIDIA Tesla V100 GPU (32 GB memory) with 400 and 100 epochs for the translation and segmentation modules, taking 40 and 8 hours on average, respectively.

2) MRM NeAt Dataset: On this external dataset, we executed only the segmentation module of MouseGAN++ to investigate how the disentanglement representation pretrained on multiple modalities facilitates downstream segmentation tasks. We conducted the 5-fold cross validation study at the subject-level, using 8 subjects for training and the remaining 2 subjects for testing in each equal-size fold. The parameter setting and training steps were kept in line with the multi-modality dataset other than the number of output channels.

C. Evaluation Metrics

These metrics were used to evaluate the quality of translation images, Learned Perceptual Image Patch Similarity (LPIPS) [42], Visual Information Fidelity (VIF) [43], Peak Signal Noise Rate (PSNR), Structural Similarity Index Measure (SSIM), and Multi-Scale Structural Similarity (MS-SSIM). Lower LPIPS, higher VIF, PSNR, SSIM, and MS-SSIM refer to a better translation. The dice coefficient and average surface distance (ASD) were used for evaluating segmentation performance at subject-level in 3D volumes for both multi-modality dataset and MRM NeAt dataset. All the results were presented with mean \pm std to exhibit both the mean performance and cross-subject variance.

D. Comparison to SOTA Methods

We first compared the translation module of MouseGAN++ with the following state-of-the-art (SOTA) image translation methods: CycleGAN [2], UNIT [19], MUNIT [20], StarGAN-v2 [21] and our prior MouseGAN [8]. As for segmentation,

we further compared our method with two atlas-based methods, aMAP [11] and Natverse [12]; two direct segmentation methods: U-Net [44] and MU-Net [13]; one disentanglement method: D^2 -Net [6]; in addition to four synthesis-and-segmentation methods: CycleGAN [2], SynSeg [45], UNIT [19], MUNIT [20]. For SOTA synthesis-and-segmentation methods, we trained a segmentation network on both real images in the original modality M_i and the synthetic ones $M_{j \rightarrow i}$ in the transferred modalities, thus the network training can be benefited from the expansion of training data both in terms of size and modality. For SynSeg revised from [45], we conducted translation and segmentation simultaneously, not separately like in CycleGAN. As for MouseGAN and MouseGAN++ (the backbone, U-Net, is identical to that utilized in comparison methods), we utilized the translation module in our framework for imputing modalities, allowing the fused multi-modality semantic information to directly pass into the segmentation network. All comparison methods are either directly taken from original implementation from the corresponding github repository (if the codes are released) or implemented followed the original paper. The impact of each module in MouseGAN++ is detailed in the ablation study. Additionally, we directly deployed U-Net [44] as a baseline with only real images input for training. Since synthesis-and-segmentation methods require multiple modalities for image synthesis training, they fail to apply to the MRM NeAt dataset, which contains only one modality (T2w). Thus, we compared U-Net, the backbone of synthesis-and-segmentation methods, as well as other SOTA methods. Then, using several metrics, we evaluated the segmentation performance of these methods on several critical brain structures, including hippocampus, superior colliculus, striatum, and thalamus.

TABLE I
TRANSLATION PERFORMANCE COMPARISON WITH DIFFERENT METHODS FROM T1w TO T2w AND VICE VERSA

Method	T1w \rightarrow T2w				
	PSNR \uparrow	SSIM \uparrow	MS-SSIM \uparrow	VIF \uparrow	LPIPS \downarrow
CycleGAN	20.51 \pm 1.36	0.685 \pm 0.021	0.726 \pm 0.046	0.178 \pm 0.010	0.239 \pm 0.012
SynSeg	20.10 \pm 1.32	0.677 \pm 0.035	0.725 \pm 0.036	0.182 \pm 0.013	0.252 \pm 0.015
UNIT	19.98 \pm 1.32	0.579 \pm 0.075	0.727 \pm 0.039	0.145 \pm 0.018	0.240 \pm 0.012
MUNIT	20.88 \pm 1.21	0.602 \pm 0.085	0.751 \pm 0.039	0.160 \pm 0.017	0.202 \pm 0.018
StarGAN-v2	21.13 \pm 1.68	0.677 \pm 0.071	0.717 \pm 0.026	0.126 \pm 0.021	0.212 \pm 0.014
MouseGAN	20.46 \pm 1.13	0.693 \pm 0.036	0.775 \pm 0.033	0.197 \pm 0.015	0.186 \pm 0.014
MouseGAN++	22.72\pm1.02	0.716\pm0.031	0.796\pm0.029	0.215\pm0.017	0.173\pm0.010
Method	T2w \rightarrow T1w				
	PSNR \uparrow	SSIM \uparrow	MS-SSIM \uparrow	VIF \uparrow	LPIPS \downarrow
CycleGAN	25.61 \pm 1.38	0.813 \pm 0.018	0.842 \pm 0.032	0.256 \pm 0.013	0.248 \pm 0.021
SynSeg	23.39 \pm 1.54	0.789 \pm 0.024	0.809 \pm 0.028	0.241 \pm 0.016	0.246 \pm 0.018
UNIT	25.17 \pm 2.02	0.754 \pm 0.059	0.871 \pm 0.019	0.261 \pm 0.020	0.178 \pm 0.014
MUNIT	24.39 \pm 1.83	0.740 \pm 0.070	0.871 \pm 0.015	0.231 \pm 0.026	0.177 \pm 0.015
StarGAN-v2	23.05 \pm 1.47	0.810 \pm 0.028	0.836 \pm 0.044	0.179 \pm 0.050	0.188 \pm 0.017
MouseGAN	24.39 \pm 1.53	0.820 \pm 0.023	0.873 \pm 0.016	0.304 \pm 0.033	0.169 \pm 0.011
MouseGAN++	26.70\pm1.30	0.841\pm0.019	0.899\pm0.013	0.349\pm0.025	0.151\pm0.012

V. RESULTS AND DISCUSSION

A. Evaluation on Two-Modality Image Translation

As recent state-of-the-art methods mostly focus on one-to-one domain translation (i.e., one-to-one modality translation in our case), we first conducted T1w-T2w modality translation

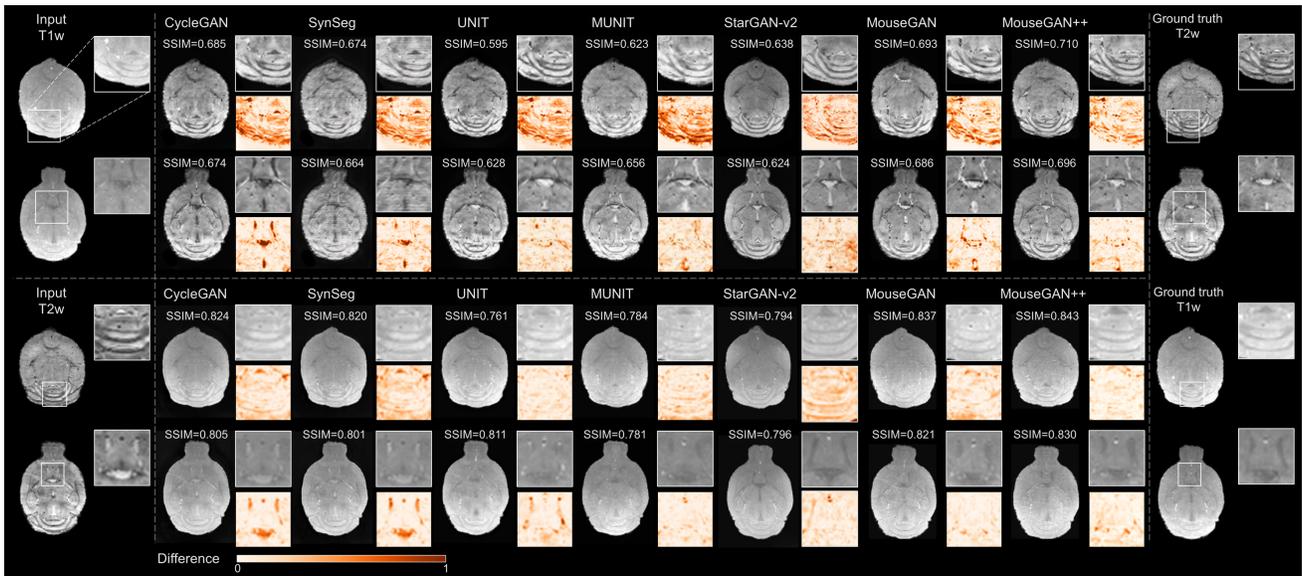


Fig. 4. Visualization of modality translation results between T1w and T2w. From left to right: raw input images (1st column), the synthetic images produced by different translation methods (2nd-7th column), and the ground truth target images (last column). For each real or synthetic image, we also show the zoomed-in regions of interest (ROI) in gray tones as well as the zoomed-in residual image between the synthetic image and the ground truth of the same ROI. The colorbar indicates the magnitude of the normalised residual.

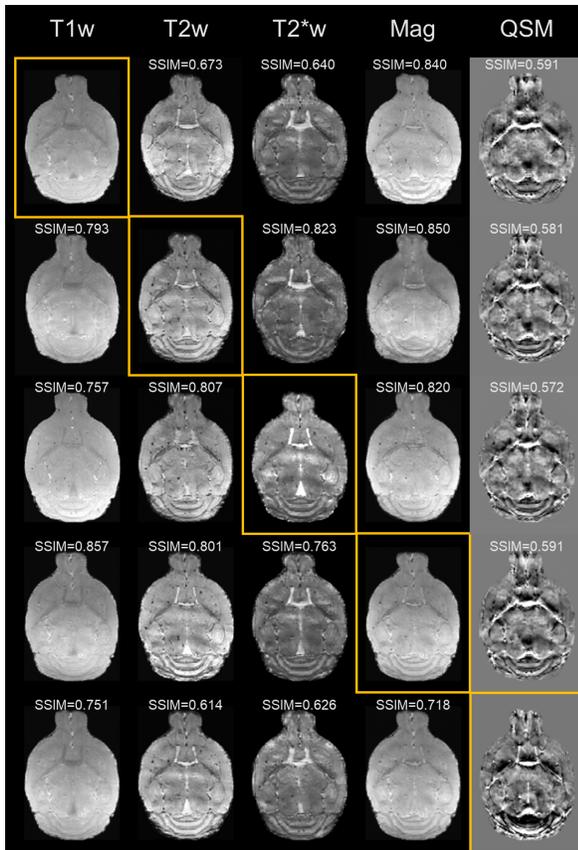


Fig. 5. Visualization of MouseGAN++ modality translation results on all five modalities. From top to bottom are the exemplary samples from the results when T1w, T2w, T2*w, Mag, or QSM is used as a single input modality, respectively. Yellow boxes represent input real images for each row, and could also serve as ground truth for each column. SSIMs between the synthetic and ground-truth real images are shown on the top of each image.

experiments in order to ensure comparability. The quantitative results of T1w-T2w translation are summarized in Table 1 and a few exemplary synthetic images are shown in Fig. 4. We can observe that for both two-direction translations, MouseGAN++ surpasses other methods in terms of all five metrics that are used to measure image synthesis quality. More specifically, for the translation from T1w to T2w, MouseGAN++ achieves the best performance with the highest SSIM of 0.716 and the lowest LPIPS of 0.173. We can also observe from Fig. 4 that MouseGAN++ produces sharper and more realistic textures, with the content better aligned with anatomy than other methods. For example, the lateral ventricles (zoomed-in ROIs from the second example) ought to have low contrast in T1w and, conversely, high contrast in T2w. Unlike MouseGAN++, synthetic images obtained from CycleGAN and SynSeg ignore these contrast features. While images synthesized by StarGAN-v2 show better contrast features, unexpected image deformation is present in the third ventricle and fourth ventricle (the second row, Fig.4). Moreover, MouseGAN++ yields significantly less pixel-level bias than others, especially in the cerebellar cortex. From our perspective, a fundamental reason for the better performance of MouseGAN++ is the attribute and content disentanglement, which enables more precise translation to different modalities. For the inverse translation task from T2w to T1w, MouseGAN++ achieves the best performance with an average SSIM of 0.841 and an average PSNR of 26.70, which is consistent with the T1w to T2w task. By contrast, the competing methods, CycleGAN, SynSeg, and StarGAN-v2, produce erroneous contrast or deformation in some brain regions during translation (zoomed-in ROIs in the last row).

B. Evaluation on Multi-Modality Image Translation

Here we further test MouseGAN++ on the five-modality translation task, where only images from a single modality were used as input to synthesize images of the other four modalities. Figure 5 reports the exemplary synthetic images along with the average SSIM for each test modality. Each row represents the test samples of the same subject from T1w, T2w, T2*w, Mag, and QSM, respectively.

We have two important observations from Fig. 5. First, some modalities are easier to synthesize than others: e.g., the averaged SSIMs of Mag and T1w being the target are 0.807 and 0.790, which are much higher than other modalities, e.g., QSM, with an averaged SSIM of 0.584. The suboptimal translation performance of QSM could be due to its immense different intensity distribution from other modalities, making this translation task particularly challenging. On the other hand, T2w, T2*w, and Mag can be good candidates for source modality, yielding average SSIMs of 0.762, 0.739, and 0.753 when being transferred into other modalities. By contrast, T1w and QSM generate lower SSIM (0.686 and 0.677, respectively). Particularly, we observe an asymmetric translation between T1w and T2w, the two most commonly used MRI modalities in practice: T2w \rightarrow T1w has a SSIM of 0.793, whilst T1w \rightarrow T2w lies only 0.673, which suggests that T2w could provide more structural information than T1w for mouse studies.

C. Evaluation on Brain Structure Segmentation

The quantitative results of brain structural segmentation on T1w and T2w, the two most important MR modalities used in mouse brain studies, are summarized in Table 2, with a few exemplary segmentation results presented in Fig. 6. In both modalities, MouseGAN++ consistently surpasses comparison methods, demonstrating the effectiveness of our proposed framework. Compared to the baseline U-Net segmentation with an average Dice of 73.6%, conventional translation methods (CycleGAN, SynSeg, UNIT, MUNIT) contribute to an improved Dice of around 76.3%-79.7%, which could be attributed to the synthetic images that increase both the amount and the diversity of training data. As one of the SOTA disentanglement methods, D^2 -Net also outperforms the baseline U-Net (yet inferior to MouseGAN and MouseGAN++), owing to the advantage of untangling during multi-modality training. In comparison, the two atlas-based segmentation methods, aMAP and Natverse, obtain average Dice scores of 82.9% and 75.5%, respectively. Their segmentation performance is largely determined by the quality of atlas registration. As we know, registration is not an easy task, as the heterogeneous contrasts existing in different image spaces might not provide enough guidance for intensity-based registration metrics, such as mutual information. Finally, MouseGAN++ achieves the best Dice scores (89.0%) and generates delicate segmentation results, e.g., in the regions of deep brain structures and cerebellum (Fig. 6), which demonstrates that the content encoder in the modality transfer model can successfully encode the structural information shared between different modalities, thus mitigating the influence of the domain gap and the

restriction of scarce information in a single modality. The additional incremented performance of MouseGAN++ as compared to MouseGAN could be due to a better disentanglement of the latent properties thanks to contrastive learning. Also, when making use of five modalities, MouseGAN++ gains a further 1% increment of Dice and 0.06 lower distance of ASD compared to the two-modality version, demonstrating different modalities do provide complementary information. Last but not least, the segmentation performance of T2w is slightly better than T1w images, which may be due to the fact that T1w images provide less structural information than T2w, consistent with our previous modality transfer analysis in Section 4.B.

D. Evaluation Over MRM NeAt Dataset

It is interesting to compare the performance of MouseGAN++ with SOTA methods on an additional dataset, MRM NeAt. As shown in Table 3 and Fig. 7, atlas-based methods produce unsatisfied results as the registration procedure causes brain structural deformation, especially for the elongated structures (e.g., corpus callosum). Note that these small regions are also challenging for deep learning methods, especially when the training set is limited in size. Encouragingly, on this challenging small dataset, MouseGAN++ achieves the best Dice and ASD, outperforming U-Net, MU-Net, D^2 -Net and the previous MouseGAN. Even without the translation module, the pretrained weight learnt from disentanglement representation in MouseGAN++ is superior to other networks due to the decoupled modality-agnostic knowledge, hence improving the performance of downstream tasks.

E. Ablation Study and Understanding of Contrastive Learning

In this section, we investigate why MouseGAN++ further improves the performance of our previous MouseGAN, i.e., how attribute and content contrastive learning improves modality translation. We first conduct an ablation experiment on the five modality translation tasks to show the effectiveness of each contrastive learning strategy. The baseline method is MouseGAN, which is deployed without attribute \mathcal{L}_{cl}^a and content \mathcal{L}_{cl}^c contrastive losses. To fully understand the advantages of our method, we exhibit the following experiment results and analyses from three aspects, the effectiveness of contrastive losses over mouse brain structural MR images, visualization of latent spaces, and comparison between Gaussian prior and contrastive learning prior.

1) *Effectiveness of Contrastive Losses*: As shown in Table 4, compared with the baseline, the attribute contrastive loss alone gains an average SSIM of 0.665 and 0.742 when taking T1w and T2w as input, respectively. A similar performance increase is observed when we add the content contrastive loss alone (0.671 to 0.751). Adding both components together makes MouseGAN++, which achieves 0.686 and 0.762 in terms of average SSIM. This incessant increase illustrates that attribute and content discrimination can be jointly implemented to produce better disentanglement representation in translation

TABLE II
COMPARISON OF SEGMENTATION PERFORMANCE OF DIFFERENT METHODS ON T1w AND T2w USING MULTI-MODALITY DATASET

Method	T2w							ASD [Voxel] ↓						
	Hippocampus		SC		Dice [%] ↑		P-value	Hippocampus		SC		ASD [Voxel] ↓		P-value
	Striatum	Thalamus	Mean	Hippocampus	SC	Striatum		Thalamus	Mean					
aMAP	84.4±2.0	82.5±3.2	83.0±1.9	81.7±2.8	82.9±2.7	1e-6	1.20±0.18	0.96±0.22	1.27±0.15	1.13±0.20	1.14±0.22	7e-9		
Natverse	78.2±4.4	75.3±3.7	75.5±2.8	73.1±4.6	75.5±4.1	6e-6	2.11±0.24	2.26±0.43	2.06±0.35	2.32±0.40	2.19±0.46	5e-8		
U-Net	77.9±7.0	70.6±6.1	76.1±5.5	69.9±6.2	73.6±6.6	2e-12	2.77±0.82	3.97±0.61	4.09±1.65	3.67±1.48	3.63±1.37	1e-12		
MU-Net	80.5±2.9	74.6±2.1	82.0±2.5	70.1±3.5	76.8±4.0	3e-8	2.14±0.43	3.14±0.51	2.64±0.47	2.14±0.59	2.52±0.48	2e-11		
D ² -Net	83.3±3.2	77.8±6.6	85.4±4.2	72.4±3.8	79.7±5.9	1e-11	1.82±0.24	2.09±0.31	1.96±0.29	2.32±0.40	2.05±0.36	5e-11		
CycleGAN	80.9±7.4	69.4±9.2	79.8±3.0	75.6±2.4	76.4±5.7	6e-11	2.28±0.39	3.53±0.56	2.68±0.62	2.47±0.43	2.74±0.54	1e-12		
SynSeg	82.5±6.4	74.3±5.8	80.2±4.2	76.2±4.6	78.3±5.5	4e-12	2.20±0.44	2.28±0.38	2.54±0.51	2.63±0.49	2.41±0.47	2e-10		
UNIT	84.2±4.4	75.6±5.8	81.8±5.3	74.6±5.6	79.1±5.8	3e-10	2.53±0.46	2.34±0.42	2.26±0.54	2.54±0.58	2.42±0.52	7e-9		
MUNIT	81.8±4.0	73.8±6.5	78.9±4.4	70.8±6.3	76.3±5.6	4e-6	2.20±0.57	2.43±0.48	2.17±0.64	2.30±0.68	2.28±0.61	2e-8		
MouseGAN (2-M)	87.2±3.6	84.6±4.7	87.6±3.3	82.3±3.1	85.4±4.2	1e-6	1.07±0.32	0.69±0.23	1.26±0.36	1.14±0.48	1.04±0.38	3e-8		
MouseGAN (5-M)	88.5±3.4	86.4±3.8	89.2±3.6	85.5±4.2	87.4±3.9	2e-7	0.94±0.28	0.66±0.12	0.91±0.22	1.04±0.50	0.89±0.33	1e-9		
MouseGAN++ (2-M)	91.4±1.2	88.7±1.2	92.0±1.9	83.9±4.0	89.0±2.2	1e-5	0.72±0.07	0.60±0.06	0.63±0.24	0.81±0.31	0.69±0.19	3e-7		
MouseGAN++ (5-M)	91.8±1.5	90.0±1.2	92.1±1.5	86.0±1.4	90.0±1.6	-	0.67±0.12	0.53±0.07	0.63±0.26	0.68±0.31	0.63±0.20	-		

Method	T1w							ASD [Voxel] ↓						
	Hippocampus		SC		Dice [%] ↑		P-value	Hippocampus		SC		ASD [Voxel] ↓		P-value
	Striatum	Thalamus	Mean	Hippocampus	SC	Striatum		Thalamus	Mean					
Amap	83.7±2.2	81.9±3.5	81.9±2.4	81.4±2.7	82.2±2.9	3e-6	1.25±0.18	1.00±0.24	1.39±0.19	1.15±0.20	1.20±0.25	2e-7		
Natverse	77.3±5.2	74.8±4.1	75.9±3.5	73.4±4.2	75.4±4.4	4e-7	2.24±0.26	2.45±0.48	1.92±0.36	2.26±0.39	2.22±0.40	1e-8		
U-Net	74.8±5.9	69.2±8.9	77.1±6.7	70.8±4.5	73.0±6.8	3e-10	3.35±0.78	4.15±0.79	3.93±0.62	2.87±1.09	3.58±0.93	4e-9		
MU-Net	78.1±3.7	74.0±4.8	79.2±4.5	74.9±4.2	76.6±4.6	1e-11	2.05±0.33	1.85±0.44	2.21±0.59	3.07±0.52	2.29±0.54	7e-9		
D ² -Net	81.7±2.1	76.5±4.1	83.5±4.6	75.2±2.7	79.2±3.6	2e-9	2.21±0.42	2.07±0.35	1.96±0.27	1.63±0.26	1.97±0.38	2e-10		
CycleGAN	79.5±4.6	68.1±3.8	78.8±4.0	74.3±3.4	75.2±4.2	6e-10	2.43±0.37	3.07±0.63	2.27±0.48	2.59±0.42	2.59±0.49	4e-10		
SynSeg	80.6±4.5	75.2±3.9	76.9±3.0	75.9±3.8	77.2±4.0	4e-10	2.68±0.64	2.87±0.54	2.79±0.54	2.48±0.43	2.71±0.58	3e-11		
UNIT	81.1±3.5	76.2±3.6	78.8±3.3	75.5±3.9	77.9±3.7	3e-9	2.32±0.38	2.18±0.56	2.34±0.40	2.40±0.52	2.31±0.51	5e-9		
MUNIT	78.2±2.6	74.0±3.0	76.1±2.8	73.6±3.1	75.5±3.1	1e-11	2.39±0.69	2.07±0.22	2.16±0.30	2.25±0.45	2.22±0.43	1e-12		
MouseGAN (2-M)	86.3±3.6	83.5±3.5	84.2±3.7	81.6±4.0	83.9±3.9	3e-10	1.31±0.36	1.24±0.27	1.46±0.37	1.29±0.34	1.78±0.36	2e-9		
MouseGAN (5-M)	87.4±2.4	84.3±3.0	89.0±2.6	83.5±3.2	86.1±3.9	1e-9	1.16±0.27	0.85±0.23	1.28±0.30	0.75±0.31	1.21±0.29	3e-9		
MouseGAN++ (2-M)	88.8±2.4	85.2±1.8	90.9±1.7	83.8±2.1	87.2±2.1	4e-8	0.73±0.23	0.79±0.32	0.68±0.27	0.53±0.27	0.68±0.28	5e-8		
MouseGAN++ (5-M)	89.7±2.0	86.3±1.6	91.4±1.5	84.1±1.8	87.9±1.9	-	0.71±0.17	0.72±0.13	0.62±0.25	0.56±0.20	0.65±0.21	-		

SC: Superior colliculus; 2-M: two-modality dataset; 5-M: five-modality dataset; P-value: paired t-test.

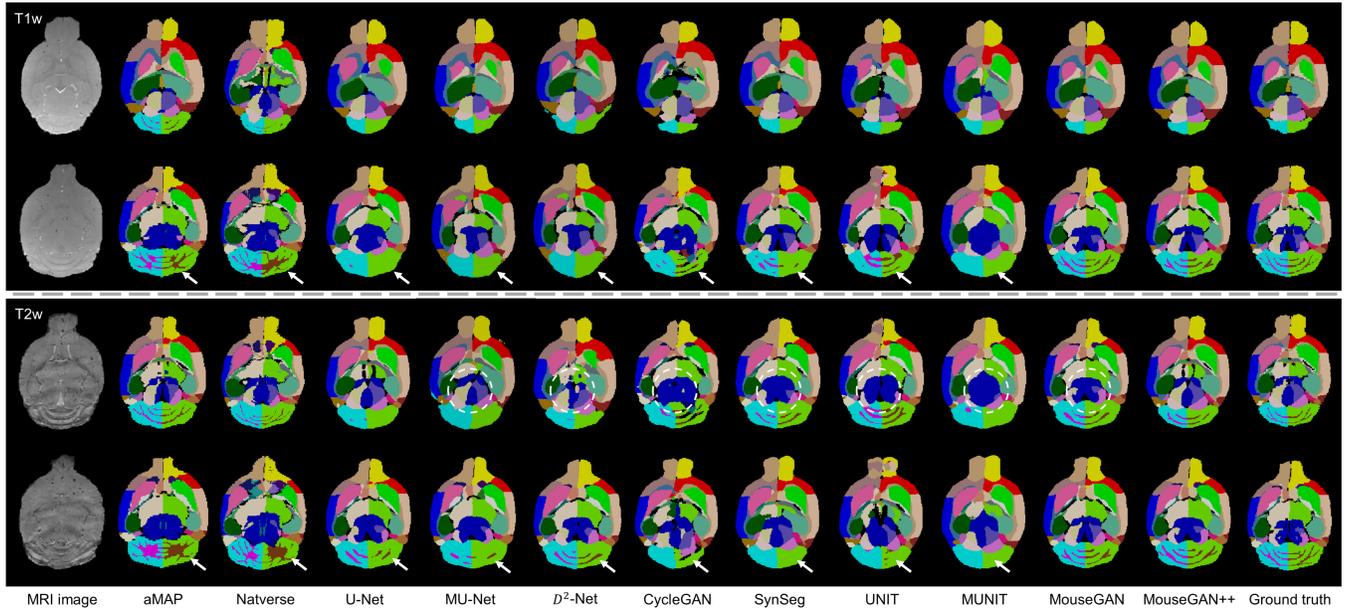


Fig. 6. Exemplary segmentation results on T1w and T2w. From left to right: images to be segmented (1st column), segmentation results of different comparative methods (2nd-11th column), results of our MouseGAN++ (12th column), and the ground truth segmentation (the last column). Our model can segment up to 50 mouse brain structures shown in different colors. White circles and arrows indicate erroneous segmentation details.

tasks. In Fig. 8, we present the training loss curve of self-modality reconstruction $\mathcal{L}_1^{self-recon}$ as well as \mathcal{L}_1^{cyc} , which is associated with the quality of translation. The MouseGAN++ with contrastive losses shows a more stable training procedure than the original MouseGAN as the adversarial losses achieve the balance between the discriminator and generator in the multi-modality dataset.

One corner case we intend to clarify here is how the network knows to deal with two very similar mouse brains in one mini-batch. Even though this case can be deemed as a hard example for network discrimination, the anatomical contents from the same subject should be located closest to each other in the latent content space despite the presence of modality translation. Additionally, some works report that attribute

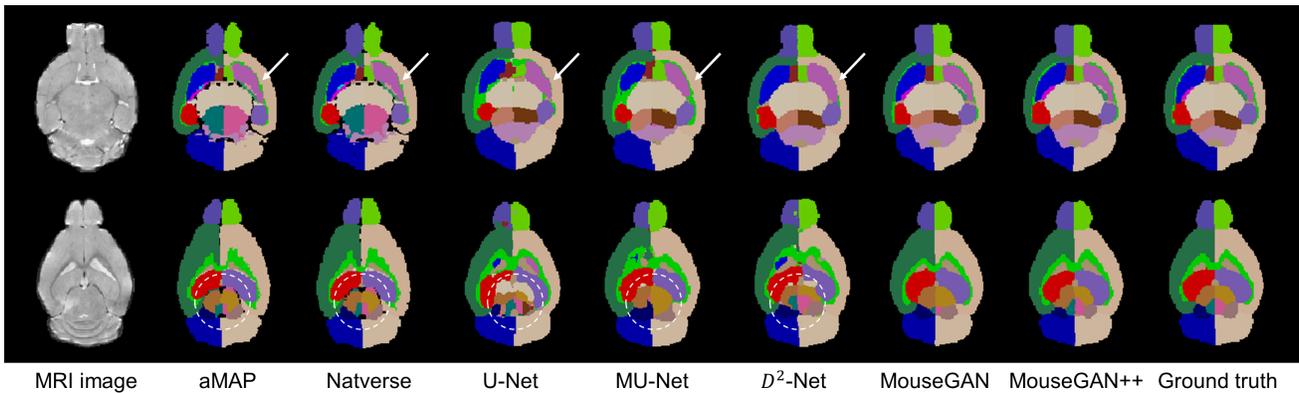


Fig. 7. Exemplary segmentation results on MRM NeAT T2w dataset. From left to right: images to be segmented(1st column), segmentation results of different comparative methods (2rd-7th column), results of our MouseGAN++ (8th column), and the ground truth segmentation (the last column). White circles and arrows indicate erroneous segmentation details.

TABLE III
COMPARISON OF SEGMENTATION PERFORMANCE OF DIFFERENT METHODS ON T2W USING MRM NEAT DATASET

Method	Dice [%] \uparrow						ASD [Voxel] \downarrow					
	Hippocampus	SC	Striatum	Thalamus	Mean	P-value	Hippocampus	SC	Striatum	Thalamus	Mean	P-value
aMAP	83.9 \pm 1.9	71.8 \pm 4.1	84.2 \pm 1.6	67.9 \pm 2.7	77.0 \pm 2.8	9e-9	1.11 \pm 0.13	1.60 \pm 0.55	0.99 \pm 0.07	2.94 \pm 0.17	1.66 \pm 0.27	3e-10
Natverse	81.3 \pm 3.4	69.4 \pm 7.2	84.4 \pm 2.3	66.5 \pm 1.2	75.4 \pm 3.9	1e-8	1.30 \pm 0.22	1.76 \pm 0.46	1.00 \pm 0.12	3.28 \pm 0.27	1.84 \pm 0.32	2e-10
U-Net	82.3 \pm 5.6	78.5 \pm 4.7	81.8 \pm 4.1	81.4 \pm 5.3	81.0 \pm 5.2	2e-6	1.82 \pm 0.75	3.24 \pm 0.54	2.64 \pm 0.69	2.37 \pm 0.68	2.52 \pm 0.72	7e-8
MU-Net	85.9 \pm 3.8	81.5 \pm 4.0	83.9 \pm 3.4	84.2 \pm 4.7	83.9 \pm 4.5	8e-5	1.34 \pm 0.35	2.05 \pm 0.48	2.18 \pm 0.59	2.19 \pm 0.47	1.94 \pm 0.51	5e-6
D^2 -Net	87.3 \pm 2.9	83.6 \pm 2.1	87.0 \pm 2.5	87.2 \pm 3.5	86.3 \pm 4.0	4e-4	1.04 \pm 0.44	1.49 \pm 0.51	1.26 \pm 0.47	2.11 \pm 0.59	1.48 \pm 0.59	1e-5
MouseGAN	89.7 \pm 1.8	85.3 \pm 2.0	90.2 \pm 2.2	91.6 \pm 2.8	89.2 \pm 4.0	6e-3	0.84 \pm 0.33	0.94 \pm 0.31	0.74 \pm 0.37	0.98 \pm 0.35	0.82 \pm 0.32	2e-4
MouseGAN++	90.6\pm1.5	87.1\pm2.5	91.7\pm1.0	92.8\pm0.7	90.6\pm2.6	-	0.64\pm0.11	0.66\pm0.08	0.57\pm0.06	0.60\pm0.06	0.62\pm0.09	-

SC: Superior colliculus; P-value: paired t-test.

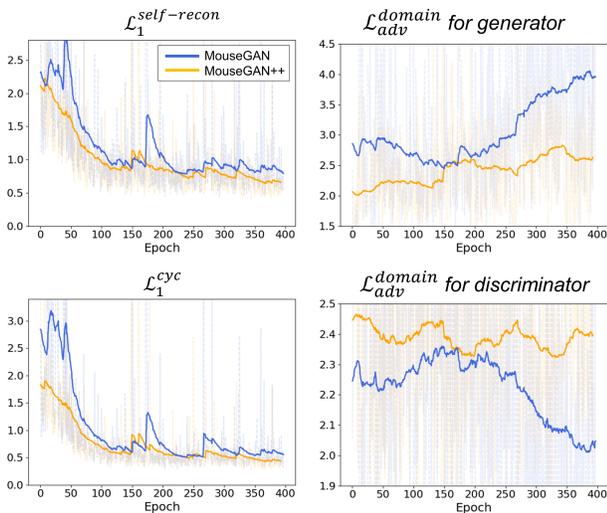


Fig. 8. Comparison of translation training loss curves in ablation study between previous MouseGAN (blue) and MouseGAN++ (orange).

features may be further divided into domain-specific style and domain-invariant style [46], [47]. In MouseGAN++, the modality-related (domain-specific) features are preserved in attribute space, while the modality-invariant (domain-invariant, or in other words, structure-related) information is condensed in the content space, which is adequate for MR image disentanglement.

Furthermore, we qualitatively analyze the contribution of each individual module of MouseGAN++ with regard to

TABLE IV
THE ABLATION STUDY FOR TRANSLATION RESULTS TAKING T2W OR T1W AS A SINGLE INPUT, RESPECTIVELY. THE AVERAGE METRICS ARE BASED ON THE TEST MODALITIES

Method	T2w				
	PSNR \uparrow	SSIM \uparrow	MS-SSIM \uparrow	VIF \uparrow	LPIPS \downarrow
MouseGAN	19.46 \pm 3.69	0.737 \pm 0.124	0.752 \pm 0.134	0.204 \pm 0.046	0.222 \pm 0.028
MouseGAN + L_{ct}^c	19.76 \pm 3.71	0.742 \pm 0.126	0.761 \pm 0.133	0.180 \pm 0.052	0.225 \pm 0.025
MouseGAN + L_{ct}^c	20.82 \pm 3.37	0.751 \pm 0.132	0.781 \pm 0.127	0.246 \pm 0.048	0.193 \pm 0.024
MouseGAN++	21.28\pm3.46	0.762\pm0.110	0.791\pm0.126	0.281\pm0.046	0.187\pm0.024
Method	T1w				
	PSNR \uparrow	SSIM \uparrow	MS-SSIM \uparrow	VIF \uparrow	LPIPS \downarrow
MouseGAN	18.16 \pm 3.77	0.645 \pm 0.104	0.683 \pm 0.119	0.167 \pm 0.089	0.250 \pm 0.028
MouseGAN + L_{ct}^c	19.01 \pm 3.37	0.665 \pm 0.111	0.692 \pm 0.121	0.188 \pm 0.076	0.246 \pm 0.025
MouseGAN + L_{ct}^c	19.68 \pm 3.42	0.671 \pm 0.110	0.721 \pm 0.114	0.191 \pm 0.066	0.214 \pm 0.024
MouseGAN++	20.42\pm3.15	0.686\pm0.100	0.722\pm0.114	0.227\pm0.087	0.205\pm0.024

segmentation. We first set up the backbone of the network as a baseline model without shared content encoder weight and translation module (baseline U-Net). Then we add each module one by one into the baseline model. We evaluate different methods on five modalities with every single modality as input successively and compute the average metrics. As shown in Table 5, both the shared content encoder and the translation module bring significant performance advancement, indicating that each module can work collaboratively in our proposed framework.

2) *Visualization of Latent Spaces*: To illustrate how contrastive learning improves attribute and content feature embedding themselves, we visualize attribute features of five modalities via t-SNE [48]. As shown in Fig. 9a, the attribute features of our proposed methods are well separated in at-

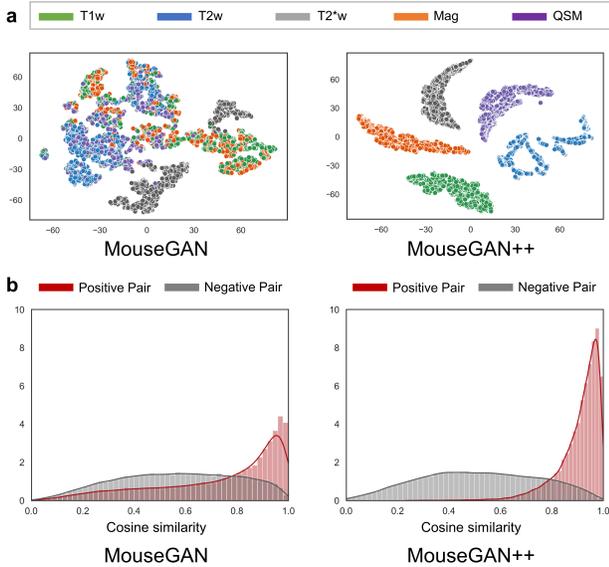


Fig. 9. Visualization of latent attribute and content embedding space without (MouseGAN) and with contrastive learning losses (MouseGAN++). (a) Scatter plot of attribute feature clusters using t-SNE. (b) The cosine similarity distributions of content features.

TABLE V

THE ABLATION STUDY FOR EXEMPLARY BRAIN STRUCTURE SEGMENTATION EVALUATED BY DICE SCORES ON THE 5-MODALITY DATASET

Method	Hippo-campus	Superior colliculus	Striatum	Thalamus	Mean
Backbone (U-Net)	75.5±6.0	69.4±5.4	74.0±4.7	70.7±6.6	72.4±6.2
Backbone + shared E^c	79.0±3.7	72.6±4.2	79.3±3.3	75.4±3.6	76.6±3.8
Backbone + translation module	86.7±2.1	83.1±3.2	85.9±2.6	81.3±3.0	84.3±2.9
Backbone + shared E^c + translation module = MouseGAN++	90.4±1.6	87.2±2.0	88.9±1.8	83.6±2.2	87.5±2.0

tribute space, while the baseline method confuses all modality features except T2*w modality, which further demonstrates the capacity of MouseGAN++ to achieve better attribute discrimination.

Then we further evaluate the instance-specific content discrimination. We calculate the cosine similarity between content features from the input modality and the modality-translated images of the same subject (positive pairs), as well as features from different subjects (negative pairs). As shown in Fig. 9b, a better feature embedding is denoted by a more separable distribution between positive and negative samples, which demonstrates that MouseGAN++ is able to learn more instance-specific and less modality-relevant features, hence conducting better semantic information retention during modality translation.

3) *Gaussian Prior vs. Contrastive Learning Prior*: By contrast with existing approaches [27]–[30] that assume the Gaussian approximation over the latent attribute space as the primary bias or prior to promote decoupling, we claim that the contrastive learning prior is more appropriate for approaching an ideal latent distribution of a multi-modality dataset. Unlike MouseGAN, which follows Gaussian priors, MouseGAN++ follows contrastive learning priors. The comprehensive results

of the above-mentioned comparisons, ablation study, and visualization of latent space demonstrate that the MouseGAN++ achieves better performance as well as translation quality compared with SOTA methods.

Nowadays, in practice, multimodal images are commonly unpaired, which can be a major obstacle to deploying such paired training methods [6]. Our efforts also extend this promising direction via fusing cross-modality information in an unpaired manner and yielding more robust performance when the modality is missing. Another point we want to emphasize is the maximum flexibility of our framework as we trained separately rather than jointly. This is reflected in the flexibility in choosing pretext tasks, which is important for disentanglement and downstream tasks. In order to impute a missing modality, we chose modality translation as our pretext task in this work. However, one limitation of our work is the dilemma caused by domain shifts, since the image quality and contrast from various centers may differ greatly. An appealing and promising solution is to convert our pretext task to cross-center image translation so that the learned center-agnostic features in the content space would alleviate the segmentation performance degradation.

VI. CONCLUSION

In summary, we propose a novel GAN-based framework, MouseGAN++, for simultaneous image synthesis and segmentation for mouse brain MRI. Based on a disentangled representation of content and style attributes strengthened by contrastive learning, MouseGAN++ is able to synthesize multiple MR modalities from single ones in a structure-preserving manner and can hence handle cases with missing modalities. Furthermore, it uses the learned modality-invariant information to improve structural segmentation. Our results demonstrate that MouseGAN++ achieves significant performance improvement over both translation and segmentation tasks and has the potential to be applied in more neuroimaging applications.

The broader implications of our work are the publicly available packaged pipeline provided by MouseGAN++ and the multiple modality MR dataset for facilitating preclinical research, especially for the communities interested in rodent brain. In addition, the unsupervised pretext task alleviates the cost of deployment and promotes the clinical use of cutting-edge machine learning techniques. In the future work, one of the attractive directions is to integrate our previous brain extraction tool [49] with MouseGAN++ as an end-to-end neuroimaging processing pipeline.

REFERENCES

- [1] T. Zhou, S. Ruan, and S. Canu, “A review: Deep learning for medical image segmentation using multi-modality fusion,” *Array*, vol. 3, p. 100004, 2019.
- [2] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, “Unpaired image-to-image translation using cycle-consistent adversarial networks,” in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2223–2232.
- [3] Z. Yi, H. Zhang, P. Tan, and M. Gong, “Dualgan: Unsupervised dual learning for image-to-image translation,” in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2849–2857.

- [4] T. Kim, M. Cha, H. Kim, J. K. Lee, and J. Kim, "Learning to discover cross-domain relations with generative adversarial networks," in *International conference on machine learning*. PMLR, 2017, pp. 1857–1865.
- [5] H.-Y. Lee, H.-Y. Tseng, Q. Mao, J.-B. Huang, Y.-D. Lu, M. Singh, and M.-H. Yang, "Drit++: Diverse image-to-image translation via disentangled representations," *International Journal of Computer Vision*, vol. 128, no. 10, pp. 2402–2417, 2020.
- [6] Q. Yang, X. Guo, Z. Chen, P. Y. Woo, and Y. Yuan, "D2-net: Dual disentanglement network for brain tumor segmentation with missing modalities," *IEEE Transactions on Medical Imaging*, 2022.
- [7] C. Chen, Q. Dou, Y. Jin, H. Chen, J. Qin, and P.-A. Heng, "Robust multimodal brain tumor segmentation via feature disentanglement and gated fusion," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2019, pp. 447–456.
- [8] Z. Yu, Y. Zhai, X. Han, T. Peng, and X.-Y. Zhang, "Mousegan: Gan-based multiple mri modalities synthesis and segmentation for mouse brain structures," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2021, pp. 442–450.
- [9] M. Cabezas, A. Oliver, X. Lladó, J. Freixenet, and M. B. Cuadra, "A review of atlas-based segmentation for magnetic resonance brain images," *Computer methods and programs in biomedicine*, vol. 104, no. 3, pp. e158–e177, 2011.
- [10] J. Bai, T. L. H. Trinh, K.-H. Chuang, and A. Qiu, "Atlas-based automatic mouse brain image segmentation revisited: model complexity vs. image registration," *Magnetic resonance imaging*, vol. 30, no. 6, pp. 789–798, 2012.
- [11] C. J. Niedworok, A. P. Brown, M. Jorge Cardoso, P. Osten, S. Ourselin, M. Modat, and T. W. Margrie, "amap is a validated pipeline for registration and segmentation of high-resolution mouse brain data," *Nature communications*, vol. 7, no. 1, pp. 1–9, 2016.
- [12] A. S. Bates, J. D. Manton, S. R. Jagannathan, M. Costa, P. Schlegel, T. Rohlfing, and G. S. Jefferis, "The natverse, a versatile toolbox for combining and analysing neuroanatomical data," *Elife*, vol. 9, 2020.
- [13] R. De Feo, A. Shatillo, A. Sierra, J. M. Valverde, O. Gröhn, F. Giove, and J. Tohka, "Automated joint skull-stripping and segmentation with multi-task u-net in large mouse brain mri databases," *NeuroImage*, vol. 229, p. 117734, 2021.
- [14] A. Chartsias, T. Joyce, G. Papanastasiou, S. Semple, M. Williams, D. E. Newby, R. Dharmakumar, and S. A. Tsafaris, "Disentangled representation learning in cardiac image analysis," *Medical image analysis*, vol. 58, p. 101535, 2019.
- [15] S. Reangamornrat, H. Sari, C. Catana, and A. Kamen, "Multi-modal image synthesis based on disentanglement representations of anatomical and modality specific features, learned using uncooperative relativistic gan," *Medical Image Analysis*, vol. 80, p. 102514, 2022.
- [16] Y. Wang, L. Zhou, B. Yu, L. Wang, C. Zu, D. S. Lalush, W. Lin, X. Wu, J. Zhou, and D. Shen, "3d auto-context-based locality adaptive multi-modality gans for pet synthesis," *IEEE transactions on medical imaging*, vol. 38, no. 6, pp. 1328–1339, 2018.
- [17] B. Zhan, L. Zhou, Z. Li, X. Wu, Y. Pu, J. Zhou, Y. Wang, and D. Shen, "D2fe-gan: Decoupled dual feature extraction based gan for mri image synthesis," *Knowledge-Based Systems*, vol. 252, p. 109362, 2022.
- [18] F. Locatello, S. Bauer, M. Lucic, G. Raetsch, S. Gelly, B. Schölkopf, and O. Bachem, "Challenging common assumptions in the unsupervised learning of disentangled representations," in *international conference on machine learning*. PMLR, 2019, pp. 4114–4124.
- [19] M.-Y. Liu, T. Breuel, and J. Kautz, "Unsupervised image-to-image translation networks," *Advances in neural information processing systems*, vol. 30, 2017.
- [20] X. Huang, M.-Y. Liu, S. Belongie, and J. Kautz, "Multimodal unsupervised image-to-image translation," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 172–189.
- [21] Y. Choi, Y. Uh, J. Yoo, and J.-W. Ha, "Stargan v2: Diverse image synthesis for multiple domains," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 8188–8197.
- [22] H.-Y. Lee, H.-Y. Tseng, J.-B. Huang, M. Singh, and M.-H. Yang, "Diverse image-to-image translation via disentangled representations," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 35–51.
- [23] C. Chen, Q. Dou, H. Chen, J. Qin, and P. A. Heng, "Unsupervised bidirectional cross-modality adaptation via deeply synergistic image and feature alignment for medical image segmentation," *IEEE transactions on medical imaging*, vol. 39, no. 7, pp. 2494–2505, 2020.
- [24] K. Yao, J. Sun, K. Huang, L. Jing, H. Liu, D. Huang, and C. Jude, "Analyzing cell-scaffold interaction through unsupervised 3d nuclei segmentation," *International journal of bioprinting*, vol. 8, no. 1, 2022.
- [25] A. Chartsias, T. Joyce, G. Papanastasiou, S. Semple, M. Williams, D. Newby, R. Dharmakumar, and S. A. Tsafaris, "Factorised spatial representation learning: Application in semi-supervised myocardial segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2018, pp. 490–498.
- [26] A. Chartsias, G. Papanastasiou, C. Wang, C. Stirrat, S. Semple, D. Newby, R. Dharmakumar, and S. A. Tsafaris, "Multimodal cardiac segmentation using disentangled representation learning," in *International Workshop on Statistical Atlases and Computational Models of the Heart*. Springer, 2019, pp. 128–137.
- [27] S. J. Wagner, N. Khalili, R. Sharma, M. Boxberg, C. Marr, W. d. Back, and T. Peng, "Structure-preserving multi-domain stain color augmentation using style-transfer with disentangled representations," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2021, pp. 257–266.
- [28] M. Havaei, X. Mao, Y. Wang, and Q. Lao, "Conditional generation of medical images via disentangled adversarial inference," *Medical Image Analysis*, vol. 72, p. 102106, 2021.
- [29] S. Wang and L. Rui, "Sgdr: Semantic-guided disentangled representation for unsupervised cross-modality medical image segmentation," *arXiv preprint arXiv:2203.14025*, 2022.
- [30] C. I. Bercea, B. Wiestler, D. Rueckert, and S. Albarqouni, "Feddis: Disentangled federated learning for unsupervised brain pathology segmentation," *arXiv preprint arXiv:2103.03705*, 2021.
- [31] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," in *International conference on machine learning*. PMLR, 2020, pp. 1597–1607.
- [32] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 9729–9738.
- [33] Z. Huang, L. Lin, P. Cheng, L. Peng, and X. Tang, "Multi-modal brain tumor segmentation via missing modality synthesis and modality-level attention fusion," *arXiv preprint arXiv:2203.04586*, 2022.
- [34] D. Zeng, Y. Wu, X. Hu, X. Xu, H. Yuan, M. Huang, J. Zhuang, J. Hu, and Y. Shi, "Positional contrastive learning for volumetric medical image segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2021, pp. 221–230.
- [35] Y. Tian, C. Sun, B. Poole, D. Krishnan, C. Schmid, and P. Isola, "What makes for good views for contrastive learning?" *Advances in Neural Information Processing Systems*, vol. 33, pp. 6827–6839, 2020.
- [36] P. Khosla, P. Teterwak, C. Wang, A. Sarna, Y. Tian, P. Isola, A. Maschinot, C. Liu, and D. Krishnan, "Supervised contrastive learning," *Advances in Neural Information Processing Systems*, vol. 33, pp. 18 661–18 673, 2020.
- [37] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 9729–9738.
- [38] M. Ye, X. Zhang, P. C. Yuen, and S.-F. Chang, "Unsupervised embedding learning via invariant and spreading instance feature," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6210–6219.
- [39] P. Khosla, P. Teterwak, C. Wang, A. Sarna, Y. Tian, P. Isola, A. Maschinot, C. Liu, and D. Krishnan, "Supervised contrastive learning," *Advances in Neural Information Processing Systems*, vol. 33, pp. 18 661–18 673, 2020.
- [40] J.-Y. Zhu, X. Zhang, D. Pathak, T. Darrell, A. A. Efros, O. Wang, and E. Shechtman, "Toward multimodal image-to-image translation," *Advances in neural information processing systems*, vol. 30, 2017.
- [41] Y. Ma, D. Smith, P. R. Hof, B. Foerster, S. Hamilton, S. J. Blackband, M. Yu, and H. Benveniste, "In vivo 3d digital atlas database of the adult c57bl/6j mouse brain by magnetic resonance microscopy," *Frontiers in neuroanatomy*, p. 1, 2008.
- [42] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 586–595.
- [43] H. R. Sheikh and A. C. Bovik, "Image information and visual quality," *IEEE Transactions on image processing*, vol. 15, no. 2, pp. 430–444, 2006.
- [44] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Confer-*

ence on *Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.

- [45] Y. Huo, Z. Xu, H. Moon, S. Bao, A. Assad, T. K. Moyo, M. R. Savona, R. G. Abramson, and B. A. Landman, “Synseg-net: Synthetic segmentation without target modality ground truth,” *IEEE transactions on medical imaging*, vol. 38, no. 4, pp. 1016–1025, 2018.
- [46] X. Liu, S. Thermos, A. O’Neil, and S. A. Tsaftaris, “Semi-supervised meta-learning with disentanglement for domain-generalised medical image segmentation,” in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2021, pp. 307–317.
- [47] M. Ilse, J. M. Tomczak, C. Louizos, and M. Welling, “Diva: Domain invariant variational autoencoders,” in *Medical Imaging with Deep Learning*. PMLR, 2020, pp. 322–348.
- [48] L. Van der Maaten and G. Hinton, “Visualizing data using t-sne.” *Journal of machine learning research*, vol. 9, no. 11, 2008.
- [49] Z. Yu, X. Han, W. Xu, J. Zhang, C. Marr, D. Shen, T. Peng, X.-Y. Zhang, and J. Feng, “Ben: a generalizable brain extraction net for multimodal mri data from rodents, nonhuman primates, and humans,” *bioRxiv*, 2022.