
GUIDED RECONSTRUCTION WITH CONDITIONED DIFFUSION MODELS FOR UNSUPERVISED ANOMALY DETECTION IN BRAIN MRIs

Finn Behrendt
Hamburg University of Technology
Hamburg, Germany

Debayan Bhattacharya
Hamburg University of Technology
Hamburg, Germany

Robin Mieling
Hamburg University of Technology
Hamburg, Germany

Lennart Maack
Hamburg University of Technology
Hamburg, Germany

Julia Krüger
Jung Diagnostics GmbH
Hamburg, Germany

Roland Opfer
Jung Diagnostics GmbH
Hamburg, Germany

Alexander Schlaefer
Hamburg University of Technology
Hamburg, Germany

ABSTRACT

Unsupervised anomaly detection in Brain MRIs aims to identify abnormalities as outliers from a healthy training distribution. Reconstruction-based approaches that use generative models to learn to reconstruct healthy brain anatomy are commonly used for this task. Diffusion models are an emerging class of deep generative models that show great potential regarding reconstruction fidelity. However, they face challenges in preserving intensity characteristics in the reconstructed images, limiting their performance in anomaly detection. To address this challenge, we propose to condition the denoising mechanism of diffusion models with additional information about the image to reconstruct coming from a latent representation of the noise-free input image. This conditioning enables high-fidelity reconstruction of healthy brain structures while aligning local intensity characteristics of input-reconstruction pairs. We evaluate our method’s reconstruction quality, domain adaptation features and finally segmentation performance on publicly available data sets with various pathologies. Using our proposed conditioning mechanism we can reduce the false-positive predictions and enable a more precise delineation of anomalies which significantly enhances the anomaly detection performance compared to established state-of-the-art approaches to unsupervised anomaly detection in brain MRI. Furthermore, our approach shows promise in domain adaptation across different MRI acquisitions and simulated contrasts, a crucial property of general anomaly detection methods.

Keywords Unsupervised Anomaly Detection · Zero-Shot Segmentation · Brain MRI · Diffusion Models

1 Introduction

The interpretation of MRI scans is a critical task in medical imaging, providing valuable diagnostic information for various neurological conditions Vernooij et al. (2007); Lundervold and Lundervold (2019). However, this process is error-prone, time-consuming, and places a significant workload on available radiologists Bruno et al. (2015); McDonald et al. (2015). To address these challenges and improve diagnostic efficiency, deep learning techniques like convolutional neural networks (CNN) have shown great promise in assisting radiologists by automating certain aspects of the analysis Lundervold and Lundervold (2019). A common task is the detection and delineation of Pathological structures in the MRI scans such as tumors Perkuhn et al. (2018), White Matter lesions Moeskops et al. (2018) or Alzheimer’s disease Islam and Zhang (2018). Supervised deep learning approaches have been proposed for these tasks, relying on

large-scale and balanced data sets for training, especially in MRI imaging, where there is considerable heterogeneity across hospitals and scanners. However, gathering such data sets is a cumbersome and costly process. Therefore, exploring alternative approaches that liberate the dependence on annotated data while being able to detect and locate anomalies, holds great potential.

Unsupervised anomaly detection (UAD) in neuroimaging is an active research area with the potential to identify abnormalities without relying on costly data annotation. In UAD methods the goal is to learn the underlying data distribution of healthy brain MRI scans and detect anomalies as outliers from that learned distribution. Reconstruction-based UAD, in particular, trains generative models to reconstruct healthy anatomy, enabling the identification of anomalies through discrepancies between (unhealthy) inputs and pseudo-healthy reconstructions at test time Baur et al. (2021); Kascenas et al. (2022); Chen et al. (2020); Pinaya et al. (2022b). This is based on the assumption that anomalous structures (e.g. tumors) in the input image are replaced by an estimation of healthy anatomical structures similar to the training distribution. This approach addresses limitations of supervised learning, such as the requirement for large-scale and balanced annotated data sets Johnson and Khoshgoftaar (2019); Karimi et al. (2020); Ellis et al. (2022), which is particularly crucial in neuroimaging where pathologies can exhibit complex and variable morphological characteristics. Additionally, the ability of UAD methods to perform zero-shot detection and segmentation of previously unseen pathologies makes them highly practical in a wide range of clinical scenarios.

Recent advancements have shown promise in utilizing denoising diffusion probabilistic models (DDPMs) Ho et al. (2020) for UAD in neuroimaging Pinaya et al. (2022a); Wyatt et al. (2022); Graham et al. (2022). DDPMs generate images by denoising images that are corrupted by artificial noise, leveraging a high-dimensional latent space to preserve spatial context and achieve high-fidelity reconstructions. However, a significant challenge remains in accurately reconstructing healthy brain anatomy that exhibits anatomical coherence (specific brain structures such as ventricles should appear at the same location in both, input and reconstruction) and aligned intensity characteristics with the input image Behrendt et al. (2023). The forward and backward processes of DDPMs do not adequately capture the highly variable local intensity distributions of MRI scans, leading to false positives in the residual map and impaired detection performance, particularly when facing domain shifts at test time.

To address this challenge, we propose context-conditioned DDPMs (cDDPMs) for UAD in brain MRI. Our approach involves training a DDPM to reconstruct healthy brain anatomy and incorporating a latent feature representation of the noise-free input image into the denoising process. Thereby, we utilize a CNN-based image encoder to obtain the feature representation. This representation is then used to linearly transform the feature maps of the denoising Unet to incorporate the information in the denoising process. While the dense feature representation is not suitable for high-fidelity reconstruction, it captures local intensity information of the input image that is partially lost during the forward process of DDPMs. Hence, our approach is designed to align the intensities of input and reconstruction, which is considered an important property for reconstruction-based UAD.

To gauge the impact of our conditioning approach on the quality of reconstruction, we analyze the intensity-based, structural, and perceptual similarity between input and reconstructions using healthy brain MRIs. Furthermore, we explore the domain adaptation capabilities of our approach by evaluating the histogram alignment of input and reconstruction for out-of-domain data sets, unseen during training, and additionally simulate different contrast levels. Finally, we analyze the unsupervised anomaly segmentation performance on diverse data sets, encompassing various pathologies that are unseen during training.

Our approach demonstrates superior or competitive UAD performance compared to recent state-of-the-art architectures on all tested data sets. It effectively addresses the domain shift inherent in different MRI data sets, showcasing its potential for identifying pathologies even in the absence of large-scale annotated data sets.

In summary, the main contributions of this work are:

- We propose conditioned DDPMs to incorporate additional information of the noise-free input image during the denoising process of DDPMs.
- We systematically analyze the effect our conditioning approach has on the reconstructed images and domain adaptation capabilities and thereby show its promising features for the UAD task.
- We demonstrate the effectiveness of our conditioning approach for the UAD task by outperforming state-of-the-art solutions and show its robustness to distribution shifts on various data sets.

This paper is organized as follows: In Section 2, we provide a review of relevant literature in the field of UAD in brain MRI. In Section 3, we introduce DDPMs and subsequently explain our conditioning approach. In Section 4, we provide details of the experimental setup. In Section 5, we present the results and subsequently discuss them in Section 6. Finally, we provide a conclusion in Section 7.

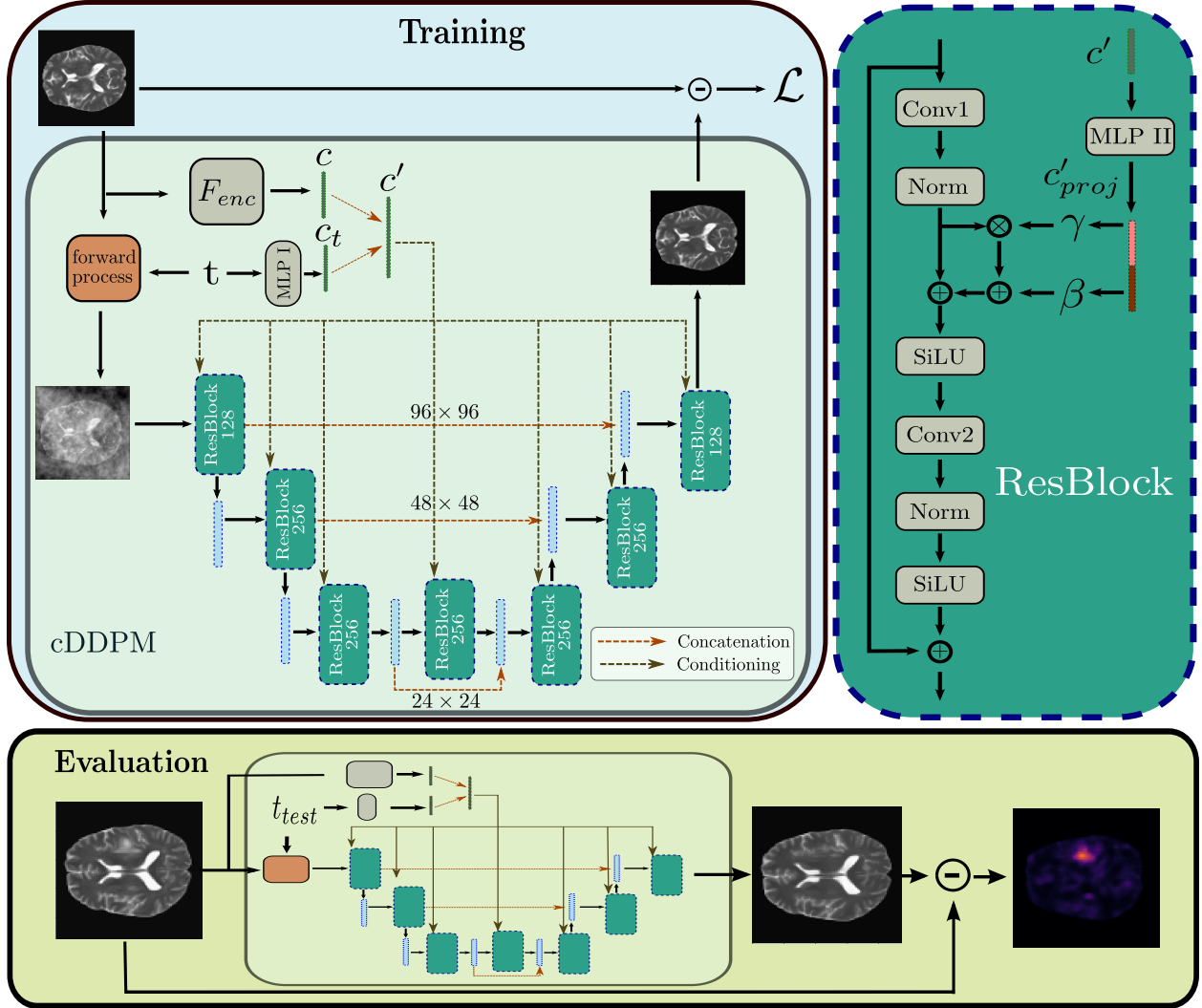


Figure 1: Overview of our proposed approach. The encoder representations are learned along with the DDPM in the main training stage to condition the denoising process. The feature maps of the denoising Unet are scaled and shifted based on the projected encoder representation of the input image in each residual block. Each residual block consists of two convolution operators (conv1, conv2), group normalization (Norm) and Sigmoid Linear Units (SiLU). During evaluation, the residual map between unhealthy brain images and their healthy reconstructions is used for anomaly detection.

2 Recent Work

Autoencoders (AE) have been the primary focus of recent research on reconstruction-based UAD in brain MRI. Although these models exhibit potential in capturing the underlying healthy distribution, their effectiveness in UAD is limited by their blurry reconstructions Baur et al. (2021). To overcome this limitation, researchers have focused on improving the representations and reconstructions by adding skip connections with dropout Baur et al. (2020a), using multi-scale features Baur et al. (2020b), or utilizing feature activation maps Silva-Rodríguez et al. (2022). Furthermore, online outlier removal has been proposed for AEs Behrendt et al. (2022b). In parallel, Variational Autoencoders (VAE) have been investigated for the UAD task Zimmerer et al. (2019a), focusing on enhancing the used context in 2D Zimmerer et al. (2019b) and 3D Bengs et al. (2021); Behrendt et al. (2022a) or utilizing restoration methods Chen et al. (2020). Also, Generative Adversarial Networks (GAN) have been proposed for UAD either as pure GAN Han et al. (2021); Schlegl et al. (2019) or in combination with VAEs Baur et al. (2018) and VQ-VAEs Pinaya et al. (2022b). While AEs with skip connections and a spatial latent space enable reconstructions of high fidelity, they tend to perform

a ‘copy task’ which enables the reconstruction of unhealthy anatomy and therefore contradicts the UAD principle Baur et al. (2021); Bercea et al. (2023b). Lately, Kascenas et al. (2022) have shown that AEs with skip connections can be effectively used for UAD in brain MRI if they are regularized by an additional denoising task. Congruently, DDPMs have shown promise in the field of UAD Wyatt et al. (2022); Pinaya et al. (2022a); Graham et al. (2022). DDPMs provide high reconstruction fidelity, but due to the noising process, important information about the input image can be lost. To address this, Behrendt et al. (2023) proposed patch-based DDPMs, that allow the use of parts of the original image content to provide information for the reconstruction of the input image. However, using this patching strategy increases complexity and computational effort and can lead to artifacts in regions of overlapping patches. A more efficient approach is seen in conditioning the denoising process of DDPMs with knowledge of the input image. Conditioned DDPMs have been successful in text-to-image synthesis tasks Rombach et al. (2022); Dhariwal and Nichol (2021) and image-guided synthetic image generation Saharia et al. (2022); Wang et al. (2022); Wolleb et al. (2022). However, in the specific case of UAD, the objective is not to generate new images or to transfer styles, but to accurately estimate a given input image while ensuring that unhealthy anatomy is absent in the estimation. Thus, directly conditioning DDPMs with information from the input image can pose a risk of reconstructing unhealthy anatomy. Therefore, to achieve the UAD task, we develop a conditioning approach for DDPMs that can effectively provide the denoising process with relevant context information of the individual input image without enabling the reconstruction of unhealthy anatomy.

3 Methods

We propose cDDPMs where we use an image encoder network and embed the input image in a context vector $\mathbf{c} \in \mathbb{R}^d$ to condition the denoising Unet on meaningful features of the input image. Our motivation is that the additional information in \mathbf{c} guides the generation process towards consistent intensity characteristics across the input image and its reconstruction. Hence, by introducing the context vector \mathbf{c} we aim to recover local intensity information that is lost during the forward (noising) process of DDPMs. We utilize an image encoder with a dense latent space to extract information regarding the coarse shape and local intensity information of the noise-free input image. This latent representation can then be used to condition the denoising process and supplement the individual context of the input image without providing detailed pixel-wise information that could be used to perform a copy task. A general depiction of our approach is shown in Fig. 1.

3.1 DDPMs

DDPMs are generative models that learn the underlying data distribution of images $\mathbf{x} \in \mathbb{R}^{H,W,C}$ with height H , width W and C channels, given a training set. Training of DDPMs consists of two steps. The forward process, where an input image \mathbf{x}_0 is gradually transformed to Gaussian noise $\mathbf{x}_T = \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and the backward process, where reversing the forward process is learned.

In the forward process, transforming \mathbf{x}_0 to \mathbf{x}_T follows a predefined schedule β_1, \dots, β_T , where intermediate versions \mathbf{x}_t are derived as

$$\mathbf{x}_t \sim q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}),$$

$$\text{with } \bar{\alpha}_t = \prod_{s=0}^t (1 - \beta_s).$$

The time step t controls the amount of added noise and is sampled from $t \sim \text{Uniform}(1, \dots, T)$. For edge cases, the image \mathbf{x}_t is transformed to pure noise ($t = T$) or no transformation is applied ($t = 0$). In the backward process, the reconstructed image \mathbf{x}_0^{rec} is recovered from \mathbf{x}_t by

$$\mathbf{x}_0^{rec} \sim p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t),$$

$$\text{with } p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t)).$$

Here, following Ho et al. (2020), $\boldsymbol{\mu}_\theta$ is estimated by a Unet Ronneberger et al. (2015) with trainable parameters θ , and $\boldsymbol{\Sigma}_\theta(t) = \boldsymbol{\Sigma}(t) = \frac{1 - \alpha_t - 1}{1 - \alpha_t} \beta_t \mathbf{I}$ is fixed. Variational inference is used to achieve a tractable loss function and the variational lower bound (VLB) is derived as

$$\mathcal{L}_{VLB} = -\log(p_\theta(\mathbf{x}_0))$$

$$+ D_{KL}(q(\mathbf{x}_{1:T} | \mathbf{x}_0) || p_\theta(\mathbf{x}_{1:T} | \mathbf{x}_0)).$$

which can be reformulated to

$$\mathcal{L}_{simple} = \|\epsilon - \epsilon_\theta(\mathbf{x}_t, t)\|^2$$

by applying simplifications and by conditioning the denoising step on \mathbf{x}_0 , as shown in Ho et al. (2020). In our work, instead of predicting the noise ϵ we perform the equivalent task of directly estimating $\mathbf{x}_0^{rec} = \mathbf{x}_t - \epsilon$. Hence, we derive our loss function as

$$\mathcal{L}_{rec} = |\mathbf{x}_0 - \mathbf{x}_0^{rec}|.$$

To generate new images with DDPMs, typically, the backward step is applied in a step-wise fashion, to gradually denoise a random noise vector. For the given UAD task, we do not aim for the generation of new images but to estimate healthy brain anatomy given an input image. Therefore, we directly estimate \mathbf{x}_0^{rec} given \mathbf{x}_t at test time as it is done in Behrendt et al. (2023). The time step $t_{test} < T$, controls the level of noise to remove from \mathbf{x}_t at test time. Optionally, to become agnostic to the noise magnitude, we use an ensemble of different values $t_{test} = [250, 500, 750]$ and average the reconstructions of each noise level, similar to Graham et al. (2022).

3.2 context-conditioned DDPMs

A general depiction of our conditioning approach is provided in Fig. 1. Formally, we condition the backward process of DDPMs on a context vector \mathbf{c} as follows

$$\mathbf{x}_0^{rec} \sim p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{c}),$$

with $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{c}) = \mathcal{N}(\boldsymbol{\mu}_\theta(\mathbf{x}_t, t, \mathbf{c}), \boldsymbol{\Sigma}(t)).$

We use an image encoder F_{enc} to achieve a latent representation $\mathbf{c} = F_{enc}(\mathbf{x}_0)$ of the input image \mathbf{x}_0 where $\mathbf{c} \in \mathbb{R}^d$ with d as conditioning dimension.

To integrate the context vector \mathbf{c} , we manipulate the denoising Unet of the DDPM. Therefore, we individually adapt the features $\mathbf{f}_i \in \mathbb{R}^{H_i, W_i, C_i}$ at each level of the denoising Unet based on \mathbf{c} where H_i , W_i and C_i are the respective feature map dimensions. To achieve this, we adapt the time step conditioning of DDPMs as follows. First, the time step is encoded using a sinusoidal position embedding. Next, t is projected to a vector $\mathbf{c}_t \in \mathbb{R}^d$, by a multi-layer perceptron (MLP I). Subsequently, we concatenate the context vector \mathbf{c} and the time step vector \mathbf{c}_t , resulting in a conditioning vector $\mathbf{c}' \in \mathbb{R}^{2 \cdot d}$. Finally, \mathbf{c}' is projected to $\mathbf{c}'_{proj} \in \mathbb{R}^{2 \cdot C_i}$ by another multi-layer perceptron (MLP II) at each feature level i . The vector \mathbf{c}'_{proj} is then split into half, where the first and last C_i elements resemble the scaling factor γ and the shift value β . Inspired by Perez et al. (2018) the variables γ and β are then used to transform the individual feature maps $\mathbf{f}'_i = \mathbf{f}_i * (\gamma + 1) + \beta$ in each residual block.

By this, can learn both, the feature extraction $\mathbf{c} = F_{enc}(\mathbf{x}_0)$ and the individual feature adjustments of the denoising Unet, based on the extracted context vector \mathbf{c} in an end-to-end fashion. Optionally, to achieve a meaningful starting point for the calculation of the context vector $\mathbf{c} = F_{enc}(\mathbf{x}_0)$, we pre-train the feature extraction of the image encoder F_{enc} which is described in the next section.

3.3 Pre-Training

We utilize a generative pre-training strategy for F_{enc} . More precisely, we utilize masked pre-training where typically transformer-based AEs are trained to reconstruct an image where a large fraction of patches are masked out He et al. (2022). We utilize the SparK framework Tian et al. (2023), where sparse convolutions and hierarchical features are used to enable the masked pre-training for CNNs. During pre-training, we utilize the same healthy training set as for the main training task to learn the general feature representations that are required to capture important information of the MRI scans. After the pre-training stage, we discard the decoder and only use F_{enc} and fine-tune it along with the denoising Unet during the main training stage of the cDDPM.

4 Experimental Setup

4.1 Data Sets

Following the principle of UAD, we train our models for the reconstruction task on healthy data only (IXI). At test time, we evaluate the models' anomaly detection ability on unhealthy test sets of various pathologies (BraTS21, ATLAS, MSLUB, WMH).

4.1.1 Training Data

We use the publicly available IXI data set¹ as our healthy reference data set for training. This data set includes 560 3D brain MRI scans, collected from three different medical facilities. Of the training data, 158 samples are set aside

¹<https://brain-development.org/ixi-data set/>

for testing, while the remaining data is divided into 5 folds, each containing 358 training samples and 44 validation samples for cross-validation.

4.1.2 Evaluation Data

For evaluation, we utilize four different publicly available data sets that contain different types of pathologies and the corresponding manual expert annotations:

1. Multimodal Brain Tumor Segmentation Challenge 2021 (BraTS21) Baid et al. (2021); Menze et al. (2014); Bakas et al. (2017)
2. multiple sclerosis data set from the University Hospital of Ljubljana (MSLUB) Lesjak et al. (2018)
3. Anatomical Tracings of Lesions After Stroke v2.0 (ATLAS) Liew et al. (2022)
4. White Matter Hyperintensity (WMH) data set Kuijf et al. (2019)

The BraTS21 data set includes 2040 3D brain routine MRI scans of patients with glioma with a pathologically confirmed diagnosis. Accompanying the MRI scans, annotations from expert neuroradiologists are provided for 1251 scans that delineate tumor sub-regions as categorical masks. In this work, fuse all sub-regions to obtain a binary segmentation mask to evaluate the anomaly detection task. All scans are available as T1-weighted volumes with and without contrast enhancement (T1-CE, T1) and T2-weighted or T2 fluid attenuated inversion recovery (T2, FLAIR) volumes. The 1251 annotated samples are split into an unhealthy validation set of 100 samples and an unhealthy test set of 1151 samples. The MSLUB data set includes 3D brain MRI scans of 30 patients with multiple sclerosis (MS) lesions. For each patient, along with the T1, T2 and FLAIR MRI scans, ground truth annotations are available derived based on multi-rater consensus. The data is split into an unhealthy validation set of 10 samples and an unhealthy test set of 20 samples.

The ATLAS data set consists of 655 T1-weighted MRI scans of stroke patients, collected from 44 research cohorts. The stroke lesions are annotated by domain experts and binary segmentation masks are provided. The data is split into an unhealthy validation set of 175 samples and an unhealthy test set of 480 samples.

The WMH data set consists of 60 MRI scans of patients with white matter hyperintensities from three different institutions and scanner types. WMH segmentation masks are derived from the consensus of two expert radiologists. The data set is split into an unhealthy validation set of 15 samples and an unhealthy test set of 45 samples.

Across the data sets, different weightings are available. For the BraTS21 and MSLUB data set multiple weightings are available (T1, T2, FLAIR), for the ATLAS data set, only T1-images and for the WMH data set T1 and FLAIR images exist. As our training data set contains T1 and T2 images of each patient, we train our models on both weightings separately and evaluate the BraTS21 and MSLUB data set on T2 images while for ATLAS and WMH, T1 images are used.

An overview of the data set sizes is provided in Table 1.

Table 1: Data set Information regarding the number of samples per data split. The IXI data set contains only healthy brain MRI scans and is considered as training set and to test the overall reconstruction quality of healthy brain anatomy. The remaining data sets are used to evaluate the domain adaptation and segmentation performance.

Data set	Healthy	Healthy	Healthy
	Training Samples	Validation Samples	Test Samples
IXI	358	44	158
		Unhealthy	Unhealthy
		Validation Samples	Test Samples
BraTS21	-	100	1151
MSLUB	-	10	20
ATLAS	-	175	480
WMH	-	15	45

4.2 Pre-Processing

We pre-process the images according to established pre-processing strategies for UAD in brain MRI Baur et al. (2021). First, we resample all MRI scans to the isotropic resolution of $1 \text{ mm} \times 1 \text{ mm} \times 1 \text{ mm}$ using cubic spline interpolation. Second, we register all MRI scans to the SRI24-Atlas. Third, we remove the skull from the MRI scans by skull stripping with HD-BET Isensee et al. (2019). Subsequently, we cut black borders and perform N4 bias field correction. Finally, we pad the images to a unified size of $192 \times 192 \times 160$. To save computational resources, we reduce volume resolution

by a factor of two and remove the 15 top and bottom slices parallel to the transverse plane leading to a final resolution of $96 \times 96 \times 50$ voxels.

4.3 Post-Processing

At test time, we derive a binary segmentation map from the residual map $\mathbf{R} = |\mathbf{x}_0 - \mathbf{x}_0^{rec}|$ where regions of high residuals indicate anomalies. To binarize \mathbf{R} , we first apply the following post-processing steps that are commonly used in the field of UAD in brain MRI Baur et al. (2021); Kascenas et al. (2022); Behrendt et al. (2023); Zimmerer et al. (2019b). First, a median filter with a kernel size of $K = 5 \times 5 \times 5$ is applied to smooth the residual map and to remove smaller false positives. Subsequently, we perform brain mask eroding for 3 iterations. This step is mainly applied to filter out residuals that occur due to poor reconstructions at sharp edges near the brain mask Baur et al. (2021). We then perform a greedy threshold search. Hereby, the test threshold is determined by searching for the best Dice score across different thresholds on the unhealthy validation set, as proposed in Zimmerer et al. (2019b). After binarization, we use connected component filtering to remove areas that include less than 7 voxels. We note that this post-processing step aims to remove false-positive predictions, much smaller than the anomalies considered in this study. We provide a systematic comparison of the post-processing strategies for each data set in the supplemental material.

4.4 Implementation Details

In our study, we compare our proposed method, called cDDPM, with multiple established baselines for UAD in brain MRI. The baselines include *AE*, *VAE* Baur et al. (2021), its sequential extension *SVAE* Behrendt et al. (2022a), and denoising AEs *DAE* Kascenas et al. (2022). We also compare with simple thresholding *Thresh* Meissen et al. (2022) and the GAN-based *AnoGAN* Schlegl et al. (2019). For a direct comparison, we also include *DDPM* Wyatt et al. (2022) and *pDDPM* Behrendt et al. (2023) as a counterpart to our proposed method.

We adapt the baseline implementations by tuning hyper-parameters based on the unhealthy validation set if required to improve training stability and performance. We set β_{VAE} to 0.001 for *VAE* and *SVAE*. For *f-AnoGAN*, we set the latent size to 128 and the learning rate to $1e - 4$.

For *DDPM*, *pDDPM* and cDDPM, the following adaptations are applied. We use structured simplex noise, as it has shown to strongly improve the UAD performance on MRI images Wyatt et al. (2022). Furthermore, we uniformly sample $t \in [1, T]$ with $T = 1000$ and either use a fixed value of $t_{test} = \frac{T}{2} = 500$, or an ensemble of different values $t_{test} = [250, 500, 750]$ and average the individual reconstructions of each noise level at test time. The denoising network for all DDPM-based approaches is a U-net similar to Dhariwal and Nichol (2021), with channel dimensions of [128, 256, 256]. As encoder network F_{enc} , we utilize a ResNet-backbone with a fully connected layer to match the target dimension of $c \in \mathbb{R}^d$ with $d = 128$ as conditioning dimension. During pre-training, a mask-out ratio of 65 % is used. For data augmentation, we utilize random -blur (p=0.25), -bias (p=0.25), -gamma (p=0.5) and -ghosting (p=0.5) from the torchio library Pérez-García et al. (2021). If not specified otherwise, all models are trained for a maximum of 1600 epochs on NVIDIA V100 (32GB) GPUs, using Adam as optimizer, a learning rate of $1e - 5$, and a batch size of 32. The best model checkpoint, as determined by performance on the healthy validation set, is used for testing. The volumes are processed in a slice-wise manner, uniformly sampling slices with replacement during training and iterating over all slices to reconstruct the full volume at test time. We implement all models in Pytorch (v0.10)².

4.5 Experiments and Evaluation

To assess the reconstruction quality, the domain adaptation and the final UAD performance, we conduct various experiments on different data sets, specified in the following.

4.5.1 Reconstruction Quality:

To evaluate the overall reconstruction quality, we utilize the held-out test set of the healthy IXI data set and calculate similarity metrics between input and reconstruction. We consider the Structural Similarity Index Measure (SSIM) Wang et al. (2004), the Peak Signal To Noise Ratio (PSNR) and the Learned Perceptual Image Patch Similarity (LPIPS) as metrics to assess the reconstruction quality. For the feature-based LPIPS metric Zhang et al. (2018), features are extracted by a resnet-based network, pre-trained on 3D medical data Chen et al. (2019). Furthermore, we report the overall $l1 - error$ for the healthy data set. As for UAD, only healthy anatomy should be estimated with high reconstruction quality, it is of interest to consider the $l1$ error of healthy and unhealthy anatomy separately, given the unhealthy evaluation data sets. Therefore, we calculate the $l1 - error$ for both healthy and unhealthy anatomy, as indicated by

²Code available at <https://github.com/FinnBehrendt/Conditioned-Diffusion-Models-UAD>

the annotation masks and calculate an $l1 - ratio$ as follows:

$$l1 - ratio = \frac{l1_{unhealthy}}{l1_{healthy}}.$$

A higher value for the $l1 - ratio$ indicates that the model successfully reconstructs the healthy anatomy while struggling to reconstruct the unhealthy parts of the input, and vice versa. This ratio serves as a metric to assess the model’s performance in capturing the distinction between healthy and unhealthy anatomical structures.

4.5.2 Domain Adaptation:

To investigate the domain adaptation ability of our proposed approach, we utilize both, healthy, in-domain data from the IXI data set and unhealthy out-of-domain data from the BraTS21 data set. To evaluate the effect of the conditioning mechanism, we utilize the IXI data set and simulate different levels of conditioning information and simulated domain shift to investigate the reconstructions qualitatively. Thereby, we alter the available information of the image that is fed to the image encoder to condition the cDDPM. To achieve this, we crop the conditioning image at a given width of 50% and 100% where 100% indicates that the full input image is used as the conditioning image. Furthermore, we simulate different contrast levels ranging from $cl \in [0.3, 0.7, 1.0, 1.5, 2.0]$. The images of different contrast levels are derived by potentiating the gray values by the respective contrast level i.e. $x_0^{cl=2} = x_0^2$. In addition to the qualitative evaluation of reconstructed, simulated data, we quantitatively assess the domain shift across input and reconstruction by comparing their intensity histograms. Therefore, we first calculate and plot the histograms. Thereby, we partition the intensity values into 500 bins and divide the raw count by the total number of counts and the bin width. For a quantitative analysis of the distance between the intensity distributions, we calculate the Kullback-Leibler Divergence (KLD) as follows:

$$\text{KLD} = \left[- \sum_i p_{input} \log(p_{input}) \right] - \left[- \sum_i p_{reconstruction} \log(p_{reconstruction}) \right]$$

where $p = [p_1, p_2, \dots, p_n]$ represents each intensity distribution.

4.5.3 Segmentation Performance:

To assess the segmentation performance for the UAD task, we utilize all unhealthy test sets. We report the average Dice score across all predicted anomaly maps (DICE). The formula for the DICE score is given by:

$$\text{DICE} = \frac{2 \cdot |A \cap B|}{|A| + |B|}$$

where A and B are the predicted anomaly map and the ground truth annotation, respectively.

To obtain a metric that is independent of the chosen threshold, we additionally calculate the Area Under Precision-Recall Curve (AUPRC) as follows:

$$\text{AUPRC} = \sum_r (R(r) - R(r - 1)) \cdot P(r).$$

Here, $R(r)$ represents the recall at a given threshold or rank r , and $P(r)$ represents the precision at the corresponding recall $R(r)$. The sum is taken over all thresholds or ranks r at which the precision and recall are computed.

4.5.4 Statistical Testing:

To conduct significance tests, we utilize the MLXtend library’s permutation test Raschka (2018) with 10,000 rounds of permutations and a significance level of $\alpha = 5\%$. This test computes the mean difference of the considered scores of two models for each permutation, and the resulting p-value is computed by counting the number of times the mean differences were equal to or greater than the sample differences, divided by the total number of permutations.

5 Results

In this section, we first compare the overall reconstruction quality. Second, we evaluate the domain adaptation properties of our approach and lastly, we report the Segmentation performance of all models.

Table 2: Comparison of the reconstruction quality of the different models with the best results highlighted in bold. The asterisk * denotes superior performance with statistical significance compared to all baselines ($p < 0.05$). For all metrics, the mean \pm standard deviation across the different folds are reported. The arrows \uparrow and \downarrow indicate that higher and lower values are favorable, respectively. The $l1 - ratio$ is derived by dividing the $l1 - error$ of unhealthy anatomy by the $l1 - error$ of healthy anatomy. DDPM-based models are evaluated by ensembling different values for $t_{test} = [250, 500, 750]$

Model	IXI (T2)				BraTS21 (T2)	MSLUB (T2)	ATLAS (T1)	WMH (T1)
	SSIM \uparrow	PSNR \uparrow	LPIPS (e-3) \downarrow	$l1 - error$ (e-3) \downarrow	$l1 - ratio$ \uparrow	$l1 - ratio$ \uparrow	$l1 - ratio$ \uparrow	$l1 - ratio$ \uparrow
VAE	74.98 \pm 0.54	23.38 \pm 0.14	4.03 \pm 0.50	32.32 \pm 0.64	3.52 \pm 0.08	2.92 \pm 0.06	4.43 \pm 0.03	2.36 \pm 0.04
SVAE	77.87 \pm 0.15	23.94 \pm 0.06	3.31 \pm 0.24	29.08 \pm 0.16	3.90 \pm 0.05	3.13 \pm 0.05	3.38 \pm 0.11	2.07 \pm 0.01
AE	76.11 \pm 0.27	23.41 \pm 0.14	3.19 \pm 0.54	31.67 \pm 0.41	3.84 \pm 0.17	3.26 \pm 0.18	4.40 \pm 0.07	2.36 \pm 0.04
DAE	98.69 \pm 0.15*	36.69 \pm 0.38*	0.14 \pm 0.01	8.14 \pm 0.17*	7.17 \pm 0.63	2.69 \pm 0.15	4.51 \pm 0.15	2.99 \pm 0.14
DDPM	93.96 \pm 0.37	31.79 \pm 0.26	0.49 \pm 0.14	14.29 \pm 0.32	6.16 \pm 0.53	3.37 \pm 0.24	5.00 \pm 0.23	3.16 \pm 0.15
pDDPM	96.62 \pm 0.25	34.58 \pm 0.39	0.09 \pm 0.04*	9.70 \pm 0.43	7.16 \pm 0.15	4.34 \pm 0.13	5.58 \pm 0.28	3.00 \pm 0.16
cDDPM (Ours)	96.80 \pm 0.19	34.87 \pm 0.23	0.11 \pm 0.05	9.68 \pm 0.16	7.43 \pm 0.17	4.49 \pm 0.18	5.69 \pm 0.27	3.12 \pm 0.08

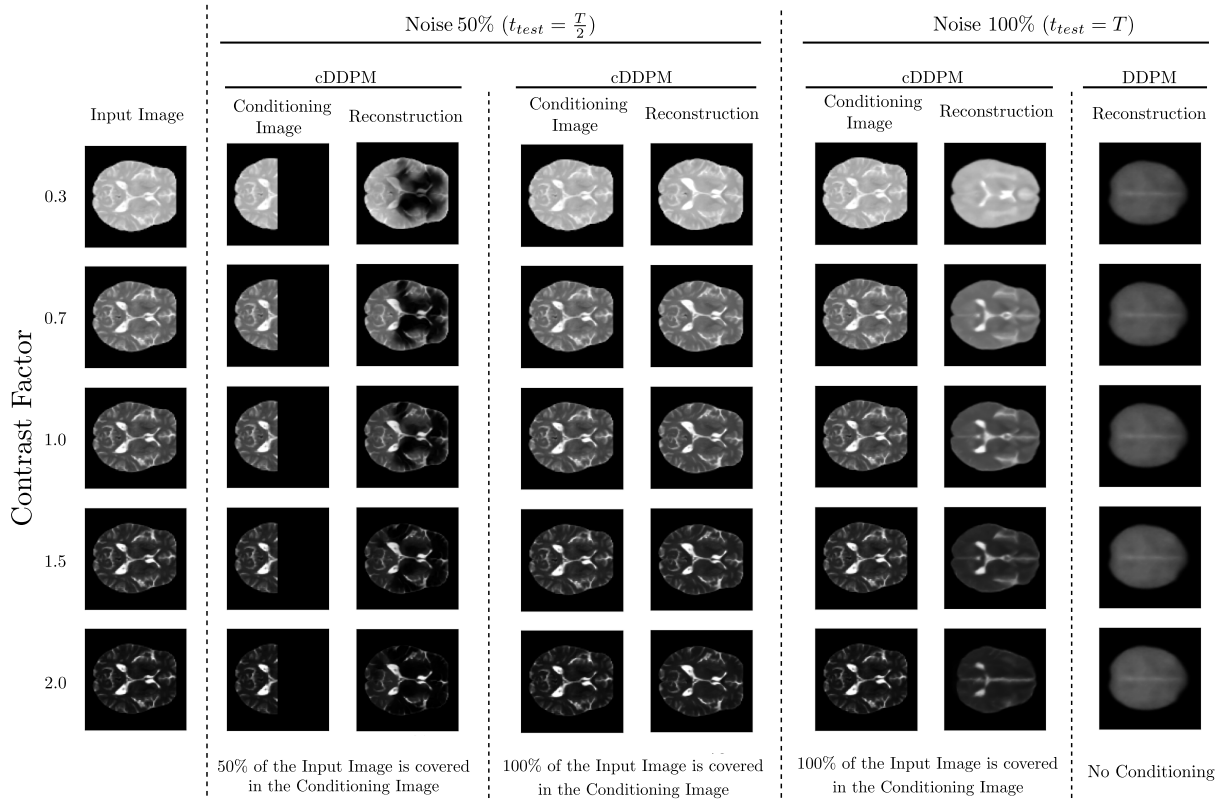


Figure 2: Simulating the conditioning effect of the cDDPM for 50 % of noise and 100% noise in the input image. In the first block the input image, that is fed to the DDPM or cDDPM is shown. In the second block, the reconstructions of cDDPM for different conditioning inputs are shown when a noise level of 50% is applied. In the third block, the reconstructions of cDDPMs and DDPMs are compared at a noise level of 100%. From top to bottom, the contrast level of the conditioning and input image is increased, respectively for all columns.

5.1 Reconstruction Quality

In Table 2 we compare baseline models regarding their ability to reconstruct the healthy anatomy, given the held-out test set of the IXI data set. Overall, DAEs, pDDPMs and cDDPMs show high performance regarding the image-based similarity metrics. In contrast, for the dense autoencoder-based baselines lower reconstruction quality is reported. Comparing the DDPM-based approaches, both, pDDPM and cDDPM outperform the baseline DDPM in terms of

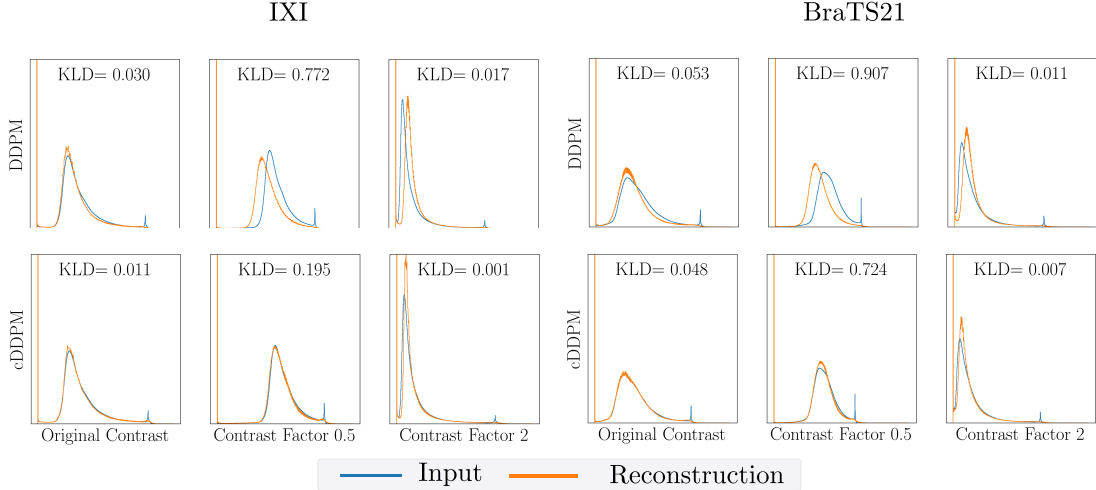


Figure 3: Comparison of the histograms for input-reconstruction pairs of the healthy IXI (left) and the unhealthy BraTS21 (right) data set with original and augmented contrast. The top row shows the baseline *DDPM* without conditioning and the bottom row our proposed *cDDPM* with conditioning. The Kullback-Leibler divergence (KLD) for both histograms is indicated within each plot, respectively (lower values are preferable). Both models are evaluated by ensembling different values for $t_{test} = [250, 500, 750]$.

Table 3: Comparison of the evaluated models with the best results highlighted in bold. The asterisk * denotes superior performance with statistical significance compared to all baselines ($p < 0.05$). For all metrics, the mean \pm standard deviation across the different folds are reported. A checkmark at SSL denotes that the encoder is pre-trained by self-supervision. A checkmark at ENS denotes the ensembling of different values for $t_{test} = [250, 500, 750]$. Otherwise, a fixed value of $t_{test} = 500$ is used for *DDPM*-based models.

Model	Modification		BraTS21 (T2)		MSLUB (T2)		ATLAS (T1)		WMH (T1)	
	ENS	SSL	DICE [%]	AUPRC [%]	DICE [%]	AUPRC [%]	DICE [%]	AUPRC [%]	DICE [%]	AUPRC [%]
<i>Thresh</i>			19.69	20.27	6.21	4.23	4.41	1.71	9.38	4.72
<i>AnoGAN</i>			23.99 \pm 6.15	21.08 \pm 6.23	4.86 \pm 2.02	3.77 \pm 1.32	9.91 \pm 1.80	9.17 \pm 1.41	6.25 \pm 1.00	3.45 \pm 0.33
<i>VAE</i>			28.81 \pm 1.26	25.72 \pm 1.54	6.16 \pm 0.58	4.46 \pm 0.20	15.08 \pm 0.28	15.27 \pm 0.30	5.23 \pm 0.82	4.16 \pm 0.13
<i>SVAE</i>			31.93 \pm 0.42	30.30 \pm 0.52	6.04 \pm 0.26	4.81 \pm 0.12	10.49 \pm 0.67	10.06 \pm 0.55	6.64 \pm 0.01	3.11 \pm 0.04
<i>AE</i>			31.51 \pm 1.94	28.80 \pm 1.92	7.23 \pm 0.90	5.71 \pm 0.90	14.91 \pm 0.33	14.76 \pm 0.46	4.53 \pm 0.36	4.16 \pm 0.19
<i>DAE</i>			45.37 \pm 4.40	49.38 \pm 4.68	3.88 \pm 1.35	4.47 \pm 0.78	8.53 \pm 0.28	12.45 \pm 0.94	7.31 \pm 1.02	6.32 \pm 0.88
<i>DDPM</i>			44.25 \pm 1.49	49.98 \pm 2.41	4.80 \pm 1.98	6.36 \pm 1.84	12.90 \pm 0.89	15.67 \pm 0.76	10.03\pm1.06	8.86 \pm 0.95
<i>DDPM</i>		✓	44.50 \pm 2.20	50.73 \pm 3.09	6.46 \pm 2.05	6.31 \pm 1.40	14.67 \pm 0.86	17.56 \pm 1.12	9.63 \pm 1.06	8.12 \pm 1.17
<i>pDDPM</i>			49.36 \pm 0.66	54.70 \pm 0.53	9.40 \pm 1.23	9.88 \pm 0.59	12.95 \pm 0.45	17.37 \pm 0.39	8.03 \pm 0.62	7.78 \pm 0.51
<i>pDDPM</i>		✓	49.78 \pm 0.85	55.10 \pm 0.57	9.21 \pm 1.35	10.11 \pm 0.45	13.24 \pm 0.93	17.76 \pm 1.17	7.97 \pm 0.95	7.65 \pm 0.91
<i>cDDPM (Ours)</i>			51.34 \pm 1.68	56.84 \pm 2.21	10.71 \pm 1.52	10.13 \pm 1.19	19.06 \pm 1.27	20.98 \pm 1.14	9.94 \pm 0.55	9.28 \pm 0.42
<i>cDDPM (Ours)</i>		✓	52.35 \pm 0.95	58.14 \pm 1.47	11.09 \pm 0.87	10.85 \pm 1.02	18.95 \pm 1.94	20.86 \pm 1.53	9.92 \pm 1.26	9.23 \pm 1.08
<i>cDDPM (Ours)</i>		✓	53.37\pm1.80*	58.84\pm1.76*	11.51\pm1.24	11.13\pm1.26	19.99\pm1.55*	22.21\pm1.47*	9.88 \pm 1.22	9.33\pm1.07

reconstruction quality with statistical significance ($p < 0.05$). We additionally provide an analysis of the unhealthy-to-healthy error ratio based on the unhealthy test sets. Notably, here the $l1$ - ratio is highest for *cDDPM* across almost all data sets, except for the WMH, where *DDPM* shows competitive performance.

5.2 Conditioning Effect

To evaluate the effect the additional conditioning input has on individual reconstructions, we simulate different conditioning inputs, varying in the amount of used image information. Furthermore, we apply artificial contrast levels for the input images to mimic strong domain shifts. For each conditioning input and contrast level, we provide the reconstructions, generated by *cDDPMs* in Fig. 2 for a noise level of $t_{test} = 500$ (50%). Moreover, we compare the reconstruction of *DDPMs* and *cDDPMs* in the extreme case of pure noise as input ($t_{test} = T = 1000$ (100%)). For a noise level of 50%, we observe that while the overall shape of the reconstruction is preserved across the different conditioning masks, only at regions that are covered in the conditioning image, local intensity information is captured in the reconstruction. At a noise level of 100%, we observe that the reconstructions of *cDDPMs* coarsely follow the

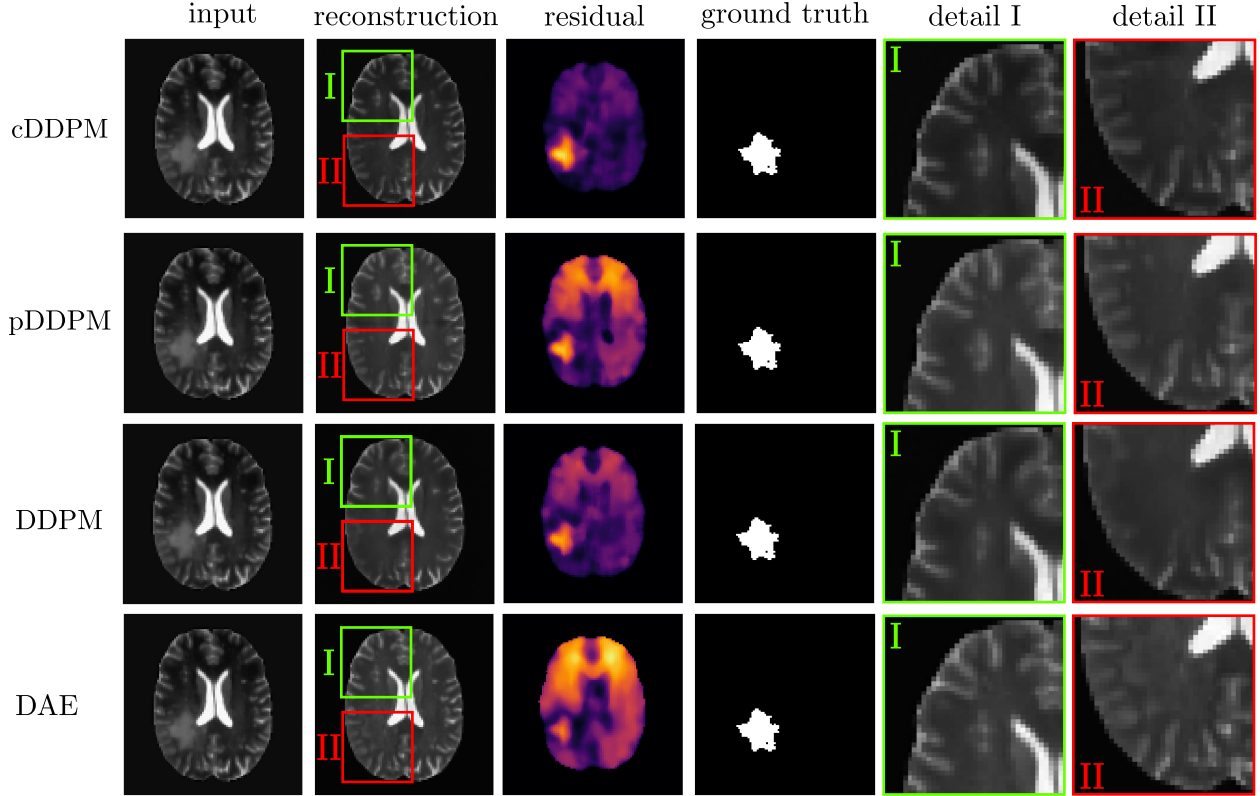


Figure 4: Exemplary reconstruction and anomaly map taken from the BraTS21 data set. From top to bottom, *cDDPM*, *pDDPM*, *DDPM* and *DAE* are compared.

content provided by the conditioning image, leading to a blurry reconstruction of the input image. In contrast, for unconditioned DDPMs, only a very generic reconstruction can be obtained that shares low similarity with the given input image.

5.3 Domain Adaptation

To evaluate the domain adaptation to different data sets. In our experiments, we consider the healthy IXI data set as in-domain data set and the unhealthy BraTS21 data set as out-of-domain data set. Note that for the BraTS21 data set, we only consider regions that have been annotated as healthy. Thereby, we ensure to evaluate domain shifts regarding scanner and brain diversity and not domain shifts that arise from unhealthy structures in brain MRIs. In Fig. 3, we provide the histograms of input and reconstructions of the healthy IXI data set (left) and the unhealthy BraTS21 data set (right). We observe that DDPMs show substantial discrepancies across the intensity distributions of input and reconstruction. Particularly for simulated contrast levels, the histograms deviate. In contrast, the intensity distribution of images reconstructed by *cDDPMs* exhibits higher similarities with the input intensity distribution for both in-domain and out-of-domain data. Considering the quantitative KLD measurement, the KLD of DDPMs is increased by a factor of 2.3, 4.0 and 17.0 for the original contrast, a contrast factor of 0.5 and 2.0, respectively compared to *cDDPMs* for the IXI data set.

5.4 Segmentation Performance

Overall, improved performance is reported for *cDDPMs* compared to all baselines across all data sets, except for the WMH data set, where the performance is on par with *DDPMs*. While the improvements for the *cDDPM* are statistically significant for the BraTS21 and ATLAS data sets ($p < 0.05$), for the MSLUB and WMH data sets, no significant difference can be observed. Furthermore, we report enhanced performance of *cDDPMs* when pre-training the encoder (SSL checkmark in Table 3) and ensembling the reconstructions of different noise levels (ENS checkmark in Table 3) for most data sets. Notably, the inference time of *cDDPMs* is reduced by $\sim 37\%$ compared to *pDDPMs* and increased

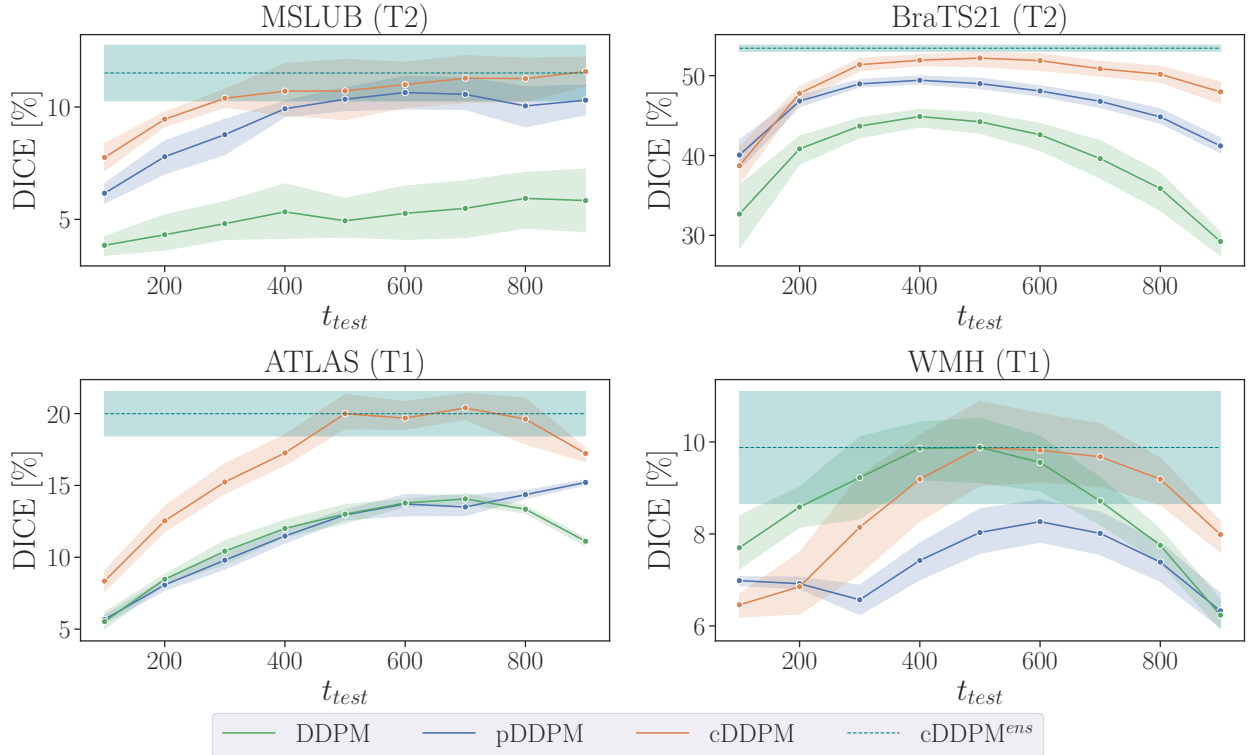


Figure 5: Comparison of different noise levels t_{test} regarding the DICE for the MSLUB (left) and BraTS21 (right) data set. The superscript *ens* denotes the ensembling of reconstructions from different noise levels t_{test} in $[250, 500, 750]$.

by $\sim 2\%$ compared to *DDPMs*.

Comparing the reconstructions and residual maps in Fig. 4, we observe crisp reconstructions for both *cDDPMs* and *pDDPMs*, whereas the reconstructions of *DDPMs* show missing details, particularly at regions where the tumor is located in the input image. We observe that while *DAEs* provide detailed reconstructions, also unhealthy anatomy is reproduced. For, *cDDPMs* we observe aligned intensity information across input and reconstructions. Hence, the residual map shows a higher contrast across normal and abnormal regions, which enables a better delineation of the present pathology.

Considering Fig. 5, it can be observed that the achieved DICE score is dependent on the noise level that is applied at test time. By applying the ensembling strategy, this dependency is reduced and consistent performance is achieved without specifying an individual noise level.

We supply a collection of exemplary residual maps for different baseline methods in Fig. 6 where across all examples, the different baselines show residual maps that are conceptually similar. Nonetheless, our suggested *cDDPM* demonstrates the lowest occurrences of false positives and the most pronounced contrast in the residual maps, particularly in comparison to *DDPMs* and *DAEs*.

We provide additional ablation studies of the applied post-processing steps in the supplementary material.

6 Discussion

Unsupervised anomaly detection in neuroimaging has gained significant attention due to its potential to identify abnormalities without the need for costly data annotation. Compared to supervised approaches that rely on annotated data sets, UAD methods take a different approach by learning the underlying data distribution of healthy brain anatomy and identifying anomalies as outliers from that distribution.

In this study, we focus on *DDPMs* for UAD in brain MRI. These models generate images by reconstructing an input that is corrupted by noise, leveraging the high-dimensional latent space to achieve high-fidelity reconstructions with preserved spatial context. However, while the overall brain structure is reconstructed well, we show that the forward and backward processes of *DDPMs* do not capture the highly variable intensity characteristics of MRI scans sufficiently,

resulting in false positives in the residual map and reduced detection performance. This becomes especially prominent in the presence of domain shifts at test time.

To address this challenge, we propose context-conditioned DDPMs (cDDPMs) for UAD in brain MRI. Here, we train a *DDPM* to reconstruct healthy brain anatomy and incorporate a latent feature representation of the noise-free input image, derived by an additional image encoder as a conditioning input to the denoising process. While the additional feature representation is not suitable for high-fidelity reconstruction Baur et al. (2021); Bercea et al. (2023b), we show that it can capture local intensity information of the image to reconstruct. We demonstrate that by incorporating the feature representation of individual input images, our proposed cDDPMs reconstruct the brain MRIs with more accurate intensity information compared to the unconditional DDPMs. Additionally, we observe enhanced domain adaptation capabilities to both, real and simulated intensity profiles with our conditioning mechanism. Finally, we demonstrate that these appealing features of our approach are crucial properties to improve the segmentation performance in reconstruction-based UAD.

Overall, we systematically evaluate the performance of our approach in terms of reconstruction quality, domain adaptation, and segmentation performance based on five different data sets.

6.1 Reconstruction Quality

We compare the reconstruction quality of our method with baseline models on the healthy IXI data set in Table 2. For the *AE* and *(S)VAE*, overall the worst reconstruction quality is reported. A reason for this is seen in the strict bottleneck enforced by the dense latent space as it inhibits information flow Baur et al. (2021). In contrast, methods like *DAE* or *DDPMs* that are not constrained by a dense latent space but by a noising strategy Kascenas et al. (2022), show improved reconstruction performance. While *pDDPMs* and cDDPMs outperform the baseline *DDPM* in terms of reconstruction quality, we observe that all models are outperformed by the baseline *DAE*. We note that the overall training objective of the compared generative models is to reconstruct the image with high accuracy and hence copying the input image would be a trivial solution. However, for the UAD task, it is crucial that the given input image is not solely copied but that unhealthy anatomy is replaced by pseudo-healthy representatives. Hence, comparing only the reconstruction quality of healthy anatomy does not necessarily reflect the usefulness for the UAD task. Therefore, we utilize the *l1 - ratio* where high values indicate a better trade-off between the reconstruction of healthy and unhealthy anatomy and vice-versa. While DDPMs and particularly the cDDPM achieve a high *l1 - ratio*, across all unhealthy data sets it becomes evident, that the *DAE* fails to generalize to different pathology types, which is a crucial property of UAD methods. A reason for that is seen in the chosen noise type in *DAE* that is highly optimized to the BraTS21 data set, mimicking the visual appearance of tumors Bercea et al. (2023b); Lagogiannis et al. (2023).

In summary, the cDDPM shows improved reconstruction quality compared to DDPMs while preserving a high *l1 - ratio*. This indicates that the conditioning mechanism effectively captures information from the input image for an improved reconstruction without providing too much detailed information that would enable the cDDPM to solely copy the input image.

6.2 Domain Adaptation

We evaluate the domain adaptation capabilities of cDDPMs by simulating different contrast levels and conditioning inputs in Fig. 2. The reconstructions show that while the overall shape is preserved across different conditioning masks, meaningful reconstructions are achieved only in regions covered by the conditioning image. Particularly, the conditioning image plays a critical role in capturing local intensity information. This demonstrates the ability of cDDPMs to adapt to different contrast levels and to capture varying intensity information effectively. This becomes even more evident when high noise levels are considered, where the only source of information concerning the given input image is the conditioning image. Here, the reconstruction becomes totally dependent on the shape and intensity characteristics of the conditioning image. The conditioning facilitates a blurred reconstruction of prominent local intensity changes, suggesting that the conditioning mechanism allows the capture and reconstruction of local intensity details from the input image. These findings indicate that cDDPMs effectively learn to balance information from the noisy input image and the conditioning encoder features during training, adapting according to the level of noise inherent in the input.

We explore the domain adaptation capabilities of cDDPMs in real-world scenarios where a different, out-of-domain data set is used for testing. To assess the domain adaptation ability, we investigate the deviation between the intensity distributions of the input and reconstruction by plotting histograms and calculating the Kullback-Leibler Divergence (KLD) as a proxy in Fig. 3. Our findings reveal that cDDPMs exhibit improved performance in capturing and estimating the intensity distribution. Particularly when simulating contrast levels considerably different from the training distribution, cDDPMs demonstrate superior alignment of the histograms and lower KLD values. This analysis highlights the potential of the conditioning mechanism in cDDPMs to effectively adapt to unseen variations in intensity

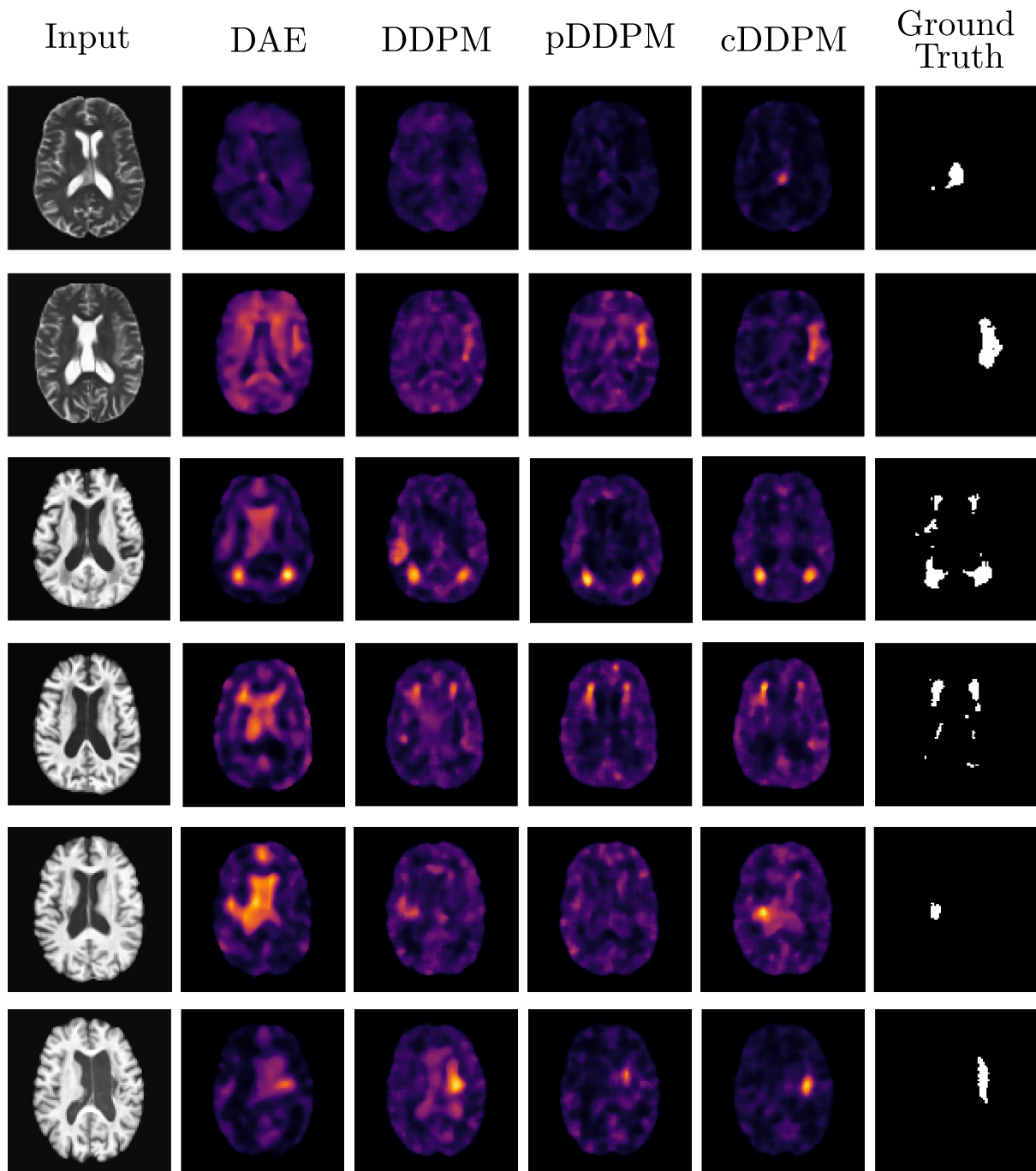


Figure 6: Exemplary residual maps from the BraTS21 data set (rows 1 and 2), the WMH data set (rows 3 and 4) and the ATLAS data set (rows 5 and 6). From left to right, the input image, the residual maps and the ground truth are shown.

profiles and improve the coherence between input and reconstruction, contributing to improved domain adaptation in real-world scenarios.

6.3 Segmentation Performance

Overall, cDDPMs demonstrate competitive or superior results compared to traditional autoencoder-based approaches, as well as the baseline *DDPM* and *pDDPM*, as presented in Table 3. We conclude that the high reconstruction quality, the accurate modeling of global and local intensity information, and the effective domain adaptation capabilities that are attributed to our conditioning approach are crucial to improve the UAD segmentation performance. Furthermore, we observe that pre-training the encoder F_{enc} slightly enhances the segmentation performance in most cases, indicating that starting from an already learned representation space has the potential to improve the overall integration of the conditioning features, compared to simultaneously training the parameters of both, *DDPM* and F_{enc} from scratch. Comparing the cDDPM to the *DAE* reveals that despite the *DAEs*' superior reconstruction quality on healthy data, they are outperformed by cDDPMs by a margin given the UAD task. A reason for this is seen in the *DAEs*' ability to reconstruct unhealthy anatomy, particularly for pathologies differing from the BraTS21 data set where the noise type is not optimized for, as discussed in Subsection 6.1.

To further analyze the effectiveness of cDDPMs, we provide visual comparisons of reconstructions and residual maps in Fig. 4. In comparison to pDDPMs, the reconstructions of cDDPMs demonstrate following the intensity information of the respective input images, resulting in improved contrast and intensity alignment between input and reconstruction pairs. This leads to a higher contrast in the residual maps, which facilitates the delineation of anomalies such as tumors. Furthermore, it becomes evident that besides adding a data-dependent hyper-parameter, the patching strategy introduces subtle artifacts at patch borders. Thus, our results indicate that cDDPMs make use of the additional information more effectively compared to pDDPMs. Furthermore, cDDPMs provide reduced complexity and inference time as there is no need for a costly patching strategy, making it a practical and efficient solution for UAD in brain MRI.

In Fig. 5, we explore the impact of noise levels on the segmentation performance. We demonstrate that cDDPMs outperform the baseline models across different noise levels for most data sets. However, we also observe that the noise level serves as a crucial hyper-parameter. In general, high noise levels tend to result in more blurry and generic reconstructions, whereas low noise levels enable sharper reconstructions, including unhealthy anatomy. Thus, selecting an appropriate noise level is essential to achieve reasonable performance. However, the optimal value for the applied noise depends on the evaluated data set as shown in Fig 5. We assume that the main reason for this dependency is the different size of pathologies in the considered data sets, as also indicated in Bercea et al. (2023a). To address this dependency, we apply different noise levels and average the resulting reconstructions. Thereby, we effectively mitigate the dependency on the noise level which enhances the model's generalization abilities, which are vital for UAD methods.

6.4 Limitations and Future Work

Overall, while our approach demonstrates promising results, several limitations should be acknowledged. Our study focuses primarily on brain MRI and may not generalize seamlessly to other imaging modalities or anatomical regions. Further investigation and adaptation of our approach to different medical imaging domains or even industrial defect detection would be beneficial.

Another potential avenue for improvement in our study is the inclusion of FLAIR (Fluid-Attenuated Inversion Recovery) data, which could improve the overall performance Meissen et al. (2022) and provide valuable insights into white matter abnormalities and lesions, especially in conditions like MS. Furthermore, it is important to acknowledge that the available data sets for white matter hyperintensities and MS lesions, such as the MSLUB and WMH data sets, are relatively small compared to the BraTS21 and ATLAS data sets. This limited sample size for WMH and MS lesions may reduce the generalizability of our findings and the availability of a larger data set would provide a more comprehensive representation of WMH and MS lesions, enabling more accurate and reliable evaluation.

Another avenue for future work is the incorporation of multi-scale image encodings into our conditioning mechanism. Currently, our study does not utilize multi-scale analysis, which could be advantageous in capturing fine-grained details and contextual information at different resolutions. By carefully integrating multi-scale image encodings without allowing to copy the conditioning image, we see potential to enhance the performance of our cDDPMs in capturing both global and local features of the input images.

To further enhance the utilized context, an additional direction for future work is to explore the use of 3D input for the image encoder. We have shown that it improves the reconstruction quality for VAEs Bengs et al. (2021); Behrendt et al. (2022a) and expect similar improvements for DDPMs. Currently, our approach operates on 2D slices of the MRI data, which may limit the preservation of 3D context and spatial relationships between slices. By incorporating 3D input into the image encoder, we can potentially capture and preserve the 3D structure and contextual information without the need to train the full 3D DDPM.

7 Conclusion

In this work, we addressed the task of reconstruction-based UAD in brain MRI. To this end, we proposed cDDPMs where we introduced a conditioning mechanism to DDPMs that incorporates feature representations of noise-free input images to the denoising process. We have shown that this conditioning mechanism effectively addresses challenges of accurate reconstruction, intensity capture, and domain adaptation and thus enables a more accurate delineation of pathologies from the generated residual maps. As a consequence, our approach outperformed state-of-the-art architectures for UAD in brain MRI, on various publicly available data sets. Our findings contribute to the development of effective UAD methods in brain MRI and have practical implications for the detection and segmentation of pathologies in clinical scenarios, where domain shifts are likely.

Acknowledgements

This work was partially funded by grant numbers KK5208101KS0 and ZF4026303TS9 and by the Free and Hanseatic City of Hamburg (Interdisciplinary Graduate School) from University Medical Center Hamburg-Eppendorf

References

- Baid, U., Ghodasara, S., Mohan, S., Bilello, M., Calabrese, E., Colak, E., Farahani, K., Kalpathy-Cramer, J., Kitamura, F.C., Pati, S., et al., 2021. The rsna-asnr-miccai brats 2021 benchmark on brain tumor segmentation and radiogenomic classification. arXiv preprint arXiv:2107.02314 .
- Bakas, S., Akbari, H., Sotiras, A., Bilello, M., Rozycki, M., Kirby, J.S., Freymann, J.B., Farahani, K., Davatzikos, C., 2017. Advancing the cancer genome atlas glioma mri collections with expert segmentation labels and radiomic features. *Scientific data* 4, 1–13.
- Baur, C., Denner, S., Wiestler, B., Navab, N., Albarqouni, S., 2021. Autoencoders for unsupervised anomaly segmentation in brain mr images: a comparative study. *Med. Image Anal.* , 101952.
- Baur, C., Wiestler, B., Albarqouni, S., Navab, N., 2018. Deep autoencoding models for unsupervised anomaly segmentation in brain mr images, in: MICCAI brainlesion workshop, Springer. pp. 161–169.
- Baur, C., Wiestler, B., Albarqouni, S., Navab, N., 2020a. Bayesian skip-autoencoders for unsupervised hyperintense anomaly detection in high resolution brain mri, in: IEEE ISBI, IEEE. pp. 1905–1909.
- Baur, C., Wiestler, B., Albarqouni, S., Navab, N., 2020b. Scale-space autoencoders for unsupervised anomaly segmentation in brain mri, in: Computer Assisted Radiology and Surgery, Springer. pp. 552–561.
- Behrendt, F., Bengs, M., Bhattacharya, D., Krüger, J., Opfer, R., Schlaefer, A., 2022a. Capturing inter-slice dependencies of 3d brain MRI-scans for unsupervised anomaly detection, in: Medical Imaging with Deep Learning. URL: <https://openreview.net/forum?id=db8wDgKH4p4>.
- Behrendt, F., Bengs, M., Rogge, F., Krüger, J., Opfer, R., Schlaefer, A., 2022b. Unsupervised anomaly detection in 3d brain mri using deep learning with impured training data, in: 2022 IEEE 19th International Symposium on Biomedical Imaging (ISBI), IEEE. pp. 1–4.
- Behrendt, F., Bhattacharya, D., Krüger, J., Opfer, R., Schlaefer, A., 2023. Patched diffusion models for unsupervised anomaly detection in brain mri. arXiv preprint arxiv.org/abs/2303.03758 .
- Bengs, M., Behrendt, F., Krüger, J., Opfer, R., Schlaefer, A., 2021. Three-dimensional deep learning with spatial erasing for unsupervised anomaly segmentation in brain mri. *Computer Assisted Radiology and Surgery* 16, 1413–1423.
- Bercea, C.I., Neumayr, M., Rueckert, D., Schnabel, J.A., 2023a. Mask, stitch, and re-sample: Enhancing robustness and generalizability in anomaly detection through automatic diffusion models, in: ICML 3rd Workshop on Interpretable Machine Learning in Healthcare (IMLH). URL: <https://openreview.net/forum?id=kTpafpXrqa>.
- Bercea, C.I., Wiestler, B., Rueckert, D., Schnabel, J.A., 2023b. Generalizing unsupervised anomaly detection: Towards unbiased pathology screening, in: Medical Imaging with Deep Learning. URL: <https://openreview.net/forum?id=8ojx-Ld3yjR>.
- Bruno, M.A., Walker, E.A., Abujudeh, H.H., 2015. Understanding and confronting our mistakes: the epidemiology of error in radiology and strategies for error reduction. *Radiographics* 35, 1668–1676.
- Chen, S., Ma, K., Zheng, Y., 2019. Med3d: Transfer learning for 3d medical image analysis. arXiv preprint arXiv:1904.00625 .
- Chen, X., You, S., Tezcan, K.C., Konukoglu, E., 2020. Unsupervised lesion detection via image restoration with a normative prior. *Med. Image Anal.* 64, 101713.

- Dhariwal, P., Nichol, A., 2021. Diffusion models beat gans on image synthesis. NIPS 34, 8780–8794.
- Ellis, R.J., Sander, R.M., Limon, A., 2022. Twelve key challenges in medical machine learning and solutions. *Intelligence-Based Medicine* 6, 100068.
- Graham, M.S., Pinaya, W.H., Tudosiu, P.D., Nachev, P., Ourselin, S., Cardoso, M.J., 2022. Denoising diffusion models for out-of-distribution detection. arXiv preprint arXiv:2211.07740 .
- Han, C., Rundo, L., Murao, K., Noguchi, T., Shimahara, Y., Milacski, Z.Á., Koshino, S., Sala, E., Nakayama, H., Satoh, S., 2021. Madgan: Unsupervised medical anomaly detection gan using multiple adjacent brain mri slice reconstruction. *BMC bioinformatics* 22, 1–20.
- He, K., Chen, X., Xie, S., Li, Y., Dollár, P., Girshick, R., 2022. Masked autoencoders are scalable vision learners, in: CVPR, pp. 16000–16009.
- Ho, J., Jain, A., Abbeel, P., 2020. Denoising diffusion probabilistic models. NIPS 33, 6840–6851.
- Isensee, F., Schell, M., Pflueger, I., Brugnara, G., Bonekamp, D., Neuberger, U., Wick, A., Schlemmer, H.P., Heiland, S., Wick, W., et al., 2019. Automated brain extraction of multisequence mri using artificial neural networks. *Human brain mapping* 40, 4952–4964.
- Islam, J., Zhang, Y., 2018. Brain mri analysis for alzheimer’s disease diagnosis using an ensemble system of deep convolutional neural networks. *Brain informatics* 5, 1–14.
- Johnson, J.M., Khoshgoftaar, T.M., 2019. Survey on deep learning with class imbalance. *Journal of Big Data* 6, 1–54.
- Karimi, D., Dou, H., Warfield, S.K., Gholipour, A., 2020. Deep learning with noisy labels: Exploring techniques and remedies in medical image analysis. *Med. Image Anal.* 65, 101759.
- Kascenas, A., Pugeault, N., O’Neil, A.Q., 2022. Denoising autoencoders for unsupervised anomaly detection in brain mri, in: *Medical Imaging with Deep Learning*, PMLR.
- Kuijf, H.J., Biesbroek, J.M., De Bresser, J., Heinen, R., Andermatt, S., Bento, M., Berseth, M., Belyaev, M., Cardoso, M.J., Casamitjana, A., et al., 2019. Standardized assessment of automatic segmentation of white matter hyperintensities and results of the wmh segmentation challenge. *IEEE transactions on medical imaging* 38, 2556–2568.
- Lagogiannis, I., Meissen, F., Kaissis, G., Rueckert, D., 2023. Unsupervised pathology detection: A deep dive into the state of the art. *IEEE Transactions on Medical Imaging* , 1–doi:doi:10.1109/TMI.2023.3298093.
- Lesjak, Ž., Galimzianova, A., Koren, A., Lukin, M., Pernuš, F., Likar, B., Špiclin, Ž., 2018. A novel public mr image dataset of multiple sclerosis patients with lesion segmentations based on multi-rater consensus. *Neuroinformatics* 16, 51–63.
- Liew, S.L., Lo, B.P., Donnelly, M.R., Zavaliangos-Petropulu, A., Jeong, J.N., Barisano, G., Hutton, A., Simon, J.P., Juliano, J.M., Suri, A., et al., 2022. A large, curated, open-source stroke neuroimaging dataset to improve lesion segmentation algorithms. *Scientific data* 9, 320.
- Lundervold, A.S., Lundervold, A., 2019. An overview of deep learning in medical imaging focusing on mri. *Zeitschrift für Medizinische Physik* 29, 102–127. URL: <https://www.sciencedirect.com/science/article/pii/S0939388918301181>, doi:doi:https://doi.org/10.1016/j.zemedi.2018.11.002. special Issue: Deep Learning in Medical Physics.
- McDonald, R.J., Schwartz, K.M., Eckel, L.J., Diehn, F.E., Hunt, C.H., Bartholmai, B.J., Erickson, B.J., Kallmes, D.F., 2015. The effects of changes in utilization and technological advancements of cross-sectional imaging on radiologist workload. *Academic radiology* 22, 1191–1198.
- Meissen, F., Kaissis, G., Rueckert, D., 2022. Challenging current semi-supervised anomaly segmentation methods for brain mri, in: *MICCAI brainlesion workshop*, Springer. pp. 63–74.
- Menze, B.H., Jakab, A., Bauer, S., Kalpathy-Cramer, J., Farahani, K., Kirby, J., Burren, Y., Porz, N., Slotboom, J., Wiest, R., et al., 2014. The multimodal brain tumor image segmentation benchmark (brats). *IEEE transactions on medical imaging* 34, 1993–2024.
- Moeskops, P., de Bresser, J., Kuijf, H.J., Mendrik, A.M., Biessels, G.J., Pluim, J.P., Išgum, I., 2018. Evaluation of a deep learning approach for the segmentation of brain tissues and white matter hyperintensities of presumed vascular origin in mri. *NeuroImage: Clinical* 17, 251–262.
- Perez, E., Strub, F., De Vries, H., Dumoulin, V., Courville, A., 2018. Film: Visual reasoning with a general conditioning layer, in: *AAAI*.
- Perkuhn, M., Stavrinou, P., Thiele, F., Shakirin, G., Mohan, M., Garmpis, D., Kabbasch, C., Borggrefe, J., 2018. Clinical evaluation of a multiparametric deep learning model for glioblastoma segmentation using heterogeneous magnetic resonance imaging data from clinical routine. *Investigative radiology* 53, 647.

- Pinaya, W.H., Graham, M.S., Gray, R., Da Costa, P.F., Tudosiu, P.D., Wright, P., Mah, Y.H., MacKinnon, A.D., Teo, J.T., Jager, R., et al., 2022a. Fast unsupervised brain anomaly detection and segmentation with diffusion models. arXiv preprint arXiv:2206.03461 .
- Pinaya, W.H., Tudosiu, P.D., Gray, R., Rees, G., Nachev, P., Ourselin, S., Cardoso, M.J., 2022b. Unsupervised brain imaging 3d anomaly detection and segmentation with transformers. *Med. Image Anal.* 79, 102475.
- Pérez-García, F., Sparks, R., Ourselin, S., 2021. Torchio: A python library for efficient loading, preprocessing, augmentation and patch-based sampling of medical images in deep learning. *Computer Methods and Programs in Biomedicine* 208, 106236.
- Raschka, S., 2018. Mlxtend: Providing machine learning and data science utilities and extensions to python’s scientific computing stack. *The Journal of Open Source Software* 3.
- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B., 2022. High-resolution image synthesis with latent diffusion models, in: *CVPR*, pp. 10684–10695.
- Ronneberger, O., Fischer, P., Brox, T., 2015. U-net: Convolutional networks for biomedical image segmentation, in: *Medical Image Computing and Computer Assisted Intervention*, Springer. pp. 234–241.
- Saharia, C., Chan, W., Chang, H., Lee, C., Ho, J., Salimans, T., Fleet, D., Norouzi, M., 2022. Palette: Image-to-image diffusion models, in: *ACM*, pp. 1–10.
- Schlegl, T., Seeböck, P., Waldstein, S.M., Langs, G., Schmidt-Erfurth, U., 2019. f-AnoGAN: Fast unsupervised anomaly detection with generative adversarial networks. *Med. Image Anal.* 54, 30–44.
- Silva-Rodríguez, J., Naranjo, V., Dolz, J., 2022. Constrained unsupervised anomaly segmentation. *Med. Image Anal.* 80, 102526.
- Tian, K., Jiang, Y., Diao, Q., Lin, C., Wang, L., Yuan, Z., 2023. Designing bert for convolutional networks: Sparse and hierarchical masked modeling. arXiv:2301.03580 .
- Vernooij, M.W., Ikram, M.A., Tanghe, H.L., Vincent, A.J., Hofman, A., Krestin, G.P., Niessen, W.J., Breteler, M.M., van der Lugt, A., 2007. Incidental findings on brain mri in the general population. *New England Journal of Medicine* 357, 1821–1828.
- Wang, T., Zhang, T., Zhang, B., Ouyang, H., Chen, D., Chen, Q., Wen, F., 2022. Pretraining is all you need for image-to-image translation. arXiv preprint arXiv:2205.12952 .
- Wang, Z., Bovik, A., Sheikh, H., Simoncelli, E., 2004. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing* 13, 600–612. doi:doi:10.1109/TIP.2003.819861.
- Wolleb, J., Bieder, F., Sandkühler, R., Cattin, P.C., 2022. Diffusion models for medical anomaly detection. arXiv preprint arXiv:2203.04306 .
- Wyatt, J., Leach, A., Schmon, S.M., Willcocks, C.G., 2022. Anoddpn: Anomaly detection with denoising diffusion probabilistic models using simplex noise, in: *CVPR*, pp. 650–656.
- Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O., 2018. The unreasonable effectiveness of deep features as a perceptual metric, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 586–595.
- Zimmerer, D., Isensee, F., Petersen, J., Kohl, S., Maier-Hein, K., 2019a. Unsupervised anomaly localization using variational auto-encoders, in: Shen, D., Liu, T., Peters, T.M., Staib, L.H., Essert, C., Zhou, S., Yap, P.T., Khan, A. (Eds.), *Medical Image Computing and Computer Assisted Intervention – MICCAI 2019*, Springer International Publishing, Cham. pp. 289–297.
- Zimmerer, D., Kohl, S., Petersen, J., Isensee, F., Maier-Hein, K., 2019b. Context-encoding variational autoencoder for unsupervised anomaly detection, in: *Medical Imaging with Deep Learning*.

Supplementary Material

Post-Processing Analysis

Model	without CC	without MF	without BE	BraTS21	MSLUB	ATLAS	WMH
<i>DAE</i>	✓	✗	✗	45.34 ± 3.94 (-0.03)	3.92 ± 1.19 (0.04)	8.53 ± 0.24 (0.00)	7.31 ± 0.91 (0.00)
<i>DAE</i>	✗	✓	✗	39.75 ± 3.37 (-5.62)	3.94 ± 0.58 (0.06)	9.21 ± 0.32 (0.68)	7.01 ± 0.68 (-0.30)
<i>DAE</i>	✗	✗	✓	44.36 ± 3.88 (-1.01)	3.51 ± 0.89 (-0.37)	8.77 ± 0.28 (0.24)	6.80 ± 0.81 (-0.51)
<i>DDPM</i>	✓	✗	✗	44.42 ± 2.03 (-0.08)	6.05 ± 2.10 (-0.41)	14.65 ± 0.30 (-0.02)	10.27 ± 0.96 (0.64)
<i>DDPM</i>	✗	✓	✗	36.56 ± 1.86 (-7.94)	8.05 ± 1.39 (1.59)	12.02 ± 0.45 (-2.65)	8.80 ± 0.71 (-0.83)
<i>DDPM</i>	✗	✗	✓	43.81 ± 2.12 (-0.69)	5.70 ± 2.06 (-0.76)	15.11 ± 0.32 (0.44)	9.89 ± 1.00 (0.26)
<i>cDDPM</i>	✓	✗	✗	53.43 ± 1.49 (0.06)	10.70 ± 1.50 (-0.81)	19.92 ± 1.45 (-0.07)	9.86 ± 1.18 (-0.02)
<i>cDDPM</i>	✗	✓	✗	40.59 ± 1.27 (-12.78)	11.05 ± 1.04 (-0.46)	16.56 ± 1.11 (-3.43)	8.34 ± 0.76 (-1.54)
<i>cDDPM</i>	✗	✗	✓	52.29 ± 1.70 (-1.08)	10.00 ± 1.48 (-1.51)	20.41 ± 1.41 (0.42)	9.36 ± 1.15 (-0.52)
<i>pDDPM</i>	✓	✗	✗	49.62 ± 0.84 (-0.16)	10.09 ± 0.76 (0.88)	13.35 ± 0.22 (0.11)	7.62 ± 0.80 (-0.35)
<i>pDDPM</i>	✗	✓	✗	36.09 ± 0.35 (-13.69)	10.40 ± 0.59 (1.19)	12.14 ± 0.21 (-1.10)	7.33 ± 0.39 (-0.64)
<i>pDDPM</i>	✗	✗	✓	47.76 ± 0.86 (-2.02)	8.13 ± 1.03 (-1.08)	13.71 ± 0.28 (0.47)	7.30 ± 0.95 (-0.67)

Table 4: Post-processing analysis for all data sets. The checkmarks indicate the exclusion of Connected Component (CC), Medianfiltering (MF) or Brain Eroding (BE) in the evaluation phase. For all models, the mean ± standard deviation are provided. Color-coded absolute differences concerning the respective baseline models are provided in the brackets.

In Table 4, we provide an analysis of the applied post-processing steps by excluding individual post-processing steps from the evaluation protocol. We show that while the median filter shows to have a large effect, the other post-processing techniques only show minor changes. Moreover, no post-processing strategy consistently works for all models or data sets, motivating further research and a systematic study about the effect of different post-processing steps for UAD in brain MRI.