# MedGen3D: A Deep Generative Framework for Paired 3D Image and Mask Generation

Kun Han<sup>1</sup>, Yifeng Xiong<sup>1</sup>, Chenyu You<sup>2</sup>, Pooya Khosravi<sup>1</sup>, Shanlin sun<sup>1</sup>, Xiangyi Yan<sup>1</sup>, James Duncan<sup>2</sup>, and Xiaohui Xie<sup>1</sup>

<sup>1</sup> University of California, Irvine <sup>2</sup> Yale University

Abstract. Acquiring and annotating sufficient labeled data is crucial in developing accurate and robust learning-based models, but obtaining such data can be challenging in many medical image segmentation tasks. One promising solution is to synthesize realistic data with ground-truth mask annotations. However, no prior studies have explored generating complete 3D volumetric images with masks. In this paper, we present MedGen3D, a deep generative framework that can generate paired 3D medical images and masks. First, we represent the 3D medical data as 2D sequences and propose the Multi-Condition Diffusion Probabilistic Model (MC-DPM) to generate multi-label mask sequences adhering to anatomical geometry. Then, we use an image sequence generator and semantic diffusion refiner conditioned on the generated mask sequences to produce realistic 3D medical images that align with the generated masks. Our proposed framework guarantees accurate alignment between synthetic images and segmentation maps. Experiments on 3D thoracic CT and brain MRI datasets show that our synthetic data is both diverse and faithful to the original data, and demonstrate the benefits for downstream segmentation tasks. We anticipate that MedGen3D's ability to synthesize paired 3D medical images and masks will prove valuable in training deep learning models for medical imaging tasks.

**Keywords:** Deep Generative Framework  $\cdot$  3D Volumetric Images with Masks  $\cdot$  Fidelity and Diversity  $\cdot$  Segmentation

# 1 Introduction

In medical image analysis, the availability of a substantial quantity of accurately annotated 3D data is a prerequisite for achieving high performance in tasks like segmentation and detection [26,17,29,9,31,34,35,32,33]. This, in turn, leads to more precise diagnoses and treatment plans. However, obtaining and annotating such data presents many challenges, including the complexity of medical images, the requirement for specialized expertise, and privacy concerns.

Generating realistic synthetic data presents a promising solution to the above challenges as it eliminates the need for manual annotation and alleviates privacy risks. However, most prior studies [16,6,7,4,3,20,28,34,32,33] have focused on 2D image synthesis, with only a few generating corresponding segmentation masks. For instance, [15] uses dual generative adversarial networks (GAN) [14,35] to

synthesize 2D labeled retina fundus images, while [12] combines a label generator [25] with an image generator [24] to generate 2D brain MRI data. More recently, [27] uses WGAN [5] to generate small 3D patches and corresponding vessel segmentations.

However, there has been no prior research on generating whole 3D volumetric images with the corresponding segmentation masks. Generating 3D volumetric images with corresponding segmentation masks faces two major obstacles. First, directly feeding entire 3D volumes to neural networks is impractical due to GPU memory constraints, and downsizing the resolution may compromise the quality of the synthetic data. Second, treating the entire 3D volume as a single data point during training is suboptimal because of the limited availability of annotated 3D data. Thus, innovative methods are required to overcome these challenges and generate high-quality synthetic 3D volumetric data with corresponding segmentation masks.

We propose MedGen3D, a novel diffusion-based deep generative framework that generates paired 3D volumetric medical images and multi-label masks. Our approach treats 3D medical data as sequences of slices and employs an autoregressive process to sequentially generate 3D masks and images. In the first stage, a Multi-Condition Diffusion Probabilistic Model (MC-DPM) generates mask sequences by combining conditional and unconditional generation processes. Specifically, the MC-DPM generates mask subsequences (i.e., several consecutive slices) at any position directly from random noise or by conditioning on existing slices to generate subsequences forward or backward. Given that medical images have similar anatomical structures, slice indices serve as additional conditions to aid the mask subsequence generation. In the second stage, we introduce a conditional image generator with a seq-to-seq model from [30] and a semantic diffusion refiner. By conditioning on the mask sequences generated in the first stage, our image generator synthesizes realistic medical images aligned with masks while preserving spatial consistency across adjacent slices.

The main contributions of our work are as follows: 1) Our proposed framework is the *first* to address the challenge of synthesizing complete 3D volumetric medical images with their corresponding masks; 2) we introduce a multicondition diffusion probabilistic model for generating 3D anatomical masks with high fidelity and diversity; 3) we leverage the generated masks to condition an image sequence generator and a semantic diffusion refiner, which produces realistic medical images that align accurately with the generated masks; and 4) we present experimental results that demonstrate the fidelity and diversity of the generated 3D multi-label medical images, highlighting their potential benefits for downstream segmentation tasks.

# 2 Preliminary

### 2.1 Diffusion Probabilistic Model

A diffusion probabilistic model (DPM) [18] is a parameterized Markov chain of length T, which is designed to learn the data distribution p(X). DPM builds

the Forward Diffusion Process (FDP) to get the diffused data point  $X_t$  at any time step t by  $q(X_t | X_{t-1}) = \mathcal{N}(X_t; \sqrt{1 - \beta_t}X_{t-1}, \beta_t I)$ , with  $X_0 \sim q(X_0)$  and  $p(X_T) = \mathcal{N}(X_T; 0, I)$ . Let  $\alpha_t = 1 - \beta_t$  and  $\bar{\alpha}_t = \prod_{s=1}^t (1 - \beta_s)$ , Reverse Diffusion Process (RDP) is trained to predict the noise added in the FDP by minimizing:

 $Loss(\theta) = \mathbb{E}_{X_0 \sim q(X_0), \epsilon \sim \mathcal{N}(0, I), t} \left[ \left\| \epsilon - \epsilon_\theta \left( \sqrt{\bar{\alpha}_t} X_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t \right) \right\|^2 \right], \quad (1)$ where  $\epsilon_\theta$  is predicted noise and  $\theta$  is the model parameters.

#### 2.2 Classifier-free Guidance

Samples from conditional diffusion models can be improved with classifier-free guidance [19] by setting the condition c as  $\emptyset$  with probability p. During sampling, the output of the model is extrapolated further in the direction of  $\epsilon_{\theta} (X_t \mid c)$  and away from  $\epsilon_{\theta} (X_t \mid \emptyset)$  as follows:

$$\hat{\epsilon}_{\theta} \left( X_t \mid c \right) = \epsilon_{\theta} \left( X_t \mid \emptyset \right) + s \cdot \left( \epsilon_{\theta} \left( X_t \mid c \right) - \epsilon_{\theta} \left( X_t \mid \emptyset \right) \right), \tag{2}$$

where  $\emptyset$  represents a null condition and  $s \ge 1$  is the guidance scale.

### 3 Methodology

We propose a sequential process to generate complex 3D volumetric images with masks, as illustrated in Figure 1. The first stage generates multi-label segmentation, and the second stage performs conditional medical image generation. The details will be presented in the following sections.



Fig. 1. Overview of the proposed MedGen3D, including a 3D mask generator to autoregressively generate the mask sequences starting from a random position z, and a conditional image generator to generate 3D images conditioned on generated masks.

#### 3.1 3D Mask Generator

Due to the limited annotated real data and GPU memory constraints, directly feeding the entire 3D volume to the network is impractical. Instead, we treat 3D medical data as a series of subsequences. To generate an entire mask sequence, an initial subsequence of m consecutive slices is **unconditionally** generated from random noise. Then the subsequence is expanded **forward** and **backward** in an autoregressive manner, conditioned on existing slices.

Inspired by classifier-free guidance in Section 2.2, we propose a general Multi-Condition Diffusion Probabilistic Model (MC-DPM) to unify all three conditional generations (unconditional, forward, and backward). As shown in Fig. 2, MC-DPM is able to generate mask sequences directly from random noise or conditioning on existing slices.

Furthermore, as 3D medical data typically have similar anatomical structures, slices with the same relative position roughly correspond to the same anatomical regions. Therefore, we can utilize the relative position of slices as



**Fig. 2.** Proposed 3D mask generator. Given target position z, MC-DPM is designed to generate mask subsequences (length of m) for specific region, unconditionally or conditioning on first or last n slices, according to the pre-defined probability  $p^C \in \{p_F, p_B, p_U\}$ . Binary indicators are assigned to slices to signify the conditional slices. We ignore the binary indicators in the inference process for clear visualization with red outline denoting the conditional slices and green outline denoting the generated slices.

conditions to guide the MC-DPM in generating subsequences of the target region and control the length of generated sequences.

**Train:** For a given 3D multi-label mask  $M \in \mathbb{R}^{D \times H \times W}$ , subsequnces of m consecutive slices are selected as  $\{M_z, M_{z+1}, \ldots, M_{z+(m-1)}\}$ , with z as the randomly selected starting indices. For each subsequence, we determine the conditional slices  $X^C \in \{\mathbb{R}^{n \times H \times W}, \emptyset\}$  by selecting either the first or the last n slices, or no slice, based on a probability  $p^C \in \{p_{Forward}, p_{Backward}, p_{Uncondition}\}$ . The objective of the MC-DPM is to generate the remaining slices, denoted as  $X^P \in \mathbb{R}^{(m-\ln(X^C)) \times H \times W}$ .

To incorporate the position condition, we utilize the relative position of the subsequence  $\tilde{z} = z/D$ , where z is the index of the subsequence's starting slice. Then we embed the position condition and concatenate it with the time embedding to aid the generation process. We also utilize a binary indicator for each slice in the subsequence to signify the existence of conditional slices.

The joint distribution of reverse diffusion process (RDP) with the conditional slices  $X^C$  can be written as:

$$p_{\theta}(X_{0:T}^{P}|X^{C},\tilde{z}) = p(X_{T}^{P})\prod_{t=1}^{I} p_{\theta}(X_{t-1}^{P} \mid X_{t}^{P}, X^{C}, \tilde{z}).$$
(3)

where  $p(X_T^P) = \mathcal{N}(X_T^P; 0, I)$ ,  $\tilde{z} = z/D$  and  $p_{\theta}$  is the distribution parameterized by the model.

Overall, the model will be trained by minimizing the following loss function, with  $X_t^P = \sqrt{\overline{\alpha}_t} X_0^P + \sqrt{1 - \overline{\alpha}_t} \epsilon$ :

$$\operatorname{Loss}(\theta) = \mathbb{E}_{X_0 \sim q(X_0), \epsilon \sim \mathcal{N}(0, I), p^C, z, t} \left[ \left\| \epsilon - \epsilon_\theta \left( X_t^P, X^C, z, t \right) \right\|^2 \right].$$
(4)

**Inference:** During inference, MC-DPM first generates a subsequence of m slices from random noise given a random location z. The entire mask sequence can then be generated autoregressively by expanding in both directions, conditioned on the existing slices, as shown in Figure 2. Please refer to the **Supplementary** for a detailed generation process and network structure.

### 3.2 Conditional Image Generator

In the second step, we employ a sequence-to-sequence method to generate medical images conditioned on masks, as shown in Figure 3.

**Image Sequence Generator:** In the sequence-to-sequence generation task, new slice is the combination of the warped previous slice and newly generated texture, weighted by a continuous mask [30]. We utilize Vid2Vid [30] as our image sequence generator. We train Vid2Vid with its original loss, which includes GAN loss on multi-scale images and video discriminators, flow estimation loss, and feature matching loss.



**Fig. 3.** Image Sequence Generator. Given the generated 3D mask, the initial image is generated by Vid2Vid model sequentially. To utilize the semantic diffusion model (SDM) to refine the initial result, we first apply small steps (10 steps) noise, and then use three SDMs to refine. The final result is the mean 3D images from 3 different views (Axial, Coronal, and Sagittal), yielding significant improvements over the initially generated image.

Semantic Diffusion Refiner: Despite the high cross-slice consistency and spatial continuity achieved by vid2vid, issues such as blocking, blurriness and suboptimal texture generation persist. Given that diffusion models have been shown to generate superior images [11], we propose a semantic diffusion refiner utilizing a diffusion probabilistic model to refine the previously generated images.

For each of the 3 different views, we train a semantic diffusion model (SDM), which takes 2D masks and noisy images as inputs to generate images aligned with input masks. During inference, we only apply small noising steps (10 steps) to the generated images so that the overall anatomical structure and spatial continuity are preserved. After that, we refine the images using the pre-trained semantic diffusion model. The final refined 3D images are the mean results from 3 views. Experimental results show an evident improvement in the quality of generated images with the help of semantic diffusion refiner.

6 K. Han et al.

# 4 Experiments and Results

#### 4.1 Datasets and Setups

**Datasets:** We conducted experiments on the thoracic site using three thoracic CT datasets and the brain site with two brain MRI datasets. For both generative models and downstream segmentation tasks, we utilized the following datasets:

- SegTHOR [22]: 3D thorax CT scans (25 training, 5 validation, 10 testing);

- OASIS [23]: 3D brain MRI T1 scans (40 training, 10 validation, 10 testing); For the downstream segmentation task only and the transfer learning, we utilized 10 fine-tuning, 5 validation, and 10 testing scans from each of the 3D thorax CT datasets of StructSeg-Thorax [2] and Public-Thor [9], as well as the 3D brain MRI T1 dataset from ADNI [1].

**Implementation:** For thoracic datasets, we crop and pad CT scans to  $(96 \times 320 \times 320)$ . The annotations of six organs (left lung, right lung, spinal cord, esophagus, heart, and trachea) are examined by an experienced radiation oncologist. We also include a body mask to aid in the image generation of body regions. For brain MRI datasets, we use Freesurfer [13] to get segmentations of four regions (cortex, subcortical gray matter, white matter, and CSF), and then crop the volume to  $(192 \times 160 \times 160)$ . We assign discrete values to masks of different regions or organs for both thoracic and brain datasets and then combine them into one 3D volume. When synthesizing mask sequences, we resize the width and height of the masks to  $128 \times 128$  and set the length of the subsequence m to 6. We use official segmentation models provided by MONAI[8] along with standard data augmentations, including spatial and color transformations.

**Setup:** We compare the synthetic image quality with DDPM [18], 3D- $\alpha$ -WGAN [21] and Vid2Vid [30], and utilize four segmentation models with different training strategies to demonstrate the benefit for the downstream task.

### 4.2 Evaluate the Quality of Synthetic Image.

**Synthetic Dataset:** To address the limited availability of annotated 3D medical data, we used only 30 CT scans from SegTHOR (25 for training and 5 for validation) and 50 MRI scans from OASIS (40 for training and 10 for validation) to generate 110 3D thoracic CT scans and 110 3D brain MRI scans, respectively.

	Thor	acic CT	Brai	in MRI
	$\mathrm{FID}\downarrow$	LPIPS $\uparrow$	$\mathrm{FID}\downarrow$	LPIPS $\uparrow$
DDPM [18] 3D-α-WGAN [21] Vid2Vid [30]	<b>35.2</b> 136.2 47.3	<b>0.316</b> 0.286 0.300	<b>34.9</b> 136.4 48.2	$\begin{array}{c} 0.298 \\ 0.289 \\ 0.324 \end{array}$
Ours	39.6	0.305	40.3	0.326

 Table 1. Synthetic image quality comparison between baselines and ours.

We compare the fidelity and diversity of our synthetic data with DDPM [18] (train 3 for different views), 3D- $\alpha$ -WGAN [21], and vid2vid [30] by calculating the mean Frèchet Inception Distance (FID) and Learned Perceptual Image Patch Similarity (LPIPS) from 3 different views.

According to Table 1, our proposed method has a slightly lower FID score but a similar LPIPS score compared to DDPM. We speculate that this is because DDPM is trained on 2D images without explicit anatomical constraints and only generates 2D images. On the other hand, 3D- $\alpha$ -WGAN [18], which uses much larger 3D training data (146 for thorax and 414 for brain), has significantly worse FID and LPIPS scores than our method. Moreover, our proposed method outperforms Vid2Vid, showing the effectiveness of our semantic diffusion refiner.



Fig. 4. Our proposed method produces more anatomically accurate images compared to  $3D-\alpha$ -WGAN and vid2vid, as demonstrated by the clearer organ boundaries and more realistic textures. Left: Qualitative comparison between different generative models. Right: Visualization of synthetic 3D brain MRI slices at different relative positions.

#### 4.3 Evaluate the Benefits for Segmentation Task.

We explore the benefits of synthetic data for downstream segmentation tasks by comparing Sørensen–Dice coefficient (DSC) of 4 segmentation models, including Unet2D [26], UNet3D [10], UNETR [17], and Swin-UNETR [29]. In Table 2 and 3, we utilize real training data (from SegTHOR and OASIS) and synthetic data to train the segmentation models with 5 different strategies, and test on all 3 thoracic CT datasets and 2 brain MRI datasets. In Table 4, we aim to demonstrate whether the synthetic data can aid transfer learning with limited real finetuning data from each of the testing datasets (StructSeg-Thorax, Public-Thor and ADNI) with four training strategies.

	SegTHOR*					Structs	Seg-Tho	rax	Public-Thor				
	Unet 2D	Unet 3D	UNETR	Swin UNETR	Unet 2D	Unet 3D	UNETR	Swin UNETR	Unet 2D	Unet 3D	UNETR	Swin UNETR	
E2-1 E2-2 E2-3 E2-4 E2-5	0.817 0.815 0.845 <b>0.855</b> 0.847	0.873 0.846 0.881 0.887 0.891	0.867 0.845 0.886 <b>0.894</b> 0.890	0.878 0.854 0.886 <b>0.899</b> 0.897	0.722 0.736 0.772 0.775 0.783	0.793 0.788 0.827 <b>0.833</b> 0.833	0.789 0.788 0.824 <b>0.825</b> 0.823	0.810 0.803 0.827 0.833 0.835	0.822 0.786 0.812 <b>0.824</b> 0.818	0.837 0.838 0.856 0.861 0.864	0.836 0.814 0.853 0.852 0.858	0.847 0.842 0.856 <b>0.867</b> 0.867	

**Table 2.** Experiment 2: DSC of different thoracic segmentation models. There are 5 training strategies, namely: **E2-1**: Training with real SegTHOR training data; **E2-2**: Training with synthetic data; **E2-3**: Training with both synthetic and real data; **E2-4**: Finetuning model from E2-2 using real training data; and **E2-5**: finetuning model from E2-3 using real training data source.)

According to Table 2 and Table 3, the significant DSC difference between 2D and 3D segmentation models underlines the crucial role of 3D annotated data.

### 8 K. Han et al.

While purely synthetic data (E2-2) fails to achieve the same performance as real training data (E2-1), the combination of real and synthetic data (E2-3) improves model performance in most cases, except for Unet2D on the Public-Thor dataset. Furthermore, fine-tuning the pre-trained model with real data (E2-4 and E2-5) consistently outperforms the model trained only with real data. Please refer to Supplementary for organ-level DSC comparisons of the Swin-UNETR model with more details.

		0	ASIS*		ADNI					
	Unet 2D	Unet 3D	UNETR	Swin UNETR	Unet 2D	Unet 3D	UNETR	Swin UNETR		
E2-1 E2-2 E2-3	$0.930 \\ 0.905 \\ 0.938$	$0.951 \\ 0.936 \\ 0.953$	$0.952 \\ 0.935 \\ 0.953$	$0.954 \\ 0.934 \\ 0.955$	$0.815 \\ 0.759 \\ 0.818$	$0.826 \\ 0.825 \\ 0.888$	$0.880 \\ 0.828 \\ 0.898$	0.894 0.854 <b>0.906</b>		
E2-4 E2-5	0.940 0.940	<b>0.955</b> 0.954	$0.954 \\ 0.954$	$0.956 \\ 0.956$	0.819 0.819	0.891 0.894	<b>0.903</b> 0.902	0.903 <b>0.906</b>		

**Table 3.** Experiment 2: DSC of brain segmentation models. Please refer to Table 2 for detailed training strategies. (\* denotes the training data source.)

According to Table 4, for transfer learning, utilizing the pre-trained model (E3-2) leads to better performance compared to training from scratch (E3-1). Additionally, pretraining the model with synthetic data (E3-3 and E3-4) can facilitate transfer learning to a new dataset with limited annotated data.

	Thoracic	CT	Brain MRI
	StructSeg-Thorax*	Public-Thor*	ADNI*
E3-1 E3-2 E3-3	$0.845 \\ 0.865 \\ 0.878$	$0.897 \\ 0.901 \\ 0.913$	0.946 0.948 <b>0.949</b>
E3-4	0.882	0.914	0.949

Table 4. Experiment 3: DSC of Swin-UNETR finetuned with real dataset. There are 4 training strategies: E3-1: Training from scratch for each dataset using limited finetuning data; E3-2 Finetuning the model E2-1 from experiment 2; E3-3 Finetuning the model E2-4 from experiment 2; and E3-4 Finetuning the model E2-5 from experiment 2. (\* denotes the finetuning data source.)

We have included video demonstrations of the generated 3D volumetric images in the **supplementary material**, which offer a more comprehensive representation of the generated image's quality.

## 5 Conclusion

This paper introduces MedGen3D, a new framework for synthesizing 3D medical mask-image pairs. Our experiments demonstrate its potential in realistic data generation and downstream segmentation tasks with limited annotated data. Future work includes merging the image sequence generator and semantic diffusion refiner for end-to-end training and extending the framework to synthesize 3D medical images across modalities. Overall, we believe that our work opens up new possibilities for generating 3D high-quality medical images paired with masks, and look forward to future developments in this field.

### References

- 1. https://adni.loni.usc.edu/
- 2. https://structseg2019.grand-challenge.org/dataset/
- Abbasi-Sureshjani, S., Amirrajab, S., Lorenz, C., Weese, J., Pluim, J., Breeuwer, M.: 4d semantic cardiac magnetic resonance image synthesis on xcat anatomical model. In: Medical Imaging with Deep Learning. pp. 6–18. PMLR (2020)
- Abhishek, K., Hamarneh, G.: Mask2lesion: Mask-constrained adversarial skin lesion image synthesis. In: Simulation and Synthesis in Medical Imaging: 4th International Workshop, SASHIMI 2019, Held in Conjunction with MICCAI 2019, Shenzhen, China, October 13, 2019, Proceedings. pp. 71–80. Springer (2019)
- Arjovsky, M., Chintala, S., Bottou, L.: Wasserstein gan. arXiv preprint arXiv: Arxiv-1701.07875 (2017)
- Baur, C., Albarqouni, S., Navab, N.: Melanogans: high resolution skin lesion synthesis with gans. arXiv preprint arXiv:1804.04338 (2018)
- Bermudez, C., Plassard, A.J., Davis, L.T., Newton, A.T., Resnick, S.M., Landman, B.A.: Learning implicit brain mri manifolds with deep learning. In: Medical Imaging 2018: Image Processing. vol. 10574, pp. 408–414. SPIE (2018)
- Cardoso, M.J., Li, W., Brown, R., Ma, N., Kerfoot, E., Wang, Y., Murrey, B., Myronenko, A., Zhao, C., Yang, D., et al.: Monai: An open-source framework for deep learning in healthcare. arXiv preprint arXiv:2211.02701 (2022)
- Chen, X., Sun, S., Bai, N., Han, K., Liu, Q., Yao, S., Tang, H., Zhang, C., Lu, Z., Huang, Q., et al.: A deep learning-based auto-segmentation system for organsat-risk on whole-body computed tomography images for radiation therapy. Radiotherapy and Oncology 160, 175–184 (2021)
- Çiçek, Ö., Abdulkadir, A., Lienkamp, S.S., Brox, T., Ronneberger, O.: 3d u-net: learning dense volumetric segmentation from sparse annotation. In: Medical Image Computing and Computer-Assisted Intervention-MICCAI 2016: 19th International Conference, Athens, Greece, October 17-21, 2016, Proceedings, Part II 19. pp. 424– 432. Springer (2016)
- Dhariwal, P., Nichol, A.: Diffusion models beat gans on image synthesis. Advances in Neural Information Processing Systems 34, 8780–8794 (2021)
- Fernandez, V., Pinaya, W.H.L., Borges, P., Tudosiu, P.D., Graham, M.S., Vercauteren, T., Cardoso, M.J.: Can segmentation models be trained with fully synthetically generated data? In: Simulation and Synthesis in Medical Imaging: 7th International Workshop, SASHIMI 2022, Held in Conjunction with MICCAI 2022, Singapore, September 18, 2022, Proceedings. pp. 79–90. Springer (2022)
- 13. Fischl, B.: Freesurfer. Neuroimage **62**(2), 774–781 (2012)
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial networks. Communications of the ACM 63(11), 139–144 (2020)
- Guibas, J.T., Virdi, T.S., Li, P.S.: Synthetic medical images from dual generative adversarial networks. arXiv preprint arXiv:1709.01872 (2017)
- Han, C., Hayashi, H., Rundo, L., Araki, R., Shimoda, W., Muramatsu, S., Furukawa, Y., Mauri, G., Nakayama, H.: Gan-based synthetic brain mr image generation. In: 2018 IEEE 15th international symposium on biomedical imaging (ISBI 2018). pp. 734–738. IEEE (2018)
- Hatamizadeh, A., Tang, Y., Nath, V., Yang, D., Myronenko, A., Landman, B., Roth, H.R., Xu, D.: Unetr: Transformers for 3d medical image segmentation. In: Proceedings of the IEEE/CVF winter conference on applications of computer vision. pp. 574–584 (2022)

- 10 K. Han et al.
- Ho, J., Jain, A., Abbeel, P.: Denoising diffusion probabilistic models. Advances in Neural Information Processing Systems 33, 6840–6851 (2020)
- Ho, J., Salimans, T.: Classifier-free diffusion guidance. arXiv preprint arXiv: Arxiv-2207.12598 (2022)
- Kim, S., Kim, B., Park, H.: Synthesis of brain tumor multicontrast mr images for improved data augmentation. Medical Physics 48(5), 2185–2198 (2021)
- Kwon, G., Han, C., Kim, D.s.: Generation of 3d brain mri using auto-encoding generative adversarial networks. In: Medical Image Computing and Computer Assisted Intervention–MICCAI 2019: 22nd International Conference, Shenzhen, China, October 13–17, 2019, Proceedings, Part III 22. pp. 118–126. Springer (2019)
- Lambert, Z., Petitjean, C., Dubray, B., Kuan, S.: Segthor: Segmentation of thoracic organs at risk in ct images. In: 2020 Tenth International Conference on Image Processing Theory, Tools and Applications (IPTA). pp. 1–6. IEEE (2020)
- Marcus, D.S., Wang, T.H., Parker, J., Csernansky, J.G., Morris, J.C., Buckner, R.L.: Open access series of imaging studies (oasis): cross-sectional mri data in young, middle aged, nondemented, and demented older adults. Journal of cognitive neuroscience 19(9), 1498–1507 (2007)
- Park, T., Liu, M.Y., Wang, T.C., Zhu, J.Y.: Semantic image synthesis with spatially-adaptive normalization. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 2337–2346 (2019)
- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 10684–10695 (2022)
- Ronneberger, O., Fischer, P., Brox, T.: U-net: Convolutional networks for biomedical image segmentation. In: Medical Image Computing and Computer-Assisted Intervention-MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18. pp. 234–241. Springer (2015)
- 27. Subramaniam, P., Kossen, T., Ritter, K., Hennemuth, A., Hildebrand, K., Hilbert, A., Sobesky, J., Livne, M., Galinovic, I., Khalil, A.A., et al.: Generating 3d tof-mra volumes and segmentation labels using generative adversarial networks. Medical Image Analysis 78, 102396 (2022)
- Sun, L., Chen, J., Xu, Y., Gong, M., Yu, K., Batmanghelich, K.: Hierarchical amortized gan for 3d high resolution medical image synthesis. IEEE journal of biomedical and health informatics 26(8), 3966–3975 (2022)
- 29. Tang, Y., Yang, D., Li, W., Roth, H.R., Landman, B., Xu, D., Nath, V., Hatamizadeh, A.: Self-supervised pre-training of swin transformers for 3d medical image analysis. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 20730–20740 (2022)
- Wang, T.C., Liu, M.Y., Zhu, J.Y., Liu, G., Tao, A., Kautz, J., Catanzaro, B.: Video-to-video synthesis. arXiv preprint arXiv:1808.06601 (2018)
- Yan, X., Tang, H., Sun, S., Ma, H., Kong, D., Xie, X.: After-unet: Axial fusion transformer unet for medical image segmentation. In: Proceedings of the IEEE/CVF winter conference on applications of computer vision. pp. 3971–3981 (2022)
- You, C., Dai, W., Liu, F., Su, H., Zhang, X., Staib, L., Duncan, J.S.: Mine your own anatomy: Revisiting medical image segmentation with extremely limited labels. arXiv preprint arXiv:2209.13476 (2022)
- 33. You, C., Dai, W., Min, Y., Liu, F., Zhang, X., Clifton, D.A., Zhou, S.K., Staib, L.H., Duncan, J.S.: Rethinking semi-supervised medical image segmentation: A variance-reduction perspective. arXiv preprint arXiv:2302.01735 (2023)

- You, C., Dai, W., Staib, L., Duncan, J.S.: Bootstrapping semi-supervised medical image segmentation with anatomical-aware contrastive distillation. arXiv preprint arXiv:2206.02307 (2022)
- You, C., Zhao, R., Liu, F., Dong, S., Chinchali, S., Topcu, U., Staib, L., Duncan, J.: Class-aware adversarial transformers for medical image segmentation. In: NeurIPS (2022)

#### Algorithm 1 3D Mask Generation (Inference)

Require: n: number of conditional slices, m: length of subsequence, L: length of the sequence to generate,  $P_{\theta}$ : the predicted probability distribution  $Z \leftarrow L - (m-1)$  $z \sim \text{Uniform} (\{0, 1, ..., Z\})$  $\triangleright$  Randomly pick one z as start position Initialize an empty mask sequence  $\mathcal{M}$  $X \leftarrow \{M_z, M_{z+1}, ..., M_{z+(m-1)}\} \sim P_{\theta}(X_0^P \mid X^C = \emptyset, z)$  $\mathcal{M} \leftarrow \mathcal{M} \cup X$ //Forward Sampling z' $\begin{array}{c} z' \leftarrow z \\ \textbf{while} \ z' <= Z \ \textbf{do} \\ z' \leftarrow z' + (m-n) \end{array}$  $X^C \leftarrow \mathcal{M}[-n:]$  $\triangleright$  Select the last n masks as condition  $\begin{array}{l} X^P \sim P_{\theta}(X^P_0 \mid X^C, z') \\ \mathcal{M} \leftarrow \mathcal{M} \cup X^P \qquad \triangleright \end{array}$  $\triangleright$  Sample the following (m-n) masks  $\triangleright$  Add the generated masks to the end of sequence end while /Backward Sampling  $z' \leftarrow z \\ \text{while } z' \ge 0 \text{ do} \\ z' \leftarrow z' - (m-n)$  $X^C \leftarrow \mathcal{M}[:n]$  $\triangleright$  Select the first n masks as condition  $\begin{array}{l} X^P \sim P_{\theta}(X_0^P \mid X^C, z') \\ \mathcal{M} \leftarrow X^P \cup \mathcal{M} \qquad \vartriangleright \end{array}$  $\triangleright$  Sample the previous (m-n) masks  $\triangleright$  Add the generated masks to the start of sequence end while return  $\mathcal{M}$  $\triangleright$  Return the generated mask sequence

		SegTHOR*						StructSeg-Thorax					Public-Thor					
	11	rl	ht	eso	$\operatorname{tra}$	$_{\rm spin}$	11	rl	ht	eso	$\operatorname{tra}$	$_{\rm spin}$	11	rl	ht	eso	$\operatorname{tra}$	spin
E2-1	0.98	0.99	0.90	0.64	0.86	0.91	0.95	0.96	0.91	0.60	0.68	0.77	0.97	0.98	0.88	0.69	0.73	0.84
E2-2	0.98	0.98	0.91	0.55	0.83	0.89	0.94	0.95	0.89	0.54	0.63	0.86	0.97	0.98	0.88	0.68	0.69	0.85
E2-3	0.98	0.98	0.92	0.64	0.87	0.93	0.95	0.95	0.90	0.64	0.65	0.87	0.97	0.98	0.88	0.73	0.72	0.86
E2-4	0.98	0.99	0.93	0.67	0.90	0.94	0.95	0.96	0.90	0.64	0.69	0.87	0.97	0.98	0.90	0.74	0.75	0.86
E2-5	0.98	0.99	0.92	0.67	0.88	0.94	0.95	0.96	0.90	0.67	0.69	0.86	0.97	0.98	0.89	0.75	0.77	0.85
	ll:	left lı	ung	rl:	right	lung	ht	: hea	rt	eso:	esopł	nagus	$\operatorname{tra}$	trac	hea	spin:	spina	l cord

**Table 1.** Experiment 2: Organ-level DSC Comparison of Swin-UNETR for thoracic site. Please refer to **Table 2** in **main** submission for detailed strategies.

		OASIS	*	ADNI					
	Cortex	Subcortical Gray	White Matter	CSF	Cortex	Subcortical Gray	White Matter	CSF	
E2-1	0.941	0.965	0.970	0.941	0.884	0.856	0.929	0.908	
E2-2	0.922	0.951	0.958	0.909	0.847	0.841	0.910	0.818	
E2-3	0.942	0.964	0.970	0.944	0.895	0.876	0.932	0.925	
E2-4	0.943	0.965	0.971	0.946	0.888	0.870	0.931	0.924	
E2-5	0.943	0.966	0.971	0.947	0.895	0.875	0.931	0.919	

Table 2. Experiment 2: Organ-level DSC Comparison of Swin-UNETR for brain site.



Fig. 1. Network Structure of MC-DPM.



Fig. 2. Experiment 2: Qualitative comparison of different training strategies using Swin-UNETR. For brain segmentation, improvements brought by synthetic data are more evident when training and testing data come from different datasets (ADNI: train on OASIS, test on ADNI) rather than same dataset (OASIS\*).



Fig. 3. Experiment 3: Qualitative comparison of different finetuning strategies using Swin-UNETR. Please refer to Table 4 in main submission for detailed strategies.