# A Recycling Training Strategy for Medical Image Segmentation with Diffusion Denoising Models

Yunguan Fu  HTTPS://ORCID.ORG/0000-0002-1184-7421
University College London, UK; InstaDeep, UK

yunguan.fu.18@ucl.ac.uk; y.fu@instadeep.com

Yiwen Li  HTTPS://ORCID.ORG/0000-0002-7794-9391
University of Oxford, UK

yiwen.li@st-annes.ox.ac.uk

Shaheer U. Saeed  HTTPS://ORCID.ORG/0000-0002-5004-0663
University College London, UK

shaheer.saeed.17@ucl.ac.uk

Matthew J. Clarkson  HTTPS://ORCID.ORG/0000-0002-5565-1252
University College London, UK

m.clarkson@ucl.ac.uk

Yipeng Hu  HTTPS://ORCID.ORG/0000-0003-4902-0486
University College London, UK

yipeng.hu@ucl.ac.uk

## Abstract

Denoising diffusion models have found applications in image segmentation by generating segmented masks conditioned on images. Existing studies predominantly focus on adjusting model architecture or improving inference, such as test-time sampling strategies. In this work, we focus on improving the training strategy and propose a novel recycling method. During each training step, a segmentation mask is first predicted given an image and a random noise. This predicted mask, which replaces the conventional ground truth mask, is used for the denoising task during training. This approach can be interpreted as aligning the training strategy with inference by eliminating the dependence on ground truth masks for generating noisy samples. Our proposed method significantly outperforms standard diffusion training, self-conditioning, and existing recycling strategies across multiple medical imaging data sets: muscle ultrasound, abdominal CT, prostate MR, and brain MR. This holds for two widely adopted sampling strategies: denoising diffusion probabilistic model and denoising diffusion implicit model. Importantly, existing diffusion models often display a declining or unstable performance during inference, whereas our novel recycling consistently enhances or maintains performance. We show for the first time that, under a fair comparison with the same network architectures and computing budget, the proposed recycling-based diffusion models achieved on-par performance with non-diffusion-based supervised training. Furthermore, by ensembling the proposed diffusion model and the non-diffusion counterpart, significant improvements to the non-diffusion models have been observed across all applications, demonstrating the value of this novel training method. This paper summarizes these quantitative results and discusses their values, with a fully reproducible JAX-based implementation, released at https://github.com/mathpluscode/ImgX-DiffSeg.

**Keywords:** Image Segmentation, Diffusion Model, Recycling, Muscle Ultrasound, Abdominal CT, Prostate MR, Brain MR

# 1. Introduction

Diffusion denoising models, first proposed by Sohl-Dickstein et al. (2015); Ho et al. (2020); Ho and Salimans (2022), are generative models that produce data samples through iterative denoising processes. These models achieved superior performance compared to generative adversarial networks (Goodfellow et al., 2020) and became the foundation for many image generation applications such as DALL·E 2 (Ramesh et al., 2022), stable diffusion, and Midjourney (Rombach et al., 2022), etc. Given the success in computer vision, diffusion models have been adapted in medical imaging in various fields, including image synthesis (Dorjsembe et al., 2022; Khader et al., 2022), image denoising (Hu et al., 2022), anomaly detection (Wolleb et al., 2022a), classification (Yang et al., 2023), segmentation (Wu et al., 2022; Rahman et al., 2023), and registration (Kim et al., 2022). Among these, segmentation is one of the most foundational tasks in medical imaging and a variety of applications have been explored, including liver CT (Xing et al., 2023), lung CT (Zbinden et al., 2023; Rahman et al., 2023), abdominal CT (Wu et al., 2023; Fu et al., 2023), brain MR (Pinaya et al., 2022a; Wolleb et al., 2022b; Wu et al., 2023; Xing et al., 2023; Bieder et al., 2023), and prostate MR (Fu et al., 2023).

For segmentation tasks, although various model architectures and training strategies (Wang et al., 2022) have been proposed, U-net equipped with attention mechanisms and trained by supervised learning consistently remains the state-of-the-art model and an important baseline. In comparison, divergent observations have emerged: some studies reported superior performance of diffusion-based segmentation models (Amit et al., 2021; Wu et al., 2022, 2023; Xing et al., 2023), while others observed the opposite trend (Pinaya et al., 2022a; Wolleb et al., 2022b; Kolbeinsson and Mikolajczyk, 2022; Fu et al., 2023). This inconsistency may result from different training schemes, network architectures, and application-specific modifications in comparisons, suggesting that challenges persist in applying diffusion models for image segmentation.

Formally, conditioning on an image, diffusion-based segmentation models operate by progressive denoising, starting with random noise and ultimately yielding the corresponding segmentation masks. In comparison to their non-diffusion counterparts, the necessity of supplementary noisy masks as input leads to increased memory demands that can pose challenges, particularly for processing 3D volumetric medical images. To address this, volume slicing (Wu et al., 2023) or patching (Xing et al., 2023; Bieder et al., 2023) has been used to manage memory limitations. However, diffusion model training still requires considerable computation due to its inherent iterative nature, since the same model needs to learn to denoise masks with varying levels of noise. Consequently, enhancing the diffusion model performance while adhering to a fixed compute budget is of significant importance. Empirically, using the reparametrisation (Kingma et al., 2021), the denoising training task has shifted from noise prediction (Wolleb et al., 2022b; Wu et al., 2022) to mask prediction (Fu et al., 2023; Zbinden et al., 2023) due to the superior performance and faster learning. Furthermore, Fu et al. (2023) highlighted a limitation of diffusion models, noting the misalignment between training and inference procedures, since training samples were generated from ground truth masks. This raises concerns of data leakage as discussed in Chen et al. (2022a). However, there have been limited studies in medical image segmentation that rigorously compare diffusion models with their non-diffusion counterparts and examine diffusion training efficiency.

In this work, we present a substantial extension to the preliminary work (Fu et al., 2023) and focus on an improvement in the diffusion denoising model training strategy that applies to 2D and 3D medical image segmentation in different modalities. First, a novel recycling approach has been introduced. Different from Fu et al. (2023), in the first step during training, the input is completely corrupted by noise instead of a partially corrupted ground truth. This seemingly minor adjustment effectively eliminates the ground truth information from model inputs, which further aligns the training strategy toward inference. The proposed diffusion models can refine or maintain segmentation accuracy throughout the inference process. On the contrary, all other diffusion models demonstrate declining or unstable performance trends. Our research showcases the superior performance of our method compared to established diffusion training strategies (Ho et al., 2020; Chen et al., 2022b; Watson et al., 2023; Fu et al., 2023) for both denoising diffusion probabilistic model-based (Ho et al., 2020) and denoising diffusion implicit model-based (Song et al., 2020a) sampling procedures. We also achieved on-par performance with non-diffusion baselines that had not been observed in the previous study (Fu et al., 2023). Second, we introduce an ensemble model that averages the predicted probabilities from the proposed diffusion-based model and non-diffusion counterpart, resulting in significant improvement to the non-diffusion baseline. Third, we extended the experiments to four large data sets – 2D muscle ultrasound with 3910 images, 3D abdominal CT with 300 images, 3D prostate MR with 589 images, and 3D brain MR with 1251 images, further demonstrating the robustness of the proposed method against different applications and data types. Lastly, we integrated a Transformer block into our network architecture. This brings our models in line with contemporary state-of-the-art approaches, rendering our findings more pertinent to real-world applications. To mitigate the increased memory consumption resulting from this addition, we employed patch-based training and inference strategies. The JAX-based framework has been released on `https://github.com/mathpluscode/ImgX-DiffSeg`.

## 2. Related Works

The diffusion process is a Markov process where data structures are gradually noise-corrupted and eventually destroyed (noising process). A reverse diffusion process (denoising process) can then be learned, where the objective is to gradually recover the data structure. Sohl-Dickstein et al. (2015) first proposed diffusion models which map the disrupted data to a noise distribution. Ho et al. (2020) have shown that such modeling is equivalent to score-matching models, a class of models that estimates the gradient of the log-density (Hyvärinen and Dayan, 2005; Vincent, 2011; Song and Ermon, 2019, 2020). This led to a simplified variational lower bound training objective and a denoising diffusion probabilistic model (DDPM) (Ho et al., 2020). DDPM achieved state-of-the-art performance for unconditional image generation on CIFAR10 at the time. In practice, DDPMs were found suboptimal on log-likelihood estimation and Nichol and Dhariwal (2021) addressed this with a learnable variance schedule, sinusoidal noise schedule, and an importance sampling for time steps. Furthermore, diffusion models were trained with hundreds or thousands of steps, inference with the same number of steps is time-consuming. Therefore, different strategies have been proposed to enable faster sampling. While Nichol and Dhariwal (2021) suggested variance resampling without modifying the probabilistic distribution, Song et al. (2020a) derived a deterministic model,

denoising diffusion implicit model (DDIM), which shares the same marginal distribution as DDPM. Liu et al. (2022) further generalized the reverse step of DDIM into an ordinary differential equation and used high-order numerical methods (e.g., Runge-Kutta method) with predicted noise to perform sampling with second-order convergence. Besides, Zheng et al. (2022); Lyu et al. (2022); Guo et al. (2022) accelerated diffusion model training by shortening the noising schedule and only considering a truncated diffusion chain with less noise. These unconditioned denoising diffusion models have been successfully applied in multiple medical imaging applications (Kazerouni et al., 2023), including brain MR image generation (Dorjsembe et al., 2022; Khader et al., 2022), optical coherence tomography denoising (Hu et al., 2022), and chest X-ray pleural effusion detection (Wolleb et al., 2022a).

Guided diffusion models have been developed to generate data in a controllable manner. Song et al. (2020b); Dhariwal and Nichol (2021) used gradients of pre-trained classifiers to bias the sampling process, without modifying the diffusion model training. Ho and Salimans (2022), on the other hand, modified the models to take additional information as input, enabling end-to-end conditional diffusion model training. For medical image synthesis, conditions can be patient biometric information (Pinaya et al., 2022b), genotypes data (Moghadam et al., 2023), or images from different modalities (Saeed et al., 2023). Conditional diffusion models have also been explored for medical image classification (Yang et al., 2023), segmentation (Wu et al., 2022; Rahman et al., 2023), and registration (Kim et al., 2022). Particularly for image segmentation, the diffusion models apply the noising and denoising on the segmentation masks, and the network takes a noisy mask and an image to perform denoising.

In contrast to the continuous spectrum of values found in natural images, image segmentation mask values are categorical and nominal. The Gaussian-based continuous diffusion processes behind DDPM and DDIM cannot be directly applied. Chen et al. (2022b) therefore encoded categories with binary bits and relaxed them to real values for continuous diffusion models. Han et al. (2022); Fu et al. (2023) encoded categories with one-hot embeddings and performed diffusion on scaled values. Li et al. (2022a); Strudel et al. (2022) encoded the discrete data and applied diffusion processes in embedding spaces directly. Alternatively, discrete diffusion models have been proposed to model the transition matrix between categories based on discrete probability distributions, including binomial distribution (Sohl-Dickstein et al., 2015), categorical distribution (Hoogeboom et al., 2021; Austin et al., 2021; Gu et al., 2022), and Bernoulli distribution (Chen et al., 2023). In this work, we follow Fu et al. (2023) to perform diffusion on scaled binary masks.

Originally, DDPM models were trained through noise prediction (Ho et al., 2020), where the loss was calculated between the predicted and sampled noises. Many diffusion-based segmentation models directly adopted this strategy (Wolleb et al., 2022b; Wu et al., 2022). Alternatively, Kingma et al. (2021) derived an equivalent formulation (often known as $\mathbf{x}_0$ reparameterization) of the variational lower bound and simplified the loss to a norm between predicted data and the corresponding ground truth. For segmentation, this is equivalent to predicting the segmentation mask for each time step. Compared to noise prediction, multiple studies found that this mask prediction strategy is more efficient (Fu et al., 2023; Wang et al., 2023; Lai et al., 2023). Furthermore, Chen et al. (2022b) suggested self-conditioning to use these predictions as additional input to improve diffusion models for image synthesis. Self-conditioning contains two steps: the first step predicts a noise-free sample given a

noise-corrupted sample only; the second step uses the same timestep and inputs the same noise-corrupted sample, as well as the prediction from the first step. This technique was later adopted for protein design (Watson et al., 2023) with an additional reverse step, where the second step performs denoising in a smaller timestep where the noise level is lower. However, in both cases, the noisy samples are directly derived from the ground truth, which is not available during inference. This risks data leakage during training and empirically leads to overfitting and lack of generalization as discussed in Chen et al. (2022a); Kolbeinsson and Mikolajczyk (2022); Lai et al. (2023). Chen et al. (2022a); Young et al. (2022) addressed this issue by controlling the signal-to-noise ratio so that less information is preserved after noising: Chen et al. (2022a) scaled the mask value ranges to implicitly amplify the noise level, and Young et al. (2022) explicitly varied the scale and standard deviation of the Gaussian noise added to the masks. On the other hand, Kolbeinsson and Mikolajczyk (2022) proposed recursive denoising that iterates through each step during training, without using ground truth as input. However, such a strategy extends the training length by a factor of hundreds or more, making it practically infeasible for larger 3D medical image data sets. Following these studies, Fu et al. (2023) concluded that the lack of generalization in diffusion-based segmentation models is due to the misalignment between training and inference processes. Fu et al. (2023) thus presented a two-step recycling training strategy: the first step ingests a partially noisied sample for mask prediction; the predicted mask is then noise-corrupted again for denoising training. Compared to recursive denoising, this method requires a limited training time increase. This method also resembles PD-DDPM (Guo et al., 2022), where a pre-segmentation is used for noising. However, PD-DDPM requires a separate pre-segmentation network and more device memory, thus not suitable for 3D image segmentation applications.

## 3. Background

### 3.1 Denoising Diffusion Probabilistic Model

$$\mathbf{x}_T \rightleftharpoons \cdots \rightleftharpoons \mathbf{x}_t \xrightleftharpoons[q(\mathbf{x}_t \mid \mathbf{x}_{t-1})]{p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t)} \mathbf{x}_{t-1} \rightleftharpoons \cdots \rightleftharpoons \mathbf{x}_0 \qquad (1)$$

**Definition** The denoising diffusion probabilistic models (DDPM) (Ho et al., 2020) consider a *forward* process (illustrated from right to left in Equation (1)): given a sample $\mathbf{x}_0 \sim q(\mathbf{x}_0)$, a noise-corrupted sample $\mathbf{x}_t$ follows a multivariate normal distribution at timestep $t \in \{1, \cdots, T\}$, $q(\mathbf{x}_t \mid \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1-\beta_t}\mathbf{x}_{t-1}, \beta_t \mathbf{I})$, where $\beta_t \in [0, 1]$. As Gaussians are closed under convolution, given $\mathbf{x}_0$, $\mathbf{x}_t$ can be directly sampled from $\mathbf{x}_0$ as follows, $q(\mathbf{x}_t \mid \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1-\bar{\alpha}_t)\mathbf{I})$. Correspondingly, a *reverse* process (illustrated from left to right in Equation (1)) denoises $\mathbf{x}_t$ at each step, for $t \in \{T, \cdots, 1\}$, $p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I})$, with a predicted mean $\boldsymbol{\mu}_\theta(\mathbf{x}_t, t)$ and variance $\sigma_t^2 \mathbf{I}$. $\sigma_t$ is a pre-defined schedule dependent on timestep $t$. In this work, $\sigma_t^2 = \tilde{\beta}_t = \frac{1-\bar{\alpha}_{t-1}}{1-\bar{\alpha}_t}\beta_t$ with $\alpha_t = 1 - \beta_t$ and $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$. The mean $\boldsymbol{\mu}_\theta(\mathbf{x}_t, t) = \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1-\bar{\alpha}_t}\hat{\mathbf{x}}_0 + \frac{1-\bar{\alpha}_{t-1}}{1-\bar{\alpha}_t}\sqrt{\alpha_t}\mathbf{x}_t$, also know as $\mathbf{x}_0$ parameterization, where $\hat{\mathbf{x}}_0$ is an estimation of $\mathbf{x}_0$ from a learned neural network $\hat{\mathbf{x}}_0 = f_\theta(\mathbf{x}_t, t)$.

**Training** For each step during training, a noise-corrupted sample $\mathbf{x}_t$ is sampled and input to the neural network $f_\theta$ to predict $\mathbf{x}_0$. The network is then trained with loss $L_{\text{denoising}}(\theta) = \mathbb{E}_{t,\mathbf{x}_0,\mathbf{x}_t} L(\mathbf{x}_0, \hat{\mathbf{x}}_0)$ with $t$ sampled from 1 to $T$. $L(\cdot, \cdot)$ is a loss function in the space of $\mathbf{x}$. In this work, importance sampling (Nichol and Dhariwal, 2021) is used for time step $t$, where the weight for $t$ is proportional to $\mathbb{E}_{\mathbf{x}_0,\mathbf{x}_t} L(\mathbf{x}_0, \hat{\mathbf{x}}_0)$.

$$\mathbf{x}_t \sim \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I}), \qquad \text{(Sampling)} \qquad \text{(2a)}$$

$$\hat{\mathbf{x}}_0 = f_\theta(t, \mathbf{x}_t), \qquad \text{(Prediction)} \qquad \text{(2b)}$$

$$L_{\text{denoising}}(\theta) = \mathbb{E}_{t,\mathbf{x}_0,\mathbf{x}_t} L(\mathbf{x}_0, \hat{\mathbf{x}}_0), \qquad \text{(loss calculation)} \qquad \text{(2c)}$$

**Inference** At inference time, the denoising starts with a randomly sampled Gaussian noise $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and the data is denoised step-by-step for $t = T, \cdots, 1$:

$$p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I})$$

$$\boldsymbol{\mu}_\theta(\mathbf{x}_t, t) = \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t}\hat{\mathbf{x}}_0 + \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}\sqrt{\alpha_t}\,\mathbf{x}_t$$

Optionally, the variance schedule $\beta_t$ can be down-sampled to reduce the number of inference steps (Nichol and Dhariwal, 2021). A detailed review of DDPM and the loss has been summarised in Appendix A and we refer the readers to Sohl-Dickstein et al. (2015); Ho et al. (2020); Nichol and Dhariwal (2021); Kingma et al. (2021) and other literature for in-depth understanding and derivations.

### 3.2 Diffusion for Segmentation

When applying diffusion models for segmentation, noising and denoising are performed on the segmentation masks. The ground-truth binary mask, where channels correspond to classes that include the background, is denoted by $\mathbf{x}_0$. For the $i$-th pixel/voxel, the value for the $j$-th channel is in 1 if it belongs to class $j$ and $-1$ otherwise. The training process (illustrated in Figure 1) is similar to Equation (2) except that the segmentation network $f_\theta(I, \mathbf{x}_t, t)$ now takes the image $I$ as an additional input for prediction $\hat{\mathbf{x}}_0$. $L(\cdot, \cdot)$ is a weighted sum of cross entropy and foreground-only Dice loss.

$$\mathbf{x}_t \sim \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I}), \qquad \text{(Sampling)} \qquad \text{(3a)}$$

$$\hat{\mathbf{x}}_0 = f_\theta(I, t, \mathbf{x}_t), \qquad \text{(Prediction)} \qquad \text{(3b)}$$

$$L_{\text{denoising}}(\theta) = \mathbb{E}_{t,\mathbf{x}_0,\mathbf{x}_t} L(\mathbf{x}_0, \hat{\mathbf{x}}_0), \qquad \text{(loss calculation)} \qquad \text{(3c)}$$

## 4. Methods

At each training step, the recycling considers a sampled time step $t < T$ and a data sample $\mathbf{x}_0$. First, a noise-corrupted sample $\mathbf{x}_T$ at time step $T$ is sampled, with $\sqrt{\bar{\alpha}_T} \approx 0$. $\mathbf{x}_T$ is fed to the network $f_\theta$ to perform a prediction $\hat{\mathbf{x}}_0 = f_\theta(I, T, \mathbf{x}_T)$. This prediction is then noise-corrupted to generate $\mathbf{x}_t$. A second prediction $\hat{\mathbf{x}}_0 = f_\theta(I, t, \mathbf{x}_t)$ (overriding the $\hat{\mathbf{x}}_0$ for simplicity) is produced and used for loss calculation (see Figure 1). Formally, at each
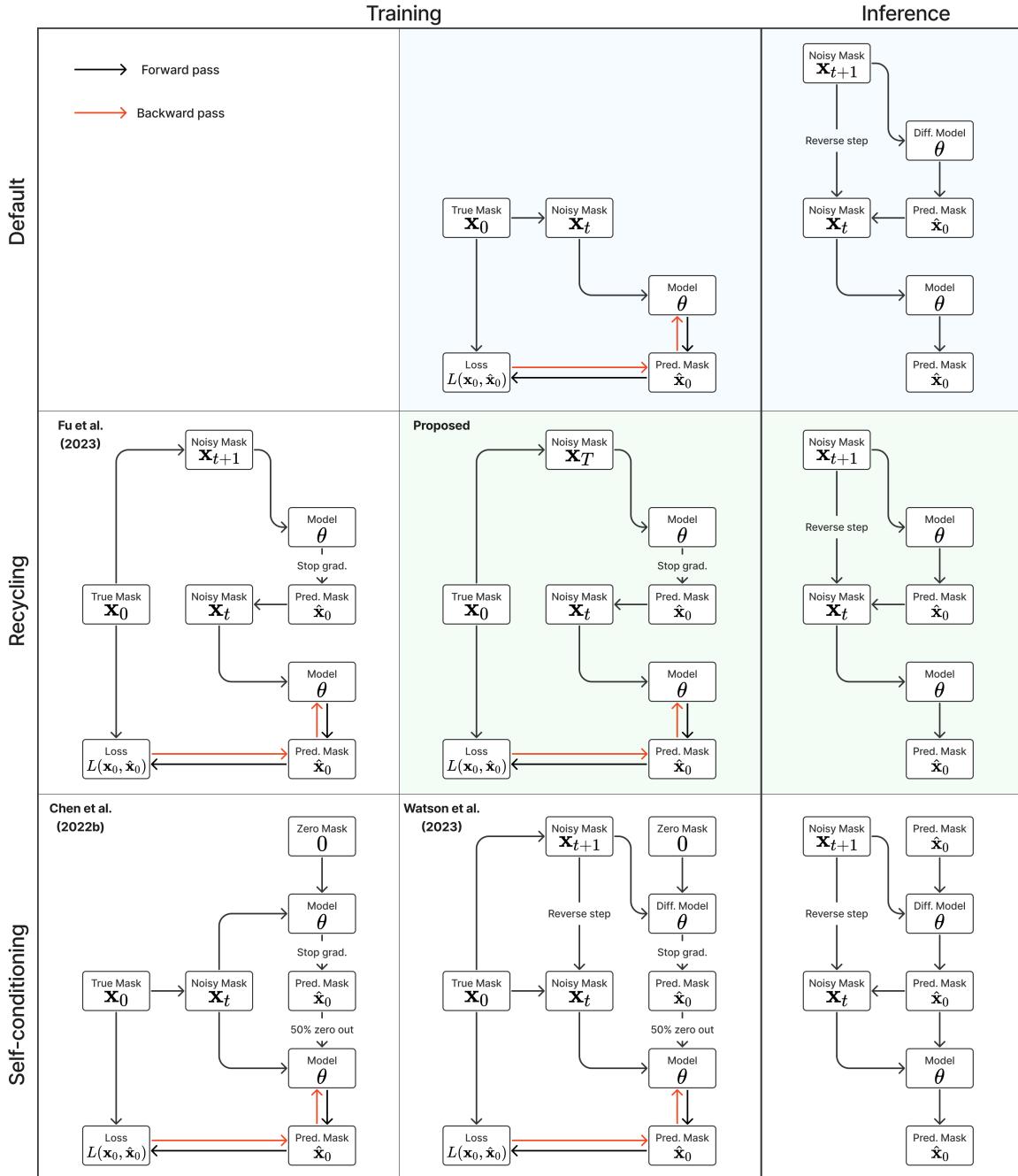
Figure 1: **Illustration of training and inference processes.** The top, middle, and bottom rows show the training and inference steps for default diffusion (highlighted in blue), diffusion with recycling, and diffusion with self-conditioning, respectively. For training, different settings are presented for recycling and self-conditioning. The proposed method is highlighted in green. Notably, recycling shares the same inference steps as default diffusion, while self-conditioning is different as a result of the additional input. "Pred." and "Diff." stands for predicted and diffusion, respectively.

timestep $t$, the proposed recycling (denoted as "Diff. rec. $\mathbf{x}_T$") has the following steps.

$$\mathbf{x}_T \sim \mathcal{N}(\mathbf{x}_T; \sqrt{\bar{\alpha}_T}\, \mathbf{x}_0, (1 - \bar{\alpha}_T)\mathbf{I}), \qquad \text{(rec. } \mathbf{x}_T, \text{ step 1, sampling)} \qquad (4a)$$
$$\hat{\mathbf{x}}_0 = \text{StopGradient}(f_\theta(I, T, \mathbf{x}_T)), \qquad \text{(rec. } \mathbf{x}_T, \text{ step 1, prediction)} \qquad (4b)$$
$$\mathbf{x}_t \sim \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\hat{\mathbf{x}}_0, (1 - \bar{\alpha}_t)\mathbf{I}), \qquad \text{(rec. } \mathbf{x}_T, \text{ step 2, sampling)} \qquad (4c)$$
$$\hat{\mathbf{x}}_0 = f_\theta(I, t, \mathbf{x}_t), \qquad \text{(rec. } \mathbf{x}_T, \text{ step 2, prediction)} \qquad (4d)$$
$$L_{\text{denoising}}(\theta) = \mathbb{E}_{t, \mathbf{x}_0, \mathbf{x}_t}\, L(\mathbf{x}_0, \hat{\mathbf{x}}_0), \qquad \text{(loss calculation)} \qquad (4e)$$

In particular, stop gradient is applied to $\hat{\mathbf{x}}_0$ in the first step to prevent the gradient calculation across two steps, to reduce training time. Optionally, a model with exponential moving averaged weights can be used, but it requires even more memory. Compared to Equation (3), recycling modification only affects training and does not change network architecture. It is independent of the sampling strategy during inference. Therefore, the DDIM sampler can also be used for inference.

The recycling strategy we propose in this work differs from the one introduced in Fu et al. (2023) (denoted as "Diff. rec. $\mathbf{x}_{t+1}$"), illustrated in Figure 1 and the equations below,

$$\mathbf{x}_{t+1} \sim \mathcal{N}(\mathbf{x}_{t+1}; \sqrt{\bar{\alpha}_{t+1}}\, \mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I}), \qquad \text{(rec. } \mathbf{x}_{t+1}, \text{ step 1, sampling)} \qquad (5a)$$
$$\hat{\mathbf{x}}_0 = \text{StopGradient}(f_\theta(I, t+1, \mathbf{x}_{t+1})), \qquad \text{(rec. } \mathbf{x}_{t+1}, \text{ step 1, prediction)} \qquad (5b)$$
$$\mathbf{x}_t \sim \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\hat{\mathbf{x}}_0, (1 - \bar{\alpha}_t)\mathbf{I}), \qquad \text{(rec. } \mathbf{x}_{t+1}, \text{ step 2, sampling)} \qquad (5c)$$
$$\hat{\mathbf{x}}_0 = f_\theta(I, t, \mathbf{x}_t), \qquad \text{(rec. } \mathbf{x}_{t+1}, \text{ step 2, prediction)} \qquad (5d)$$
$$L_{\text{denoising}}(\theta) = \mathbb{E}_{t, \mathbf{x}_0, \mathbf{x}_t}\, L(\mathbf{x}_0, \hat{\mathbf{x}}_0), \qquad \text{(loss calculation)} \qquad (5e)$$

In the new approach ("Diff. rec. $\mathbf{x}_T$"), the first step is consistently executed at the time step $T$ instead of $t+1$ as shown in Equation (4). Compared to $\mathbf{x}_{t+1}$ in Equation (5), $x_T$ is fully noised and contains even less ground truth information during the initial step. Specifically, for a given time step $t$, $\mathbf{x}_t \sim \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\, \mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I})$, which can be reparameterized as $\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\, \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\, \boldsymbol{\epsilon}_t$ with $\boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $\bar{\alpha}_t = \prod_{s=1}^{t} \alpha_s$. In this work, $\alpha_t$ is a monotonically decreasing noise schedule ranging from 0.999 to 0.98 for $t = 1$ to $T$. Correspondingly, $\sqrt{\bar{\alpha}_t}$ monotonically decreases from 0.99995 to 0.00632. $\mathbf{x}_T = \sqrt{\bar{\alpha}_T}\, \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_T}\, \boldsymbol{\epsilon}_T$ with $\sqrt{\bar{\alpha}_T} = 0.00632$ can be considered to contain almost no ground truth information. The information can also be empirically measured by cross entropy and Dice score, and an example is presented in Figure 7 in Appendix D. This seemingly minor modification removes the ground truth information from model inputs, essentially reducing the risk of data leakage and training overfitting. This adaptation guides the model to learn the denoising task based on its initial prediction, rather than ground truth. Consequently, the model can effectively denoise and refine the provided noisy mask, ultimately predicting the ground truth.

Recycling also differs from the self-conditioning methods proposed in Chen et al. (2022b) ("Diff. sc. $\mathbf{x}_t$") and Watson et al. (2023) ("Diff. sc. $\mathbf{x}_{t+1}$"). Although self-conditioning also requests two forward loops during training, it differs from recycling in multiple aspects. First, noisy samples in self-conditioning are always generated based on ground truth $\mathbf{x}_0$, while the second forward step of recycling does not rely on ground truth for noisy sample generation. Second, self-conditioning provides an additional input $\hat{\mathbf{x}}_0$, while recycling does not. Lastly, in self-conditioning, $\hat{\mathbf{x}}_0$ is replaced by zeros with 50% probabilities, while recycling is applied

constantly. The training strategy has been detailed in Figure 1 and Appendix C. For further details, we refer the reader to the reference papers (Chen et al., 2022b; Watson et al., 2023).

## 5. Experiments

### 5.1 Experiment Setting

A range of experiments have been performed in four data sets (Section 5.2) to evaluate the proposed method and the trained models from different aspects.

#### 5.1.1 DIFFUSION TRAINING STRATEGY COMPARISON

First, the proposed recycling training strategy ("Diff. rec. $\mathbf{x}_T$") was compared with standard diffusion models ("Diff.") and other diffusion training strategies that require two forward steps to evaluate the training efficiency with identical network architectures and compute budget. The compared diffusion training strategies include the previously proposed recycling method Fu et al. (2023) ("Diff. rec. $\mathbf{x}_{t+1}$") and two self-conditioning techniques from Chen et al. (2022b) ("Diff. sc. $\mathbf{x}_t$") and Watson et al. (2023) ("Diff. sc. $\mathbf{x}_{t+1}$"). For each trained model using a different strategy, both DDPM and DDIM samplers were evaluated. Importantly, the predictions at each inference step were assessed to study the variation of performance along the inference process.

#### 5.1.2 COMPARISON TO NON-DIFFUSION MODELS

The proposed methods were compared with non-diffusion-based models using identical architectures and the same compute budget. An ensemble model was also evaluated, where the predicted probabilities from the diffusion model and non-diffusion model were averaged. Models' segmentation accuracy was assessed with different granularities: per foreground class or averaged across foreground classes. Balnd-altmann plots were used to analyze the differences between models.

#### 5.1.3 ABLATION STUDIES FOR RECYCLING

Ablation studies were performed, including assessing the performance with different lengths of inference and evaluating the stochasticity across different seeds during inference. Compared to the previous work (Fu et al., 2023), the effectiveness of the Transformer architecture and the change of training noise schedule was evaluated.

#### 5.1.4 EVALUATION METRICS

Different methods were evaluated using binary Dice score (DS) and 95% Hausdorff distance (HD), averaging over foreground classes on the test sets. Dice score is reported in percentage, between 0% and 100%. For Hausdorff distance, the values are in mm for 3D volumes and pixels for 2D images. Paired Student's t-tests with a significance level of $\alpha = 0.05$ were performed on the Dice score to test statistical significance between model performance.

## 5.2 Data

### 5.2.1 Muscle Ultrasound

The data set[1] (Marzola et al., 2021) provides 3910 labeled transverse musculoskeletal ultrasound images, which were split into 2531, 666, and 713 images for training, validation, and test sets, respectively. Images had the shape $480 \times 512$. The predicted masks were post-processed, following Marzola et al. (2021). After filling the holes, multiple morphological operations were performed, including an erosion with a disk of radius 3 pixels, a dilation with a disk of radius 5 pixels, and an opening with a disk of radius 10 pixels. Afterward, only the largest connected component was preserved if the second largest structure was smaller than 75% of the largest one; otherwise, the most superficial (i.e., towards the top of the image) one between the two largest components was preserved. Finally, holes were filled if there were any.

### 5.2.2 Abdominal CT (AMOS)

The data set[2] (Ji et al., 2022) provides 200 and 100 CT image-mask pairs for 15 abdominal organs in training and validation sets. The validation set was randomly split into non-overlapping validation and test sets, with 10 and 90 images, respectively. The images were first resampled with a voxel dimension of $1.5 \times 1.5 \times 5.0$ (mm). HU values were clipped to $[-991, 362]$ and images were normalized so that the intensity had zero mean and unit variance. Lastly, images were center-cropped to shape $192 \times 128 \times 128$. During training, the patch size was $128 \times 128 \times 128$. During inference, the overlap between patches is $64 \times 0 \times 0$, and the predictions on the overlap were averaged.

### 5.2.3 Prostate MR

The data set[3] (Li et al., 2022b) contains 589 T2-weighted image-mask pairs for 8 anatomical structures from 7 institutions. The images were randomly split into non-overlapping training, validation, and test sets, with 411, 14, and 164 images in each split, respectively. The validation split has two images of each institution. The images were resampled with a voxel dimension of $0.75 \times 0.75 \times 2.5$ (mm). Afterward, images were normalized so that the intensity had zero mean and unit variance. Lastly, the images were center-cropped to shape $256 \times 256 \times 48$. During training, the patch size was $256 \times 256 \times 32$. During inference, the overlap between patches was $0 \times 0 \times 16$, and the predictions on the overlap were averaged.

### 5.2.4 Brain MR (BraTS 2021)

The data set[4] (Baid et al., 2021) provides 1251MR segmented mpMRI scans for brain tumour. The data set was randomly split into non-overlapping training, validation, and test sets, with 938, 31, and 282 samples, respectively. The whole tumor mask was generated as foreground class, including GD-enhancing tumor, the peritumoral edematous/invaded tissue, and the necrotic tumor core. Therefore, the task was a binary segmentation. Four modalities

---

1. https://data.mendeley.com/datasets/3jykz7wz8d/1
2. https://zenodo.org/record/7155725#.ZAkbe-zP2rO
3. https://zenodo.org/record/7013610#.ZAkaXuzP2rM
4. https://www.kaggle.com/datasets/dschettler8845/brats-2021-task1

are available, including T1-weighted (T1), post-contrast T1-weighted (T1Gd), T2-weighted (T2), and T2 Fluid Attenuated Inversion Recovery (T2-FLAIR). The voxel dimension was $1.0 \times 1.0 \times 1.0$ (mm). Images were firstly normalized so that the intensity has zero mean and unit variance. Lastly, images were center-cropped to shape $179 \times 219 \times 155$ to remove the common background. During training, the patch size was $128 \times 128 \times 128$. During inference, the overlap between patches was $77 \times 37 \times 101$, and the predictions on the overlap were averaged.
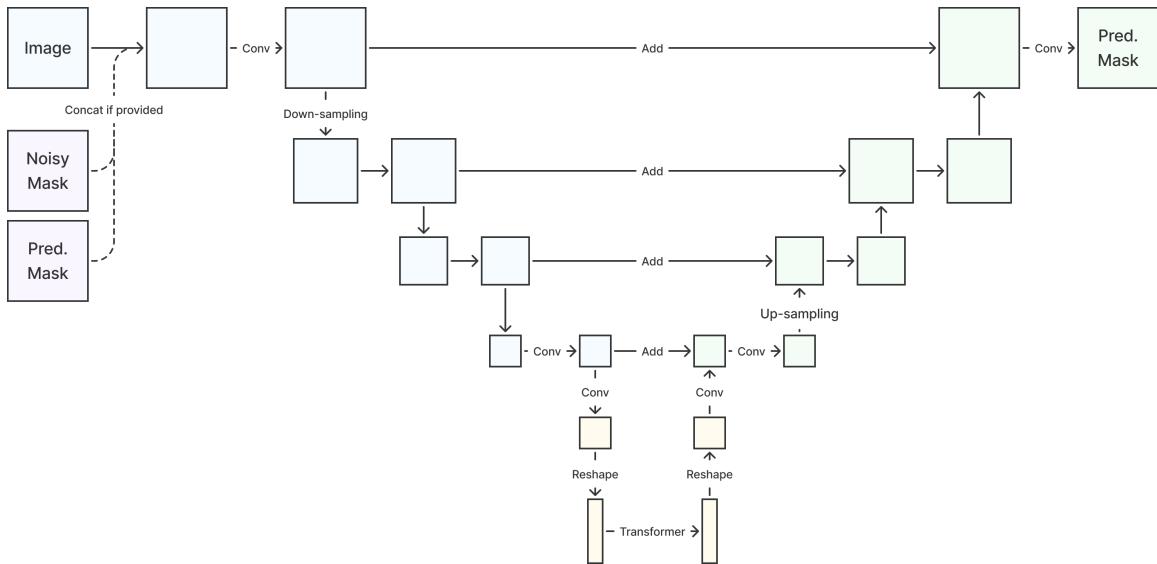
## 5.3 Implementation Details



Figure 2: **Unet architecture for diffusion and non-diffusion models.** The inputs are concatenated when a noisy mask (from diffusion models) or predicted mask (from self-conditioning) is provided. The tensor is enriched with convolution (time-conditioned for diffusion models) and down-sampling layers, then passed into a Transformer with positional encoding, the output is then enriched with convolution and up-sampling layers, and finally, prediction is performed with an additional $1 \times 1$ convolutional layer. "Pred." stands for predicted.

2D and 3D U-net variants with attention mechanisms were used for benchmarking the reference performance from cross-data-set non-diffusion models. The architecture is illustrated in Figure 2. U-nets have four layers with 32, 64, 128, and 256 channels, respectively. The numbers of learnable parameters are summarized in Table 7 in Appendix E. For diffusion-based models, the noise-corrupted masks were concatenated. Time was encoded using sinusoidal positional embedding (Rombach et al., 2022) and used in the convolution layers.

For denoising training, a linear $\beta$ schedule between 0.0001 and 0.02 was used for $T = 1001$ (illustrated in Figure 7 in Appendix D). The segmentation-specific loss function is a weighted sum of cross-entropy and foreground-only Dice loss, with weight 20 and 1 respectively (Kirillov et al., 2023). Random rotation, translation, and scaling were adopted for data augmentation

during training. Training hyper-parameters are listed in Table 6 in Appendix E. Hyper-parameters were configured empirically without extensive tuning.

Models were trained once and checkpoints were saved every 500 step. The checkpoint that had the best mean binary Dice score (without background class) in the validation set was used for the testing. For DDIM, the training was the same as DDPM while both validation and testing were performed using DDIM. The variance schedule was down-sampled to 5 steps (Nichol and Dhariwal, 2021). Experiments were carried out using bfloat16 mixed precision on TPU v3-8, which has $16 \times 8$ GB device memory. However, each device has only 16 GB memory, meaning that the model and data have to fit into 16 GB memory. The JAX-based framework has been released on `https://github.com/mathpluscode/ImgX-DiffSeg`.

## 6. Results and Discussion

### 6.1 Diffusion Training Strategy Comparison

Our proposed recycling method (Diff. rec. $\mathbf{x}_T$) achieved mean Dice scores of 88.23%, 87.45%, 85.54%, and 92.29% on muscle ultrasound, abdominal CT, prostate MR, and brain MR data sets, respectively. These scores marked absolute improvements of 1.63%, 2.20%, 1.93%, and 2.00% over standard diffusion models, respectively. The relative improvements are 1.88%, 2.58%, 2.31%, and 2.22% respectively. Impressively, this novel strategy consistently outperformed the other three training approaches in terms of both Dice score and Hausdorff distance. The observed differences were significant for all data sets in terms of Dice score ($p = 0.003$ for muscle ultrasound and $p < 0.001$ for other data sets). These findings held for both the DDPM and the DDIM samplers, underscoring the wide applicability of the proposed training strategy.

As depicted in Figure 8 in Appendix F.1, standard diffusion models often produce segmentation masks in the last step that are less accurate than the initial prediction. Similar challenges were observed with self-conditioning strategies and previously proposed recycling methods. The newly introduced recycling method was the only approach that improved initial segmentation predictions for more than half of the test images. Moreover, the average performance per step has been visualized in Figure 3, where diffusion models frequently exhibit gradually declining or unstable performance during inference, in terms of both Dice score and Hausdorff distance. It is interesting to observe that often the optimal prediction emerges not at the final step but rather at an intermediate stage. This has been observed in all diffusion models except the newly proposed diffusion model with the innovative recycling method. In the latter case, the quality of segmentation consistently improved or remained stable throughout the inference process, distinguishing it from the observed trend. A qualitative comparison on an example muscle ultrasound image has been illustrated in Figure 4, where the proposed diffusion model was able to refine the segmentation mask progressively. Similar observations have been noted with the DDIM sampler as well, as shown in Figure 9 and Figure 10. This finding aligns with the discussions from Kolbeinsson and Mikolajczyk (2022); Lai et al. (2023) that the diffusion-based segmentation model performance is strongly influenced by the prediction of the initial step. For self-conditioning or the previously proposed recycling, the denoising training relies on the ground truth to varying degrees therefore the diffusion models are trained with ground truth-like initial predictions. However, no ground truth is available during inference, and the distributions of

Table 1: **Diffusion training strategies comparison.** "Diff." represents standard diffusion. "Diff. sc. $\mathbf{x}_t$" and "Diff. sc. $\mathbf{x}_{t+1}$" represents self-conditioning from Chen et al. (2022b) and Watson et al. (2023), respectively. "Diff. rec. $\mathbf{x}_{t+1}$" and "Diff. rec. $\mathbf{x}_T$" represents recycling from Fu et al. (2023) and the proposed recycling in this work, respectively. The best results are in bold and underline indicates the difference to the second best is significant with p-value $< 0.05$.

| Method | DDPM | | DDIM | |
|---|---|---|---|---|
| | DS ↑ | HD ↓ | DS ↑ | HD ↓ |
| Diff. | $86.60 \pm 12.38$ | $41.11 \pm 35.48$ | $86.18 \pm 12.41$ | $42.31 \pm 35.82$ |
| Diff. sc. $\mathbf{x}_t$ | $86.35 \pm 14.14$ | $40.42 \pm 37.53$ | $85.96 \pm 13.78$ | $42.00 \pm 36.76$ |
| Diff. sc. $\mathbf{x}_{t+1}$ | $87.14 \pm 11.48$ | $39.24 \pm 32.83$ | $86.30 \pm 11.49$ | $41.89 \pm 32.72$ |
| Diff. rec. $\mathbf{x}_{t+1}$ | $87.44 \pm 12.39$ | $39.68 \pm 36.21$ | $87.43 \pm 12.25$ | $39.82 \pm 35.39$ |
| Diff. rec. $\mathbf{x}_T$ | $\underline{\mathbf{88.23 \pm 11.69}}$ | $\underline{\mathbf{35.37 \pm 31.79}}$ | $\underline{\mathbf{88.21 \pm 11.70}}$ | $\underline{\mathbf{35.52 \pm 31.91}}$ |

(a) Muscle Ultrasound

| Method | DDPM | | DDIM | |
|---|---|---|---|---|
| | DS ↑ | HD ↓ | DS ↑ | HD ↓ |
| Diff. | $85.25 \pm 5.36$ | $7.12 \pm 3.83$ | $85.59 \pm 5.24$ | $7.13 \pm 3.98$ |
| Diff. sc. $\mathbf{x}_t$ | $86.04 \pm 5.12$ | $7.06 \pm 4.20$ | $85.50 \pm 5.14$ | $7.21 \pm 4.16$ |
| Diff. sc. $\mathbf{x}_{t+1}$ | $85.86 \pm 5.27$ | $6.98 \pm 3.54$ | $85.25 \pm 5.42$ | $7.28 \pm 3.72$ |
| Diff. rec. $\mathbf{x}_{t+1}$ | $86.48 \pm 5.24$ | $6.69 \pm 4.59$ | $86.35 \pm 5.31$ | $6.75 \pm 4.55$ |
| Diff. rec. $\mathbf{x}_T$ | $\underline{\mathbf{87.45 \pm 5.43}}$ | $\mathbf{6.56 \pm 5.44}$ | $\underline{\mathbf{87.45 \pm 5.43}}$ | $\mathbf{6.55 \pm 5.43}$ |

(b) Abdominal CT

| Method | DDPM | | DDIM | |
|---|---|---|---|---|
| | DS ↑ | HD ↓ | DS ↑ | HD ↓ |
| Diff. | $83.61 \pm 4.87$ | $5.10 \pm 2.40$ | $83.11 \pm 4.81$ | $5.00 \pm 2.35$ |
| Diff. sc. $\mathbf{x}_t$ | $83.47 \pm 4.85$ | $5.17 \pm 2.65$ | $82.49 \pm 4.88$ | $5.42 \pm 2.70$ |
| Diff. sc. $\mathbf{x}_{t+1}$ | $83.97 \pm 4.85$ | $4.93 \pm 2.66$ | $83.00 \pm 4.89$ | $5.10 \pm 2.64$ |
| Diff. rec. $\mathbf{x}_{t+1}$ | $84.29 \pm 5.12$ | $4.59 \pm 2.21$ | $84.21 \pm 4.89$ | $4.96 \pm 2.92$ |
| Diff. rec. $\mathbf{x}_T$ | $\mathbf{85.54 \pm 5.20}$ | $\mathbf{4.40 \pm 1.96}$ | $\mathbf{85.54 \pm 5.20}$ | $\underline{\mathbf{4.41 \pm 1.96}}$ |

(c) Prostate MR

| Method | DDPM | | DDIM | |
|---|---|---|---|---|
| | DS ↑ | HD ↓ | DS ↑ | HD ↓ |
| Diff. | $90.29 \pm 12.98$ | $8.46 \pm 15.55$ | $89.94 \pm 13.00$ | $8.55 \pm 15.50$ |
| Diff. sc. $\mathbf{x}_t$ | $90.12 \pm 12.39$ | $9.55 \pm 17.18$ | $89.73 \pm 12.61$ | $9.67 \pm 16.86$ |
| Diff. sc. $\mathbf{x}_{t+1}$ | $89.11 \pm 14.70$ | $9.63 \pm 17.47$ | $88.75 \pm 14.77$ | $9.62 \pm 16.97$ |
| Diff. rec. $\mathbf{x}_{t+1}$ | $86.97 \pm 10.94$ | $9.83 \pm 12.62$ | $84.76 \pm 13.42$ | $12.52 \pm 15.55$ |
| Diff. rec. $\mathbf{x}_T$ | $\underline{\mathbf{92.29 \pm 8.55}}$ | $\mathbf{7.03 \pm 13.48}$ | $\underline{\mathbf{92.29 \pm 8.55}}$ | $\mathbf{7.03 \pm 13.48}$ |

(d) Brain MR

(a) Dice Score for muscle ultrasound

(b) Hausdorff distance for muscle ultrasound

(c) Dice score for abdominal CT

(d) Hausdorff distance for abdominal CT

(e) Dice score for prostate MR

(f) Hausdorff distance for prostate MR

(g) Dice score for brain MR

(h) Hausdorff distance for brain MR

Figure 3: **Segmentation performance per step.** "Diff." represents standard diffusion. "Diff. sc. $\mathbf{x}_t$" and "Diff. sc. $\mathbf{x}_{t+1}$" represents self-conditioning from Chen et al. (2022b) and Watson et al. (2023), respectively. "Diff. rec. $\mathbf{x}_{t+1}$" and "Diff. rec. $\mathbf{x}_T$" represents recycling from Fu et al. (2023) and the proposed recycling in this work, respectively. The sampler is DDPM.
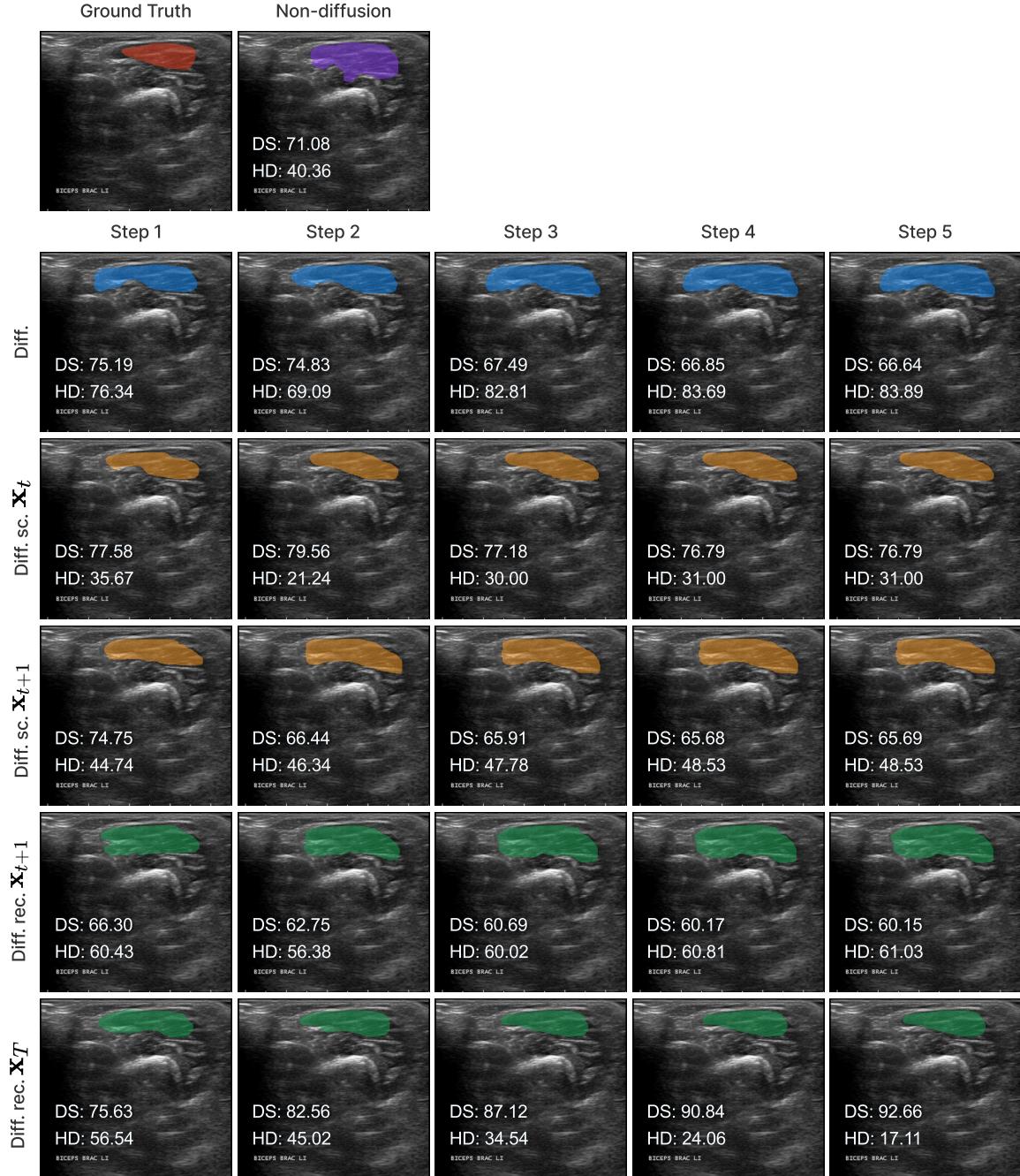
Figure 4: **Diffusion training strategies comparison on muscle ultrasound example.** "Diff." represents standard diffusion. "Diff. sc. $\mathbf{x}_t$" and "Diff. sc. $\mathbf{x}_{t+1}$" represents self-conditioning from Chen et al. (2022b) and Watson et al. (2023), respectively. "Diff. rec. $\mathbf{x}_{t+1}$" and "Diff. rec. $\mathbf{x}_T$" represents recycling from Fu et al. (2023) and the proposed recycling in this work, respectively. The Dice score (DS) and Hausdorff distance (HD) for each sample are labeled at the bottom. While different diffusion models have similar performance on the first step, the proposed method (last row) can refine the segmentation mask.

initial predictions from the trained models are dissimilar from ground truths. This results in an out-of-sample inference and therefore a declining performance. In contrast, the proposed method ingests model predictions for both the training and inference phases without the bias toward ground truth. These observations reaffirm the importance and benefits of harmonizing the training and inference processes. This alignment is crucial to mitigate data leakage, prevent overfitting, and help generalization.
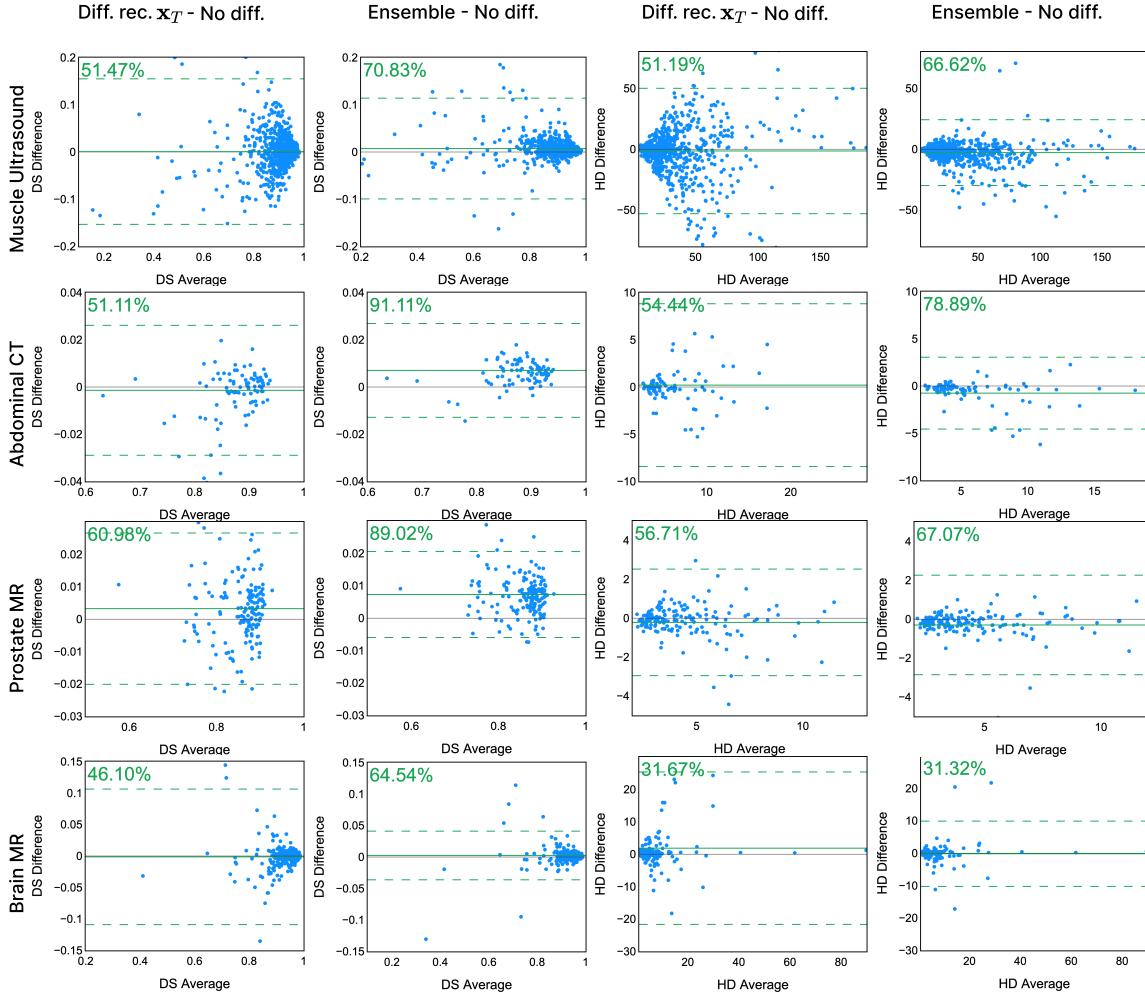
## 6.2 Comparison to Non-diffusion Models

Table 2: **Segmentation performance comparison to non-diffusion models.** "No diff." represents non-diffusion model. "Diff. rec. $\mathbf{x}_T$" represents the diffusion model with proposed recycling. "Ensemble" represents the model averaging the probabilities from "No diff." and "Diff. rec. $\mathbf{x}_T$". The inference sampler is DDPM. The best results are in bold and underline indicates the difference to non-diffusion model is significant with p-value $< 0.05$.
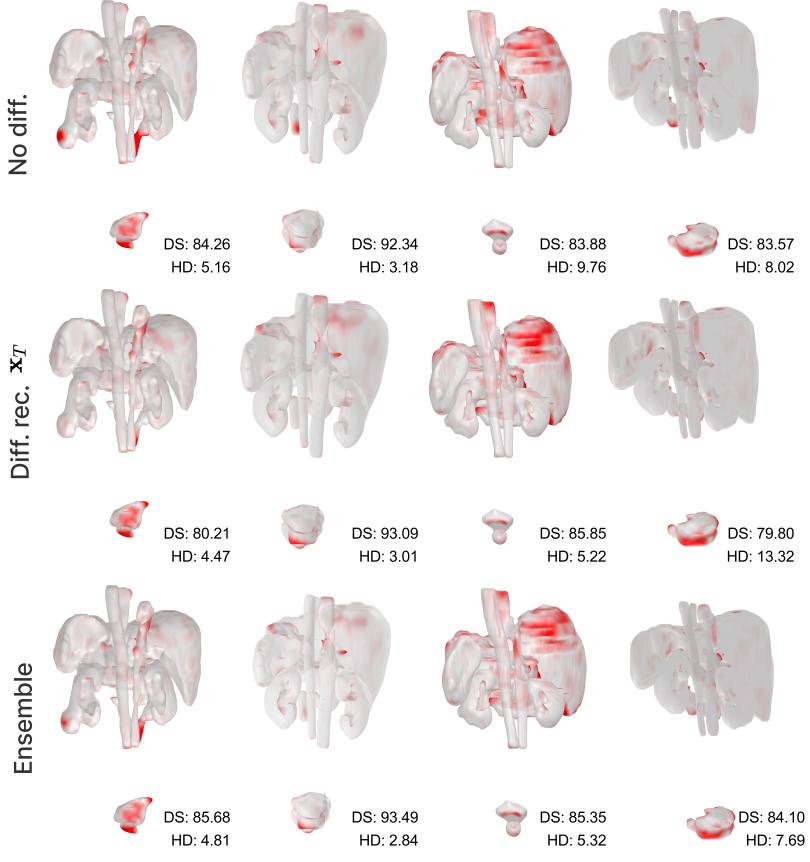
| Data Set | Method | DS $\uparrow$ | HD $\downarrow$ |
|---|---|---|---|
| Muscle Ultrasound | No diff. | $88.15 \pm 10.77$ | $36.86 \pm 30.04$ |
| | Diff. rec. $\mathbf{x}_T$ | $88.23 \pm 11.69$ | $35.37 \pm 31.79$ |
| | Ensemble | $\mathbf{88.88 \pm 10.59}$ | $\mathbf{34.01 \pm 28.75}$ |
| Abdominal CT | No diff. | $87.59 \pm 5.10$ | $6.36 \pm 3.86$ |
| | Diff. rec. $\mathbf{x}_T$ | $87.45 \pm 5.43$ | $6.56 \pm 5.44$ |
| | Ensemble | $\underline{\mathbf{88.29 \pm 5.21}}$ | $\mathbf{5.60 \pm 3.13}$ |
| Prostate MR | No diff. | $85.22 \pm 5.18$ | $4.62 \pm 2.37$ |
| | Diff. rec. $\mathbf{x}_T$ | $\underline{85.54 \pm 5.20}$ | $4.40 \pm 1.96$ |
| | Ensemble | $\underline{\mathbf{85.95 \pm 5.12}}$ | $\mathbf{4.32 \pm 2.01}$ |
| Brain MR | No diff. | $92.43 \pm 9.10$ | $5.20 \pm 9.56$ |
| | Diff. rec. $\mathbf{x}_T$ | $92.29 \pm 8.55$ | $\underline{7.03 \pm 13.48}$ |
| | Ensemble | $\underline{\mathbf{92.67 \pm 8.60}}$ | $\mathbf{5.03 \pm 8.41}$ |

The proposed diffusion models ("Diff. rec. $\mathbf{x}_T$") were compared with their non-diffusion counterparts ("No diff."), where models with identical architectures were trained under the same scheme with the same compute budget. This provides a fair comparison without application-specific adjustments. For diffusion models, the performance with DDPM was selected. As shown in Table 2, The diffusion models yielded similar performance across all data sets. The difference in Dice score is not significant for muscle ultrasound, abdominal CT, and brain MR, but the diffusion model had a higher Dice score for prostate MR ($p = 0.001$). Furthermore, Figure 5 shows that the proposed diffusion model achieved a higher Dice score on more than 50% samples for muscle ultrasound, abdominal CT, and prostate MR data sets. To the best of our knowledge, this is the first time that diffusion models achieved comparable performance against standard non-diffusion-based models with the same architecture and compute budget.
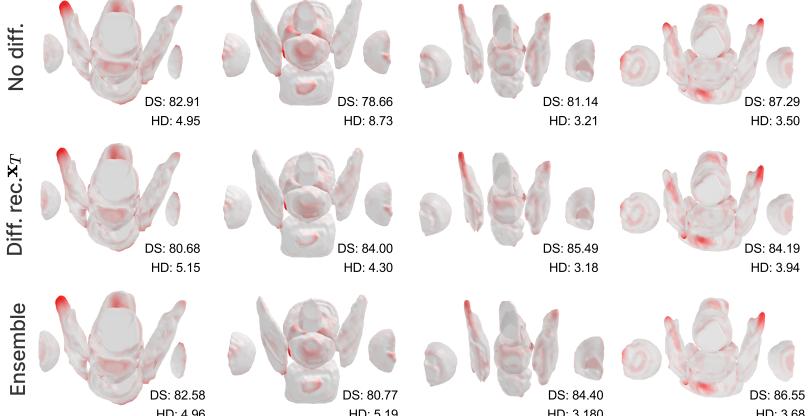
By ensembling these two models via averaging the probabilities, we achieved mean Dice scores of 88.88%, 88.29%, 85.95%, and 92.67% on muscle ultrasound, abdominal CT, prostate MR, and brain MR data sets, respectively. The improvements in Dice score were

Figure 5: **Balnd-altmann plot for comparison of diffusion and ensemble models against non-diffusion models.** "No diff." represents non-diffusion model. "Diff. rec. $\mathbf{x}_T$" represents the diffusion model with proposed recycling. "Ensemble" represents ensembled model by averaging predicted probabilities. The inference sampler is DDPM. DS and HD represents Dice score and Hausdorff distance, respectively. The differences are calculated against non-diffusion models. Positive dice score difference and negative Hausdorff distances indicate improvements. The green solid lines indicates the average difference and the dash lines are mean $\pm 1.96$ standard deviation of the difference. The percentage indicates the number of samples having better performance against non-diffusion baseline. Ensemble models brings an improvement of Dice score for $18.44\% - 40.00\%$ samples across applications.

No diff.

DS: 84.26
HD: 5.16

DS: 92.34
HD: 3.18

DS: 83.88
HD: 9.76

DS: 83.57
HD: 8.02

Diff. rec. $\mathbf{x}_T$

DS: 80.21
HD: 4.47

DS: 93.09
HD: 3.01

DS: 85.85
HD: 5.22

DS: 79.80
HD: 13.32

Ensemble

DS: 85.68
HD: 4.81

DS: 93.49
HD: 2.84

DS: 85.35
HD: 5.32

DS: 84.10
HD: 7.69

(a) Structures in abdominal CT.



No diff.

DS: 82.91
HD: 4.95

DS: 78.66
HD: 8.73

DS: 81.14
HD: 3.21

DS: 87.29
HD: 3.50

Diff. rec. $\mathbf{x}_T$

DS: 80.68
HD: 5.15

DS: 84.00
HD: 4.30

DS: 85.49
HD: 3.18

DS: 84.19
HD: 3.94

Ensemble

DS: 82.58
HD: 4.96

DS: 80.77
HD: 5.19

DS: 84.40
HD: 3.180

DS: 86.55
HD: 3.68

(b) Structures in prostate MR.

Figure 6: **Segmentation error of non-diffusion-based and diffusion-based models.** "No diff." represents non-diffusion model. "Diff. rec. $\mathbf{x}_T$" represents the diffusion model with proposed recycling. "Ensemble" represents ensembled model by averaging predicted probabilities. The ground truth segmentation is visualised. For each point on the surface, the distance to the surface of predicted segmentation is calculated and displayed with red color. The Dice score (DS) and Hausdorff distance (HD) for each sample are labeled at bottom.

significant across all four data sets ($p = 0.037$ for brain MR and $p < 0.001$ for other data sets). Especially, Figure 5 shows that the ensemble model reached a higher Dice score compared to non-diffusion models on 70.83%, 91.11%, 89.02%, and 64.54% samples in the test set for muscle ultrasound, abdominal CT, prostate MR data and brain MR, respectively. These scores marked an absolute increase of 19.36%, 40.00%, 28.04%, and 19.44% compared to the diffusion model alone. Moreover, Abdominal CT and prostate MR are two data sets with multiple classes and their per-class segmentation performances are summarised in Table 8 and Table 9 in Appendix F.1, respectively. Upon comparing diffusion models and non-diffusion models, neither consistently outperformed the other across all classes. However, the ensemble model reached the best performance across all classes and the improvement of Dice score is significant for 13 out of 15 classes in Abdominal CT data and all classes in prostate MR data (all p-values $<= 0.01$, excluding Spleen $p = 0.06$ and Gall bladder $p = 0.876$). Multiple examples have also been visualized in Figure 6 and Figure 11 for the segmentation error.

We highlight that the value of the competitive performance from alternative methods, in particular a different class of generative model-based approaches, is beyond the replacement of current segmentation algorithms for specific potential applications. Our results demonstrate a consistent improvement by combining diffusion and non-diffusion models across applications, even when they yielded a similar performance individually. This is one of the possible potential uses of the proposed improved diffusion models in addition to the well-established non-diffusion baseline. Future research could explore application-specific tuning for further performance improvements.

## 6.3 Ablation Studies

### 6.3.1 Number of sampling steps

Diffusion models were trained using a thousand steps, yet employing the same number of steps for inference can be cost-prohibitive, particularly for processing 3D image volumes. As a result, practical inference commonly utilizes a condensed schedule with a limited number of steps. While this approach reduces computational expenses, the resulting sample quality might be compromised. An ablation study of the numbers of timesteps during inference has therefore been performed across data sets with the proposed recycling-based diffusion model. DDPM sampler was used. The results have been summarised in Table 3. Notably, increasing the number of steps yielded a higher Dice score for the muscle ultrasound dataset but the difference is not significant ($p >= 0.05$). For prostate MR and brain MR data sets, the models maintained almost the same performance regardless of the inference length ($p >= 0.05$). Given that longer inference times and increased device memory usage are associated with more timesteps (e.g. out-of-memory errors were encountered with Abdominal CT at 11 steps), the trade-off between computational resources and performance suggests that a five-step sampling schedule provides the optimal balance.

### 6.3.2 Inference Variance

Different from deterministic models, the inference process of the diffusion model inherently incorporates stochasticity and models a distribution of the segmentation masks. Using the DDPM sampler with the proposed recycling-based diffusion model, the inference on each

Table 3: **Diffusion with different number of sampling steps.** Sampler is DDPM. Diffusion models were trained using the proposed recycling method (Diff. rec. $\mathbf{x}_T$). OOM indicates that out of memory errors were encountered. Best results are in bold.

| Data Set | # Sampling Steps | Dice Score | Hausdorff Distance |
|---|---|---|---|
| Muscle Ultrasound | 2 | 88.01 ± 12.07 | 36.55 ± 32.66 |
| | 5 | 88.23 ± 11.69 | 35.37 ± 31.79 |
| | 11 | **88.30 ± 11.29** | **35.25 ± 30.64** |
| Abdominal CT | 2 | 87.44 ± 5.43 | **6.56 ± 5.42** |
| | 5 | **87.45 ± 5.43** | 6.56 ± 5.44 |
| | 11 | OOM | OOM |
| Prostate MR | 2 | 85.54 ± 5.19 | 4.40 ± 1.96 |
| | 5 | **85.54 ± 5.20** | **4.40 ± 1.96** |
| | 11 | **85.54 ± 5.20** | **4.40 ± 1.96** |
| Brain MR | 2 | 92.29 ± 8.54 | 7.03 ± 13.47 |
| | 5 | **92.29 ± 8.55** | 7.03 ± 13.48 |
| | 11 | 92.29 ± 8.57 | **7.02 ± 13.48** |

Table 4: **Diffusion model performance across different inference seeds.** For each sample, the maximum difference ($\Delta$) across five random seeds is calculated. The average across all samples is reported.

| Data Set | Mean $\Delta$ Dice Score | | | | |
|---|---|---|---|---|---|
| | Step 1 | Step 2 | Step 3 | Step 4 | Step 5 |
| Muscle Ultrasound | 0.0212 | 0.0165 | 0.0122 | 0.0081 | 0.0051 |
| Abdominal CT | 0.0009 | 0.0010 | 0.0009 | 0.0008 | 0.0004 |
| Prostate MR | 0.0004 | 0.0004 | 0.0004 | 0.0004 | 0.0002 |
| Brain MR | 0.0005 | 0.0005 | 0.0005 | 0.0003 | 0.0001 |

| Data Set | Mean $\Delta$ Hausdorff Distance | | | | |
|---|---|---|---|---|---|
| | Step 1 | Step 2 | Step 3 | Step 4 | Step 5 |
| Muscle Ultrasound | 10.0582 | 7.0020 | 4.7440 | 3.1758 | 1.8164 |
| Abdominal CT | 0.1481 | 0.1339 | 0.1221 | 0.0751 | 0.0673 |
| Prostate MR | 0.0447 | 0.0426 | 0.0499 | 0.0431 | 0.0209 |
| Brain MR | 0.0758 | 0.0779 | 0.0678 | 0.0616 | 0.0197 |

data set has been repeated with five different random seeds. Consequently, each sample has five distinct predicted masks. The maximum differences across five predictions were computed for the Dice score and Hausdorff distance, denoted by $\Delta$ Dice score and $\Delta$ Hausdorff distance, respectively. The average of this performance difference across all samples in the test set has been reported in Table 4 for all data sets. While the magnitude of the average difference (mean $\Delta$) varies across data sets, a common trend was observed where mean $\Delta$

diminished during the sampling process for both metrics. In other words, despite different initial predictions, the model's predictions gradually converge as the difference across seeds decreases. Moreover, the relative magnitude of the mean $\Delta$ Hausdorff distance (e.g. 1.82 at the last step for muscle ultrasound represents around 5% fluctuation compared to 35.37, the mean Hausdorff distance to ground truth) was larger than the relative magnitude for Dice score (e.g. 0.0051 at the last step for muscle ultrasound was around 0.006% fluctuation compared to 88.23 the mean Hausdorff distance to ground truth). We hypothesize that the variation among predictions may predominantly revolve around local refinements in mask boundaries, as opposed to significant alterations like expansion or contraction of foreground areas. This may open a direction for further improving diffusion training: instead of performing independent noising per pixel/voxel results in fragmented and disjointed masks, the noising can be morphology-informed such that the noise-corrupted masks expand or contract the foreground with continuous boundaries.

### 6.3.3 Transformer

Table 5: **Segmentation performance without Transformer.** "No diff." represents non-diffusion model. "Diff. rec. $\mathbf{x}_T$" represents the diffusion model with proposed recycling. The inference sampler is DDPM. The best results are in bold and underline indicates the difference to non-diffusion model is significant with p-value < 0.05.

| Data Set | Method | Transformer | DS ↑ | HD ↓ |
|---|---|---|---|---|
| Muscle US | No diff. | | $86.66 \pm 13.16$ | $45.01 \pm 38.86$ |
| | | ✓ | $\mathbf{88.15 \pm 10.77}$ | $\mathbf{36.86 \pm 30.04}$ |
| | Diff. rec. $\mathbf{x}_T$ | | $\mathbf{88.36 \pm 12.60}$ | $35.67 \pm 34.12$ |
| | | ✓ | $88.23 \pm 11.69$ | $\mathbf{35.37 \pm 31.79}$ |
| Abdominal CT | No diff. | | $87.48 \pm 5.02$ | $6.63 \pm 4.03$ |
| | | ✓ | $\mathbf{87.59 \pm 5.10}$ | $\mathbf{6.36 \pm 3.86}$ |
| | Diff. rec. $\mathbf{x}_T$ | | $86.89 \pm 5.49$ | $6.91 \pm 4.35$ |
| | | ✓ | $\underline{\mathbf{87.45 \pm 5.43}}$ | $\mathbf{6.56 \pm 5.44}$ |
| Prostate MR | No diff. | | $84.82 \pm 5.69$ | $\mathbf{4.55 \pm 2.17}$ |
| | | ✓ | $\underline{85.22 \pm 5.18}$ | $4.62 \pm 2.37$ |
| | Diff. rec. $\mathbf{x}_T$ | | $\mathbf{85.63 \pm 5.19}$ | $4.59 \pm 2.71$ |
| | | ✓ | $85.54 \pm 5.20$ | $\mathbf{4.40 \pm 1.96}$ |
| Brain MR | No diff. | | $92.03 \pm 9.67$ | $5.29 \pm 8.53$ |
| | | ✓ | $\underline{\mathbf{92.43 \pm 9.10}}$ | $\mathbf{5.20 \pm 9.56}$ |
| | Diff. rec. $\mathbf{x}_T$ | | $92.04 \pm 9.47$ | $7.25 \pm 13.76$ |
| | | ✓ | $\mathbf{92.29 \pm 8.55}$ | $\mathbf{7.03 \pm 13.48}$ |

Compared to Fu et al. (2023), the model includes a Transformer layer at the bottom encoder of U-net. This component has one layer representing 16% and 6% of the trainable parameters for 2D and 3D networks, correspondingly (see Table 7 in Appendix E). An ablation study has been performed for the proposed recycling approach and non-diffusion models. The results have been summarised in Table 5. For non-diffusion models, the addition

of the Transformer component brought improvement in Dice score across all applications ($p < 0.001$ for muscle ultrasound; $p >= 0.05$ for abdominal CT; $p = 0.001$ for prostate MR; and $p = 0.0178$ for brain MR), making this architecture the stronger reference model. For diffusion, significantly higher Dice scores have been observed for abdominal CT data ($p < 0.001$), and the differences were not significant for other applications ($p >= 0.05$).

### 6.3.4 LENGTH OF TRAINING NOISE SCHEDULE

It's worth noting that Fu et al. (2023) recommended incorporating a shortened variance schedule during training, mirroring that used during inference, in addition to the recycling technique. This modification resulted in enhanced performance for every training strategy on the muscle ultrasound data set (as detailed in Table 10a). However, this adaptation did not yield enhancements for the proposed training strategies ("Diff. rec. $\mathbf{x}_T$") in the abdominal CT data set (as depicted in Table 10b). Moreover, not all differences observed were statistically significant. This may suggest that the advantage of the modified training variance schedule may be application-dependent and sensitive to the change of model architectures and hyper-parameters. In this work, the variance schedule was maintained at 1001 steps.

## 7. Conclusion

In this research, we have proposed a novel training strategy for diffusion-based segmentation models. The aim is to remove the dependency on ground truth masks during denoising training. In contrast to the standard diffusion-based segmentation models and those employing self-conditioning or alternative recycling techniques, our approach consistently maintains or enhances segmentation performance throughout progressive inference processes. Through extensive experiments across four medical imaging data sets with different dimensionalities and modalities, we demonstrated statistically significant improvement against all diffusion baseline models for both DDPM and DDIM samplers. Our analysis for the first time identified a common limitation of existing diffusion model training for segmentation tasks. The use of ground truth data for denoising training leads to data leakage. By utilizing the model's prediction at the initial step instead, we align the training process with inference procedures, effectively reducing over-fitting and promoting better generalization. While existing diffusion models underperformed non-diffusion-based segmentation model baselines, our innovative recycling training strategies effectively bridged the performance gap. This enhancement allowed diffusion models to attain comparable performance levels. To the best of our knowledge, this is the first time diffusion models have achieved such parity in performance while maintaining identical architecture and compute budget. By ensembling the diffusion and non-diffusion models, constant and significant improvements have been observed across all data sets, demonstrating one of its potential values. Nevertheless, challenges remain on the road to advancing diffusion-based segmentation models further. Future work could explore discrete diffusion models that are tailored for categorical data or implement diffusion in latent space to further reduce compute costs. Although the presented experimental results primarily demonstrated methodological development, the fact that these were obtained on four large clinical data sets represents a promising step toward real-world applications. We would like to argue the potential importance of the reported development, which may lead to better clinical outcomes and improved patient care in respective applications. For example,

avoiding surrounding healthy structures may be sensitive to their localization in planning imaging, in both the abdominal CT and prostate MR tasks. This sensitivity can be high and nonlinear therefore arguably a perceived marginal improvement might benefit those with smaller targets, such as those in liver resection and focal therapy of prostate cancer, or highly variable ultrasound imaging guidance.

## Acknowledgments

## Ethical Standards

The work follows appropriate ethical standards in conducting research and writing the manuscript, following all applicable laws and regulations regarding treatment of animals or human subjects.

## Conflicts of Interest

We declare we do not have conflicts of interest.

## References

Tomer Amit, Eliya Nachmani, Tal Shaharbany, and Lior Wolf. Segdiff: Image segmentation with diffusion probabilistic models. arXiv preprint arXiv:2112.00390, 2021.

Jacob Austin, Daniel D Johnson, Jonathan Ho, Daniel Tarlow, and Rianne Van Den Berg. Structured denoising diffusion models in discrete state-spaces. Advances in Neural Information Processing Systems, 34:17981–17993, 2021.

Ujjwal Baid, Satyam Ghodasara, Suyash Mohan, Michel Bilello, Evan Calabrese, Errol Colak, Keyvan Farahani, Jayashree Kalpathy-Cramer, Felipe C Kitamura, Sarthak Pati, et al. The rsna-asnr-miccai brats 2021 benchmark on brain tumor segmentation and radiogenomic classification. arXiv preprint arXiv:2107.02314, 2021.

Florentin Bieder, Julia Wolleb, Alicia Durrer, Robin Sandkühler, and Philippe C Cattin. Diffusion models for memory-efficient processing of 3d medical images. arXiv preprint arXiv:2303.15288, 2023.

Tao Chen, Chenhui Wang, and Hongming Shan. Berdiff: Conditional bernoulli diffusion model for medical image segmentation. arXiv preprint arXiv:2304.04429, 2023.

Ting Chen, Lala Li, Saurabh Saxena, Geoffrey Hinton, and David J Fleet. A generalist framework for panoptic segmentation of images and videos. arXiv preprint arXiv:2210.06366, 2022a.

Ting Chen, Ruixiang Zhang, and Geoffrey Hinton. Analog bits: Generating discrete data using diffusion models with self-conditioning. arXiv preprint arXiv:2208.04202, 2022b.

Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. Advances in Neural Information Processing Systems, 34:8780–8794, 2021.

Zolnamar Dorjsembe, Sodtavilan Odonchimed, and Furen Xiao. Three-dimensional medical image synthesis with denoising diffusion probabilistic models. In Medical Imaging with Deep Learning, 2022.

Yunguan Fu, Yiwen Li, Shaheer U Saeed, Matthew J Clarkson, and Yipeng Hu. Importance of aligning training strategy with evaluation for diffusion models in 3d multiclass segmentation. arXiv preprint arXiv:2303.06040, 2023.

Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. Communications of the ACM, 63(11):139–144, 2020.

Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, Bo Zhang, Dongdong Chen, Lu Yuan, and Baining Guo. Vector quantized diffusion model for text-to-image synthesis. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 10696–10706, 2022.

Xutao Guo, Yanwu Yang, Chenfei Ye, Shang Lu, Yang Xiang, and Ting Ma. Accelerating diffusion models via pre-segmentation diffusion sampling for medical image segmentation. arXiv preprint arXiv:2210.17408, 2022.

Xiaochuang Han, Sachin Kumar, and Yulia Tsvetkov. Ssd-lm: Semi-autoregressive simplex-based diffusion language model for text generation and modular control. arXiv preprint arXiv:2210.17432, 2022.

Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. arXiv preprint arXiv:2207.12598, 2022.

Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. Advances in Neural Information Processing Systems, 33:6840–6851, 2020.

Emiel Hoogeboom, Didrik Nielsen, Priyank Jaini, Patrick Forré, and Max Welling. Argmax flows and multinomial diffusion: Learning categorical distributions. Advances in Neural Information Processing Systems, 34:12454–12465, 2021.

Dewei Hu, Yuankai K Tao, and Ipek Oguz. Unsupervised denoising of retinal oct with diffusion probabilistic model. In Medical Imaging 2022: Image Processing, volume 12032, pages 25–34. SPIE, 2022.

Aapo Hyvärinen and Peter Dayan. Estimation of non-normalized statistical models by score matching. Journal of Machine Learning Research, 6(4), 2005.

Yuanfeng Ji, Haotian Bai, Jie Yang, Chongjian Ge, Ye Zhu, Ruimao Zhang, Zhen Li, Lingyan Zhang, Wanling Ma, Xiang Wan, et al. Amos: A large-scale abdominal multi-organ benchmark for versatile medical image segmentation. arXiv preprint arXiv:2206.08023, 2022.

Amirhossein Kazerouni, Ehsan Khodapanah Aghdam, Moein Heidari, Reza Azad, Mohsen Fayyaz, Ilker Hacihaliloglu, and Dorit Merhof. Diffusion models in medical imaging: A comprehensive survey. Medical Image Analysis, page 102846, 2023.

Firas Khader, Gustav Mueller-Franzes, Soroosh Tayebi Arasteh, Tianyu Han, Christoph Haarburger, Maximilian Schulze-Hagen, Philipp Schad, Sandy Engelhardt, Bettina Baessler, Sebastian Foersch, et al. Medical diffusion–denoising diffusion probabilistic models for 3d medical image generation. arXiv preprint arXiv:2211.03364, 2022.

Boah Kim, Inhwa Han, and Jong Chul Ye. Diffusemorph: unsupervised deformable image registration using diffusion model. In European Conference on Computer Vision, pages 347–364. Springer, 2022.

Diederik Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. Advances in neural information processing systems, 34:21696–21707, 2021.

Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. arXiv preprint arXiv:2304.02643, 2023.

Benedikt Kolbeinsson and Krystian Mikolajczyk. Multi-class segmentation from aerial views using recursive noise diffusion. arXiv preprint arXiv:2212.00787, 2022.

Zeqiang Lai, Yuchen Duan, Jifeng Dai, Ziheng Li, Ying Fu, Hongsheng Li, Yu Qiao, and Wenhai Wang. Denoising diffusion semantic segmentation with mask prior modeling. arXiv preprint arXiv:2306.01721, 2023.

Xiang Li, John Thickstun, Ishaan Gulrajani, Percy S Liang, and Tatsunori B Hashimoto. Diffusion-lm improves controllable text generation. Advances in Neural Information Processing Systems, 35:4328–4343, 2022a.

Yiwen Li, Yunguan Fu, Iani Gayo, Qianye Yang, Zhe Min, Shaheer Saeed, Wen Yan, Yipei Wang, J Alison Noble, Mark Emberton, et al. Prototypical few-shot segmentation for cross-institution male pelvic structures with spatial registration. arXiv preprint arXiv:2209.05160, 2022b.

Luping Liu, Yi Ren, Zhijie Lin, and Zhou Zhao. Pseudo numerical methods for diffusion models on manifolds. arXiv preprint arXiv:2202.09778, 2022.

Zhaoyang Lyu, Xudong Xu, Ceyuan Yang, Dahua Lin, and Bo Dai. Accelerating diffusion models via early stop of the diffusion process. arXiv preprint arXiv:2205.12524, 2022.

Francesco Marzola, Nens van Alfen, Jonne Doorduin, and Kristen M Meiburger. Deep learning segmentation of transverse musculoskeletal ultrasound images for neuromuscular disease assessment. Computers in Biology and Medicine, 135:104623, 2021.

Puria Azadi Moghadam, Sanne Van Dalen, Karina C Martin, Jochen Lennerz, Stephen Yip, Hossein Farahani, and Ali Bashashati. A morphology focused diffusion probabilistic model for synthesis of histopathology images. In Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision, pages 2000–2009, 2023.

Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In International Conference on Machine Learning, pages 8162–8171. PMLR, 2021.

Walter HL Pinaya, Mark S Graham, Robert Gray, Pedro F Da Costa, Petru-Daniel Tudosiu, Paul Wright, Yee H Mah, Andrew D MacKinnon, James T Teo, Rolf Jager, et al. Fast unsupervised brain anomaly detection and segmentation with diffusion models. In International Conference on Medical Image Computing and Computer-Assisted Intervention, pages 705–714. Springer, 2022a.

Walter HL Pinaya, Petru-Daniel Tudosiu, Jessica Dafflon, Pedro F da Costa, Virginia Fernandez, Parashkev Nachev, Sebastien Ourselin, and M Jorge Cardoso. Brain imaging generation with latent diffusion models. arXiv preprint arXiv:2209.07162, 2022b.

Aimon Rahman, Jeya Maria Jose Valanarasu, Ilker Hacihaliloglu, and Vishal M Patel. Ambiguous medical image segmentation using diffusion models. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 11536–11546, 2023.

Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. arXiv preprint arXiv:2204.06125, 2022.

Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 10684–10695, 2022.

Shaheer U. Saeed, Tom Syer, Wen Yan, Qianye Yang, Mark Emberton, Shonit Punwani, Matthew J. Clarkson, Dean C. Barratt, and Yipeng Hu. Bi-parametric prostate mr image synthesis using pathology and sequence-conditioned stable diffusion. arXiv preprint arXiv:2303.02094, 2023.

Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In International Conference on Machine Learning, pages 2256–2265. PMLR, 2015.

Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. arXiv preprint arXiv:2010.02502, 2020a.

Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. Advances in neural information processing systems, 32, 2019.

Yang Song and Stefano Ermon. Improved techniques for training score-based generative models. Advances in neural information processing systems, 33:12438–12448, 2020.

Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. arXiv preprint arXiv:2011.13456, 2020b.

Robin Strudel, Corentin Tallec, Florent Altché, Yilun Du, Yaroslav Ganin, Arthur Mensch, Will Grathwohl, Nikolay Savinov, Sander Dieleman, Laurent Sifre, et al. Self-conditioned embedding diffusion for text generation. arXiv preprint arXiv:2211.04236, 2022.

Pascal Vincent. A connection between score matching and denoising autoencoders. Neural computation, 23(7):1661–1674, 2011.

Hefeng Wang, Jiale Cao, Rao Muhammad Anwer, Jin Xie, Fahad Shahbaz Khan, and Yanwei Pang. Dformer: Diffusion-guided transformer for universal image segmentation. arXiv preprint arXiv:2306.03437, 2023.

Risheng Wang, Tao Lei, Ruixia Cui, Bingtao Zhang, Hongying Meng, and Asoke K Nandi. Medical image segmentation using deep learning: A survey. IET Image Processing, 16(5): 1243–1267, 2022.

Joseph L Watson, David Juergens, Nathaniel R Bennett, Brian L Trippe, Jason Yim, Helen E Eisenach, Woody Ahern, Andrew J Borst, Robert J Ragotte, Lukas F Milles, et al. De novo design of protein structure and function with rfdiffusion. Nature, pages 1–3, 2023.

Julia Wolleb, Florentin Bieder, Robin Sandkühler, and Philippe C Cattin. Diffusion models for medical anomaly detection. In International Conference on Medical image computing and computer-assisted intervention, pages 35–45. Springer, 2022a.

Julia Wolleb, Robin Sandkühler, Florentin Bieder, Philippe Valmaggia, and Philippe C Cattin. Diffusion models for implicit image segmentation ensembles. In International Conference on Medical Imaging with Deep Learning, pages 1336–1348. PMLR, 2022b.

Junde Wu, Huihui Fang, Yu Zhang, Yehui Yang, and Yanwu Xu. Medsegdiff: Medical image segmentation with diffusion probabilistic model. arXiv preprint arXiv:2211.00611, 2022.

Junde Wu, Rao Fu, Huihui Fang, Yu Zhang, and Yanwu Xu. Medsegdiff-v2: Diffusion based medical image segmentation with transformer. arXiv preprint arXiv:2301.11798, 2023.

Zhaohu Xing, Liang Wan, Huazhu Fu, Guang Yang, and Lei Zhu. Diff-unet: A diffusion embedded network for volumetric segmentation. arXiv preprint arXiv:2303.10326, 2023.

Yijun Yang, Huazhu Fu, Angelica Aviles-Rivero, Carola-Bibiane Schönlieb, and Lei Zhu. Diffmic: Dual-guidance diffusion network for medical image classification. arXiv preprint arXiv:2303.10610, 2023.

Sean I Young, Adrian V Dalca, Enzo Ferrante, Polina Golland, Bruce Fischl, and Juan Eugenio Iglesias. Sud: Supervision by denoising for medical image segmentation. arXiv preprint arXiv:2202.02952, 2022.

Lukas Zbinden, Lars Doorenbos, Theodoros Pissas, Raphael Sznitman, and Pablo Márquez-Neila. Stochastic segmentation with conditional categorical diffusion models. arXiv preprint arXiv:2303.08888, 2023.

Huangjie Zheng, Pengcheng He, Weizhu Chen, and Mingyuan Zhou. Truncated diffusion probabilistic models. stat, 1050:7, 2022.

## Appendix A. Denoising Diffusion Probabilistic Model

We review the formulation of denoising diffusion probabilistic models (DDPM) from Sohl-Dickstein et al. (2015); Ho et al. (2020); Nichol and Dhariwal (2021).

### A.1 Definition

$$\mathbf{x}_T \rightleftharpoons \cdots \rightleftharpoons \mathbf{x}_t \underset{q(\mathbf{x}_t \mid \mathbf{x}_{t-1})}{\overset{p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t)}{\rightleftharpoons}} \mathbf{x}_{t-1} \rightleftharpoons \cdots \rightleftharpoons \mathbf{x}_0$$

Consider a continuous *diffusion* process (also named *forward* process or *noising* process): given a data point $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ in $\mathbb{R}^D$, we add noise to $\mathbf{x}_t$ for $t = 1, \cdots, T$ with the following multivariate normal distribution:

$$q(\mathbf{x}_t \mid \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\, \mathbf{x}_{t-1}, \beta_t \mathbf{I})$$

where $\beta_t \in [0, 1]$ is a variance schedule. Given sufficiently large $T$ and a well-defined variance schedule, the distribution of $\mathbf{x}_T$ approximates an isotropic multivariate normal distribution.

$$q(\mathbf{x}_t \mid \mathbf{x}_0) \to \mathcal{N}(\mathbf{x}_t; \mathbf{0}, \mathbf{I})$$

Therefore, we can define a *reverse* process (also named *denoising* process): given a sample $\mathbf{x}_T \sim \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$, we denoise the data using neural networks $\boldsymbol{\mu}_\theta : \mathbb{R}^D \to \mathbb{R}^D$ and $\boldsymbol{\Sigma}_\theta : \mathbb{R}^D \to \mathbb{R}^{D \times D}$ as follows:

$$p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t))$$

In this work, an isotropic variance is assumed with $\boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t) = \sigma_t^2 \mathbf{I}$, such that

$$p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \sigma_t^2 \mathbf{I})$$

### A.2 Variational Lower Bound

Consider $\mathbf{z} = \mathbf{x}_{1:T} \mid \mathbf{x}_0$ as latent variables for $\mathbf{x}_0$, we can derive the variational lower bound (VLB) as follows:

$$
\begin{aligned}
\log p_\theta(\mathbf{x}_0) =& D_{\mathrm{KL}}(q(\mathbf{z}) \parallel p_\theta(\mathbf{z} \mid \mathbf{x}_0)) + \mathbb{E}_{q(\mathbf{z})}\left[\log \frac{p_\theta(\mathbf{x}_0, \mathbf{z})}{q(\mathbf{z})}\right] \\
\geq& \mathbb{E}_{q(\mathbf{z})}\left[\log \frac{p_\theta(\mathbf{x}_0, \mathbf{z})}{q(\mathbf{z})}\right] \\
=& -\left(\mathbb{E}_{q(\mathbf{x}_1 \mid \mathbf{x}_0)} L_0 + \sum_{t=2}^{T} \mathbb{E}_{q(\mathbf{x}_t \mid \mathbf{x}_0)} L_{t-1} + L_T\right)
\end{aligned}
$$

where

$$
\begin{array}{ll}
L_0 = -\log p_\theta(\mathbf{x}_0 \mid \mathbf{x}_1) & \text{(reconstruction loss)} \\
L_{t-1} = D_{\mathrm{KL}}(q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) \| p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t)) & \text{(diffussion loss)} \\
L_T = D_{\mathrm{KL}}(q(\mathbf{x}_T \mid \mathbf{x}_0)) \| p_\theta(\mathbf{x}_T)). & \text{(prior loss)}
\end{array}
$$

## A.3 Diffusion Loss

In particular, we can derive the closed form $L_{t-1}$ with

$$q(\mathbf{x}_t \mid \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\,\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I})$$
$$q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\boldsymbol{\mu}}(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t \mathbf{I})$$

where

$$\alpha_t = 1 - \beta_t$$
$$\bar{\alpha}_t = \prod_{s=1}^{t} \alpha_s$$
$$\tilde{\boldsymbol{\mu}}(\mathbf{x}_t, \mathbf{x}_0) = \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t}\,\mathbf{x}_0 + \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}\sqrt{\alpha_t}\,\mathbf{x}_t, \qquad (6)$$
$$\tilde{\beta}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}\beta_t.$$

### A.3.1 NOISE PREDICTION LOSS ($\epsilon$-PARAMETERIZATION)S

Consider the reparameterization in Ho et al. (2020),

$$\mathbf{x}_t(\mathbf{x}_0, \boldsymbol{\epsilon}) = \sqrt{\bar{\alpha}_t}\,\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\,\boldsymbol{\epsilon}$$
$$\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{1 - \bar{\alpha}_t}}\,\mathbf{x}_t - \frac{\sqrt{\bar{\alpha}_t}}{\sqrt{1 - \bar{\alpha}_t}}\,\mathbf{x}_0$$
$$\boldsymbol{\mu}_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}}(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}}\,\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t))$$

We can derive a closed form of $L_{t-1}$

$$L_{t-1}(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{2\sigma_t^2}\frac{\beta_t^2}{\alpha_t(1 - \bar{\alpha}_t)}\|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta\|_2^2 + C$$

If $\sigma_t^2 = \tilde{\beta}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}\beta_t$, using the signal-to-noise ratio (SNR) defined in Kingma et al. (2021), $\mathrm{SNR}(t) = \frac{\bar{\alpha}_t}{1 - \bar{\alpha}_t}$, the loss can be derived as

$$L_{t-1}(\mathbf{x}_t, \mathbf{x}_0) = (\frac{\mathrm{SNR}(t-1)}{\mathrm{SNR}(t)} - 1)\|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta\|_2^2 + C$$

### A.3.2 SAMPLE PREDICTION LOSS ($\mathbf{x}_0$-PARAMETERIZATION)

Similar to Eq. (6), consider the parameterization (Kingma et al., 2021),

$$\boldsymbol{\mu}_\theta(\mathbf{x}_t, t) = \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t}\,\mathbf{x}_{0,\theta} + \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}\sqrt{\alpha_t}\,\mathbf{x}_t$$

We can derive a closed form of $L_{t-1}$

$$L_{t-1}(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{2\sigma_t^2}\frac{\bar{\alpha}_{t-1}\beta_t^2}{(1 - \bar{\alpha}_t)^2}\|\mathbf{x}_{0,\theta} - \mathbf{x}_0\|_2^2 + C.$$

If $\sigma_t^2 = \tilde{\beta}_t = \frac{1-\bar{\alpha}_{t-1}}{1-\bar{\alpha}_t}\beta_t$, using the signal-to-noise ratio (SNR) defined in Kingma et al. (2021), $\mathrm{SNR}(t) = \frac{\bar{\alpha}_t}{1-\bar{\alpha}_t}$, the loss can be derived as

$$L_{t-1}(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{2}(\mathrm{SNR}(t-1) - \mathrm{SNR}(t))\|\mathbf{x}_{0,\theta} - \mathbf{x}_0\|_2^2 + C.$$

## A.4 Training

Empirically, instead of using the variational lower bound, the neural network can be trained on one of the following simplified loss (Ho et al., 2020)

$$L_{\mathrm{simple},\boldsymbol{\epsilon}_t}(\theta) = \mathbb{E}_{t,\mathbf{x}_0,\boldsymbol{\epsilon}_t}\|\boldsymbol{\epsilon}_t - \boldsymbol{\epsilon}_{t,\theta}\|_2^2 = \mathbb{E}_{t,\mathbf{x}_0,\boldsymbol{\epsilon}_t} L(\boldsymbol{\epsilon}_t, \boldsymbol{\epsilon}_{t,\theta}), \quad (\boldsymbol{\epsilon}\text{-parameterization})$$

$$L_{\mathrm{simple},\mathbf{x}_0}(\theta) = \mathbb{E}_{t,\mathbf{x}_0,\boldsymbol{\epsilon}_t}\|\mathbf{x}_0 - \mathbf{x}_{0,\theta}\|_2^2 = \mathbb{E}_{t,\mathbf{x}_0,\boldsymbol{\epsilon}_t} L(\mathbf{x}_0, \mathbf{x}_{0,\theta}). \quad (\mathbf{x}_0\text{-parameterization})$$

with $t$ uniformly sampled from 1 to $T$ and $\boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. $L(\cdot, \cdot)$ is a loss function in the space of $\mathbf{x}$. With the importance sampling proposed in Nichol and Dhariwal (2021), $t$ can be sampled with a probability proportional to $\mathbb{E}_{\mathbf{x}_0,\mathbf{x}_t} L(\mathbf{x}_0, \hat{\mathbf{x}}_0)$. In other words, a time step $t$ is sampled more often if the loss is larger.

As the previous work (Fu et al., 2023) has extensively compared the $\boldsymbol{\epsilon}$-parameterization and $\mathbf{x}_0$-parameterization, as well as the benefits of including Dice loss, in this work, we use $\mathbf{x}_0$-parameterization with a weighted sum of cross-entropy and foreground-only Dice loss Kirillov et al. (2023).

## A.5 Variance Resampling

Given a variance schedule $\{\beta_t\}_{t=1}^T$ (e.g. $T = 1001$), a subsequence $\{\beta_k\}_{k=1}^K$ (e.g. $K = 5$) can be sampled with $\{t_k\}_{k=1}^K$. Following Nichol and Dhariwal (2021), we can define $\beta_k = 1 - \frac{\bar{\alpha}_{t_k}}{\bar{\alpha}_{t_{k-1}}}$ then $\alpha_k = 1 - \beta_k$ and $\bar{\alpha}_k = \prod_{s=1}^k \alpha_s$ can be recalculated correspondingly. In this work, $t_k$ is uniformly downsampled. For instance, if $T = 1001$ and $K = 5$, then $\{t_k\}_{k=1}^K = \{1, 251, 501, 751, 1001\}$.

## Appendix B. Denoising Diffusion Implicit Model

**Definition**    Song et al. (2020a) parameterize $q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0)$ as follows, with $\boldsymbol{\epsilon} = \frac{\mathbf{x}_t - \sqrt{\bar{\alpha}_t}\mathbf{x}_0}{\sqrt{1-\bar{\alpha}_t}}$,

$$q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t}\,\boldsymbol{\epsilon}, \sigma_t^2\mathbf{I}).$$

For any variance schedule $\sigma_t$, this formulation ensures $q(\mathbf{x}_t \mid \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I})$. Particularly, if $\sigma_t^2 = \tilde{\beta}_t$, this represents DDPM. If $\sigma_t = 0$ for $t > 1$ and $\sigma_1 = \sqrt{\tilde{\beta}_1}$, the model is deterministic and named as denoising diffusion implicit model (DDIM).

**Inference**    For DDIM, at inference time, the denoising starts with a Gaussian noise $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \boldsymbol{I})$ and the data is denoised step-by-step for $t = T, \cdots, 1$:

$$p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t) = \begin{cases} \mathcal{N}(\hat{\mathbf{x}}_0, \sigma_1^2\mathbf{I}) & t = 1 \\ q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_{0,\theta}(\mathbf{x}_t, t)) & t > 1 \end{cases}$$

## Appendix C. Self-conditioning

The self-conditioning methods proposed in Chen et al. (2022b) ("Diff. sc. $\mathbf{x}_t$" in Equation (7)) and Watson et al. (2023) ("sc. $\mathbf{x}_{t+1}$" in Equation (8)) are illustrated below.

$$\mathbf{x}_t \sim \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1-\bar{\alpha}_t)\mathbf{I}), \qquad \text{(sc. } \mathbf{x}_t, \text{ step 1, sampling)} \qquad (7a)$$

$$\hat{\mathbf{x}}_0 = \text{StopGradient}(f_\theta(I, t, \mathbf{x}_t, \mathbf{0})), \qquad \text{(sc. } \mathbf{x}_t, \text{ step 1, prediction)} \qquad (7b)$$

$$\hat{\mathbf{x}}_0 = \text{Dropout}_{p=50\%}(\hat{\mathbf{x}}_0), \qquad \text{(sc. } \mathbf{x}_t, \text{ step 2, dropout)} \qquad (7c)$$

$$\hat{\mathbf{x}}_0 = f_\theta(I, t, \mathbf{x}_t, \hat{\mathbf{x}}_0), \qquad \text{(sc. } \mathbf{x}_t, \text{ step 2, prediction)} \qquad (7d)$$

$$L_{\text{denoising}}(\theta) = \mathbb{E}_{t,\mathbf{x}_0,\mathbf{x}_t} L(\mathbf{x}_0, \hat{\mathbf{x}}_0), \qquad \text{(loss calculation)} \qquad (7e)$$

$$\mathbf{x}_{t+1} \sim \mathcal{N}(\mathbf{x}_{t+1}; \sqrt{\bar{\alpha}_{t+1}}\,\mathbf{x}_0, (1-\bar{\alpha}_{t+1})\mathbf{I}), \qquad \text{(sc. } \mathbf{x}_{t+1}, \text{ step 1, sampling)} \qquad (8a)$$

$$\hat{\mathbf{x}}_0 = \text{StopGradient}(f_\theta(I, t+1, \mathbf{x}_{t+1}, \mathbf{0})), \qquad \text{(sc. } \mathbf{x}_{t+1}, \text{ step 1, prediction)} \qquad (8b)$$

$$\hat{\mathbf{x}}_0 = \text{Dropout}_{p=50\%}(\hat{\mathbf{x}}_0), \qquad \text{(sc. } \mathbf{x}_{t+1}, \text{ step 2, dropout)} \qquad (8c)$$

$$\mathbf{x}_t \sim \mathcal{N}(\mathbf{x}_t; \tilde{\boldsymbol{\mu}}, \tilde{\beta}_{t+1}\mathbf{I}), \qquad \text{(sc. } \mathbf{x}_{t+1}, \text{ step 2, sampling)} \qquad (8d)$$

$$\tilde{\boldsymbol{\mu}} = \frac{\sqrt{\bar{\alpha}_t}\beta_{t+1}}{1-\bar{\alpha}_{t+1}}\mathbf{x}_0 + \frac{1-\bar{\alpha}_t}{1-\bar{\alpha}_{t+1}}\sqrt{\alpha_{t+1}}\,\mathbf{x}_{t+1}$$

$$\hat{\mathbf{x}}_0 = f_\theta(I, t, \mathbf{x}_t, \hat{\mathbf{x}}_0), \qquad \text{(sc. } \mathbf{x}_{t+1}, \text{ step 2, prediction)} \qquad (8e)$$

$$L_{\text{denoising}}(\theta) = \mathbb{E}_{t,\mathbf{x}_0,\mathbf{x}_t} L(\mathbf{x}_0, \hat{\mathbf{x}}_0), \qquad \text{(loss calculation)} \qquad (8f)$$

## Appendix D. Diffusion Noise Schedule

The noise schedule $\beta_t$ and $\sqrt{\bar{\alpha}_t}$ have been visualised in Figure 7. The cross entropy and dice score between $\mathbf{x}_t$ and ground truth $\mathbf{x}_0$ have also been visualized to empirically measure the amount of information of ground truth $\mathbf{x}_0$ contained in $\mathbf{x}_t$.

## Appendix E. Implementation Details

Table 6: **Training Hyper-parameters**

| Parameter | Value |
|---|---|
| Optimiser | AdamW (b1=0.9, b2=0.999, weight_decay=1E-8) |
| Learning Rate Warmup | 100 steps |
| Learning Rate Decay | 10,000 steps |
| Learning Rate Values | Initial = 1E-5, Peak = 8E-4, End = 5E-5 |
| Batch size | 256 for Muscle Ultrasound and 8 for other data sets |
| Number of samples | 320K for Muscle Ultrasound and 100K for other data sets |

(a) $\beta_t$        (b) $\sqrt{\bar{\alpha}_t}$
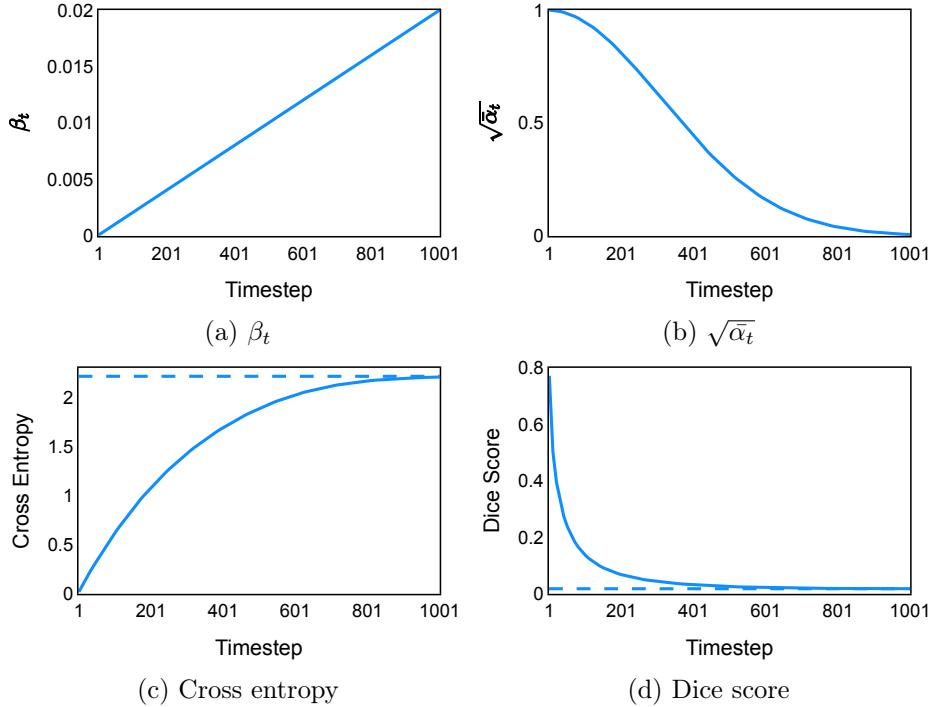
(c) Cross entropy        (d) Dice score

Figure 7: **Information contained in $\mathbf{x}_t$.** Cross entropy and dice score between $\mathbf{x}_t$ and ground truth $\mathbf{x}_0$ are used to empirically measure the amount of information of ground truth $\mathbf{x}_0$ contained in $\mathbf{x}_t$. The dashed line represents the information contained in the sampled noise (between noise and ground truth $\mathbf{x}_0$), which is considered to be the limit. The values are calculated using the sample "005095" in prostate MR data set.

Table 7: **Network Size**

| Dimension | Method | Transformer | |
|---|---|---|---|
| | | ✓ | |
| 2D | No diff. | 12,586,594 | 10,550,370 |
| | Diff. | 13,335,554 | 11,299,330 |
| 3D | No diff. | 33,385,154 | 31,283,394 |
| | Diff. | 34,135,266 | 32,033,506 |

539

## Appendix F. Results

### F.1 Diffusion Training Strategy Comparison

Table 8: **Per class Dice score comparison.** "No diff." represents non-diffusion model. "Diff. rec. $\mathbf{x}_T$" represents the diffusion model with proposed recycling. "Ensemble" represents the model averaging the probabilities from "No diff." and "Diff. rec. $\mathbf{x}_T$". The inference sampler is DDPM. The best results are in bold and underline indicates the difference to non-diffusion model is significant with p-value $< 0.05$.

| Method | Spleen | RT kidney | LT kidney | Gall bladder |
|---|---|---|---|---|
| No diff. | $96.62 \pm 1.87$ | $95.08 \pm 10.74$ | $96.29 \pm 1.73$ | $78.83 \pm 27.82$ |
| Diff. rec. $\mathbf{x}_T$ | $96.40 \pm 2.42$ | $96.24 \pm 1.90$ | $96.27 \pm 1.53$ | $76.68 \pm 29.25$ |
| Ensemble | $\mathbf{96.78 \pm 1.75}$ | $\underline{\mathbf{96.47 \pm 2.44}}$ | $\underline{\mathbf{96.50 \pm 1.51}}$ | $\mathbf{79.65 \pm 27.29}$ |

| Method | Esophagus | Liver | Stomach | Arota |
|---|---|---|---|---|
| No diff. | $83.22 \pm 11.08$ | $97.36 \pm 1.17$ | $90.53 \pm 14.78$ | $94.65 \pm 4.22$ |
| Diff. rec. $\mathbf{x}_T$ | $83.60 \pm 10.32$ | $97.33 \pm 1.13$ | $90.77 \pm 14.46$ | $94.66 \pm 4.66$ |
| Ensemble | $\underline{\mathbf{84.10 \pm 11.15}}$ | $\mathbf{97.54 \pm 1.05}$ | $\mathbf{91.07 \pm 14.91}$ | $\mathbf{94.96 \pm 4.39}$ |

| Method | Postcava | Pancreas | Right adrenal gland | Left adrenal gland |
|---|---|---|---|---|
| No diff. | $90.45 \pm 4.68$ | $84.88 \pm 11.40$ | $77.80 \pm 9.46$ | $77.98 \pm 11.95$ |
| Diff. rec. $\mathbf{x}_T$ | $90.55 \pm 4.19$ | $84.86 \pm 11.15$ | $\underline{76.63 \pm 12.84}$ | $78.01 \pm 11.60$ |
| Ensemble | $\mathbf{91.18 \pm 4.12}$ | $\mathbf{85.85 \pm 11.12}$ | $\underline{\mathbf{78.51 \pm 10.58}}$ | $\underline{\mathbf{78.95 \pm 11.45}}$ |

| Method | Duodenum | Bladder | Prostate/uterus | |
|---|---|---|---|---|
| No diff. | $79.57 \pm 14.89$ | $88.09 \pm 16.25$ | $82.35 \pm 18.90$ | |
| Diff. rec. $\mathbf{x}_T$ | $79.80 \pm 15.14$ | $87.90 \pm 16.65$ | $81.90 \pm 18.86$ | |
| Ensemble | $\underline{\mathbf{80.99 \pm 15.07}}$ | $\mathbf{88.61 \pm 16.59}$ | $\mathbf{83.06 \pm 18.68}$ | |

(a) Abdominal CT: LT and RT stand for left and right, respectively.

| Method | Bladder | Bone | Obturator internus | Transition zone |
|---|---|---|---|---|
| No diff. | $93.28 \pm 9.90$ | $93.12 \pm 5.68$ | $88.95 \pm 3.53$ | $79.61 \pm 8.37$ |
| Diff. rec. $\mathbf{x}_T$ | $93.57 \pm 9.61$ | $\mathbf{93.84 \pm 5.85}$ | $\underline{89.15 \pm 3.62}$ | $79.79 \pm 8.36$ |
| Ensemble | $\underline{\mathbf{93.66 \pm 9.84}}$ | $\underline{93.77 \pm 5.52}$ | $\mathbf{89.52 \pm 3.51}$ | $\underline{\mathbf{80.57 \pm 8.20}}$ |

| Method | Central gland | Rectum | Seminal vesicle | NV bundle |
|---|---|---|---|---|
| No diff. | $88.75 \pm 5.60$ | $93.30 \pm 3.48$ | $77.55 \pm 10.99$ | $67.17 \pm 14.34$ |
| Diff. rec. $\mathbf{x}_T$ | $\underline{89.13 \pm 5.78}$ | $93.42 \pm 3.51$ | $\underline{78.39 \pm 9.71}$ | $67.07 \pm 15.50$ |
| Ensemble | $\mathbf{89.45 \pm 5.56}$ | $\mathbf{93.70 \pm 3.37}$ | $\underline{\mathbf{78.91 \pm 10.28}}$ | $\mathbf{68.01 \pm 14.85}$ |

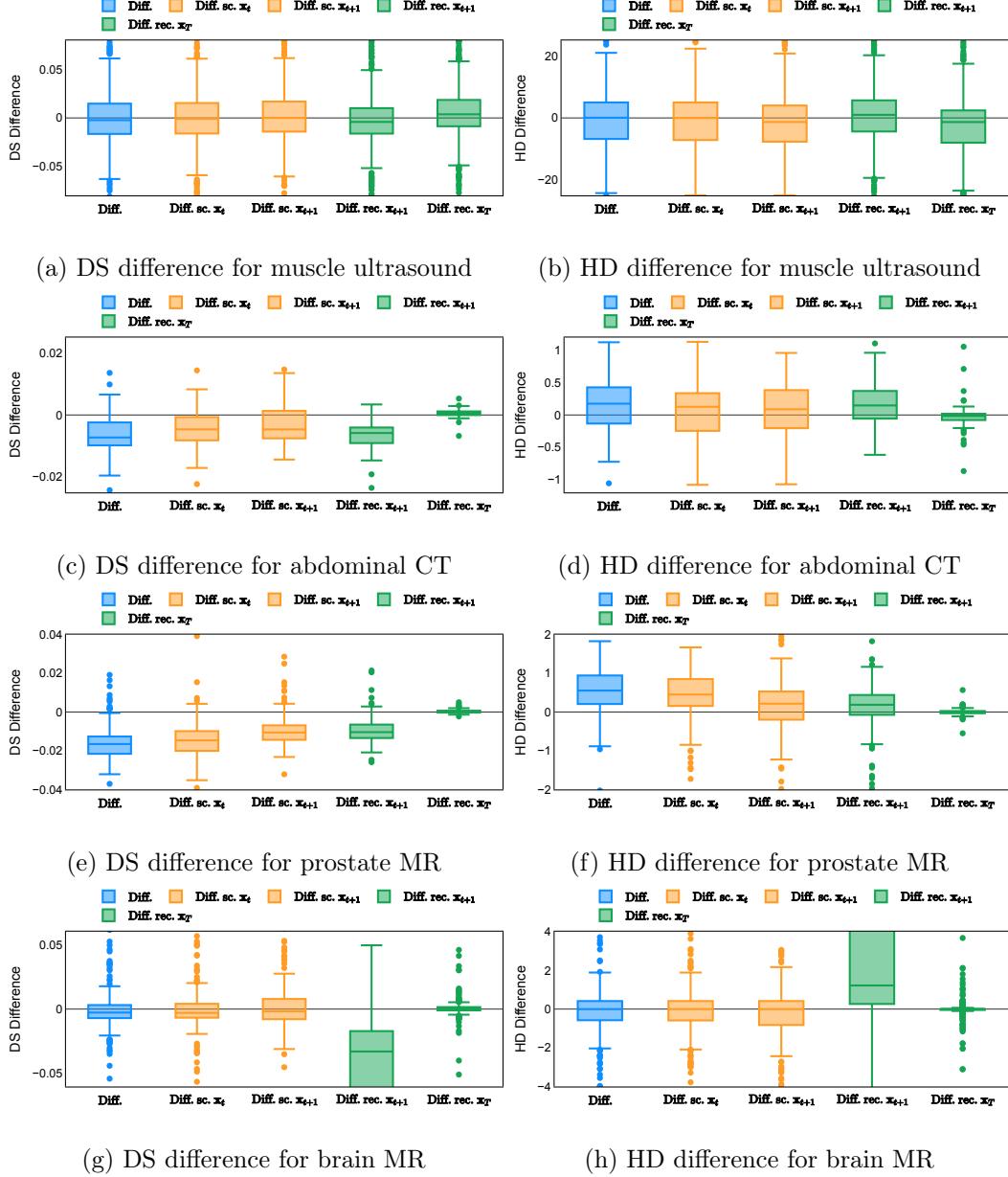(b) Prostate MR: Dice score per class. NV stands for neurovascular.

Table 9: **Per class Hausdorff distance comparison** "No diff." represents non-diffusion model. "Diff. rec. $\mathbf{x}_T$" represents the diffusion model with proposed recycling. "Ensemble" represents the model averaging the probabilities from "No diff." and "Diff. rec. $\mathbf{x}_T$". The inference sampler is DDPM. The best results are in bold and underline indicates the difference to non-diffusion model is significant with p-value $< 0.05$.

| Method | Spleen | Right kidney | Left kidney | Gall bladder |
|---|---|---|---|---|
| No diff. | $3.22 \pm 4.91$ | $1.97 \pm 1.46$ | $4.13 \pm 10.83$ | $\mathbf{9.23 \pm 16.71}$ |
| Diff. rec. $\mathbf{x}_T$ | $\mathbf{2.86 \pm 3.84}$ | $1.93 \pm 0.83$ | $3.13 \pm 8.35$ | $12.65 \pm 21.74$ |
| Ensemble | $\underline{2.89 \pm 4.28}$ | $\mathbf{1.84 \pm 1.11}$ | $\mathbf{2.70 \pm 5.86}$ | $9.57 \pm 18.86$ |

| Method | Esophagus | Liver | Stomach | Arota |
|---|---|---|---|---|
| No diff. | $5.50 \pm 6.81$ | $3.50 \pm 2.50$ | $8.96 \pm 13.99$ | $6.62 \pm 14.52$ |
| Diff. rec. $\mathbf{x}_T$ | $5.30 \pm 6.41$ | $3.79 \pm 4.16$ | $9.00 \pm 14.03$ | $\mathbf{5.41 \pm 11.20}$ |
| Ensemble | $\mathbf{5.22 \pm 6.63}$ | $\underline{\mathbf{3.06 \pm 1.63}}$ | $\mathbf{8.04 \pm 12.87}$ | $5.47 \pm 11.29$ |

| Method | Postcava | Pancreas | Right adrenal gland | Left adrenal gland |
|---|---|---|---|---|
| No diff. | $4.80 \pm 4.55$ | $7.57 \pm 8.62$ | $\mathbf{4.39 \pm 2.39}$ | $5.15 \pm 5.40$ |
| Diff. rec. $\mathbf{x}_T$ | $4.62 \pm 3.09$ | $7.50 \pm 8.62$ | $4.66 \pm 3.14$ | $4.87 \pm 4.64$ |
| Ensemble | $\mathbf{4.41 \pm 3.25}$ | $\mathbf{6.96 \pm 8.40}$ | $4.41 \pm 2.79$ | $\underline{\mathbf{4.82 \pm 4.92}}$ |

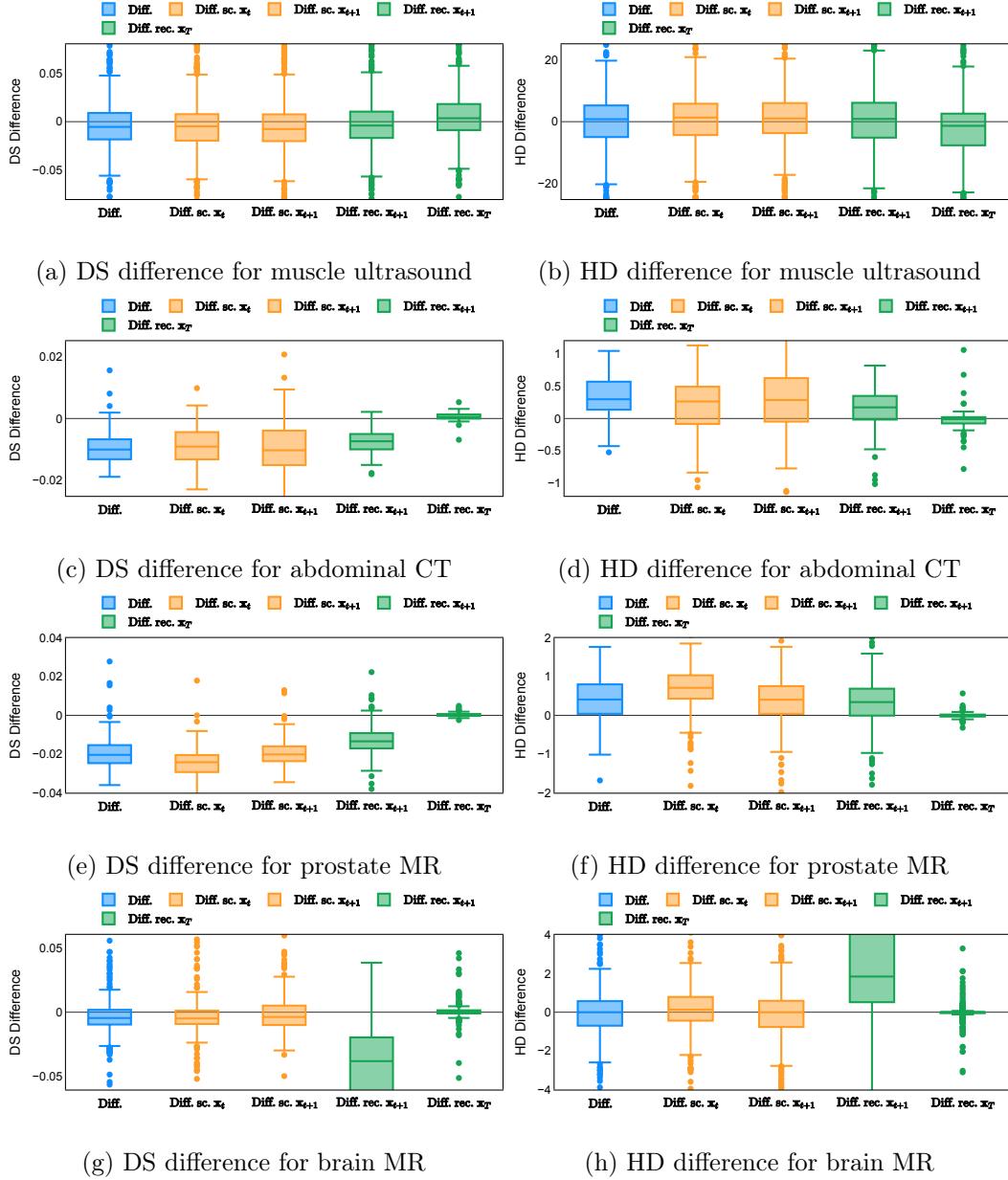| Method | Duodenum | Bladder | Prostate/uterus | |
|---|---|---|---|---|
| No diff. | $10.54 \pm 8.44$ | $9.10 \pm 23.07$ | $10.97 \pm 19.01$ | |
| Diff. rec. $\mathbf{x}_T$ | $\underline{9.31 \pm 7.13}$ | $10.70 \pm 31.83$ | $13.35 \pm 32.75$ | |
| Ensemble | $\underline{\mathbf{9.29 \pm 7.37}}$ | $\mathbf{6.52 \pm 10.34}$ | $\mathbf{9.14 \pm 13.11}$ | |

(a) Abdominal CT: LT and RT stand for left and right, respectively.

| Method | Bladder | Bone | Obturator internus | Transition zone |
|---|---|---|---|---|
| No diff. | $3.30 \pm 4.54$ | $3.18 \pm 9.77$ | $4.60 \pm 3.36$ | $5.97 \pm 4.97$ |
| Diff. rec. $\mathbf{x}_T$ | $3.20 \pm 4.12$ | $\mathbf{2.21 \pm 1.62}$ | $4.50 \pm 3.46$ | $6.25 \pm 4.96$ |
| Ensemble | $\mathbf{2.95 \pm 3.48}$ | $2.32 \pm 1.46$ | $\underline{\mathbf{4.34 \pm 3.29}}$ | $\mathbf{6.18 \pm 5.11}$ |

| Method | Central gland | Rectum | Seminal vesicle | NV bundle |
|---|---|---|---|---|
| No diff. | $3.94 \pm 2.28$ | $4.46 \pm 5.69$ | $4.82 \pm 3.85$ | $6.68 \pm 6.33$ |
| Diff. rec. $\mathbf{x}_T$ | $\underline{3.70 \pm 1.93}$ | $4.25 \pm 4.75$ | $4.57 \pm 2.66$ | $6.55 \pm 6.28$ |
| Ensemble | $\underline{\mathbf{3.66 \pm 1.93}}$ | $\underline{\mathbf{4.16 \pm 5.25}}$ | $\mathbf{4.52 \pm 2.83}$ | $\mathbf{6.45 \pm 6.34}$ |

(b) Prostate MR: Dice score per class. NV stands for neurovascular.

(a) DS difference for muscle ultrasound

(b) HD difference for muscle ultrasound

(c) DS difference for abdominal CT

(d) HD difference for abdominal CT

(e) DS difference for prostate MR

(f) HD difference for prostate MR

(g) DS difference for brain MR

(h) HD difference for brain MR

Figure 8: **Segmentation performance difference between the last step and first step using DDPM.** "Diff." represents standard diffusion. "Diff. sc. $\mathbf{x}_t$" and "Diff. sc. $\mathbf{x}_{t+1}$" represents self-conditioning from Chen et al. (2022b) and Watson et al. (2023), respectively. "Diff. rec. $\mathbf{x}_{t+1}$" and "Diff. rec. $\mathbf{x}_T$" represents recycling from Fu et al. (2023) and the proposed recycling in this work, respectively. The sampler is DDPM. DS and HD represents Dice score and Hausdorff distance, respectively. The difference is the value at the last step subtracted by the one at the first step. A positive value for Dice score difference or a negative value for Hausdorff distance means improvement.

(a) DS difference for muscle ultrasound

(b) HD difference for muscle ultrasound

(c) DS difference for abdominal CT

(d) HD difference for abdominal CT

(e) DS difference for prostate MR

(f) HD difference for prostate MR

(g) DS difference for brain MR

(h) HD difference for brain MR

Figure 9: **Segmentation performance difference between the last step and first step using DDIM.** "Diff." represents standard diffusion. "Diff. sc. $\mathbf{x}_t$" and "Diff. sc. $\mathbf{x}_{t+1}$" represents self-conditioning from Chen et al. (2022b) and Watson et al. (2023), respectively. "Diff. rec. $\mathbf{x}_{t+1}$" and "Diff. rec. $\mathbf{x}_T$" represents recycling from Fu et al. (2023) and the proposed recycling in this work, respectively. The sampler is DDIM. DS and HD represents Dice score and Hausdorff distance, respectively. The difference is the value at the last step subtracted by the one at the first step. A positive value for Dice score difference or a negative value for Hausdorff distance means improvement.
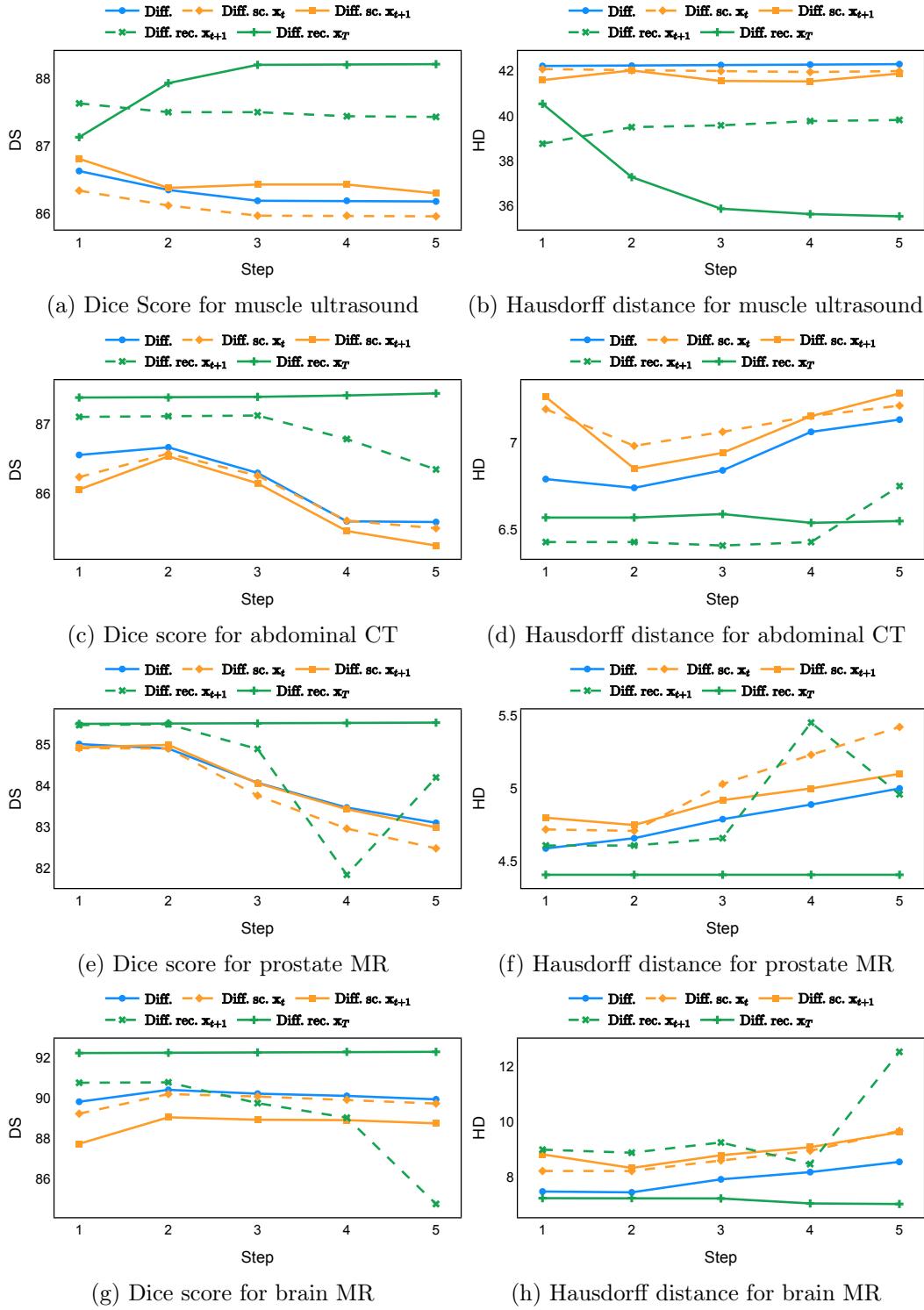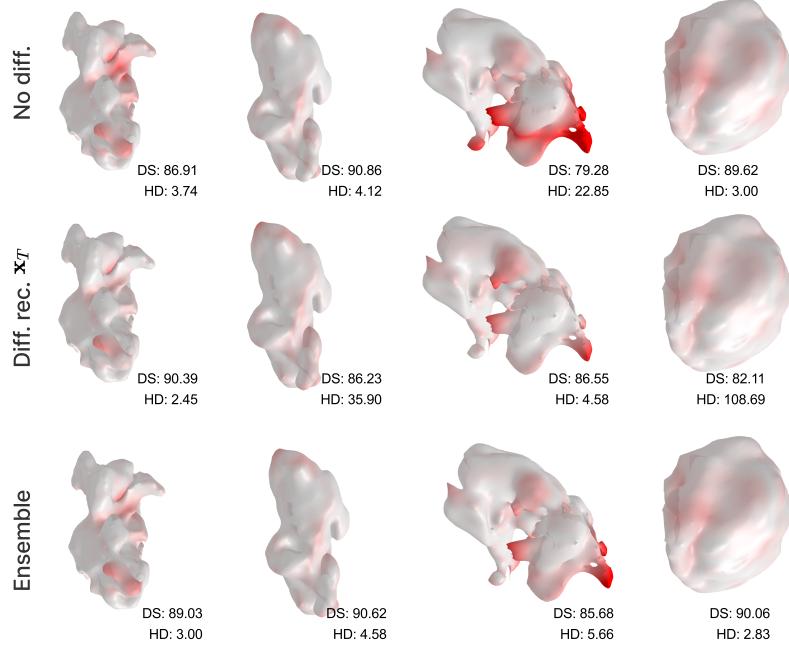
(a) Dice Score for muscle ultrasound

(b) Hausdorff distance for muscle ultrasound

(c) Dice score for abdominal CT

(d) Hausdorff distance for abdominal CT

(e) Dice score for prostate MR

(f) Hausdorff distance for prostate MR

(g) Dice score for brain MR

(h) Hausdorff distance for brain MR

Figure 10: **Segmentation performance per step.** "Diff." represents standard diffusion. "Diff. sc. $\mathbf{x}_t$" and "Diff. sc. $\mathbf{x}_{t+1}$" represents self-conditioning from Chen et al. (2022b) and Watson et al. (2023), respectively. "Diff. rec. $\mathbf{x}_{t+1}$" and "Diff. rec. $\mathbf{x}_T$" represents recycling from Fu et al. (2023) and the proposed recycling in this work, respectively. The sampler is DDIM.

## F.2 Comparison to Non-diffusion Models



Figure 11: **Segmentation error of non-diffusion-based and diffusion-based models for tumour in brain MR.** "No diff." represents non-diffusion model. "Diff. rec. $\mathbf{x}_T$" represents the diffusion model with proposed recycling. The ground truth segmentation is visualised. For each point on the surface, the distance to the surface of predicted segmentation is calculated and displayed with red color. The Dice score (DS) and Hausdorff distance (HD) for each sample are labeled at bottom.

## F.3 Ablation Studies

Table 10: **Diffusion with different training variance schedule.** "Diff." represents standard diffusion. "Diff. sc. $\mathbf{x}_t$" and "Diff. sc. $\mathbf{x}_{t+1}$" represents self-conditioning from Chen et al. (2022b) and Watson et al. (2023), respectively. "Diff. rec. $\mathbf{x}_{t+1}$" and "Diff. rec. $\mathbf{x}_T$" represents recycling from Fu et al. (2023) and the proposed recycling in this work, respectively. "T" represents the length of variance schedule during training. The best results are in bold and a underline indicates the difference to the second best is significant with p-value < 0.05.

| T | Method | DDPM | | DDIM | |
|---|---|---|---|---|---|
| | | DS ↑ | HD ↓ | DS ↑ | HD ↓ |
| 1001 | Diff. | 86.60 ± 12.38 | 41.11 ± 35.48 | 86.18 ± 12.41 | 42.31 ± 35.82 |
| | Diff. sc. $\mathbf{x}_t$ | 86.35 ± 14.14 | 40.42 ± 37.53 | 85.96 ± 13.78 | 42.00 ± 36.76 |
| | Diff. sc. $\mathbf{x}_{t+1}$ | 87.14 ± 11.48 | 39.24 ± 32.83 | 86.30 ± 11.49 | 41.89 ± 32.72 |
| | Diff. rec. $\mathbf{x}_{t+1}$ | 87.44 ± 12.39 | 39.68 ± 36.21 | 87.43 ± 12.25 | 39.82 ± 35.39 |
| | Diff. rec. $\mathbf{x}_T$ | 88.23 ± 11.69 | 35.37 ± 31.79 | 88.21 ± 11.70 | 35.52 ± 31.91 |
| 5 | Diff. | 87.81 ± 10.98 | 37.39 ± 31.17 | 87.76 ± 11.00 | 37.56 ± 31.34 |
| | Diff. sc. $\mathbf{x}_t$ | 88.11 ± 11.06 | 35.94 ± 30.13 | 88.20 ± 10.73 | 35.57 ± 29.68 |
| | Diff. sc. $\mathbf{x}_{t+1}$ | 87.61 ± 10.88 | 37.76 ± 29.91 | 88.09 ± 10.66 | 35.73 ± 29.26 |
| | Diff. rec. $\mathbf{x}_{t+1}$ | 88.19 ± 10.60 | 36.10 ± 30.38 | 87.83 ± 11.01 | 37.22 ± 30.55 |
| | Diff. rec. $\mathbf{x}_T$ | **89.01 ± 10.79** | **33.70 ± 30.29** | **88.80 ± 11.54** | **34.26 ± 31.88** |

(a) Muscle Ultrasound

| T | Method | DDPM | | DDIM | |
|---|---|---|---|---|---|
| | | DS ↑ | HD ↓ | DS ↑ | HD ↓ |
| 1001 | Diff. | 85.25 ± 5.36 | 7.12 ± 3.83 | 85.59 ± 5.24 | 7.13 ± 3.98 |
| | Diff. sc. $\mathbf{x}_t$ | 86.04 ± 5.12 | 7.06 ± 4.20 | 85.50 ± 5.14 | 7.21 ± 4.16 |
| | Diff. sc. $\mathbf{x}_{t+1}$ | 85.86 ± 5.27 | 6.98 ± 3.54 | 85.25 ± 5.42 | 7.28 ± 3.72 |
| | Diff. rec. $\mathbf{x}_{t+1}$ | 86.48 ± 5.24 | 6.69 ± 4.59 | 86.35 ± 5.31 | 6.75 ± 4.55 |
| | Diff. rec. $\mathbf{x}_T$ | **87.45 ± 5.43** | **6.56 ± 5.44** | **87.45 ± 5.43** | **6.55 ± 5.43** |
| 5 | Diff. | 86.42 ± 5.00 | 7.09 ± 4.40 | 86.52 ± 5.18 | 6.65 ± 3.88 |
| | Diff. sc. $\mathbf{x}_t$ | 86.68 ± 4.96 | 7.06 ± 6.98 | 86.39 ± 4.87 | 7.12 ± 6.95 |
| | Diff. sc. $\mathbf{x}_{t+1}$ | 86.34 ± 5.33 | 6.69 ± 3.46 | 86.13 ± 5.27 | 6.74 ± 3.55 |
| | Diff. rec. $\mathbf{x}_{t+1}$ | 87.27 ± 5.20 | 6.64 ± 4.69 | 87.27 ± 5.20 | 6.63 ± 4.69 |
| | Diff. rec. $\mathbf{x}_T$ | 87.38 ± 5.46 | 6.71 ± 4.46 | 87.37 ± 5.45 | 6.74 ± 4.49 |

(b) Abdominal CT