# Exploring Robustness of Cortical Morphometry in the presence of white matter lesions, using Diffusion Models for Lesion Filling

Vinzenz Uhr, Ivan Diaz, Christian Rummel, and Richard McKinley

*Abstract*—Cortical thickness measurements from magnetic resonance imaging, an important biomarker in many neurodegenerative and neurological disorders, are derived by many tools from an initial voxel-wise tissue segmentation. White matter (WM) hypointensities in T1-weighted imaging, such as those arising from multiple sclerosis or small vessel disease, are known to affect the output of brain segmentation methods and therefore bias cortical thickness measurements. These effects are well-documented among traditional brain segmentation tools but have not been studied extensively in tools based on deep-learning segmentations, which promise to be more robust. In this paper, we explore the potential of deep learning to enhance the accuracy and efficiency of cortical thickness measurement in the presence of WM lesions, using a high-quality lesion filling algorithm leveraging denoising diffusion networks.

A pseudo-3D U-Net architecture trained on the OASIS dataset to generate synthetic healthy tissue, conditioned on binary lesion masks derived from the MSSEG dataset, allows realistic removal of white matter lesions in multiple sclerosis patients. By applying morphometry methods to patient images before and after lesion filling, we analysed robustness of global and regional cortical thickness measurements in the presence of white matter lesions. Methods based on a deep learning-based segmentation of the brain (Fastsurfer, DL+DiReCT, ANTsPyNet) exhibited greater robustness than those using classical segmentation methods (Freesurfer, ANTs).

*Index Terms*—multiple sclerosis, lesions, inpainting, diffusion models, deep learning, brain morphometry, cortical thickness

## I. INTRODUCTION

**W**HITE MATTER (WM) lesions, often associated with neurological conditions like multiple sclerosis (MS), can significantly perturb tissue segmentation algorithms operating on magnetic resonance imaging (MRI), causing misclassification of tissue types. The misclassification varies considerably with lesion size and intensity, especially when the lesion intensity is similar to that of the Gray Matter (GM)/WM interface. As well as causing inaccuracies in volumetric gray mater measurements in brains with lesions, these perturbations can cause downstream biases in cortical thickness calculations. In the past different inpainting algorithms (also known as *lesion filling*) have been proposed which can replace voxels within a lesion mask with white matter tissue intensities, leading to more robust measurements [1], [2], [3], [4], [5]. In recent years, deep learning has emerged as a powerful tool in medical image analysis, revolutionizing the field with its ability to automatically learn and extract meaningful features from large datasets. Deep learning techniques, especially convolutional neural networks (CNNs), have shown remarkable success in various medical imaging applications, including both brain segmentation [?] and lesion filling [5]. Deep-learning-based image segmentation models have been shown to dramatically outperform previous generations of brain segmentation tools. Meanwhile denoising diffusion probabilistic models (DDPM) [6] have shown an impressive performance and experienced increasing popularity in medical image analysis [7].

In this article we use DDPMs to examine the robustness of DL and non-DL approaches to cortical morphometry. Our contributions are

- Evaluating different methods using noise diffusion models to inpaint WM lesions in MR images.
- Evaluating the impact of lesion filling on cortical thickness measurements using existing morphological tools, to determine their robustness to the presence of WM lesions.

## II. BACKGROUND

### A. Lesion Filling

The presence of white matter lesions can significantly impact MR-based measurements like cortical thickness due to misclassification of different tissue types [1] [3] [8] [9] [10]. This misclassification is particularly problematic for WM lesions with size and intensity similar to the GM/WM interface and leads to overestimation of GM atrophy [11]. Lesion filling algorithms have been developed to address this issue and improve measurements such as cortical volume, thickness and surface area estimation [3].

Early lesion filling approaches employed various strategies. For instance, [1] utilized lesion filling to enhance brain volume measurements, including normalized brain volume (NBV), normalized white matter volume (NWMV), normalized gray matter volume (NGMV), and percentage brain volume change (PBVC). Their method involved calculating intensity distributions of cerebrospinal fluid (CSF), CSF/GM, GM, and GM/WM from existing brain MR-images and filling WM lesions with pixel intensities randomly sampled from these distributions.

Another approach, proposed by [2], involved refilling WM lesions by replacing lesion voxel intensities with random values drawn from a normal distribution based on the WM signal intensity of each two-dimensional slice. Segmentation of the slices was achieved using the fuzzy c-means algorithm.

Graph theoretical network analysis, a technique used to assess brain connectivity patterns, can also benefit from lesion filling. To reduce the variability in network analysis caused by

WM lesions, [12] applied lesion filling by substituting lesion voxel intensities with intensities from nearby voxels. Their study suggests that lesion filling might improve the detection of network alterations in MS patients, but also highlights the potential for introducing artifacts. Therefore, caution is advised, especially for individuals with high lesion loads or lesions located at the WM/CSF or WM/GM interface.

More recent advancements leverage machine learning for lesion inpainting. [4] employed a total variation model to improve registration performance with brain atlases. Inspired by Gated Convolution, [5] introduced a user-guided deep adversarial inpainting model capable of filling irregularly shaped holes in high-resolution T1w MR brain images. Training data generation involved synthesizing lesion masks by sampling and deforming random circles. Additional data augmentation techniques included rotation, cropping, flipping, noise addition, and varying brightness levels.

The emergence of DDPMs [6] offers a novel approach for high-quality image generation. DDPMs exhibit superior distribution coverage and training stability compared to adversarial loss-trained models, achieving state-of-the-art performance in various image synthesis tasks.

The International Brain Tumor Segmentation (BraTS) challenge in 2023 incorporated an inpainting challenge focused on synthesizing healthy brain tissue in glioma-affected regions [13]. Due to the high computational cost of 3D processing, [14] opted for a 2D diffusion model conditioned on glioma masks. While achieving comparable results to other participants, their approach resulted in stripe artifacts due to stacking of the 2D slices. Gaussian filtering was subsequently employed to mitigate these effects at the slice borders.

[5] addressed the high computational demands of 3D diffusion models by proposing several resource-reduction strategies. Notably, they introduced PatchDDM, a memory-efficient patch-based diffusion model that allows for inference on the entire volume while training solely on patches. Additional approaches included reducing self-attention layers, incorporating additive skip connections, and training on downsampled data.

In pursuit of improved inpainting quality for 3D MR-images, [15] evaluated and modified various diffusion models, including 2D, pseudo-3D, and 3D models operating in image space, 3D wavelet or 3D latent space. Their findings suggest that the pseudo-3D model proposed by [16] achieved the best performance in terms of structural similarity index measure (SSIM), peak signal noise ratio (PSNR), and mean squared error (MSE).

## B. Denoising Diffusion Probabilistic Models

Diffusion models are a generative deep learning technique that leverage an approach for data synthesis. The core idea lies in progressively transforming a data sample $x_0$ from its original distribution into a sample $x_T$ from a normally distributed noise. The model then learns to progressively reverse this transformation process [17].

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1-\beta_t}x_{t-1}, \beta_t I)$$



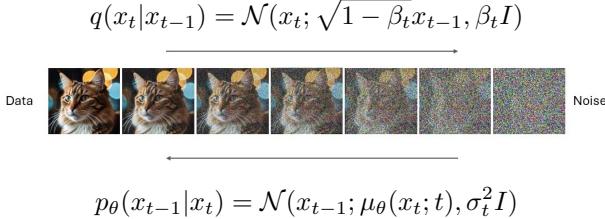$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t; t), \sigma_t^2 I)$$

Fig. 1. Diffusion process from Data to Noise and reverse process from Noise to Data.

The term forward process is the progressive transition from a clean image $x_0$ to pure noise $x_T$, via a series of steps, each adding additional random noise. At each step, zero-mean Gaussian noise is added, gradually increasing its strength until a maximum level is reached at a predefined endpoint $t = T$. This process is a *Markov chain*, since the noisy image $x_t$ only depends on the immediately previous one $x_{t-1}$ (and not the whole sequence). The mathematical formulation behind this forward noising process $q$ is denoted by

$$q(x_t|x_{t-1}) := \mathcal{N}(x_t; \sqrt{1-\beta_t}x_{t-1}, \beta_t I) \qquad (1)$$

Where $I$ represents the identity matrix and $\beta_t$ the variance schedule, which controls the amount of noise added at each step based on the current step $t$. The noisy image $x_t$ can be written

$$x_t = \sqrt{\bar{\alpha_t}}x_0 + \sqrt{1-\bar{\alpha_t}}\epsilon$$

where $\epsilon \sim \mathcal{N}(0, I)$ and $\alpha_t := 1 - \beta_t$ and $\bar{\alpha}_t = \prod_{s=1}^{t} \alpha_s$ (2)

The noise schedule $\beta_t$ is designed such that $\alpha_T \to 0$ and $q(x_T|x_0) = \mathcal{N}(0, I)$. Note that the Gaussian noise is used here for its mathematical properties (in particular, that the sum of two Gaussians is also a Gaussian), and does not imply any noise structure in the image $x_0$ (in particular, the use of Gaussian noise in the diffusion process is not incompatible with the fact that noise in MRI signals are not Gaussian)

The forward process provides training examples for the denoising process $p_\theta$, whose goal is to predict a less noisy sample $x_{t-1}$ from a noisy sample $x_t$. In general, the equation $q(x_{t-1}|x_t) \propto q(x_t|x_{t-1})q(x_{t-1})$ is intractable, but it can be approximated with a Gaussian for small transitions (small $\beta_t$). The equation for the reverse can be written as

$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \sigma_t^2 I) \qquad (3)$$

The variance $\sigma_t^2$ can be fixed, eliminating the need to learn it explicitly. In practice it is easier to learn a model $\epsilon_\theta(x_t, t)$ which predicts the noise that needs to be removed at each step [6]. A sample $x_{t-1}$ can then be generated from $x_t$ by

$$x_{t-1} = \frac{1}{\sqrt{\alpha_t}}(x_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}}\epsilon_\theta(x_t, t)) + \sigma_t z, \text{ with } z \sim \mathcal{N}(0, I)$$

$$(4)$$

Given pure noise $x_T = \mathcal{N}(0, I)$, an image can be synthesized by iteratively applying equation 4 for all timesteps

$t \in \{T, ..., 1\}$ to obtain the final prediction $x_0$. The diffusion model $\mu_\theta$ usually uses a U-Net architecture [18]. Training focuses on minimizing the MSE loss

$$\mathbb{E}_{x_0 \sim q(x_0), \epsilon \sim \mathcal{N}(0,I)}[||\epsilon - \epsilon_\theta(x_t, t)||^2] \qquad (5)$$

## III. METHODS

### A. DDPMs for lesion filling

CNNs operating on 3D volumes require significant GPU memory. To address this challenge, we treat the 3D volume as batches of 2D transversal slices and employ a 2D Unet for processing, stacking the results to yield a 3D image.
We explore two approaches in MR-images using diffusion models: conditional and unconditional. Both approaches utilize the ground truth MR-image $x$, a binary mask $m$ defining the lesion region, and the masked ground truth image $\hat{x}$.

*1) Conditional Model:* The conditional approach trains a diffusion model conditioned on the masked ground truth image and the binary mask. The conditioning information is incorporated through channel-wise concatenation. At each timestep $t$ during reverse diffusion, the model receives the concatenated input of the noisy image $x_t$, the masked ground truth image $\hat{x}$ and the binary mask $m$. The objective is to predict the noise term for calculating a less noisy image. This leads to the loss function,

$$\mathbb{E}_{x_0 \sim q(x_0), \epsilon \sim \mathcal{N}(0,I)}[||\epsilon - \epsilon_\theta(((\sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon) \oplus \hat{x} \oplus m), t)||^2] \qquad (6)$$

For sampling, we employ Denoising Diffusion Implicit Models (DDIM) [19], a computationally efficient class of iterative probabilistic models that share the training procedure of DDPM. DDIM utilizes a non-Markovian sampling process, which is deterministic. The sequence of a training step for the conditional model is described later in Figure 3.

*2) Unconditional Model:* The unconditional approach does not use conditioning information during training. We train an unconditional DDPM as a generative prior, as described in Section II-B. This essentially creates a model that can produce random 2D brain MRI samples. To condition the generation process, we modify the reverse diffusion iterations by sampling masked regions using the provided image information, as proposed in the RePaint paper [20]. This technique does not modify the original DDPM network and is applicable to any inpainting mask distribution.

Inpainting aims to predict missing pixels within a masked region based on surrounding image information. Each reverse step from image $x_t$ (noisy) to $x_{t-1}$ (less noisy) depends solely on the noisy image $x_t$, which consists of unknown pixels within the mask $m \odot x_t$ and known pixels outside the mask $(1 - m) \odot x_t$. The known pixels can be calculated for each timestep based on the forward process (Equation 2). The RePaint approach proposes separate sampling processes for unknown and known pixels during the reverse step, resulting in the following expression:

$$x_{t-1}^{known} \sim \mathcal{N}(\sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I) \qquad (7)$$

$$x_{t-1}^{unknown} \sim \mathcal{N}(\mu_\theta(x_t, t), \Sigma_\theta(x_t, t)) \qquad (8)$$

$$x_{t-1} = (1 - m) \odot x_{t-1}^{known} + m \odot x_{t-1}^{unknown} \qquad (9)$$

Here, $x_{t-1}^{known}$ is sampled using the known pixels in the given image $(1 - m) \odot x_t$, while $x_{t-1}^{unknown}$ is sampled from the model given the previous iteration $x_t$. These are then combined to form the new sample $x_{t-1}$ using the mask $m$.

*3) U-Net Architectures:* As model we used a 2D U-Net and a pseudo-3D U-Net [16] which achieved high scores in another study comparing different diffusion models architectures for 3D healthy brain inpainting [15]. The 2D U-Net has an architecture similar to [21]. It uses six feature map resolutions with two convolutional residual blocks per resolution level and one self-attention block at the 16x16 resolution after each convolutional block. From highest to lowest resolution the U-Net stages use (128, 128, 256, 256, 512, 512) channels.

The pseudo-3D U-Net has, in addition to the 2D U-Net, a volumetric layer inside the residual block after each 2D convolution. Pseudo-3D convolutions result from 1D convolutions in the z-axis, requiring the batch to be rearranged before and after. Following [15] we apply the model in the image space and directly use the pseudo-3D convolutions without the proposed fine-tuning strategy used by the original paper [16]. To setup the U-Net models and the training environment, we used the python library diffusers from huggingface [22].

### B. Datasets

To train the models, a dataset was created by combining healthy subjects from the OASIS project [23] with MS patients from the MICCAI challenge [24].

The OASIS dataset consisted of T1w MRI scans from 20 healthy subjects. It was divided into training and validation sets, with 16 samples used for training and 4 for validation.

The MICCAI challenge dataset comprised of MRI scans of 15 patients diagnosed with MS. Each scan included both T1w and FLAIR images, with the FLAIR images containing manually segmented lesion masks. This dataset was also split into training and validation sets, with 13 samples for training and 2 for validation.

Section III-B2 explores the use of additional synthetic masks in the form of random circle masks. To evaluate their effectiveness on larger datasets, the BraTS Inpainting Challenge 2023 dataset [13] was employed. The training set of this dataset comprises 1251 brains. Since the challenge has concluded and online analysis of the validation set is no longer possible, the training set was divided into 90% training data and 10% test data.

The impact of lesion filling on cortical thickness measurements was conducted using a test set composed of 65 patients diagnosed with RR-MS. This data originated from an internal longitudinal study conducted at the Insel hospital. All patients had been undergoing Natalizumab treatment for over two years and had at least four MRI scans performed over a period of approximately six months each, with corresponding clinical evaluations. MRI scans included a combination of 1.5T and 3T datasets with a slice thickness of 1mm or less in the T1w sequences. For each patient, the T1w and FLAIR images from their final visit, typically containing the highest lesion burden, were used for testing.

*1) Preprocessing:* All T1w images undergo resampling to a standardized size of 256x256x256 voxels with a 1.0x1.0x1.0 mm voxel size and are reoriented to RAS orientation. FLAIR images are resampled to 160x256x256 voxels. The resampling process is carried out using nibabel.processing.conform [25]. Values below 0.01 are discarded as noise and the remaining data is scaled to the range [-1, 1]. A deep learning-based tissue segmentation is performed on the T1w images for each patient using the DL+DiReCT model [26]. To accelerate this process, the parallelization program GNU Parallel [27] is employed. For datasets with existing lesion masks, these are registered from FLAIR to T1w images using NiftyReg [28]. In the absence of lesion masks in the test set, a separate segmentation model DeepSCAN [29] is utilized to identify MS lesions. Only 2D slices containing WM, based on the DL+DiReCT segmentation, were incorporated for training.



Fig. 2. Creation of the training dataset

*2) Mask Generation:* The distribution of masks employed for training a conditional model can have an impact on the performance of the model [30] [31] [32] [33]. Conditional models were trained on the healthy subject images using lesion masks obtained from the MS patients. To achieve this, each lesion mask from the MS patients was registered to every T1w image from the healthy subjects. This resulted in 15 registered lesion masks for each of the 20 healthy patients.

Each lesion mask was restricted to WM tissue by multiplying it with a binary WM mask derived from the DL+DiReCT segmentation. To augment the diversity of the masks, the set of connected lesions was computed for each mask. During training, a different set of connected lesions was sampled and used as the lesion mask.

Given the limited dataset of masks, a secondary approach was explored that utilized a second mask distribution consisting of random circle masks with varying locations and sizes. This resulted in three models being trained: One model trained on the distribution of real lesion masks (conditional lesions model), one model trained on the distribution of random circle masks (conditional circles model) and one model trained with a combination of 50% real lesion masks and 50% random circle masks (conditional mixture model).
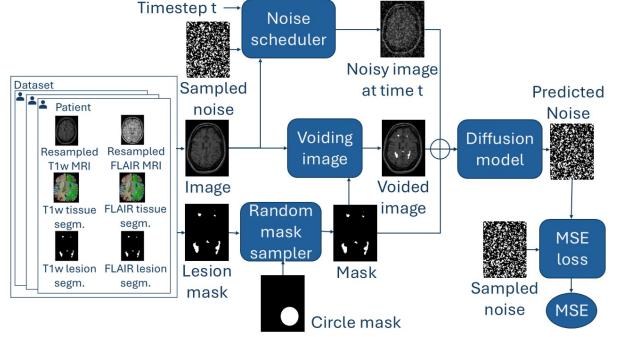


Fig. 3. Training step involving the conditional mixture model. An MR-image and its corresponding lesion mask are sampled from the dataset. Alternatively, with a 50% probability, a random circle mask is sampled instead. This mask is used to void the portion of the image, which requires inpainting. Additionally, a random timestep $t$ and random noise matching the image shape are sampled. These are used to generate a noisy image as described in Section II-B. The mask, the noisy image, and the voided image are concatenated and fed into the diffusion model, which aims to predict the sampled noise. The predicted and sampled noise are used to calculate the MSE.

*3) Min-SNR Loss weighting:* During the training of the unconditional model, the validation loss indicated faster overfitting with smaller timesteps (e.g., 200 steps) compared to larger ones (e.g., 1000 steps). This could be due to the different levels of difficulty inherent in the time steps of the diffusion models. Predicting added noise becomes progressively easier as the image approaches pure noise. Consequently, bigger timesteps naturally result in lower MSE and correspondingly weaker gradients compared to smaller timesteps. This imbalance leads the training process to prioritize optimization for smaller timesteps.

To address this and achieve a more balanced loss function, we explored the min-SNR weighting strategy proposed in [34]. This approach advocates for adapting loss weights assigned to individual timesteps based on clamped signal noise ratios.

### C. Evaluation of the DDPM

*1) Metrics:* During the training, the model was evaluated at regular intervals using the validation dataset. MSE, PSNR, and SSIM [35] are calculated inside the masks, outside the masks, and across the entire image. LPIPS [36] is solely evaluated over the whole image

All metrics are measured on 2D images. The model version with the highest SSIM score is periodically saved to disc.

*2) Mask Dilation:* The evaluation revealed the presence of artifacts at the boundaries of inpainted lesions. This occurs because RePaint replaces all areas outside the designated mask with the original image. If the annotated masks don't fully encompass the entire lesion, these small residual areas of the original lesion can lead to border artifacts. Conditional models exhibited similar, though less pronounced, artifacts. To address this issue, we implemented a minor, one-pixel dilation restricted to WM regions. This dilation strongly minimizes the artifacts.

### D. Evaluation of the robustness of cortical thickness methods

Having identified the best performing lesion filling method, further evaluations centers on the influence of lesion filling

on cortical thickness measurements. Cortical thickness assessments are performed on sixty-five patients from the test set, both before and after lesion filling, using five methods: ANTs [37], ANTsPyNet [38], FreeSurfer [39], FastSurfer [40] and DL+DiReCT [26]. FreeSurfer calculates cortical thickness by modeling the cortical band as a surface mesh. FastSurfer, a more recent method, replicates FreeSurfer's anatomical segmentation, including surface reconstruction, but leverages deep learning techniques to accelerate the process. On the other hand, ANTs,ANTSPyNet and DL+DiReCT all build on cortical thickness method based on diffeomorphic registration-based cortical thickness (DiReCT) applied to an atlas-based segmentation. ANTsPyNet extends this approach by incorporating deep learning for segmentation: ANTs applys this to a Bayesian EM-based segmentation, while ANTsPyNet and DL+DiReCT both use a deep-learning algorithm to provide the segmentation. To assess robustness, we apply each method before and after lesion filling. [41]. The average absolute changes relative to the mean (%) are calculated using the following formula:

$$\epsilon_\mu = \frac{100}{N} \sum_{i=1}^{N} \frac{1}{2} \sum_{t=1}^{2} \frac{|m_{i,t} - \mu_i|}{\mu_i}$$

Where N is the number of patients, $m_1$ the measurement before lesion filling, $m_2$ the measurement after lesion filling and $\mu_i = \frac{1}{2} \sum_{t=1}^{2} m_{i,t}$ the within-patient mean. This calculation is performed both for the global mean thickness (averaging the mean thickness of the left and right hemispheres) and across anatomical regions defined by the Desikan-Killiany (DK) atlas [42]. Since ANTs and ANTsPyNet do not deliver regional statistics, for those methods we average over the parcellation derived from DL+DiReCT. In a subanalysis, we excluded cases with MS lesions located near the cortical surface (juxtacortical lesions) which might lead to lesion-filling errors: to identify patients with juxtacortical lesions, the binary lesion masks are dilated by one pixel and multiplied with the tissue segmentations. Patients with lesions outside WM are excluded in the second analysis.

## IV. RESULTS

### A. Performance of the lesion-filling models

The 3D conditional model trained with a balanced mixture of lesion masks and random circle masks emerges as the top-performing model, attaining a SSIM of 0.96 and LPIPS of 2e-4 on the evaluation set. Metrics measured during training can be viewed in Appendix A. A comparative analysis between 2D and pseudo-3D models reveals that the latter consistently outperforms the former across all metrics. Furthermore, within the realm of conditional models, the architecture trained with random circle masks demonstrates superior performance compared to its lesion mask-trained counterpart.

|  | SSIM | PSNR | MSE | LPIPS |
|---|---|---|---|---|
| 2D unconditional RePaint | 0.83 | 28 | 8.2e-3 | 2.0e-3 |
| 2D conditional circles | 0.9 | 32 | 4e-3 | 2e-3 |
| 2D conditional lesions | 0.85 | 28 | 0.01 | 5e-3 |
| 2D conditional mixture | 0.9 | 33 | 4e-3 | 1e-3 |
| 3D unconditional RePaint | 0.90 | 32 | 3e-3 | 9e-4 |
| 3D conditional circles | 0.95 | 38 | 1e-3 | 3e-4 |
| 3D conditional lesions | 0.93 | 34 | 3e-3 | 4e-4 |
| **3D conditional mixture** | **0.96** | **39** | **8e-4** | **2e-4** |

TABLE I
METRICS MEASURED WITH VALIDATION DATASET. SSIM, PSNR AND MSE ARE MEASURED INSIDE THE MASK AND LPIPS OVER THE FULL IMAGE.
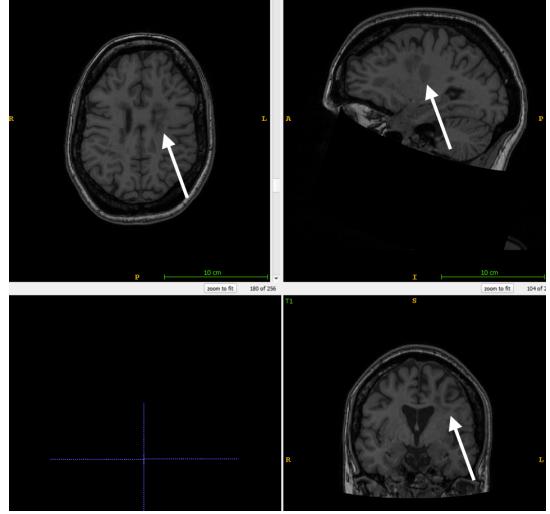


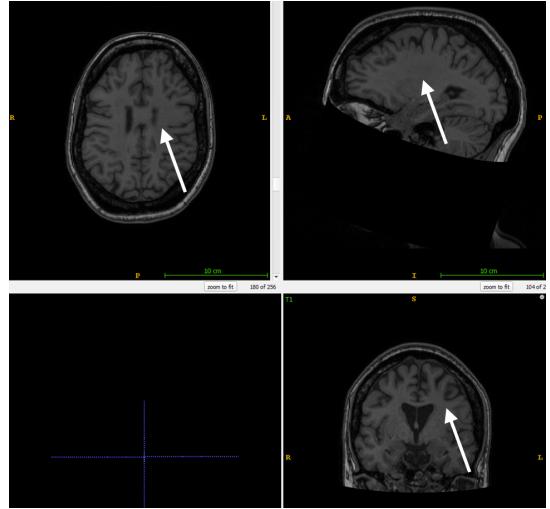Fig. 4. T1w before lesion filling with conditional mixture model



Fig. 5. T1w after lesion filling with conditional mixture model

A significant difference exists in terms of inference time. Due to the resampling approach, the inference time for the unconditional RePaint model is substantially longer compared to the conditional models. For a batch of 16 samples on a single Nvidia RTX A6000 40GB GPU, the inference time is 45 seconds for the conditional mixture model and 350 seconds for the unconditional RePaint model.

## B. Robustness of Cortical Thickness methods

Table II presents mean reproducibility errors for both global mean thickness and the average across all 68 ROIs, calculated using data from 65 patients. To account for potential influences of juxtacortical lesions, we excluded patients with such lesions and recalculated the same measurements for the remaining 17 patients, with results displayed in Table III. Lesion filling was performed using the 3D conditional mixture model. ANTsPyNet and FastSurfer, which incorporate deep learning, show significantly improved robustness compared to their predecessors. Furthermore, the DL+DiReCT approach yields a substantial additional reduction in error.

Comparing robustness across different regions reveals consistent superiority of the newer deep learning-based methods (see Figure 6). ANTsPyNet's least robust region is the left frontal pole with a 1.4% error. In contrast, FastSurfer's least robust regions are the left and right pericalcarine regions, with errors of 1.1% and 1.0% respectively. Similarly, DL+DiReCT exhibits the lowest robustness in the right and left pericalcarine regions, with error rates of 0.6% and 0.5%, respectively.

|  | Global mean thickness (%) | ROI-average (%) |
|---|---|---|
| ANTs | 1.31 | 1.68 |
| ANTsPyNet | 0.52 | 0.84 |
| FreeSurfer | 0.51 | 0.92 |
| FastSurfer | 0.14 | 0.45 |
| **DL+DiReCT** | **0.05** | **0.14** |

TABLE II
MEAN REPRODUCIBILITY ERRORS

|  | Global mean thickness (%) | ROI-average (%) |
|---|---|---|
| ANTs | 1.34 | 1.68 |
| ANTsPyNet | 0.38 | 0.81 |
| FreeSurfer | 0.61 | 0.90 |
| FastSurfer | 0.13 | 0.43 |
| **DL+DiReCT** | **0.04** | **0.12** |

TABLE III
MEAN REPRODUCIBILITY ERRORS WITHOUT PATIENTS WITH
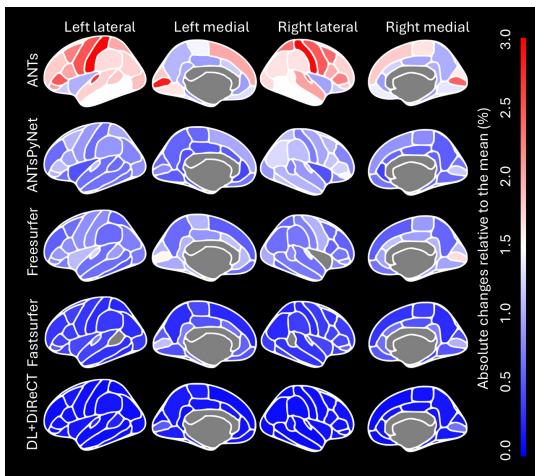JUXTACORTICAL LESIONS



Fig. 6. Color-coded reproducibility errors of the ROI-wise average cortical thicknesses evaluated on all samples.

## V. DISCUSSION

Across the whole cortex and regionally, the impact of lesion filling on cortical thickness measurements varies significantly depending on the morphometric tool employed. While the original ANTs model is strongly affected, the newer ANTsPyNet model, incorporating deep learning, shows a much smaller impact. Similarly, the newer FastSurfer model, also incorporating deep learning, improves the error rate compared to its predecessor, FreeSurfer. DL+DiReCT exhibits the smallest differences among all models. These findings suggest that deep learning models are more robust to WM lesions than classical methods.

When considering the lesion filling models, we can observe some general trends: conditional models perform better than unconditional models, and 3D models perform better than 2D models. The better performance is intuitive, as conditional models are explicitly trained for inpainting, while unconditional models are reused through the RePaint sample approach. The conditional models have significantly longer training times (days versus hours), but this is more than offset by faster image generation at test time.

Comparing different training regimes, we found that training with additional circle masks enhanced performance. Surprisingly, training exclusively with random circle masks yielded better results than using only lesion masks, suggesting that a broader range of unrealistic masks distributed across MR-images is beneficial within the given dataset, rather than a smaller set of masks sampled from the true mask distribution.

## VI. CONCLUSION

We successfully developed a deep learning model for filling MS lesions in MR-images, using it to observe the superior robustness of morphometric pipelines based on deep learning segmentations. This raises the possibility that lesion filling might become obsolete with the increasing adoption of more modern tools: meanwhile, researchers preferring to remain with more established tools may find it useful to perform lesion filling before analysis.

### A. Limitations

The dataset used in this study is relatively small, which might limit the representativeness of the ethnic groups included. The current study focused on filling multiple lesions simultaneously. The performance of the model for inpainting single lesions while preserving others remains unexplored and could potentially differ. We utilized lesion masks created by doctors based on their interpretation of MR-images. It's important to note that lesions can also influence the surrounding brain tissue, which may not be readily identifiable by humans on current MR-images. The extent of this influence and its relevance for lesion filling is a separate research question and may vary depending on the specific use case.

### B. Outlook

The models developed for filling MS lesions, could be applied to other inpainting tasks. However, performance might

vary, especially considering the training objectives. It would be interesting to determine if the performance advantage of conditional models over unconditional ones persists in these new applications. Improving the unconditional model to match the performance of the conditional model is another potential area of research. This is desirable due to the unconditional model's shorter training time and independence from mask distribution.

Although we demonstrated that deep learning-based tools are more robust to MS lesions than older methods when measuring cortical thickness, a larger population study is necessary to definitively establish the obsolescence of lesion filling.

## REFERENCES

[1] M. Battaglini, M. Jenkinson, and N. De Stefano, "Evaluating and reducing the impact of white matter lesions on brain volume measurements," *Human Brain Mapping*, vol. 33, no. 9, 2012.

[2] S. Valverde, A. Oliver, and X. Lladó, "A white matter lesion-filling approach to improve brain tissue volume measurements," *NeuroImage: Clinical*, vol. 6, 2014.

[3] H. Amiri, A. de Sitter, K. Bendfeldt, M. Battaglini, C. A. Gandini Wheeler-Kingshott, M. Calabrese, J. J. Geurts, M. A. Rocca, J. Sastre-Garriga, C. Enzinger, N. de Stefano, M. Filippi, A. Rovira, F. Barkhof, and H. Vrenken, "Urgent challenges in quantification and interpretation of brain grey matter atrophy in individual MS patients using MRI," 2018.

[4] M. R. Farazi, F. Faisal, Z. Zaman, and S. Farhan, "Inpainting multiple sclerosis lesions for improving registration performance with brain atlas," in *1st International Conference on Medical Engineering, Health Informatics and Technology, MediTec 2016*, 2017.

[5] M. Almansour, N. M. Ghanem, and S. Bassiouny, "High-resolution MRI brain inpainting," in *BHI 2021 - 2021 IEEE EMBS International Conference on Biomedical and Health Informatics, Proceedings*, 2021.

[6] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," in *Advances in Neural Information Processing Systems*, vol. 2020-December, 2020.

[7] A. Kazerouni, E. K. Aghdam, M. Heidari, R. Azad, M. Fayyaz, I. Hacihaliloglu, and D. Merhof, "Diffusion Models for Medical Image Analysis: A Comprehensive Survey," *arXiv*, 2022.

[8] V. E. Tiu, I. Enache, C. A. Panea, C. Tiu, and B. O. Popescu, "Predictive MRI Biomarkers in MS—A Critical Review," 2022.

[9] F. Bieder, J. Wolleb, A. Durrer, R. Sandkühler, and P. C. Cattin, "Memory-Efficient 3D Denoising Diffusion Models for Medical Image Processing," in *Proceedings of Machine Learning Research*, vol. 227, 2023.

[10] S. Magon, L. Gaetano, M. M. Chakravarty, J. P. Lerch, Y. Naegelin, C. Stippich, L. Kappos, E. W. Radue, and T. Sprenger, "White matter lesion filling improves the accuracy of cortical thickness measurements in multiple sclerosis patients: A longitudinal study," *BMC Neuroscience*, vol. 15, no. 1, 2014.

[11] R. Gelineau-Morel, V. Tomassini, M. Jenkinson, H. Johansen-Berg, P. M. Matthews, and J. Palace, "The effect of hypointense white matter lesions on automated gray matter segmentation in multiple sclerosis." *Human brain mapping*, vol. 33, no. 12, pp. 2802–14, 12 2012.

[12] C. W. van der Weijden, M. S. Pitombeira, Y. R. Haveman, C. A. Sanchez-Catasus, K. R. Campanholo, G. D. Kolinger, C. M. Rimkus, C. A. Buchpiguel, R. A. Dierckx, R. J. Renken, J. F. Meilof, E. F. de Vries, and D. de Paula Faria, "The effect of lesion filling on brain network analysis in multiple sclerosis using structural magnetic resonance imaging," *Insights into Imaging*, vol. 13, no. 1, 2022.

[13] F. Kofler, F. Meissen, F. Steinbauer, R. Graf, E. Oswald, E. de da Rosa, H. B. Li, U. Baid, F. Hoelzl, O. Turgut, I. Horvath, D. Waldmannstetter, C. Bukas, M. Adewole, S. M. Anwar, A. Janas, A. F. Kazerooni, D. LaBella, A. W. Moawad, K. Farahani, J. Eddy, T. Bergquist, V. Chung, R. T. Shinohara, F. Dako, W. Wiggins, Z. Reitman, C. Wang, X. Liu, Z. Jiang, A. Familiar, G.-M. Conte, E. Johanson, Z. Meier, C. Davatzikos, J. Freymann, J. Kirby, M. Bilello, H. M. Fathallah-Shaykh, R. Wiest, J. Kirschke, R. R. Colen, A. Kotrotsou, P. Lamontagne, D. Marcus, M. Milchenko, A. Nazeri, M.-A. Weber, A. Mahajan, S. Mohan, J. Mongan, C. Hess, S. Cha, J. Villanueva-Meyer, E. Colak, P. Crivellaro, A. Jakab, J. Albrecht, U. Anazodo, M. Aboian,

J. E. Iglesias, K. Van Leemput, S. Bakas, D. Rueckert, B. Wiestler, I. Ezhov, M. Piraud, and B. Menze, "The Brain Tumor Segmentation (BraTS) Challenge 2023: Local Synthesis of Healthy Brain Tissue via Inpainting," 5 2023.

[14] A. Durrer, P. C. Cattin, and J. Wolleb, "Denoising Diffusion Models for Inpainting of Healthy Brain Tissue," 2 2024.

[15] A. Durrer, J. Wolleb, F. Bieder, P. Friedrich, L. Melie-Garcia, M. Ocampo-Pineda, C. I. Bercea, I. E. Hamamci, B. Wiestler, M. Piraud, O. Yaldizli, C. Granziera, B. H. Menze, P. C. Cattin, and F. Kofler, "Denoising Diffusion Models for 3D Healthy Brain Tissue Inpainting," 3 2024.

[16] L. Zhu, Z. Xue, Z. Jin, X. Liu, J. He, Z. Liu, and L. Yu, "Make-A-Volume: Leveraging Latent Diffusion Models for Cross-Modality 3D Brain MRI Synthesis," 7 2023.

[17] J. Sohl-Dickstein, E. A. Weiss, N. Maheswaranathan, and S. Ganguli, "Deep unsupervised learning using nonequilibrium thermodynamics," in *32nd International Conference on Machine Learning, ICML 2015*, vol. 3, 2015.

[18] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional Networks for Biomedical Image Segmentation," 5 2015.

[19] J. Song, C. Meng, and S. Ermon, "DENOISING DIFFUSION IMPLICIT MODELS," in *ICLR 2021 - 9th International Conference on Learning Representations*, 2021.

[20] A. Lugmayr, M. Danelljan, A. Romero, F. Yu, R. Timofte, and L. Van Gool, "RePaint: Inpainting using Denoising Diffusion Probabilistic Models," in *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, vol. 2022-June, 2022.

[21] A. Nichol and P. Dhariwal, "Improved Denoising Diffusion Probabilistic Models," 2 2021.

[22] P. von Platen, S. Patil, A. Lozhkov, P. Cuenca, N. Lambert, K. Rasul, M. Davaadorj, D. Nair, S. Paul, S. Liu, W. Berman, X. Xu, and T. Wolf, "Diffusers: State-of-the-art diffusion models." [Online]. Available: https://github.com/huggingface/diffusers

[23] D. S. Marcus, T. H. Wang, J. Parker, J. G. Csernansky, J. C. Morris, and R. L. Buckner, "Open Access Series of Imaging Studies (OASIS): Cross-sectional MRI data in young, middle aged, nondemented, and demented older adults," *Journal of Cognitive Neuroscience*, vol. 19, no. 9, 2007.

[24] O. Commowick, M. Kain, R. Casey, R. Ameli, J. C. Ferré, A. Kerbrat, T. Tourdias, F. Cervenansky, S. Camarasu-Pop, T. Glatard, S. Vukusic, G. Edan, C. Barillot, M. Dojat, and F. Cotton, "Multiple sclerosis lesions segmentation from multiple experts: The MICCAI 2016 challenge dataset," *NeuroImage*, vol. 244, 2021.

[25] M. Brett, C. J. Markiewicz, M. Hanke, M.-A. Côté, B. Cipollini, P. McCarthy, D. Jarecka, C. P. Cheng, E. Larson, Y. O. Halchenko, M. Cottaar, S. Ghosh, D. Wassermann, S. Gerhard, G. R. Lee, Z. Baratz, H.-T. Wang, D. Papadopoulos Orfanos, E. Kastman, J. Kaczmarzyk, R. Guidotti, J. Daniel, O. Duek, A. Rokem, M. Scheltienne, C. Madison, A. Sólon, B. Moloney, F. C. Morency, M. Goncalves, R. Markello, C. Riddell, C. Burns, J. Millman, A. Gramfort, J. Leppäkangas, J. J. F. den Bosch, R. D. Vincent, H. Braun, K. Subramaniam, A. Van, K. J. Gorgolewski, P. R. Raamana, J. Klug, R. de Wael, B. N. Nichols, E. M. Baker, S. Hayashi, B. Pinsard, C. Haselgrove, M. Hymers, O. Esteban, S. Koudoro, F. Pérez-García, J. Dockès, N. N. Oosterhof, B. Amirbekian, H. Christian, I. Nimmo-Smith, L. Nguyen, P. Suter, S. Reddigari, S. St-Jean, E. Panfilov, E. Garyfallidis, G. Varoquaux, J. H. Legarreta, K. S. Hahn, L. Waller, O. P. Hinds, B. Fauber, B. Dewey, F. Perez, J. Roberts, J.-B. Poline, J. Stutters, K. Jordan, M. Cieslak, M. E. Moreno, T. Hrnčiar, V. Haenel, Y. Schwartz, B. C. Darwin, B. Thirion, C. Gauthier, I. Solovey, I. Gonzalez, J. Palasubramaniam, J. Lecher, K. Leinweber, K. Raktivan, M. Calábková, P. Fischer, P. Gervais, S. Gadde, T. Ballinger, T. Roos, V. R. Reddam, and freec84, "nipy/nibabel: 5.2.0," 12 2023. [Online]. Available: https://doi.org/10.5281/zenodo.10363247

[26] M. Rebsamen, C. Rummel, M. Reyes, R. Wiest, and R. McKinley, "Direct cortical thickness estimation using deep learning-based anatomy segmentation and cortex parcellation," *Human Brain Mapping*, vol. 41, no. 17, 2020.

[27] O. Tange, *GNU Parallel 2018*. Ole Tange, 3 2018. [Online]. Available: https://doi.org/10.5281/zenodo.1146014

[28] M. Modat, G. R. Ridgway, Z. A. Taylor, M. Lehmann, J. Barnes, D. J. Hawkes, N. C. Fox, and S. Ourselin, "Fast free-form deformation using graphics processing units," *Computer Methods and Programs in Biomedicine*, vol. 98, no. 3, pp. 278–284, 6 2010.

[29] R. McKinley, R. Wepfer, F. Aschwanden, L. Grunder, R. Muri, C. Rummel, R. Verma, C. Weisstanner, M. Reyes, A. Salmen, A. Chan, F. Wagner, and R. Wiest, "Simultaneous lesion and brain segmentation

in multiple sclerosis using deep neural networks," *Scientific Reports*, vol. 11, no. 1, 2021.

[30] R. Suvorov, E. Logacheva, A. Mashikhin, A. Remizova, A. Ashukha, A. Silvestrov, N. Kong, H. Goka, K. Park, and V. Lempitsky, "Resolution-robust Large Mask Inpainting with Fourier Convolutions," in *Proceedings - 2022 IEEE/CVF Winter Conference on Applications of Computer Vision, WACV 2022*, 2022.

[31] J. Yu, Z. Lin, J. Yang, X. Shen, X. Lu, and T. Huang, "Free-form image inpainting with gated convolution," in *Proceedings of the IEEE International Conference on Computer Vision*, vol. 2019-October, 2019.

[32] Y. Zeng, Z. Lin, J. Yang, J. Zhang, E. Shechtman, and H. Lu, "High-Resolution Image Inpainting with Iterative Confidence Feedback and Guided Upsampling," in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 12364 LNCS, 2020.

[33] G. Liu, F. A. Reda, K. J. Shih, T. C. Wang, A. Tao, and B. Catanzaro, "Image Inpainting for Irregular Holes Using Partial Convolutions," in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, vol. 11215 LNCS, 2018.

[34] T. Hang, S. Gu, C. Li, J. Bao, D. Chen, H. Hu, X. Geng, and B. Guo, "Efficient Diffusion Training via Min-SNR Weighting Strategy," in *Proceedings of the IEEE International Conference on Computer Vision*, 2023.

[35] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: From error visibility to structural similarity," *IEEE Transactions on Image Processing*, vol. 13, no. 4, 2004.

[36] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, "The Unreasonable Effectiveness of Deep Features as a Perceptual Metric," in *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2018.

[37] N. J. Tustison, P. A. Cook, A. Klein, G. Song, S. R. Das, J. T. Duda, B. M. Kandel, N. van Strien, J. R. Stone, J. C. Gee, and B. B. Avants, "Large-scale evaluation of ANTs and FreeSurfer cortical thickness measurements," *NeuroImage*, vol. 99, 2014.

[38] N. J. Tustison, P. A. Cook, A. J. Holbrook, H. J. Johnson, J. Muschelli, G. A. Devenyi, J. T. Duda, S. R. Das, N. C. Cullen, D. L. Gillen, M. A. Yassa, J. R. Stone, J. C. Gee, and B. B. Avants, "The ANTsX ecosystem for quantitative biological and medical imaging," *Scientific reports*, vol. 11, no. 1, 2021.

[39] B. Fischl, "FreeSurfer," *NeuroImage*, vol. 62, no. 2, pp. 774–781, 8 2012.

[40] L. Henschel, S. Conjeti, S. Estrada, K. Diers, B. Fischl, and M. Reuter, "FastSurfer - A fast and accurate deep learning based neuroimaging pipeline," *NeuroImage*, vol. 219, p. 117012, 10 2020.

[41] J. Jovicich, M. Marizzoni, R. Sala-Llonch, B. Bosch, D. Bartrés-Faz, J. Arnold, J. Benninghoff, J. Wiltfang, L. Roccatagliata, F. Nobili, T. Hensch, A. Tränkner, P. Schönknecht, M. Leroy, R. Lopes, R. Bordet, V. Chanoine, J. P. Ranjeva, M. Didic, H. Gros-Dagnac, P. Payoux, G. Zoccatelli, F. Alessandrini, A. Beltramello, N. Bargalló, O. Blin, and G. B. Frisoni, "Brain morphometry reproducibility in multi-center 3T MRI studies: A comparison of cross-sectional and longitudinal segmentations," *NeuroImage*, vol. 83, 2013.

[42] R. S. Desikan, F. Ségonne, B. Fischl, B. T. Quinn, B. C. Dickerson, D. Blacker, R. L. Buckner, A. M. Dale, R. P. Maguire, B. T. Hyman, M. S. Albert, and R. J. Killiany, "An automated labeling system for subdividing the human cerebral cortex on MRI scans into gyral based regions of interest," *NeuroImage*, vol. 31, no. 3, 2006.

# APPENDIX A
## QUANTITATIVE AND QUALITATIVE TRAINING PROGRESSION

This chapter presents both quantitative and qualitative visualizations of the different model's training progression.
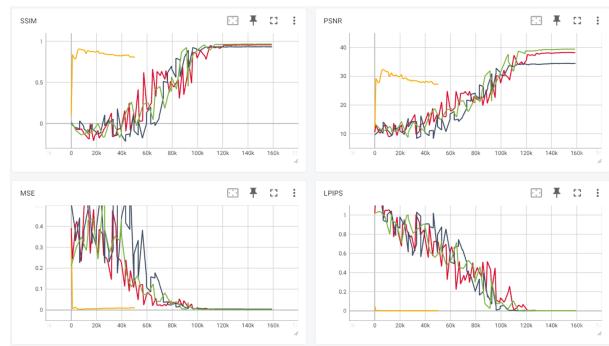


Fig. 7. Training metrics of the 3D model's conditional mixture (green), conditional circles (red), conditional lesions (black) and unconditional RePaint (orange).
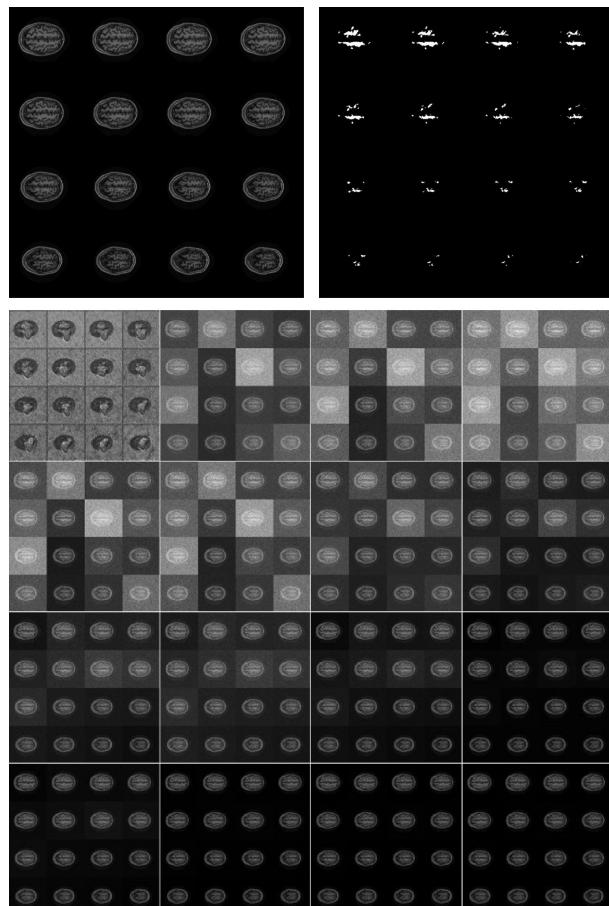


Fig. 8. 16 image-mask pairs for evaluation at the top right and left corners (Top). Below, the results of the 3D conditional mixture model training in a grid ordered from left to right and top to bottom at specific timesteps: 000001, 011'914, 017'210, 026'476, 034'419, 039'715, 051'628, 056'924, 066'190, 074'133, 079'429, 091'342, 096'638, 105'904, 113'847, 119'143
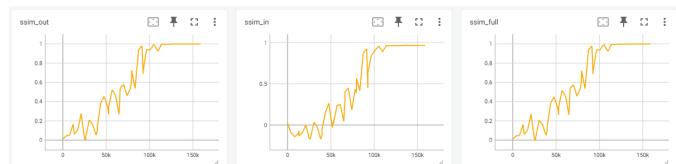


Fig. 10. SSIM metric of the 3D conditional mixture model outside and inside the mask and over the entire image.
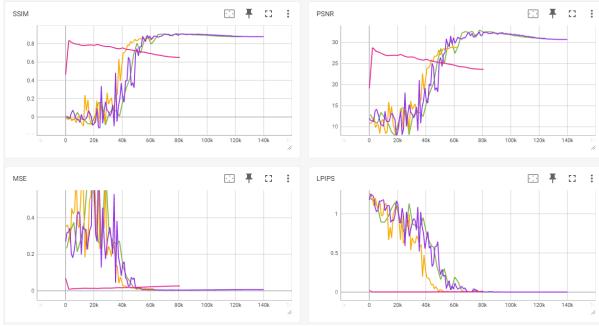
Fig. 9. Training metrics of the 2D model's conditional mixture (green), conditional circles (purple), conditional lesions (orange) and unconditional RePaint (red).

## APPENDIX B
## TRAINING ENVIRONMENT

| Number of training diffusion steps | 1000 |
|---|---|
| Number of inference steps | 50 |
| Batch size | 16 |
| Learning rate | 1e-4 |
| Optimizer | AdamW |
| Learning rate scheduler | Cosine with 500 steps warmup |
| RePaint Jump length | 8 |
| RePaint Resample | 10 |

TABLE IV
HYPERPARAMETERS FOR TRAINING AND EVALUATION

| GPU | 3x Nvidia RTX A6000 40GB |
|---|---|
| CPU | 64x Intel Xeon Gold 6226R @ 2.9Ghz |
| RAM | 196 GB |

TABLE V
HARDWARE

### A. Artifacts XXXappendix?

Incomplete lesion masking leads to recognizable artifacts in the form of residual borders at the lesion edges. This phenomenon is particularly pronounced in the unconditional RePaint approach, which replaces regions outside the mask with original image content. Although less obvious, conditional models also exhibit similar artifacts. A small one-pixel dilation limited to the WM effectively mitigates this problem.

The intrinsic two-dimensional nature of 2D models leads to another artifact: inconsistencies along the z-axis, manifesting as visible stripes. The pseudo-3D models successfully mitigate these z-axis irregularities.
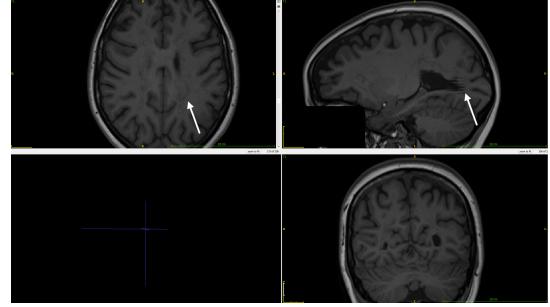


Fig. 11. Border (left arrow) and stripe artifacts (right arrow)

### B. Training Duration XXXapendix?

The unconditional RePaint model exhibits significantly faster convergence compared to conditional models, achieving a peak SSIM of 0.9 after only 6000 training steps, while conditional models require approximately 90,000 steps to reach comparable performance.



Fig. 12. SSIM score during training of the 4 3D models unconditional RePaint (violet), conditional mixture (red), conditional circles (yellow) and conditional lesions (blue).

Interestingly, RePaint achieves optimal performance when the underlying unconditional model remains unconverged. This phenomenon is evident when sampling random 2D images using a DDIM sampler instead of the RePaint sampler, resulting in highly noisy outputs. While RePaint's strong guidance produces high-quality inpainting results, DDIM sampling reveals the underlying unconditional model's immaturity. This raises the question of whether preventing overfitting and refining the unconditional RePaint model can match or surpass the performance of conditional models while requiring substantially less training time.
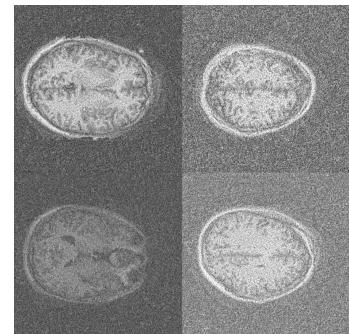


Fig. 13. Samples of unconditional model with DDIM sampler at a timestep in training, where the model achieves its best scores with the RePaint sampler.

Examining the validation loss per timestep reveals that smaller timesteps begin to overfit while larger timesteps con-

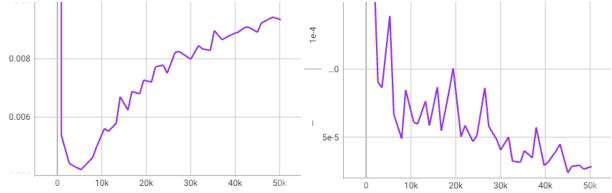tinue learning. To counteract this imbalance, we adopted the min-SNR weighting strategy outlined in Section III-B3.



Fig. 14. Comparison of the validation loss of timestep 200 (left) and 800 (right).

|  | SSIM | PSNR | MSE | LPIPS |
|---|---|---|---|---|
| 2D Unconditional RePaint | 0.83 | 28 | 8.2e-3 | **2.0e-3** |
| 2D Uncon. RePaint min-SNR | 0.83 | **29** | **7.3e-3** | 2.8e-3 |

TABLE VI
MIN-SNR METRICS MEASURED WITH VALIDATION DATASET. SSIM, PSNR AND MSE ARE MEASURED INSIDE THE MASK AND LPIPS OVER THE FULL IMAGE.

Min-SNR loss mitigated overfitting, reducing the overall validation loss across all timesteps to 0.008 compared to 0.013 with unweighted MSE loss. However, metrics such as SSIM, PSNR, MSE, and LPIPS showed no significant improvement. To simplify the training process, min-SNR loss was excluded from the model training.