

# HARMONIZING MULTI-SITE MULTI-SEQUENCE BRAIN MRI VIA SEMANTIC-GUIDED CONDITIONAL DIFFUSION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Training robust AI models for brain MRI analysis typically requires large datasets, prompting many studies to aggregate multi-site data. However, this introduces unwanted variations due to differences in scanners and/or acquisition protocols. These non-biological variations (known as site effects) can significantly compromise the performance and generalizability of downstream deep learning models. While image-level harmonization has emerged as a promising solution, existing methods frequently demand paired data (e.g., scans of the same subject at different sites) or costly encoder-decoder networks to disentangle anatomical content from pre-defined imaging style (e.g., intensity and contrast), which struggle to comprehensively capture diverse image styles. Moreover, existing methods cannot adapt well across different MRI sequences, limiting their scalability. This paper proposes a semantic-guided conditional diffusion (SGCD) framework for unpaired 3D multi-sequence MRI harmonization. The SGCD first trains a conditional diffusion model (CDM) to align multi-site multi-sequence MRIs into a unified, sequence-specific domain, reducing global site-related variations. It then fine-tunes the CDM for target-specific harmonization using a style loss derived from BiomedCLIP trained on medical imaging data. By capturing differences in disentangled semantic image style between the harmonized and target MRIs, this loss enables effective harmonization that preserves anatomical structure and does not require paired training data. We evaluate SGCD on 4,163 T1/T2-weighted MRIs from three multi-site datasets, with results suggesting its superiority over several state-of-the-art methods across voxel-level comparison, downstream classification, and brain tissue segmentation tasks.

## 1 INTRODUCTION

Recent advancements in machine learning (ML) and deep learning (DL) have led to powerful models for neuroimaging analysis, tackling tasks such as brain tissue segmentation, disease classification, and longitudinal studies from MRI scans. The statistical power and robustness of these models depend on access to large-scale training data, which often necessitates the aggregation of multi-site MRI data (An et al., 2022; Tofts & Collins, 2011; Schnack et al., 2010). However, this strategy introduces site effects—deeply embedded non-biological variations from differences in scanner hardware, imaging protocols, and software—that can confound ML/DL models and undermine their training and generalization (Gadewar et al., 2024; Parida et al., 2024).

To address site effects in multi-site studies, various harmonization techniques have been proposed, which can be broadly categorized as feature-level or image-level methods. Feature-level methods (Fortin et al., 2018; Pomponio et al., 2020) typically utilize statistical methods to correct site-specific variance in pre-extracted radiomic features. This approach is primarily limited by its reliance on feature quality and lacks generalizability across diverse downstream tasks (An et al., 2022; Cackowski et al., 2023). Harmonizing raw image data is increasingly recognized as an effective strategy to improve generalizability. Existing image-level harmonization methods that use generative models, such as generative adversarial networks (GANs) for cross-domain image translation (Liu et al., 2021; Chang et al., 2022; Modanwal et al., 2020), often suffer from training

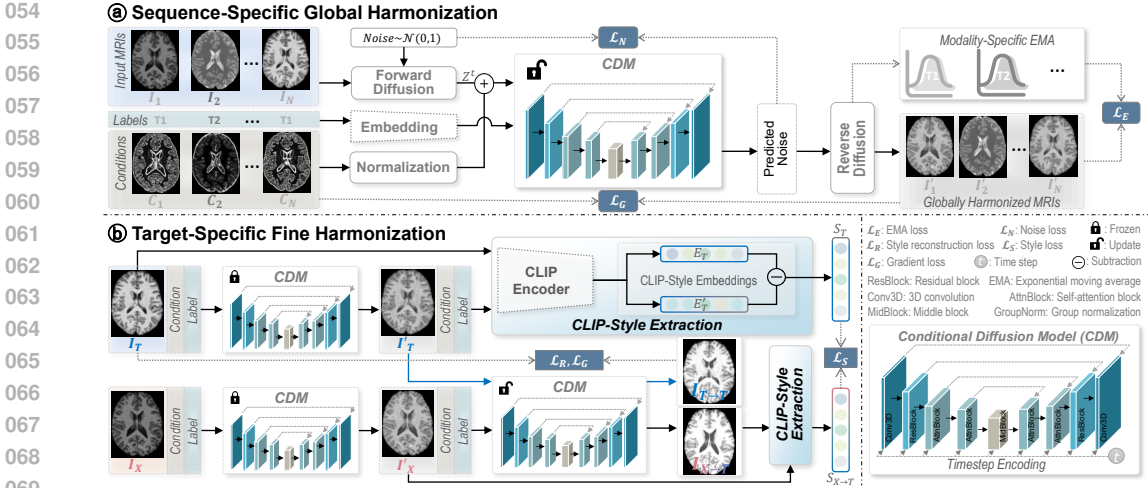


Figure 1: Overview of SGCD with two-stage training: (a) A conditional diffusion model (CDM) is first trained to globally harmonize multi-site multi-sequence MRIs into a sequence-specific unified domain, removing site-related global variations while preserving anatomical content; and (b) CDM is then fine-tuned to harmonize these globally aligned MRIs into the target style using a disentangled CLIP-style loss. This dual-stage process achieves efficient multi-site multi-sequence MRI harmonization without explicit content-style disentanglement.

instability and mode collapse. Some methods using variational autoencoders (VAEs) or multiple encoder-decoder networks aim to learn disentangled representations of anatomical and style information for harmonization (Zuo et al., 2021; Dewey et al., 2020; Cackowski et al., 2023), but they are typically computationally intensive. Another limitation of existing methods is their reliance on paired training data, such as scans from traveling subjects or multiple sequences (e.g., T1- and T2-weighted MRIs) of the same subject. This requirement is often challenging to meet in large-scale retrospective studies, posing a practical barrier to their application.

To this end, we propose a new semantic-guided conditional diffusion framework (SGCD) for multi-sequence 3D brain MRI harmonization. As illustrated in Fig. 1, our SGCD adopts a two-stage progressive training scheme. The first *global harmonization* stage trains a site-agnostic conditional diffusion model (CDM) to align multi-sequence MRIs from all sites to a sequence-specific global domain, mitigating global differences (e.g., intensity and contrast). This is achieved through style-free gradient map conditioning and losses based on exponential moving average (EMA) record, supporting multi-sequence input without requiring paired samples. The second *target-specific harmonization* stage fine-tunes this CDM to harmonize the globally aligned MRIs into a pre-defined target space using a semantic style loss. A pre-trained BiomedCLIP encoder (Zhang et al., 2024) is used to extract image styles by capturing the difference in semantic embeddings between the original target MRIs and their globally aligned counterparts. As these pairs share the same anatomical information, the difference in their semantic embeddings represents the target style, effectively disentangled from the content information. This dual-stage training scheme leverages the conditional generative power of diffusion models and the semantic-rich embeddings from BiomedCLIP, achieving content-style disentanglement during multi-sequence harmonization without requiring auxiliary encoder-decoder networks or paired training samples. The SGCD is trained and evaluated on three multi-site datasets with a total of 4,163 T1-weighted (T1w) or T2-weighted (T2w) MRIs through three tasks. Experimental results demonstrate the superiority of SGCD over several state-of-the-art (SOTA) methods in aligning multi-site and multi-sequence MRI styles while preserving critical biological and anatomical features.

## 2 RELATED WORK

Learning-based neuroimaging studies require large-scale datasets to enable robust training and statistical power (Dufumier et al., 2022; Zhu et al., 2023; Hawco et al., 2022). This is often achieved

108 by retrospectively pooling data across studies, which typically lack standardized scanning proto-  
109 cols. MRI data, however, are highly sensitive to site-related factors such as field strength, sequence  
110 design, and reconstruction software (Parida et al., 2024). These differences yield variations in inten-  
111 sity, contrast, and signal-to-noise ratio that, while negligible to radiologists, confound model training  
112 by entangling biological and site-specific features, resulting in unstable training and poor cross-site  
113 generalization.

114 Existing brain MRI harmonization methods fall into *feature-level* and *image-level* approaches.  
115 Feature-level methods, such as ComBat (Fortin et al., 2018), model site effects as additive and  
116 multiplicative biases within an empirical Bayes framework. While effective for specific radiomic  
117 features (e.g., gray matter volumes, cortical thickness), they are constrained by feature quality and  
118 lack flexibility for diverse downstream tasks (An et al., 2022; Cackowski et al., 2023). Image-level  
119 methods operate directly on raw MRIs, offering broader applicability. Basic techniques such as  
120 min-max, z-score, WhiteStripe, and histogram matching standardize intensity distributions but only  
121 capture global differences, leaving spatially varying contrast and noise uncorrected.

122 Recent ML/DL-based approaches treat image-level harmonization as a generative task akin to image  
123 translation or style transfer. Early methods relied on paired data for supervision, using *traveling sub-*  
124 *jects* scanned across sites (Xu et al., 2024) or paired multi-sequence MRIs (e.g., T1, T2, FLAIR) to  
125 disentangle anatomy from contrast (Zuo et al., 2021; Dewey et al., 2020; Zuo et al., 2023). However,  
126 such paired data are costly and impractical for large-scale retrospective studies.

127 Recent research emphasizes unpaired harmonization, often using CycleGAN (Modanwal et al.,  
128 2020; Liu et al., 2021) or normalizing flows (Beizaee et al., 2023) to map scanner domains. Other  
129 methods disentangle anatomy and style, either by unlearning site/scanner effects (Cackowski et al.,  
130 2023) or by exploiting weak supervision from multi-view 2D slices (Zuo et al., 2022). How-  
131 ever, these approaches face key limitations: (1) disentanglement typically requires multiple  
132 encoder-decoder networks and latent code swapping, which is computationally costly (Ouyang et al.,  
133 2021); (2) reliance on 2D slices neglects 3D spatial context, causing artifacts; and (3) Many are  
134 sequence-specific and require retraining for new sequences or sites, reducing practicality. To ad-  
135 dress these issues, we propose a conditional diffusion model with semantic-rich MRI embeddings  
136 from a pre-trained BiomedCLIP encoder, enabling effective multi-sequence harmonization without  
137 paired data or auxiliary networks.

### 138 3 METHODOLOGY

#### 139 3.0.1 PROBLEM FORMULATION.

140  
141 The goal of image-level harmonization is typically to either transform multi-site MRIs into a unified  
142 virtual domain (Xu et al., 2024) or harmonize all source MRIs into a pre-selected target domain  
143 derived from one or more reference scans (Cackowski et al., 2023; Zuo et al., 2022; Wu et al.,  
144 2025). Our SGCD framework adopts a progressive two-stage training scheme capable of achieving  
145 both objectives. In the first global harmonization stage, all input MRIs, across multiple sequences  
146 and sites, are harmonized into a sequence-specific unified domain, eliminating global site-related  
147 intensity variations. In the second target-specific harmonization stage, all globally aligned MRIs are  
148 further transformed into a pre-selected target domain.

#### 149 3.1 STAGE 1: SEQUENCE-SPECIFIC GLOBAL HARMONIZATION

150  
151 To reduce site effect from multi-site and multi-sequence acquisitions, the first stage aims to normal-  
152 ize global intensity variations while preserving anatomical structure. Unlike traditional site-specific  
153 training, we propose a *sequence-specific global harmonization approach* by training a conditional  
154 diffusion model (CDM) across all sites and sequences simultaneously. The CDM is trained by two  
155 key mechanisms: (1) the sequence-specific, style-free conditions, obtained from sequence labels  
156 and normalized gradient maps; and (2) the exponential moving average (EMA)-based record update,  
157 which serves as a dynamic harmonization target, ensuring the model aligns MRIs of each sequence  
158 to a stable unified domain. This formulation enables better generalization to unseen domains, elim-  
159 inates the need for site-specific models, and establishes a unified intermediate representation that  
160 aids the subsequent target-specific harmonization.  
161

### 3.1.1 1) CONDITIONAL DIFFUSION MODEL:

As shown in Fig. 1 (a), this stage trains a conditional diffusion model  $\Phi$  that takes  $N$  multi-site MRIs  $\{I_i\}_{i=1}^N \in \mathbb{R}^{1 \times W \times H \times D}$  as input. Each MRI  $I_i$  is paired with a sequence label  $m_i \in \{1, \dots, M\}$  drawn from a set of  $M$  classes to differentiate the MRI sequence type (e.g., 1 for T1w, 2 for T2w). Here  $W$ ,  $H$ , and  $D$  denote the width, height, and depth of the 3D MRI volume, respectively. The input then goes through a forward diffusion process (FDP) governed by a Markov chain with a total of  $T$  timesteps. During FDP, noise is sampled from a standard Gaussian distribution and gradually added to  $I_i$  to create a noisy image  $I_i^t$  at each timestep  $t$ :

$$I_i^t = \sqrt{\bar{\alpha}_t} I_i^0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (1)$$

where  $\epsilon$  is the sampled noise,  $\bar{\alpha}_t := \prod_{i=1}^t \alpha_i$ ,  $\alpha_t := 1 - \beta_t$ , and  $\beta_t$  follows a predefined variance schedule (Ho et al., 2020). To preserve anatomical structure of the brain, we use each MRI’s gradient map  $G_{i=1:N} \in \mathbb{R}^{1 \times W \times H \times D}$  as the input condition to the CDM, computed as:

$$G(I_i) = \text{Pad}\left(\frac{1}{3}(\nabla_H I_i + \nabla_W I_i + \nabla_D I_i)\right), \quad (2)$$

where  $\nabla$  is the forward-difference operator along each spatial axis, and  $\text{Pad}(\cdot)$  denotes zero padding to restore the original input size. Each gradient map is normalized to  $[-1, 1]$ . This style-free anatomical condition is concatenated with the noisy image  $I_i^t$  and fed into CDM, while the embedded sequence label  $m_i$  serves as the class condition. CDM is implemented as a time-conditioned 3D U-Net and trained to predict the noise  $\epsilon_\theta$  by minimizing the noise-level loss:

$$\mathcal{L}_N = \|\epsilon - \epsilon_\theta(I_i^t, t, G_i, m_i)\|_2^2. \quad (3)$$

During training, we also derive  $\hat{I}'_i$ , which is an intermediate, one-step estimate of the final globally harmonized MRIs  $I'_i$  through one reverse diffusion process (RDP):

$$\hat{I}'_i \approx I'_i = \frac{1}{\sqrt{\alpha_t}}(I_i^t - \sqrt{1 - \bar{\alpha}_t} \epsilon_\theta(I_i^t, t, G_i, m_i)), \quad (4)$$

### 3.1.2 2) EMA-BASED RECORD UPDATE:

We then use  $\hat{I}'_i$  to update a sequence-specific exponential moving average (EMA) record to guide the global alignment across different sites. The EMA record for each sequence contains a fully differentiable soft histogram, mean, and standard deviation (std.) of intensity values. Let  $x \in \mathbb{R}^F$  be a 1D flattened tensor of  $\hat{I}'_i$ ; we compute a soft-histogram over a fixed value range  $[v_{\min}, v_{\max}]$  with  $K$  bins. The bin centers can be defined as  $c_k = v_{\min} + \frac{k-1}{K-1}(v_{\max} - v_{\min})$  for  $k = 1 : K$ . For each voxel value  $x_i$ , its contribution to the  $k$ -th bin is computed using a Gaussian kernel:

$$w_{ik} = \exp\left(-\frac{1}{2\sigma^2}(x_i - c_k)^2\right), \quad (5)$$

where  $\sigma$  is the kernel bandwidth controlling smoothness. The normalized soft-histogram  $\mathcal{H}(x) \in \mathbb{R}^K$  is defined by its  $k$ -th components as:

$$\mathcal{H}_k(x) = \frac{\sum_{i=1}^F w_{ik}}{\sum_{j=1}^K (\sum_{i=1}^F w_{ij}) + \delta}, \quad \text{for } k = 1, \dots, K, \quad (6)$$

where  $\delta$  is a small constant to avoid division by zero. We denote the differentiable soft-histogram of  $\hat{I}'_i$  as  $\mathcal{H}(\hat{I}'_i)$ . We compute the soft-histogram for each  $\hat{I}'_i$  and update the corresponding EMA record for the  $m$ -th sequence as follows:

$$EMA_m^{\{\mathcal{H}\}} = \gamma \cdot EMA_m^{\{\mathcal{H}\}} + (1 - \gamma)\mathcal{H}(\hat{I}'_i), \quad (7)$$

where  $\gamma \in [0, 1)$  is the EMA decay factor controlling the update rate. We update the  $EMA_m^{\{\mu\}}$  and  $EMA_m^{\{\sigma\}}$  for the voxel mean and std. intensities similarly. This EMA record serves as the running estimate of the sequence-specific intensity statistics. After every EMA update, we calculate an EMA loss to guide the global harmonization, defined as:

$$\mathcal{L}_E = WD(EMA_m^{\{\mathcal{H}\}}, \mathcal{H}(\hat{I}'_i)) + \|EMA_m^{\{\mu\}} - \mu(\hat{I}'_i)\|_2^2 + \|EMA_m^{\{\sigma\}} - \sigma(\hat{I}'_i)\|_2^2, \quad (8)$$

where the  $WD(\cdot)$  is the differentiable Wasserstein distance, which is computed as the mean absolute difference between the cumulative distribution functions of two soft-histograms.

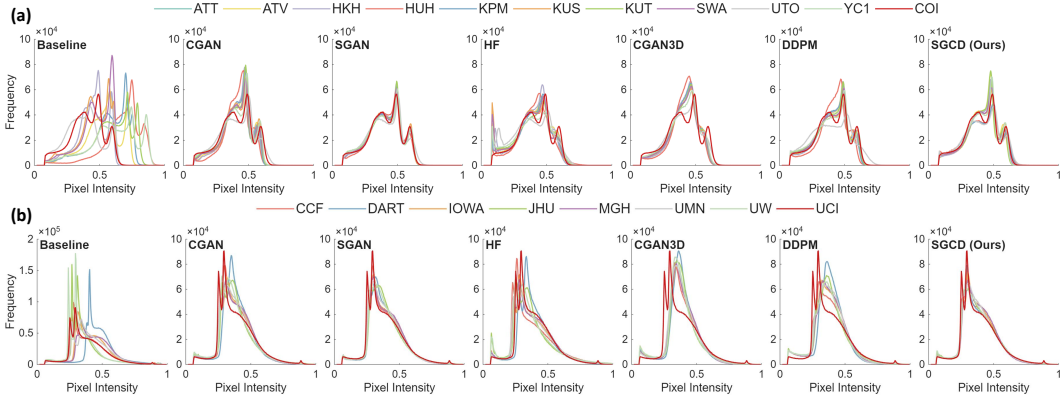


Figure 2: Histograms of (a) 22 T1w MRIs from SRPBS test set across 11 sites, with COI as the target domain; (b) 16 T2w MRIs from DWI-THP test set across 8 sites, with UCI as the target domain.

### 3.1.3 ANATOMY PRESERVATION CONSTRAINT:

To preserve brain anatomical structures during harmonization, we compute the normalized gradient map  $G_i$  of  $I'_i$ , and define a gradient loss conditioned on  $G_i$  as:  $\mathcal{L}_G = \|G_i - G(\hat{I}'_i)\|_2^2$ . The hybrid loss for Stage 1 is defined as:  $\mathcal{L}_1 = \mathcal{L}_N + \mathcal{L}_G + \mathcal{L}_E$ . This first stage training creates a global harmonizer capable of aligning multi-site, multi-sequence MRIs simultaneously to unified domains, achieving target-free harmonization.

### 3.2 STAGE 2: TARGET-SPECIFIC FINE HARMONIZATION

To adapt globally harmonized MRIs to a specific target style, the second stage fine-tunes the pre-trained CDM from Stage 1 using unpaired source and target data of the same sequence. This stage builds upon the site-agnostic intermediate representations generated earlier, introducing a novel style guidance derived from a pre-trained vision-language model BiomedCLIP (Zhang et al., 2024). The fine-tuning is guided by two key principles: (1) a diffusion-based translation process, which iteratively transforms globally harmonized MRIs to the target style; and (2) a disentangled CLIP-style loss, which measures style differences in a semantic embedding space, removing the need for paired data or explicit style-content separation. This formulation not only achieves efficient target-specific harmonization while preserving anatomy and sequence specificity but also offers improved generalizability on unseen domains.

As shown in Fig. 1 (b), given multi-sequence source MRIs  $I_X$  and unpaired target MRIs  $I_T$  of the same sequence, we first globally harmonize them into  $I'_X$  and  $I'_T$  using the pre-trained CDM from the first stage. We then fine-tune the CDM to translate  $I'_X$  into a target-style MRI  $I_{X \rightarrow T}$  guided by a pre-trained BiomedCLIP encoder (Zhang et al., 2024) that implicitly extracts style embeddings. We introduce a *hybrid disentangled CLIP-style loss* to achieve this translation without requiring paired data or explicit style definitions. The globally harmonized  $I'_X$  and  $I'_T$  are obtained through CDM inference that employs the DDIM sampling strategy (Song et al., 2020) on RDP. Instead of the one-step estimate  $\hat{I}'_i$  used in training (see Eq. 4), we obtain  $I'_i$  by iteratively denoising over  $T_r$  steps ( $t = T_r : 0$ ):

$$I_i^{t-1} = \sqrt{\bar{\alpha}_{t-1}} \hat{I}'_i + \sqrt{1 - \bar{\alpha}_{t-1}} \epsilon_\theta(I_i^t, t, G_i, m_i), \tag{9}$$

where  $\hat{I}'_i$  is the intermediate estimate used during training and  $\epsilon_\theta$  is the trained CDM function.

#### 3.2.1 TARGET-SPECIFIC FINE-TUNING:

The CDM is then fine-tuned to adapt  $I'_X$  to match the style of  $I_T$  through the DDIM sampling strategy, containing an FDP followed by an RDP. Unlike Eq. 1, where random noise is added to  $I_i$ , we now iteratively add CDM-generated noise to  $I'_X$  over  $T_f$  forward iterations ( $t = 0 : T_f$ ) to get the noisy image  $I_X^{T_f}$ :

$$I_X^{t+1} = \sqrt{\bar{\alpha}_{t+1}} I'_X + \sqrt{1 - \bar{\alpha}_{t+1}} \epsilon_\theta(I_X^t, t, G_X, m_X), \tag{10}$$

We then apply RDP to denoise  $I_X^{T_f}$  back to the final translated MRI over  $t = T_r : 0$  steps:

$$I_X^{t-1} = \sqrt{\bar{\alpha}_{t-1}} I_X' + \sqrt{1 - \bar{\alpha}_{t-1}} \epsilon_\theta(I_X^t, t, G_X, m_X), \quad (11)$$

After  $T_r$  timesteps, we obtain the harmonized MRI:  $I_{X \rightarrow T}$ .

### 3.2.2 DISENTANGLED CLIP-STYLE GUIDANCE:

To guide style translation without requiring paired data or explicit style learning, we incorporate a pre-trained BiomedCLIP encoder to extract implicit style representations. Specifically, we compute style embeddings as the difference in CLIP-space between each MRI and its globally-aligned counterpart:

$$S_T = \Psi(I_T) - \Psi(I_T'), \quad S_{X \rightarrow T} = \Psi(I_{X \rightarrow T}) - \Psi(I_X'), \quad (12)$$

where  $\Psi(\cdot)$  is the BiomedCLIP image encoder. Since  $I_T'$  shares content with  $I_T$ , the difference  $S_T$  captures target-specific style, disentangled from content information. Similarly,  $S_{X \rightarrow T}$  reflects the style of the harmonized source. We then define the *style translation loss* as:

$$\mathcal{L}_S = \|S_T - S_{X \rightarrow T}\|_1 + \left(1 - \frac{S_T \cdot S_{X \rightarrow T}}{\|S_T\| \|S_{X \rightarrow T}\|}\right), \quad (13)$$

where the 1st term is the  $l_1$  distance in the CLIP-embedding space and the 2nd term quantifies the directional discrepancy between two style embeddings.

To ensure style consistency, we further design a *style reconstruction loss* by minimizing style embeddings of each target MRI and its harmonized counterpart:  $\mathcal{L}_R = \|S_T - S_{T \rightarrow T}\|_1$ . The hybrid disentangled CLIP-style loss is defined as:  $\mathcal{L}_2 = \mathcal{L}_S + \mathcal{L}_R$ .

By leveraging BiomedCLIP’s semantic-rich embeddings, the CDM effectively translates source MRIs to the target style without requiring paired training data or explicit image style and content disentanglement learning, ensuring that anatomical content remains unchanged.

### 3.2.3 ADAPTATION TO NEW DATA.

The proposed two-stage training strategy improves the generalizable potential of SGCD. When MRIs from an unseen site are used as the source domain, they can be harmonized without retraining the model. As the global harmonizer, trained across multiple sites and sequences, generalizes to new domains by aligning input MRIs to a learned site-agnostic intermediate representation. The target-specific harmonizer then maps these globally harmonized MRIs into a learned target style. If an unseen site is used as the target domain, only the second-stage model requires fine-tuning to learn the new domain-specific style, while the first-stage harmonizer remains unchanged.

### 3.2.4 IMPLEMENTATION.

We use the MONAI (Cardoso et al., 2022) framework to implement the proposed SGCD approach. CDM is implemented as a time-conditioned 3D U-Net with a symmetric architecture comprising 2 upsampling/downsampling layers, 2 residual blocks, 4 self-attention blocks, and a middle block, with channels {32, 64, 256, 256}, respectively. Stage 1 training uses the default Adam optimizer with an initial learning rate (LR) of  $1 \times 10^{-4}$ , while Stage 2 fine-tuning uses an LR of  $1 \times 10^{-6}$ . For EMA-based histogram/record updates, we set the decay factor  $\gamma = 0.2$  and compute soft histograms with  $K = 100$  bins over the range  $[v_{\min} = 0, v_{\max} = 1]$ . The variance scheduler  $\beta$  is empirically set to increase linearly from 0.0015 to 0.0195. We apply  $T = 1,000$  noise steps for the global harmonization stage and empirically set  $T_f = 35$  and  $T_r = 25$  for the target-specific harmonization stage via grid search. See Appendix A for detailed implementations.

## 4 EXPERIMENT

### 4.0.1 STUDIED MATERIALS.

Three brain MRI datasets are utilized: (1) OpenBHB (Dufumier et al., 2022), with T1w MRIs from 3,984 healthy subjects across 58 sites; (2) SRPBS (Tanaka et al., 2021), with T1w MRIs from 9 traveling subjects across 11 sites; and (3) DWI-THP (Magnotta et al., 2012) with T1w and T2w

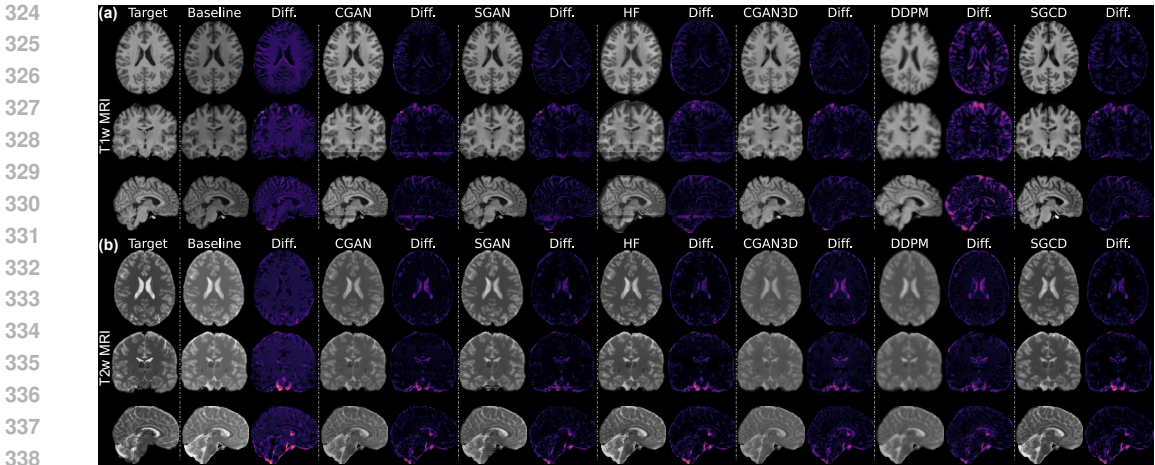


Figure 3: Visualization of harmonization results for (a) a T1-weighted MRI sample from SRPBS and (b) a T2-weighted sample from DWI-THP. Each panel shows the original source MRI (Baseline), harmonized source MRIs from six different methods, and the corresponding difference (Diff.) maps between the harmonized source and target MRIs.

Table 1: Comparison between source site MRIs and corresponding target site (COI) MRIs with matching subjects on the SRPBS test set (2 subjects across 11 sites).

Method	SSIM $\uparrow$	PSNR $\uparrow$	PCC $\uparrow$	WD $\downarrow$
Baseline	0.881 <sup>0.03</sup>	22.03 <sup>4.29</sup>	0.980 <sup>0.01</sup>	0.041 <sup>0.03</sup>
CGAN	0.888 <sup>0.03</sup>	27.64 <sup>1.64</sup>	0.976 <sup>0.01</sup>	0.006 <sup>0.002</sup>
SGAN	0.894 <sup>0.02</sup>	27.33 <sup>1.51</sup>	0.978 <sup>0.01</sup>	<b>0.005</b> <sup>0.002</sup>
HF	0.850 <sup>0.02</sup>	26.38 <sup>1.55</sup>	0.969 <sup>0.01</sup>	0.008 <sup>0.002</sup>
CGAN3D	0.866 <sup>0.02</sup>	27.55 <sup>1.48</sup>	0.976 <sup>0.01</sup>	0.009 <sup>0.003</sup>
DDPM	0.800 <sup>0.10</sup>	25.39 <sup>3.50</sup>	0.950 <sup>0.03</sup>	0.006 <sup>0.001</sup>
SGCD (Ours)	<b>0.896</b> <sup>0.02</sup>	<b>28.04</b> <sup>1.71</sup>	<b>0.988</b> <sup>0.01</sup>	<b>0.005</b> <sup>0.001</sup>

MRIs from 5 subjects scanned across 8 sites. The OpenBHB is split into a training set (3,227 MRIs) and a validation set (757 MRIs). For SRPBS, we use 66 MRIs from 6 subjects across 11 sites for training, 11 MRIs from 1 subject for validation, and 22 MRIs from 2 subjects for testing. For DWI-THP, 48 MRIs from 3 subjects across 8 sites (both modalities) are used for training, and the remaining 32 MRIs for testing. Given the lack of standardized criteria for target site selection in brain MRI harmonization, we follow (Tian et al., 2022) and select the site exhibiting the *lowest intra-site style variability* for each dataset, measured by the mean Wasserstein Distance (WD) across samples.

#### 4.0.2 COMPETING METHODS.

We compare our SGCD with five SOTA image-level MRI harmonization methods: CycleGAN (CGAN) (Chang et al., 2022), StyleGAN (SGAN) (Liu et al., 2021), 3D CycleGAN (CGAN3D) (Zhu et al., 2017), Harmonizing Flow (HF) (Beizae et al., 2023), and DDPM (Durrer et al., 2023). The 3D methods (*i.e.*, CGAN3D and DDPM) are trained using the same volumetric dataset as our method, while the 2D methods (*i.e.*, CGAN, SGAN, and HF) use axial MRI slices derived from the same dataset. We ensured consistent training hyperparameters across all methods for fair comparisons.

#### 4.0.3 TASK 1: HISTOGRAM & VOXEL-LEVEL COMPARISON.

This experiment utilizes SRPBS and DWI-THP with traveling subjects. For SRPBS, each method harmonizes MRIs from 10 source sites to a target site (ID: COI). For DWI-THP, harmonization is performed from 7 source sites to a target site (ID: UCI). Performance is assessed via histogram comparisons and voxel-level metrics: PSNR, structural similarity index (SSIM), Wasserstein distance (WD), and Pearson correlation coefficient (PCC) with raw source MRIs as **Baseline**. Our method is trained end-to-end across all T1w and T2w MRIs in DWI-THP simultaneously, while

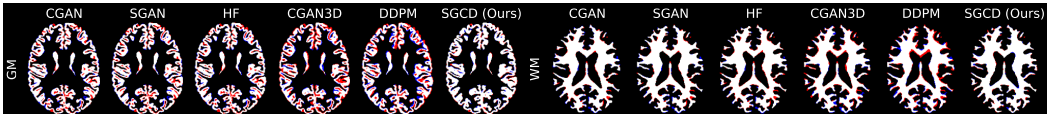


Figure 4: Segmentation maps with White: accurate segmentation; Red: under-segmentation; Blue: over-segmentation.

Table 2: Results (%) achieved by different methods in terms of AP and DSC metrics for gray matter (GM) and white matter (WM) segmentation on DWI-THP.

Method	AP $\uparrow$			DSC $\uparrow$		
	GM	WM	Mean	GM	WM	Mean
CGAN	96.1 <sup>1.3</sup>	96.0 <sup>0.6</sup>	96.1 <sup>0.9</sup>	88.0 <sup>0.9</sup>	92.0 <sup>0.3</sup>	90.0 <sup>0.6</sup>
SGAN	95.7 <sup>0.7</sup>	97.1 <sup>0.6</sup>	96.4 <sup>0.4</sup>	87.9 <sup>0.8</sup>	92.0 <sup>0.3</sup>	90.0 <sup>0.5</sup>
HF	94.5 <sup>1.3</sup>	95.2 <sup>0.8</sup>	94.8 <sup>0.9</sup>	87.8 <sup>1.0</sup>	91.7 <sup>0.4</sup>	89.7 <sup>0.6</sup>
CGAN3D	86.2 <sup>1.5</sup>	96.0 <sup>0.6</sup>	91.1 <sup>0.7</sup>	82.8 <sup>0.9</sup>	90.1 <sup>0.2</sup>	86.4 <sup>0.5</sup>
DDPM	92.9 <sup>1.4</sup>	91.9 <sup>1.4</sup>	92.4 <sup>1.3</sup>	74.4 <sup>8.3</sup>	79.0 <sup>9.8</sup>	76.7 <sup>9.0</sup>
SGCD (Ours)	<b>97.8<sup>0.4</sup></b>	<b>99.6<sup>0.3</sup></b>	<b>98.7<sup>0.2</sup></b>	<b>93.0<sup>0.7</sup></b>	<b>94.7<sup>0.2</sup></b>	<b>93.8<sup>0.4</sup></b>

the competing deep models require separate training for each sequence. As shown in Fig. 2, both datasets exhibit substantial site-wise intensity variations (Baseline). Without relying on explicit style learning or retraining for each sequence, our method effectively aligns source MRI histograms to those of the target site across both sequences. Additional results are presented in Appendix A.

As shown in Table 1, SGCD achieves the highest SSIM, PSNR, and PCC, indicating superior voxel-level agreement and anatomy preservation. It also achieves the best WD score (tied with SGAN), confirming effective style alignment. The visualization results in Fig. 3 show that SGCD-harmonized MRIs more closely resemble the target domain for T1w and T2w sequences, while the 2D methods (CGAN, SGAN, and HF) generate strip artifacts in different views. These results highlight SGCD’s ability to harmonize MRIs while maintaining high image quality and anatomy fidelity.

#### 4.0.4 TASK 2: BRAIN TISSUE SEGMENTATION.

We further evaluate anatomy preservation via a brain tissue segmentation task on T1 MRIs in DWI-THP. FreeSurfer (Billot et al., 2023) is used to generate gray matter (GM) and white matter (WM) segmentation maps for original and harmonized MRIs. Segmentation quality is assessed using the Anatomical Preservation (AP) score (Parida et al., 2024), which measures the relative absolute difference in tissue volumes, and Dice Similarity Coefficient (DSC), which quantifies the spatial overlap between segmentation maps. Table 2 shows that SGCD achieves the highest AP and DSC for both GM and WM, and mean scores. Figure 4 further shows that SGCD yields fewer segmentation errors in the WM and GM tissue boundaries. This superior anatomical fidelity may be attributed to our gradient-based anatomy conditioning and constraint, and the implicit content-style disentanglement in CLIP-style guidance during the target-specific fine-tuning stage.

Table 3: Site classification and age prediction results on original (Baseline) and harmonized MRIs on OpenBHB.

Method	Site Classification (%)				Age Prediction	
	BAC $\downarrow$	F1 $\downarrow$	PRE $\downarrow$	Recall $\downarrow$	MAE $\downarrow$	MSE $\downarrow$
Baseline	34.3 <sup>2.40</sup>	66.3 <sup>2.30</sup>	75.7 <sup>1.80</sup>	73.2 <sup>1.90</sup>	5.30 <sup>0.260</sup>	47.4 <sup>1.41</sup>
CGAN	42.5 <sup>1.60</sup>	69.5 <sup>2.70</sup>	77.0 <sup>3.00</sup>	73.9 <sup>2.00</sup>	6.63 <sup>0.264</sup>	79.0 <sup>10.5</sup>
SGAN	25.8 <sup>2.20</sup>	59.3 <sup>1.20</sup>	66.2 <sup>1.50</sup>	65.1 <sup>1.50</sup>	7.31 <sup>0.494</sup>	85.7 <sup>12.9</sup>
HF	34.2 <sup>1.10</sup>	66.5 <sup>2.00</sup>	73.6 <sup>2.10</sup>	72.3 <sup>2.10</sup>	5.84 <sup>0.221</sup>	57.0 <sup>3.85</sup>
CGAN3D	32.4 <sup>2.90</sup>	65.6 <sup>1.90</sup>	75.1 <sup>1.90</sup>	72.3 <sup>1.70</sup>	5.90 <sup>0.360</sup>	<b>33.3<sup>9.67</sup></b>
DDPM	16.6 <sup>1.60</sup>	56.0 <sup>1.30</sup>	<b>54.5<sup>0.70</sup></b>	53.5 <sup>2.00</sup>	5.33 <sup>0.258</sup>	48.5 <sup>6.80</sup>
SGCD	<b>14.5<sup>1.30</sup></b>	<b>54.3<sup>2.50</sup></b>	56.0 <sup>1.80</sup>	<b>52.5<sup>3.14</sup></b>	<b>5.24<sup>0.141</sup></b>	54.0 <sup>2.02</sup>

#### 4.0.5 TASK 3: SITE CLASSIFICATION & BRAIN AGE PREDICTION.

We evaluate SGCD in reducing site-related variations through two downstream tasks (*i.e.*, site classification and brain age prediction). Each method is trained on OpenBHB training data and applied

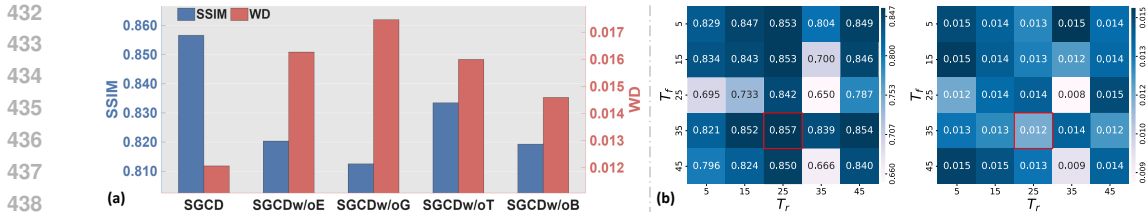


Figure 5: Results of (a) SGCD with 4 variants on DWI-THP and (b) with different parameters.

to harmonize the validation data with Site 17 as the target site. We extract deep features from these harmonized MRIs using ResNet18 (He et al., 2016) (with its final layer removed). A logistic regression, trained on 70% of these deep features and tested on 30%, performs multiclass site classification, and a ridge regressor predicts brain age. Both tasks are repeated 5 times for random data partition, with mean and standard deviation results reported in Table 3. Lower site classification results indicate better removal of site-related variations, while lower age prediction error suggests superior anatomical feature preservation. Table 3 shows that SGCD yields the lowest site classification performance in BAC, F1, and Recall, effectively removing site-related features and maintaining faithful anatomical integrity with lower mean absolute error (MAE).

## 5 DISCUSSION

### 5.0.1 ABLATION STUDY.

We assess four key components of SGCD by comparing it to its four variants: **SGCDw/oE** (without EMA-based constraint), **SGCDw/oG** (without gradient-based anatomical condition and constraint), **SGCDw/oT** (without Stage 2), and **SGCDw/oB** that uses style embeddings from CLIP pretrained on natural images (Radford et al., 2021) (rather than from BiomedCLIP pretrained on medical data). As shown in Fig. 5 (a), SGCDw/oE shows degraded performance (lower SSIM, higher WD), indicating the EMA constraint is critical for globally aligning multi-sequence inputs into their unified domains. SGCDw/oG performs the worst across all metrics, highlighting the necessity of gradient-based anatomical conditioning for maintaining structural fidelity and facilitating content-style disentanglement. SGCDw/oT aligns global intensity but fails to capture local style features such as tissue contrast and texture, confirming the importance of the target-specific stage. SGCDw/oB shows moderate WD but compromises anatomical integrity. This implies that CLIP trained on natural images lacks domain-specific knowledge to capture subtle variations in medical imaging, such as site-related style differences. In contrast, SGCD employs BiomedCLIP, pretrained on the PMC-15M dataset of medical image-text pairs with diverse modalities, like MRI, CT, and X-ray (Zhang et al., 2024), enabling domain-aware harmonization that preserves anatomy while adapting site- and sequence-specific styles.

### 5.0.2 INFLUENCE OF DIFFUSION PARAMETERS.

We perform a grid search over the forward ( $T_f$ ) and reverse ( $T_r$ ) diffusion steps, evaluating volume-level metrics on the DWI-THP test set. As shown in Fig. 5 (b), the model performs best when  $T_f$  and  $T_r$  are comparable. While some combinations yield lower WD scores, SGCD with  $T_f = 35$  and  $T_r = 25$  offers good image quality by prioritizing anatomical fidelity.

## 6 CONCLUSION

We present a semantic-guided conditional diffusion (SGCD) framework for multi-site multi-sequence MRI harmonization. SGCD first aligns MRIs into a unified, sequence-specific space via style-free gradient conditioning, then performs target-specific harmonization using CLIP-based semantic style embeddings, enabling effective volume-level harmonization without paired data or explicit style learning. Evaluated on three multi-site datasets, SGCD outperforms SOTA methods in removing site-related variations while preserving anatomical fidelity across T1w and T2w MRIs. Future work will extend SGCD to MRIs with pathological features such as lesions and tumors.

## REFERENCES

- 486  
487  
488 Lijun An, Jianzhong Chen, Pansheng Chen, Chen Zhang, Tong He, Christopher Chen, Juan Helen  
489 Zhou, BT Thomas Yeo, Lifestyle Study of Aging, Alzheimer’s Disease Neuroimaging Initiative,  
490 et al. Goal-specific brain MRI harmonization. *Neuroimage*, 263:119570, 2022.
- 491 Farzad Beizae, Christian Desrosiers, Gregory A Lodygensky, and Jose Dolz. Harmonizing Flows:  
492 Unsupervised MR harmonization based on normalizing flows. In *International Conference on*  
493 *Information Processing in Medical Imaging*, pp. 347–359. Springer, 2023.
- 494 Benjamin Billot, Colin Magdamo, You Cheng, Steven E Arnold, Sudeshna Das, and Juan Eugenio  
495 Iglesias. Robust machine learning segmentation for large-scale analysis of heterogeneous clinical  
496 brain mri datasets. *Proceedings of the National Academy of Sciences*, 120(9):e2216399120, 2023.
- 497 Stenzel Cackowski, Emmanuel L Barbier, Michel Dojat, and Thomas Christen. Imunity: A gener-  
498 alizable VAE-GAN solution for multicenter MR image harmonization. *Medical Image Analysis*,  
499 88:102799, 2023.
- 500 M Jorge Cardoso, Wenqi Li, Richard Brown, Nic Ma, Eric Kerfoot, Yiheng Wang, Benjamin Murrey,  
501 Andriy Myronenko, Can Zhao, Dong Yang, et al. Monai: An open-source framework for deep  
502 learning in healthcare. *arXiv preprint arXiv:2211.02701*, 2022.
- 503 Xiao Chang, Xin Cai, Yibo Dan, Yang Song, Qing Lu, Guang Yang, and Shengdong Nie. Self-  
504 supervised learning for multi-center magnetic resonance imaging harmonization without traveling  
505 phantoms. *Physics in Medicine & Biology*, 67(14):145004, 2022.
- 506 Blake E Dewey, Lianrui Zuo, Aaron Carass, Yufan He, Yihao Liu, Ellen M Mowry, Scott Newsome,  
507 Jiwon Oh, Peter A Calabresi, and Jerry L Prince. A disentangled latent space for cross-site MRI  
508 harmonization. In *Medical Image Computing and Computer-Assisted Intervention*, pp. 720–729.  
509 Springer, 2020.
- 510 Benoit Dufumier, Antoine Grigis, Julie Victor, Corentin Ambroise, Vincent Frouin, and Edouard  
511 Duchesnay. OpenBHB: A large-scale multi-site brain MRI data-set for age prediction and debi-  
512 asing. *NeuroImage*, 263:119637, 2022.
- 513 Alicia Durrer, Julia Wolleb, Florentin Bieder, Tim Sinnecker, Matthias Weigel, Robin Sandkühler,  
514 Cristina Granziera, Özgür Yaldizli, and Philippe C Cattin. Diffusion models for contrast harmo-  
515 nization of magnetic resonance images. *arXiv preprint arXiv:2303.08189*, 2023.
- 516 Jean-Philippe Fortin, Nicholas Cullen, Yvette I Sheline, Warren D Taylor, Irem Aselcioglu, Philip A  
517 Cook, Phil Adams, Crystal Cooper, Maurizio Fava, Patrick J McGrath, et al. Harmonization of  
518 cortical thickness measurements across scanners and sites. *NeuroImage*, 167:104–120, 2018.
- 519 Shruti P Gadewar, Alyssa H Zhu, Iyad Ba Gari, Sunanda Somu, Sophia I Thomopoulos, Paul M  
520 Thompson, Talia M Nir, and Neda Jahanshad. Synthesizing study-specific controls using gener-  
521 ative models on open access datasets for harmonized multi-study analyses. *arXiv preprint*  
522 *arXiv:2403.00093*, 2024.
- 523 Colin Hawco, Erin W Dickie, Gabrielle Herman, Jessica A Turner, Miklos Argyelan, Anil K Mal-  
524 hotra, Robert W Buchanan, and Aristotle N Voineskos. A longitudinal multi-scanner multimodal  
525 human neuroimaging dataset. *Scientific Data*, 9(1):332, 2022.
- 526 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-  
527 nition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp.  
528 770–778, 2016.
- 529 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*  
530 *Neural Information Processing Systems*, 33:6840–6851, 2020.
- 531 Mengting Liu, Piyush Maiti, Sophia Thomopoulos, Alyssa Zhu, Yaqiong Chai, Hosung Kim, and  
532 Neda Jahanshad. Style transfer using generative adversarial networks for multi-site MRI har-  
533 monization. In *Medical Image Computing and Computer Assisted Intervention, Part III 24*, pp.  
534 313–322. Springer, 2021.

- 540 Vincent A Magnotta, Joy T Matsui, Dawei Liu, Hans J Johnson, Jeffrey D Long, Bradley D Bol-  
541 ster Jr, Bryon A Mueller, Kelvin Lim, Susumu Mori, Karl G Helmer, et al. Multicenter reliability  
542 of diffusion tensor imaging. *Brain Connectivity*, 2(6):345–355, 2012.
- 543 Gourav Modanwal, Adithya Vellal, Mateusz Buda, and Maciej A Mazurowski. MRI image harmo-  
544 nization using cycle-consistent generative adversarial network. In *Computer-Aided Diagnosis*,  
545 volume 11314, pp. 259–264. SPIE, 2020.
- 547 Jiahong Ouyang, Ehsan Adeli, Kilian M Pohl, Qingyu Zhao, and Greg Zaharchuk. Representation  
548 disentanglement for multi-modal brain mri analysis. In *International conference on information*  
549 *processing in medical imaging*, pp. 321–333. Springer, 2021.
- 550 Abhijeet Parida, Zhifan Jiang, Roger J Packer, Robert A Avery, Syed M Anwar, and Marius G  
551 Linguraru. Quantitative metrics for benchmarking medical image harmonization. In *2024 IEEE*  
552 *International Symposium on Biomedical Imaging (ISBI)*, pp. 1–5. IEEE, 2024.
- 554 R Pomponio, G Erus, M Habes, J Doshi, D Srinivasan, E Mamourian, V Bashyam, IM Nasrallah,  
555 TD Satterthwaite, Y Fan, et al. Harmonization of large MRI datasets for the analysis of brain  
556 imaging patterns throughout the lifespan, 2020.
- 557 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,  
558 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual  
559 models from natural language supervision. In *International conference on machine learning*, pp.  
560 8748–8763. PmLR, 2021.
- 561 Hugo G Schnack, Neeltje EM van Haren, Rachel M Brouwer, G Caroline M van Baal, Marco Pic-  
562 chioni, Matthias Weisbrod, Heinrich Sauer, Tyrone D Cannon, Matti Huttunen, Claude Lepage,  
563 et al. Mapping reliability in multicenter MRI: Voxel-based morphometry and cortical thickness.  
564 *Human Brain Mapping*, 31(12):1967–1982, 2010.
- 565 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv*  
566 *preprint arXiv:2010.02502*, 2020.
- 568 Saori Tanaka, Ayumu Yamashita, Noriaki Yahata, Takashi Itahashi, Giuseppe Lisi, Takashi Ya-  
569 mada, Naho Ichikawa, Masahiro Takamura, Yujiro Yoshihara, Akira Kunimatsu, Naohiro Okada,  
570 Ryuichiro Hashimoto, Go Okada, Yuki Sakai, Jun Morimoto, Jin Narumoto, Yasuhiro Shimada,  
571 Hiroaki Mano, Wako Yoshida, and Hiroshi Imamizu. A multi-site, multi-disorder resting-state  
572 magnetic resonance image database. *Scientific Data*, 8(1):227, 2021.
- 573 Dezheng Tian, Zilong Zeng, Xiaoyi Sun, Qiqi Tong, Huanjie Li, Hongjian He, Jia-Hong Gao, Yong  
574 He, and Mingrui Xia. A deep learning-based multisite neuroimage harmonization framework  
575 established with a traveling-subject dataset. *NeuroImage*, 257:119297, 2022.
- 576 PS Tofts and DJ Collins. Multicentre imaging measurements for oncology and in the brain. *The*  
577 *British Journal of Radiology*, 84:S213–S226, 2011.
- 579 Mengqi Wu, Lintao Zhang, Pew-Thian Yap, Hongtu Zhu, and Mingxia Liu. Disentangled latent  
580 energy-based style translation: An image-level structural mri harmonization framework. *Neural*  
581 *Networks*, 184:107039, 2025.
- 582 Chundan Xu, Jie Li, Yakui Wang, Lixue Wang, Yizhe Wang, Xiaofeng Zhang, Weiqi Liu, Jingang  
583 Chen, Aleksandra Vatian, Natalia Gusarova, et al. Simix: A domain generalization method for  
584 cross-site brain mri harmonization via site mixing. *NeuroImage*, 299:120812, 2024.
- 586 Sheng Zhang, Yanbo Xu, Naoto Usuyama, Hanwen Xu, Jaspreet Bagga, Robert Tinn, Sam Pre-  
587 ston, Rajesh Rao, Mu Wei, Naveen Valluri, Cliff Wong, Andrea Tupini, Yu Wang, Matt Mazzola,  
588 Swadheen Shukla, Lars Liden, Jianfeng Gao, Angela Crabtree, Brian Piening, Carlo Bifulco,  
589 Matthew P. Lungren, Tristan Naumann, Sheng Wang, and Hoifung Poon. A multimodal biomed-  
590 ical foundation model trained from fifteen million image–text pairs. *NEJM AI*, 2(1), 2024.
- 591 Jun-Ding Zhu, Yung-Fu Wu, Shih-Jen Tsai, Ching-Po Lin, and Albert C Yang. Investigating brain  
592 aging trajectory deviations in different brain regions of individuals with schizophrenia using mul-  
593 timodal magnetic resonance imaging and brain-age prediction: A multicenter study. *Translational*  
*Psychiatry*, 13(1):82, 2023.

- 594 Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation  
595 using cycle-consistent adversarial networks. In *Proceedings of the IEEE International Conference*  
596 *on Computer Vision*, pp. 2223–2232, 2017.
- 597 Lianrui Zuo, Blake E Dewey, Yihao Liu, Yufan He, Scott D Newsome, Ellen M Mowry, Susan M  
598 Resnick, Jerry L Prince, and Aaron Carass. Unsupervised MR harmonization by learning disen-  
599 tangled representations using information bottleneck theory. *NeuroImage*, 243:118569, 2021.
- 600 Lianrui Zuo, Yihao Liu, Yuan Xue, Shuo Han, Murat Bilgel, Susan M Resnick, Jerry L Prince, and  
601 Aaron Carass. Disentangling a single mr modality. In *MICCAI workshop on data augmentation,*  
602 *labelling, and imperfections*, pp. 54–63. Springer, 2022.
- 603 Lianrui Zuo, Yihao Liu, Yuan Xue, Blake E Dewey, Samuel W Remedios, Savannah P Hays, Murat  
604 Bilgel, Ellen M Mowry, Scott D Newsome, Peter A Calabresi, et al. Haca3: A unified approach for  
605 multi-site mr image harmonization. *Computerized Medical Imaging and Graphics*, 109:102285,  
606 2023.
- 607  
608

## 609 A APPENDIX

610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647

# HARMONIZING MULTI-SITE MULTI-SEQUENCE BRAIN MRI VIA SEMANTIC-GUIDED CONDITIONAL DIFFUSION

– *Supplementary Materials*

**Anonymous authors**

Paper under double-blind review

## 1 SGCD ALGORITHM IMPLEMENTATION

Here we include a brief description of the algorithm of the first sequence-specific global harmonization stage and the second target-specific fine harmonization stage. The global harmonization is trained for  $N_1 = 300$  epochs and the target-specific harmonization is fine-tuned for  $N_2 = 30$  epochs.

---

### Algorithm 1 Sequence-Specific Global Harmonization

---

**Input:** Multi-site MRIs with sequence label  $\{I_i, m_i\}_{i=1}^N$

**Parameter:** Number of epochs  $n_{\text{epochs}}$ , EMA decay  $\gamma$ , histogram bins  $K$ , value range  $[v_{\min}, v_{\max}]$ , diffusion steps  $T$

**Output:** Trained conditional diffusion model (CDM)  $\Phi$

- 1: Initialize  $\Phi$  (3D U-Net), optimizer, and EMA records for each sequence  $m$ :  
 $\text{EMA}_m^{\{\mathcal{H}\}}, \text{EMA}_m^{\{\mu\}}, \text{EMA}_m^{\{\sigma\}}$
  - 2: **for** epoch = 1 to  $N_1$  **do**
  - 3:   Compute normalized gradient map  $G_i = \text{Norm}(\nabla(I_i))$  for each  $I_i$
  - 4:   Sample random timestep  $t \sim \{1, \dots, T\}$
  - 5:   Sample noise  $\epsilon \sim \mathcal{N}(0, I)$
  - 6:   Generate noisy image  $I_i^t = \sqrt{\alpha_t}I_i + \sqrt{1 - \alpha_t}\epsilon$
  - 7:   Concatenate  $I_i^t$  and  $G_i$  as input, embed  $m_i$  as class condition
  - 8:   Predict noise:  $\hat{\epsilon} = \Phi(I_i^t, t, G_i, m_i)$
  - 9:   **Noise loss:**  $\mathcal{L}_N = \|\epsilon - \hat{\epsilon}\|_2^2$
  - 10:   Estimate harmonized MRI:  $I_i' = \frac{1}{\sqrt{\alpha_t}}(I_i^t - \sqrt{1 - \alpha_t}\hat{\epsilon})$
  - 11:   Compute soft-histogram  $\mathcal{H}(I_i')$ , mean  $\mu(I_i')$ , std  $\sigma(I_i')$  over  $I_i'$
  - 12:   **EMA update for sequence  $m_i$ :**
  - 13:      $\text{EMA}_{m_i}^{\{\mathcal{H}\}} \leftarrow \gamma \cdot \text{EMA}_{m_i}^{\{\mathcal{H}\}} + (1 - \gamma)\mathcal{H}(I_i')$
  - 14:      $\text{EMA}_{m_i}^{\{\mu\}} \leftarrow \gamma \cdot \text{EMA}_{m_i}^{\{\mu\}} + (1 - \gamma)\mu(I_i')$
  - 15:      $\text{EMA}_{m_i}^{\{\sigma\}} \leftarrow \gamma \cdot \text{EMA}_{m_i}^{\{\sigma\}} + (1 - \gamma)\sigma(I_i')$
  - 16:   **EMA loss:**  $\mathcal{L}_E = WD(\text{EMA}_{m_i}^{\{\mathcal{H}\}}, \mathcal{H}(I_i')) + \|\text{EMA}_{m_i}^{\{\mu\}} - \mu(I_i')\|_2^2 + \|\text{EMA}_{m_i}^{\{\sigma\}} - \sigma(I_i')\|_2^2$
  - 17:   Compute normalized gradient map  $G(I_i')$
  - 18:   **Gradient loss:**  $\mathcal{L}_G = \|G_i - G(I_i')\|_2^2$
  - 19:   **Total loss:**  $\mathcal{L}_1 = \mathcal{L}_N + \mathcal{L}_G + \mathcal{L}_E$
  - 20:   Backpropagate  $\mathcal{L}_1$  and update  $\Phi$  parameters
  - 21: **end for**
  - 22: **return** trained CDM model  $\Phi$
-

**Algorithm 2** Target-Specific Fine Harmonization**Input:** Source MRIs  $\{I_X\}$ , Target MRIs  $\{I_T\}$  (unpaired, containing same set of sequences)**Parameter:** Pre-trained CDM  $\Phi$  from Stage 1, BiomedCLIP encoder  $\Psi$ , DDIM steps  $T_f, T_r$ **Output:** Fine-tuned CDM  $\Phi^*$  for target-specific harmonization

```

1: # Pre-processing (performed once):
2:   for each MRI  $I$  in  $\{I_X\} \cup \{I_T\}$  do
3:     Global Harmonization: Use pre-trained  $\Phi$  to harmonize  $I$  to  $I'$  via DDIM inference (FDP+RDP)
4:   end for
5:   for each globally harmonized image  $I'$  do
6:      $I^0 = I'$ 
7:     for  $t = 0 : T_f - 1$  do
8:        $I^{t+1} = \sqrt{\bar{\alpha}_{t+1}}I' + \sqrt{1 - \bar{\alpha}_{t+1}}\epsilon_\theta(I^t, t, G, m)$ 
9:     end for
10:    Save  $I^{T_f}$  as the noisy image for fine-tuning
11:  end for
12: #Target-sepcific Fine-tuning:
13: for epoch = 1 to  $N_2$  do
14:   for each source noisy image  $I_X^{T_f}$ , Globally harmonized source  $I'_X$ , Globally harmonized target  $I'_T$ , and original target  $I_T$  do
15:     # DDIM Reverse Diffusion:
16:      $I_{X \rightarrow T}^{T_r} = I_X^{T_f}$ 
17:     for  $t = T_r : 1$  do
18:        $I_{X \rightarrow T}^{t-1} = \sqrt{\bar{\alpha}_{t-1}}I'_X$ 
19:          $+ \sqrt{1 - \bar{\alpha}_{t-1}}\epsilon_\theta(I_{X \rightarrow T}^t, t, G_X, m_X)$ 
20:     end for
21:     Obtain translated MRI  $I_{X \rightarrow T} = I_{X \rightarrow T}^0$ 
22:     # CLIP-Style Embedding:
23:      $S_T = \Psi(I_T) - \Psi(I'_T)$  {Target style embedding}
24:      $S_{X \rightarrow T} = \Psi(I_{X \rightarrow T}) - \Psi(I'_X)$  {Translated style embedding}
25:     # Disentangled CLIP-Style Loss:
26:      $\mathcal{L}_S = \|S_T - S_{X \rightarrow T}\|_1 + \left(1 - \frac{S_T \cdot S_{X \rightarrow T}}{\|S_T\| \|S_{X \rightarrow T}\|}\right)$ 
27:     # Style Reconstruction Loss:
28:      $I_T^{T_f}$  from  $I'_T$  via DDIM FDP,  $I_{T \rightarrow T}^0$  via DDIM RDP as above
29:      $S_{T \rightarrow T} = \Psi(I_{T \rightarrow T}) - \Psi(I'_T)$ 
30:      $\mathcal{L}_R = \|S_T - S_{T \rightarrow T}\|_1$ 
31:     Total Loss:  $\mathcal{L}_2 = \mathcal{L}_S + \mathcal{L}_R$ 
32:     Backpropagate and update  $\Phi$  parameters using  $\mathcal{L}_2$ 
33:   end for
34: return fine-tuned model  $\Phi^*$ 

```

## 2 ADDITIONAL RESULTS

### 2.0.1 ADDITIONAL VISUALIZATION ON SRPBS

We present additional results on T1 MRIs from the SRPBS test set, showing visualizations of two subjects across 11 sites and difference maps between each method’s harmonized MRIs and the ground truth target (site COI). The Baseline denotes the unharmonized raw MRIs.

768  
769  
750  
759  
762  
763  
762  
765  
764  
763  
766  
769  
780  
789  
720

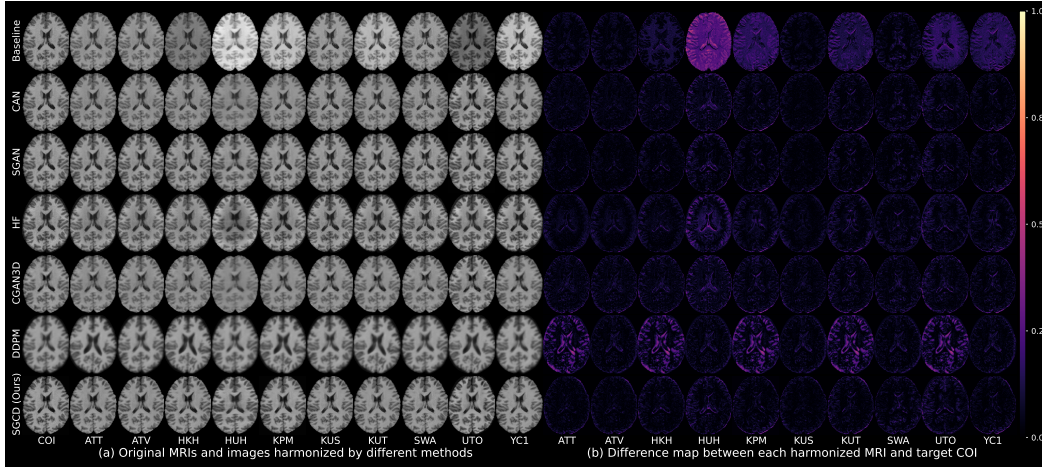


Figure 1: Axial view sample visualization results (a) and difference maps (b) of Subject 8 across 11 sites in SRPBS.

723  
722  
723  
726  
723  
726  
729  
730  
739  
780  
783

782  
785  
786  
783  
786  
789  
780  
789  
700  
703  
702  
705  
706  
703  
706

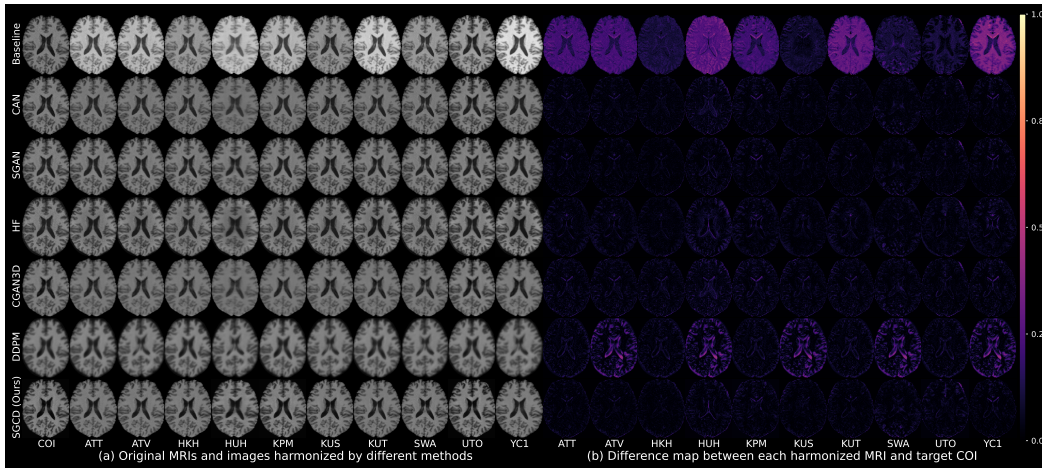


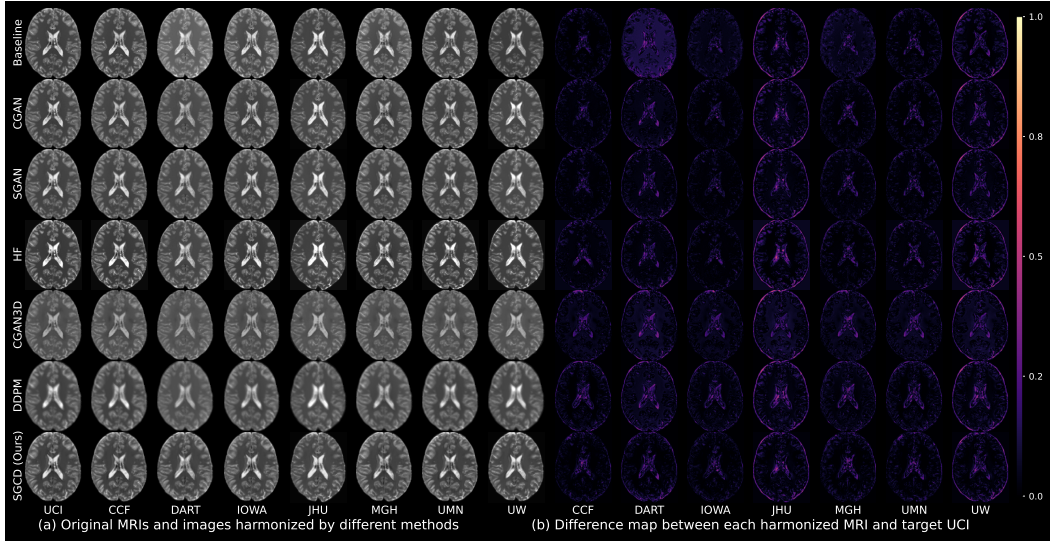
Figure 2: Axial view sample visualization results (a) and difference maps (b) of Subject 9 across 11 sites in SRPBS.

790  
799  
800  
803  
802  
805  
806  
803  
806  
809  
800  
809

### 2.0.2 ADDITIONAL VISUALIZATION ON DWI-THP

We present additional results on T2 MRIs from the DWI-THP test set, showing visualizations of two subjects across 8 sites and difference maps between each method’s harmonized MRIs and the ground truth target (site UCI). The Baseline denotes the unharmonized raw MRIs.

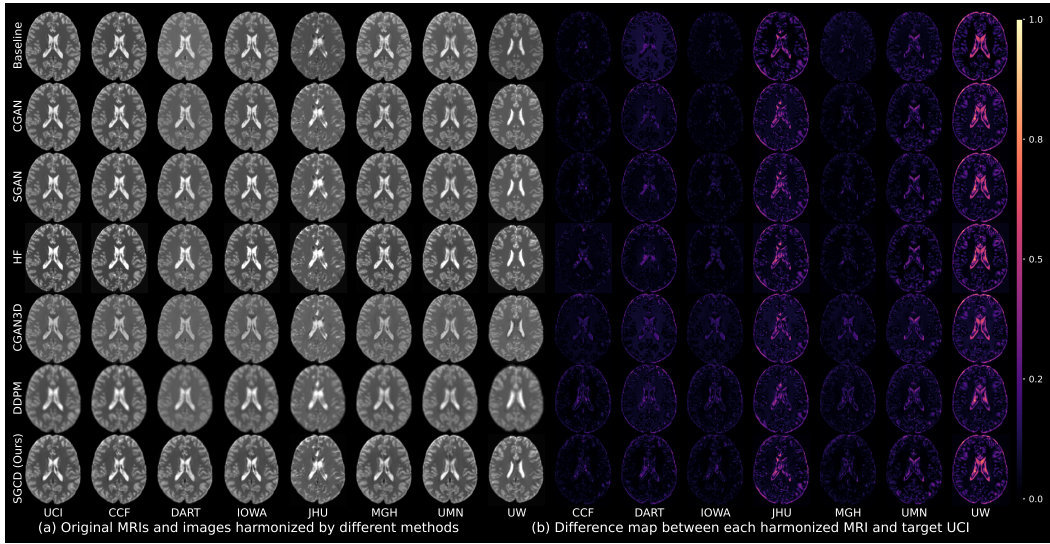
860  
863  
862  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878



879  
880  
881  
882

Figure 3: Axial view sample visualization results (a) and difference maps (b) of Subject THP0004 across 8 sites in DWI-THP.

883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908  
909  
910



911  
912  
913  
914  
915  
916  
917  
918  
919  
920  
921  
922  
923  
924  
925  
926  
927  
928  
929  
930  
931  
932  
933  
934  
935

Figure 4: Axial view sample visualization results (a) and difference maps (b) of Subject THP0005 across 8 sites in DWI-THP.