



# Harmonizing MR Images Across 100+ Scanners: Multi-site Validation with Traveling Subjects and Real-world Protocols

**Savannah P. Hays**<sup>1</sup> 


SHAYS6@JHU.EDU

<sup>1</sup> *Image Analysis and Communications Laboratory, Department of Electrical and Computer Engineering, Johns Hopkins University, USA*

**Lianrui Zuo**<sup>2</sup> 

LIANRUI.ZUO@VANDERBILT.EDU

<sup>2</sup> *Department of Electrical and Computer Engineering, Vanderbilt University, USA*

**Muhammad Faizyab Ali Chaudhary**<sup>1</sup> 

MCHAUD25@JHU.EDU

**Kathleen M. Bartz**<sup>1</sup> 


KBARTZ2@JHU.EDU

**Samuel W. Remedios**<sup>1</sup> 

SREMEDI1@JHU.EDU

**Jinwei Zhang**<sup>1</sup> 

JWZHANG@JHU.EDU

**Jiachen Zhuo**<sup>3</sup> 

JZHUO@SOM.UMARYLAND.EDU

<sup>3</sup> *Department of Diagnostic Radiology and Nuclear Medicine, University of Maryland School of Medicine, USA*

**Murat Bilgel**<sup>4</sup> 


MURAT.BILGEL@NIH.GOV

<sup>4</sup> *Laboratory of Behavioral Neuroscience, National Institute on Aging, USA*

**Shiv Saidha**<sup>5</sup> 

SSAIDHA2@JHMI.EDU

<sup>5</sup> *Department of Neurology, Johns Hopkins School of Medicine, USA*

**Ellen M. Mowry**<sup>5</sup> 


EMOWRY1@JHMI.EDU

**Scott D. Newsome**<sup>5</sup> 


SNEWSOM2@JHMI.EDU

**Jerry L. Prince**<sup>1</sup> 

PRINCE@JHU.EDU

**Blake E. Dewey**<sup>5</sup> 

BLAKE.DEWEY@JHU.EDU

**Aaron Carass**<sup>1</sup> 

AARON\_CARASS@JHU.EDU

**Editors:** Under Review for MIDL 2026

## Abstract

Reliable harmonization of heterogeneous magnetic resonance (MR) image datasets, especially those acquired in pragmatic clinical trials, is critical to advance multi-center neuroimaging studies and translational machine learning in healthcare. We present an enhanced and rigorously validated version of the HACA3 harmonization algorithm, which we refer to as HACA3<sup>+</sup>, incorporating key methodological enhancements: (1) an improved artifact encoder to better isolate and mitigate image artifacts, (2) background and foreground-sensitive attention mechanisms to increase harmonization specificity, and (3) extensive training using data spanning 100+ scanners from 64 independent sites, providing a broader diversity of scanners than other harmonization methods. Our study focuses on four commonly acquired MR image contrasts (T1-weighted, T2-weighted, proton density, & fluid-attenuated inversion recovery), reflecting realistic clinical protocols. We perform inter-site harmonization experiments using traveling subjects to assess the generalization and robustness of the harmonization model. We compare the results of the publicly available version of HACA3 and our implementation, HACA3<sup>+</sup>. Downstream relevance is further established through whole brain segmentation and image imputation. Pre-trained weights and code for HACA3<sup>+</sup> are made publicly available at <https://github.com/UponAcceptance>.

**Keywords:** MRI, Image Harmonization, Image Synthesis

## 1. Introduction

Variability in magnetic resonance (MR) image acquisition due to different scanner parameters and manufacturers leads to inconsistent image contrasts and intensities across datasets, which can severely confound biomarker development and analysis. Multi-site neuroimaging studies and the adoption of deep learning in clinical MR image processing workflows depend critically on the ability to harmonize MR imaging data to overcome this variability. However, despite advances in algorithm development, few harmonization methods undergo rigorous, large-scale validation with multiple out-of-domain datasets, which limits the algorithms translational impact while also undermining reliability.

A broad spectrum of harmonization methods for neuroimaging have been proposed, which can be broadly grouped into statistical and image-based approaches. Statistical harmonization methods, such as neuroComBat (Fortin et al., 2018) and its extensions (Beer et al., 2020; Horng et al., 2022), operate on derived measurements (e.g., cortical thickness and diffusion metrics) to adjust for scanner or site effects while preserving biological variability and have been widely applied (Fortin et al., 2018; Horng et al., 2022). In contrast, image-based harmonization aims to correct scanner-dependent variation directly at the voxel level using generative models. This enables harmonization before downstream analyzes and allows a single corrected image to support multiple tasks (Dewey et al., 2019, 2020; Liu et al., 2021). Although image-based methods offer greater flexibility, they often require substantially larger training data and may be sensitive to domain shifts.

Earlier image-based harmonization methods, such as DeepHarmony (Dewey et al., 2019), used supervised training on paired images, requiring the same subject to be scanned across different scanners and possibly sites. However, reliance on paired data limits their broader applicability as (inter-site) traveling subject data is almost never available. To relax the assumption of paired data, unpaired image-to-image translation using cycle-consistency was proposed (Modanwal et al., 2020). Although these models use distribution matching for synthesizing samples across domains, they are limited in two pertinent ways: their inability to preserve anatomical structure across domains and limited scalability across sites and contrasts. To address this, a major shift was observed towards anatomy-contrast disentanglement approaches (Liu et al., 2017; Huang et al., 2018).

Recent methods such as HACA3 (Zuo et al., 2023) and MURD (Liu and Yap, 2024) have introduced attention mechanisms, flexible contrast handling, and unified multi-domain architectures. HACA3 relaxes assumptions about contrast similarity and incorporates artifact-aware modules to robustly harmonize images across diverse contrasts and protocols, even in the presence of motion or noise. Meanwhile, MURD addresses scalability by learning a single encoder-decoder framework that generalizes across sites without direct supervision, enabling realistic harmonization while preserving subject-specific features.

Recent work by Lu et al. (2025) demonstrated that harmonization performance can vary substantially between scanners and sequences, motivating the critical importance of robust, large-scale validation of MR image harmonization methods to ensure their effectiveness and reliability for subsequent analyzes. Although all previously reported harmonization methods have shown a benefit over using unharmonized data, HACA3 (Zuo et al., 2023) showed the most consistent results in comparison to DeepHarmony (Dewey et al., 2019) and neuroCombat (Fortin et al., 2018) in Lu et al. (2025). Supervised methods that use paired

training data may outperform HACA3, but when using real-world, large-scale datasets, HACA3’s ability to handle unpaired data and missing contrasts gives it a practical edge.

In this work, we focus on enhancements to HACA3 that address three of its limitations. First, HACA3’s artifact encoder was trained using a contrastive learning framework to detect images with a high level of artifacts. This encoder is particularly important when there are multiple input images. It allows HACA3 to focus on using the images with a low level of artifacts for the synthesis task. Hays et al. (2025b) recently used the margin loss in a contrastive learning framework to score an MR image’s artifact level. We expand on that work to a 2D score to replace the HACA3 artifact encoder and train this new artifact encoder on a larger range of simulated artifacts. The margin loss within the encoder allows for the stratification by different artifact levels, which helps improve its sensitivity.

Our second enhancement to HACA3 introduces an attention mechanism that operates on a 2D slice-wise level, in all three cardinal orientations, covering the volume. HACA3’s current scalar attention vector is constant over an entire 2D slice regardless of the spatial location within the slice. When MR images are acquired with a full field-of-view (FOV), this does not raise concerns. However, when images are acquired with a limited FOV, such as the axial plane missing a portion of the superior skull, HACA3 will fail to recover the missing region even if it was present in other input source images. Recent preliminary work (Hays et al., 2025a) focused on this modification of HACA3’s attention mechanism to be background and foreground aware, allowing for variation across a 2D slice.

Our last enhancement to HACA3 is the inclusion of considerably more data in the training of the contrast encoder. HACA3 was originally trained on MR images from 21 different sites consisting of only 21 scanners, with each site contributing ten subjects. Our contribution is to train across 64 different sites covering 132 scanners, with a total of 996 subjects; this represents a six-fold increase in the number of scanners and a quadrupling in the number of subjects. We refer to our new version of HACA3, with the outlined three enhancements, as HACA3<sup>+</sup>. This paper provides a comprehensive, clinically-relevant validation of the HACA3<sup>+</sup> algorithm for MR image harmonization. Throughout the paper, we compare HACA3<sup>+</sup> with the publicly available version of HACA3. The pre-trained weights and code for HACA3<sup>+</sup> are made publicly available at <https://github.com/UponAcceptance>.

## 2. Methods

### 2.1. Technical Contributions of HACA3<sup>+</sup>

**Enhanced Artifact Encoder** The first technical contribution is the modification of the artifact encoder. We trained the artifact encoder using 297 structural MR volumes from the TRaditional vs. Early Aggressive Therapy for Multiple Sclerosis (TREAT-MS) pragmatic, clinical trial (NCT03500328) (Mowry et al., 2025). These scans were acquired from seven different imaging sites and included four structural MR image contrasts: T<sub>1</sub>-weighted (T<sub>1</sub>-w), T<sub>2</sub>-weighted (T<sub>2</sub>-w), fluid-attenuated inversion recovery (FLAIR), and proton density (PD) images. Only high-quality images were included in the training dataset. Prior to training, the images were N4 bias field corrected (Tustison et al., 2010) and 2D acquisitions were super-resolved (Remedios et al., 2023). Similarly to HACA3, we simulated common MR image artifacts using the TorchIO library (Pérez-García et al., 2021), including random noise, random ghosting, random bias field, and random anisotropy. Unlike HACA3, we mapped

the artifact simulation parameters to a normalized score, where a low artifact level mapped to a low score and a high artifact level mapped to a higher score. We incorporated this score into training through the triplet loss. The triplet loss enforces a ranking such that the clean anchor image receives a lower severity score than the artifact-degraded negative image:

$$\mathcal{L}_{\text{triplet}} = \sum_{i=1}^N \max(0, S_i^{\text{anchor}} - S_i^{\text{positive}} + m) + \max(0, S_i^{\text{negative}} - S_i^{\text{anchor}} + m) \quad (1)$$

where  $m$  is a dynamic margin based on artifact severity.

**Enhanced Attention** The second technical contribution is the modification of the attention module to handle limited FOV source images. HACA3 uses the background mask of only the first input source image. This mask is used for all synthetic images independent of the harmonization target and attention for each source image. Our modification uses the union of the background mask from all of the source images. This modification only makes a difference if there are multiple source images. For each pixel within the union background region, the attention will be equally distributed across source images. For each pixel within the union foreground region, the attention will be distributed across source images according to the similarity between the source image and the target image as computed in HACA3. For other pixels not in the union regions, the background source pixels will be forced to 0 and the attention between the foreground source pixels, as determined by the attention module, will be normalized to sum to 1.

Our approach for limited FOV imputation relies on multiple contrast source images to accurately impute missing regions. Imputation only occurs when a region is visible in at least one source image. We chose not to impute regions without any information across the source images. To validate this modification, we simulated limited FOV images using MR images (N=61) from 7 sites included in Table 1. Limited FOV images were simulated by cropping regions from full FOV images. We simulated two types of limited FOV acquisitions: anterior degradation and left/right degradation. We tested HACA3<sup>+</sup> and HACA3 using simulated degradations on T<sub>1</sub>-w, T<sub>2</sub>-w, and FLAIR images. Evaluation metrics include peak-signal-to-noise ratio (PSNR) and structural similarity index measure (SSIM) between the imputed, harmonized image and the acquired full FOV image.

### 3. Data: Training and Validation

#### 3.1. Training Dataset

The training datasets used for HACA3<sup>+</sup> are from 64 imaging sites and 132 scanners. These datasets are summarized in Table 1. The four MR image contrasts used in training are: T<sub>1</sub>-weighted (T<sub>1</sub>-w), T<sub>2</sub>-weighted (T<sub>2</sub>-w), fluid-attenuated inversion recovery (FLAIR), and proton density (PD). Not every site acquired all four contrasts, which is indicated in Table 1. All subjects included in the training have at least a T<sub>1</sub>-w image. In total, 996 subjects were used to train HACA3<sup>+</sup>. The open sites were from the following datasets: Open Access Series of Imaging Studies-3 (OASIS-3) (LaMontagne et al., 2019), Baltimore Longitudinal Study of Aging (BLSA) (Resnick et al., 2000), Human Connectome Project (HCP)<sup>1</sup>, and

1. <https://www.humanconnectome.org/study/hcp-young-adult>

Information eXtraction from Images (IXI)<sup>2</sup>. We included data from multiple imaging sites within the OASIS-3, BLSA, and IXI datasets. These are labeled as different sites in Table 1. In total, there were 11 open sites used in training. The private data included 46 sites from the TREAT-MS pragmatic clinical trial (Mowry et al., 2025). The remaining seven sites are from other private sources.

For preprocessing, all images were background removed followed by N4 bias field correction (Tustison et al., 2010). Super-resolution (Remedios et al., 2023) was applied to 2D acquisitions. Finally, 3D acquisitions and super-resolved 2D acquisitions were rigidly registered to the MNI152 atlas using ANTs (Avants et al., 2009). The 50 middle slices of each orientation were extracted from each MR image volume for training, resulting in 150 slices per volume.

Table 1: Training site characteristics and MR contrast availability for HACA3<sup>+</sup>. S1-S21 are the sites used to train the publicly available HACA3. Open refers to the public availability of the site data. The manufacturers (Manu.) are: G - GE, H - Hitachi, P - Philips, S- Siemens, T - Toshiba. The field (Field) strength column lists used field strengths in Telsa (T) at a site. The population (Pop.) column codes are: HC - healthy controls, MS - multiple sclerosis, TB - mild traumatic brain injury. For T1w, T2w, FLAIR, and PD, the number denotes the number of used scans for training HACA<sup>+</sup>.

Site ID	Open	Manu.	Field (T)	Pop.	T1w	T2w	FLAIR	PD
S1	✓	P	1.5	HC	10	10	10	10
S2	✓	P	3.0	HC	10	10	10	10
S3	✓	S	3.0	HC	10	10	10	10
S4	✓	S	3.0	HC	10	10	10	10
S5	✓	S	3.0	HC	10	10	10	10
S6	✓	S	1.5	HC	10	10	10	10
S7	✓	P	1.5	HC	10	10	10	10
S8	✓	P	3.0	HC	10	10	10	10
S9	✓	P	3.0	HC	10	10	10	10
S10	✓	P	3.0	HC	10	10	10	10
S11	✗	P	3.0	MS	10	10	10	10
S12	✗	P	3.0	MS	10	10	10	10
S13	✗	S	3.0	TB	10	10	10	0
S14	✗	S	3.0	HC	10	10	10	0
S15	✗	G	3.0	MS	10	10	10	10
S16	✗	G	3.0	MS	10	10	10	10
S17	✗	S	1.5	MS	10	0	10	0
S18	✗	S	3.0	MS	10	10	10	10
S19	✗	G,H,P,S,T	1.1,1.5,3.0	MS	125	125	125	65
S20	✗	G,P,S	1.5,3.0	MS	22	22	20	0

Continued on the next page.

2. <https://brain-development.org/ixi-dataset/>

Site ID	Open	Manu.	Field (T)	Pop.	T1w	T2w	FLAIR	PD
S21	✗	G,P,S,T	3.0	MS	26	26	26	1
S22	✗	G,H,S	1.1,1.5,3.0	MS	42	42	38	2
S23	✗	G,P,S	1.5,3.0	MS	26	26	26	2
S24	✗	G,P,S	1.5,3.0	MS	23	23	18	3
S25	✗	G,P,S,T	1.1,1.5,3.0	MS	14	14	14	0
S26	✗	G,P,S	1.5,3.0	MS	10	10	9	0
S27	✗	G,P,S	1.5,3.0	MS	5	5	5	2
S28	✗	G,P,S	1.5,3.0	MS	8	8	8	0
S29	✗	G,S	1.5,3.0	MS	13	13	13	0
S30	✗	G,H,P,S	1.1,1.5,3.0	MS	14	14	15	5
S31	✗	G,H,S	1.1,1.5,3.0	MS	19	19	19	2
S32	✗	G,P,S	1.5,3.0	MS	37	37	36	4
S33	✗	S	1.5	MS	3	3	3	0
S34	✗	G,H,S	1.1,1.5	MS	7	7	7	0
S35	✗	G,S	1.5,3.0	MS	11	11	11	4
S36	✗	G,S	1.5,3.0	MS	3	3	3	0
S37	✗	G,P,S,T	1.5,3.0	MS	21	21	20	0
S38	✗	G,P,S	1.5,3.0	MS	23	23	21	3
S39	✗	P,S	1.5,3.0	MS	3	3	3	0
S40	✗	G,S	1.5,3.0	MS	19	19	16	10
S41	✗	G,P,S	1.5,3.0	MS	20	20	20	1
S42	✗	S	1.5	MS	14	14	14	0
S43	✗	G,P,S	1.5,3.0	MS	29	29	28	0
S44	✗	G,P,S	1.5,3.0	MS	6	6	6	0
S45	✗	G,P,S	1.5,3.0	MS	21	21	18	4
S46	✗	G,P,S,T	1.5,3.0	MS	35	35	33	1
S47	✗	G,P,S,T	1.5,3.0	MS	33	33	28	1
S48	✗	G,P,S	1.5,3.0	MS	12	12	12	1
S49	✗	G,P,S	1.5,3.0	MS	20	20	20	1
S50	✗	G,H,P,S	1.5,3.0	MS	25	25	25	3
S51	✗	G,S	1.5,3.0	MS	5	5	4	0
S52	✗	G,H,S,T	1.1,1.5,3.0	MS	20	20	20	0
S53	✗	G,H,S,T	1.1,1.5,3.0	MS	8	8	8	0
S54	✗	G,S,T	1.5,3.0	MS	12	12	11	1
S55	✗	S	1.5	MS	3	3	3	3
S56	✗	G,S	1.5,3.0	MS	18	18	18	4
S57	✗	G,S,T	1.5,3.0	MS	24	24	24	1
S58	✗	G,S	1.5,3.0	MS	13	13	13	1
S59	✗	G,P	1.5,3.0	MS	5	5	5	1
S60	✗	G,P,S	1.5,3.0	MS	12	12	12	0
S61	✗	G,P,S,T	1.5,3.0	MS	21	21	21	3

Continued on the next page.

Site ID	Open	Manu.	Field (T)	Pop.	T1w	T2w	FLAIR	PD
S62	✗	G,P,S	1.5,3.0	MS	7	7	7	1
S63	✗	G,S	1.5,3.0	MS	5	5	4	0
S64	✗	T	1.5	MS	1	1	1	0

### 3.2. Validation Using Travel Subjects

A summary of the three datasets we use to validate HACA<sup>+</sup> are in Table 2.

Table 2: Summary of the Validation Datasets. We did not use data from repeated sessions for FTHP and MASiVar datasets, as not every subject had a repeated scan.

Dataset	# Subjects	# Sites	# Sessions
#1: TREAT-MS Traveling Subjects	15	9 <sup>†</sup>	126
#2: FTHP (Opfer et al., 2023)	1	116	116
#3: MASiVar (Cai et al., 2021)	5	3 <sup>*</sup>	19

<sup>†</sup>: Each subject only visited 3–5 sites.

<sup>\*</sup>: Images were acquired at three different institutions across 3–4 scanners.

**Dataset #1: TREAT-MS Traveling Subjects** For inter-site harmonization validation, we use a private traveling subjects dataset from the TREAT-MS pragmatic, clinical trial (NCT0350032) (Mowry et al., 2025). MR images were collected from 14 participants (ages 18–60), including five healthy controls with no known neurological conditions and nine people with multiple sclerosis (PwMS) with stable MS status. Scans followed a 3-month protocol: initial session at Johns Hopkins, 3–5 sessions at other eastern United States sites, and a final session at Johns Hopkins. Each session followed the same acquisition protocol, including a scan and re-scan procedure with T<sub>1</sub>-w, T<sub>2</sub>-w, FLAIR, and PD sequences. For reasons beyond our control, some subjects did not complete all scans. We used the initial scan at Johns Hopkins as the harmonization target for HACA3 and HACA3<sup>+</sup> for all harmonization tasks. For input into HACA3 and HACA3<sup>+</sup>, we use the acquired T<sub>1</sub>-w, T<sub>2</sub>-w, and FLAIR images. The PD image was withheld to demonstrate the imputation ability of HACA3<sup>+</sup>. We compute the evaluation metrics of peak-signal-to-noise ratio (PSNR) and structural similarity index measure (SSIM) for each image contrast between the harmonized image and the acquired image at the Johns Hopkins site.

**Dataset #2: Frequently Traveling Human Phantom** The Frequently Traveling Human Phantom (FTHP) (Opfer et al., 2023)<sup>3</sup> provides a single subject, multi-scanner validation dataset. It contains T<sub>1</sub>-w MR images of a single healthy male volunteer, acquired across 116 different scanners. To ensure one scan per scanner, we excluded repeat scans acquired on the same scanner, focusing on single subject, multi-scanner validation. For

3. <https://www.nitrc.org/projects/fthp>



harmonization, we used the  $T_1$ -w image from the initial scan as the harmonization target for both HACA3 and HACA3<sup>+</sup>. Whole brain segmentation was performed on both the unharmonized and harmonized images using SLANT (Huo et al., 2019), and the resulting segmentation labels were aggregated into nine brain regions for analysis. The Dice similarity coefficient (DSC) was calculated for each region, using the segmentation result from the initial scan as the reference. Region-wise brain volumes were computed as the total number of voxels in each label. To quantify inter-scanner variability, the coefficient of variation (CV) was computed for both DSC and volume for each brain region. Lower CV values indicate reduced variability and, thus, better harmonization consistency across scanners. Since FTHP contains only one subject, the CV values reported represent dispersion across different scanners for that single subject.

**Dataset #3: MASiVar** The Multisite, Multiscanner, and Multisubject Acquisitions for Studying Variability in Diffusion Weighted Magnetic Resonance Imaging (MASiVar) (Cai et al., 2021)<sup>4</sup> study provides a multi-subject, multi-scanner validation dataset. We utilized  $T_1$ -w MR images from Cohort II, comprising of five adult subjects, each scanned on three to four different scanners across three institutions. To focus on inter-scanner and inter-subject variability, repeat scans on the same scanner were excluded. For harmonization, the  $T_1$ -w image from each subject’s initial scan at site #1 served as the harmonization target for both HACA3 and HACA3<sup>+</sup>. Whole brain segmentation was performed using SLANT (Huo et al., 2019) for both unharmonized and harmonized images. For each subject, the CV for DSC and regional brain volume was calculated across their multi-site scans, with results reported as the mean and standard deviation over subjects. This approach allows assessment of harmonization consistency within subjects.

## 4. Experiments and Results

### 4.1. Region Imputation Using Enhanced Attention

Figure 1 shows qualitative examples for the limited FOV simulations using each image contrast from seven in-domain sites. The ground truth is the harmonized full FOV image. Quantitative results for the two types of limited FOV simulations using each image contrast are shown in Fig. 2. The PSNR and SSIM when using HACA3<sup>+</sup> are significantly higher for all image contrasts than when using HACA3. Statistical significance was computed using a paired Wilcoxon test with Bonferroni correction ( $p$ -value < 0.0001). Each image contrast was tested separately. For example, when the  $T_1$ -w image with simulated limited FOV was input, the other input source images included the full FOV. Qualitative examples of clinically acquired limited FOV images are shown in Fig. 3. These images do not have corresponding full FOV ground truth images and are only shown for qualitative purposes.

### 4.2. Inter-site Harmonization Using TREAT-MS Traveling Subjects

Figure 4 illustrates a representative healthy traveling subject from the TREAT-MS Traveling Subjects dataset. The input source images were obtained from a single non-Johns Hopkins imaging site and harmonized to the images acquired at the Johns Hopkins site. We quantitatively evaluate inter-site harmonization performance across the entire TREAT-MS

4. <https://openneuro.org/datasets/ds003416/versions/2.0.2>



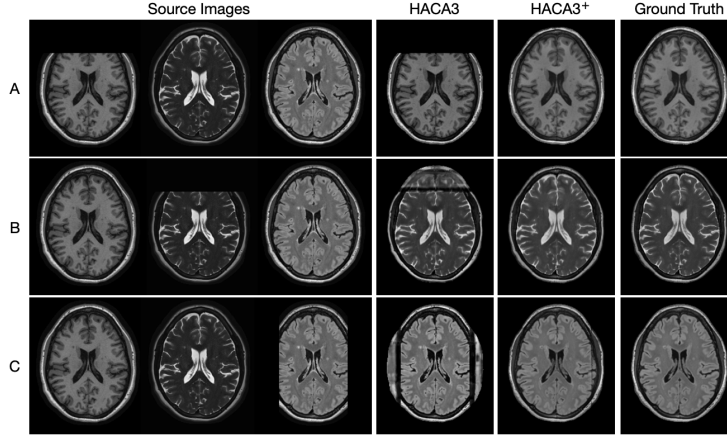


Figure 1: Qualitative results from the (A)  $T_1$ -w degraded anterior simulation, (B)  $T_2$ -w degraded anterior simulation, and (C) FLAIR degraded left/right simulation. Each row shows the different source images input to the model. We compare between HACA3 and HACA3<sup>+</sup>. HACA3 uses the background mask from the first source image, which, in this scenario, is the  $T_1$ -w image. As a result, HACA3 does not attempt to impute in (A), but it does attempt to impute in (B) and (C).

Traveling Subjects dataset, reporting both PSNR and SSIM in Fig 5. Both HACA3 and HACA3<sup>+</sup> demonstrate comparable performance on this in-domain, full FOV dataset, with no statistically significant differences observed between the methods. This outcome indicates that the methodological enhancements of HACA3<sup>+</sup> do not compromise harmonization quality in standard scenarios, thus preserving in-domain performance.

#### 4.3. Inter-site Harmonization Using FTHP Dataset

Figure 6 presents the DSC variability across brain regions in the FTHP dataset. Harmonization leads to statistically significant improvement in DSC for eight out of nine regions compared to the unharmonized result. Significant differences between results were calculated using a paired Wilcoxon test with Benjamini–Hochberg correction. There were no statistically significant differences in DSC between the results using HACA3 and HACA3<sup>+</sup>. Table 3 reports the CV for both DSC and volume, calculated across scans at different sites from the same subject. Since the FTHP dataset consists of a single subject scanned across multiple sites, only one CV value is provided per region. Harmonization consistently reduced CV for both DSC and volume in all regions. HACA3<sup>+</sup> modestly outperformed HACA3 in terms of DSC CV and in most regions for volume CV; however, these differences were not statistically significant.

#### 4.4. Inter-site Harmonization Using MASiVar Dataset

Table 4 reports the CV for DSC and volume across brain regions in the MASiVar dataset. For each region, the CV is calculated across scans for each subject, and the values in the table

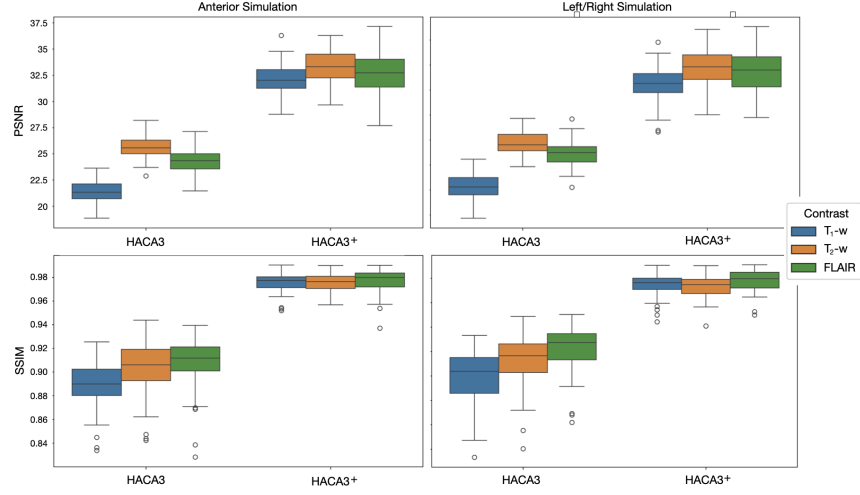


Figure 2: The PSNR and SSIM of the harmonized limited FOV images compared with the acquired full FOV image for each image contrast. We used two limited FOV simulations: anterior and left/right. Each contrast was tested separately. For example, when the T<sub>1</sub>-w image had limited FOV, the T<sub>2</sub>-w and FLAIR were full FOV. For each image contrast, HACA3<sup>+</sup> significantly ( $p$ -value < 0.0001) outperformed HACA3.

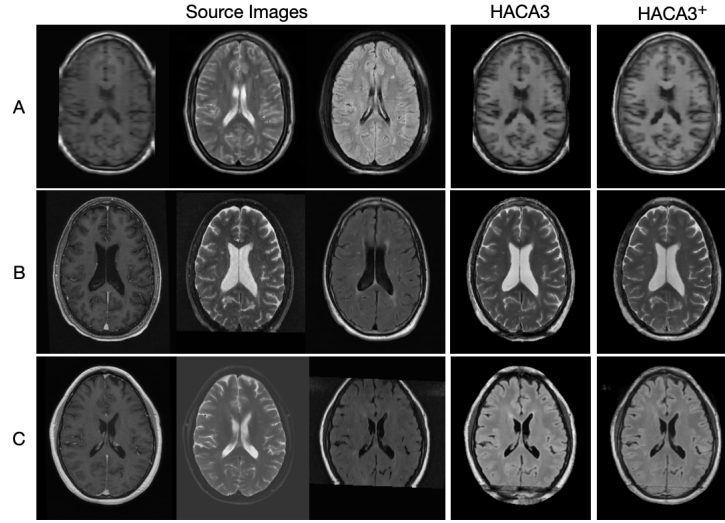


Figure 3: Qualitative results on real, clinically acquired MR images. Each row corresponds to a different person with MS. (A) shows the harmonized T<sub>1</sub>-w result using a limited FOV T<sub>1</sub>-w input image. (B) shows the harmonized T<sub>2</sub>-w result using a limited FOV T<sub>2</sub>-w input image. (C) shows the harmonized FLAIR result using a limited FOV FLAIR input image.

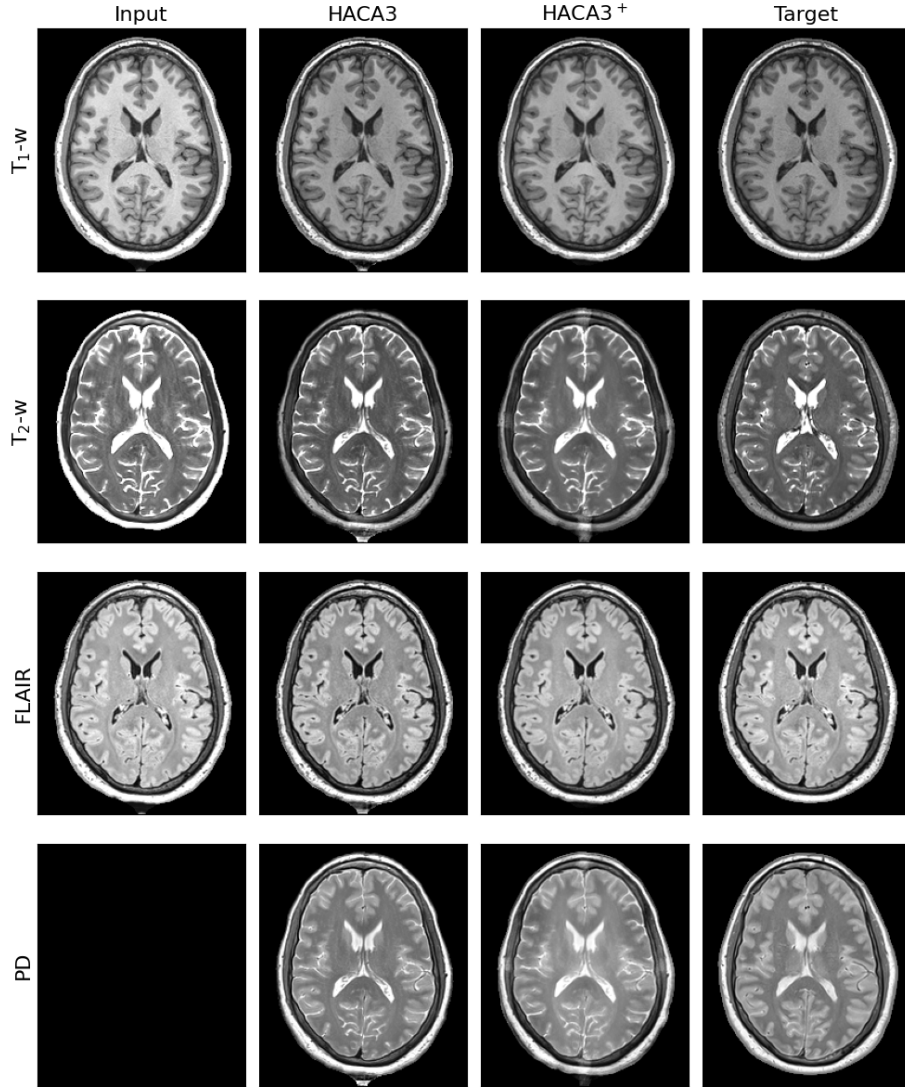


Figure 4: Inter-site harmonization results using HACA3 and HACA3<sup>+</sup>. The input source images were acquired at a single non-Johns Hopkins site. The input PD was not included to demonstrate the imputation ability. The target images were acquired at the Johns Hopkins site.

represent the mean and standard deviation across the five subjects in Cohort II. Consistent with results in the FTHP dataset, harmonization reduced CV for both DSC and volume in all regions. HACA3 and HACA3<sup>+</sup> performed similarly, with only minimal differences between them. There were no statistically significant results in this dataset most likely due to a smaller number of subjects.

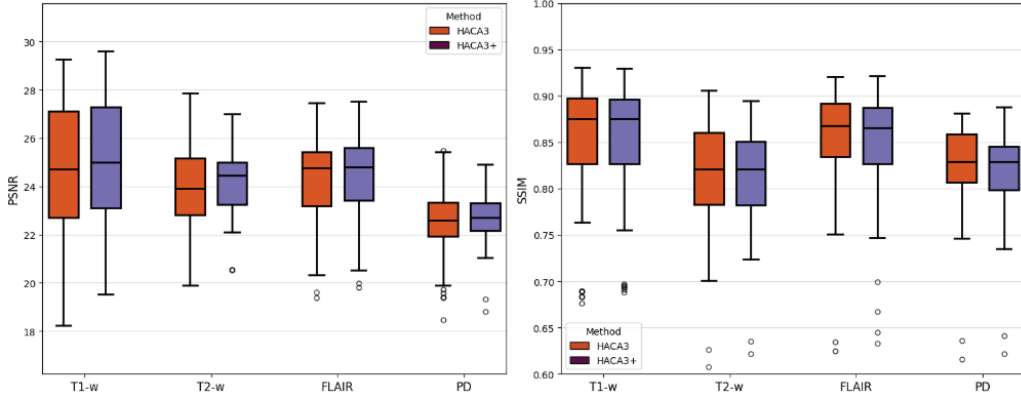


Figure 5: PSNR and SSIM of each harmonized MR image contrast over the TREAT-MS Traveling Subjects dataset. Acquired images were harmonized to the Johns Hopkins site. PSNR and SSIM were calculated using the acquired Johns Hopkins image as the reference image.

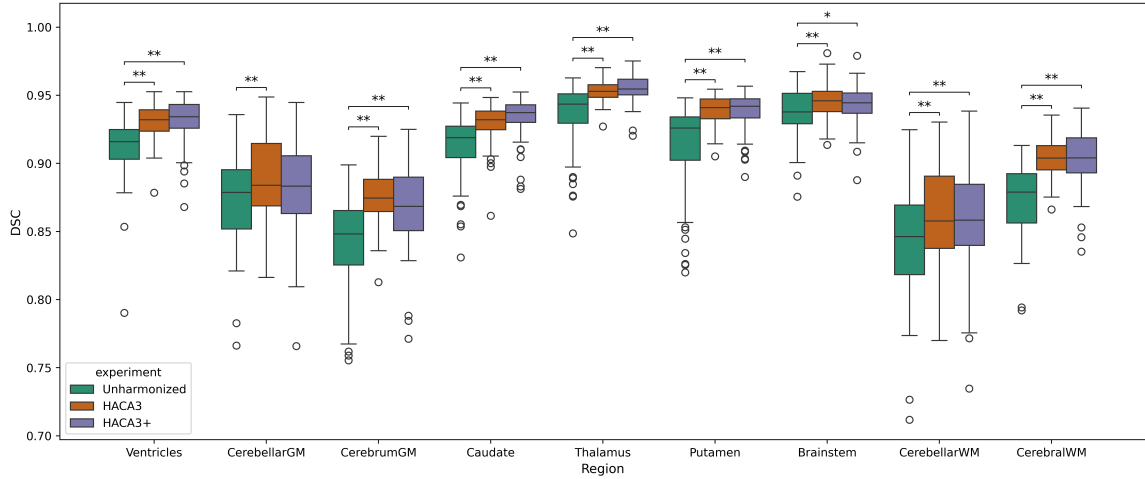


Figure 6: The DSC for each brain region computed using the segmentation result on the unharmonized, harmonized using HACA3, and harmonized using HACA3<sup>+</sup> T<sub>1</sub>-w images. Significant differences between results are indicated using a paired Wilcoxon test with Benjamini–Hochberg correction (symbols indicate significance level (\*\*:  $p$ -value < 0.01, \*:  $p$ -value < 0.05)).

## 5. Discussion and Conclusion

In this validation paper, we performed comprehensive validation of the publicly available HACA3 MR harmonization model and an enhanced variant, HACA3<sup>+</sup>. HACA3<sup>+</sup> incorporates three specific enhancements to address the limitations of HACA3. The first

Table 3: The CV for DSC and volume (Vol) of each brain region computed using the segmentation result on the unharmonized, harmonized using HACA3, and harmonized using HACA3<sup>+</sup> T<sub>1</sub>-w images from the FTHP dataset. The lowest DSC and Vol CV for each region are marked in **bold**.

Region	Unharmonized		HACA3		HACA3 <sup>+</sup>	
	DSC <sub>CV</sub> ↓	Vol <sub>CV</sub> ↓	DSC <sub>CV</sub> ↓	Vol <sub>CV</sub> ↓	DSC <sub>CV</sub> ↓	Vol <sub>CV</sub> ↓
Ventricles	0.0243	0.0408	0.0166	0.0292	<b>0.0135</b>	<b>0.0262</b>
Cerebellar GM	0.0382	0.0425	0.0379	<b>0.0269</b>	<b>0.0346</b>	0.0289
Cerebrum GM	0.0387	0.0498	0.0331	<b>0.0161</b>	<b>0.0232</b>	0.0217
Caudate	0.0234	0.0518	0.0153	0.0196	<b>0.0147</b>	<b>0.0176</b>
Thalamus	0.0231	0.0388	0.0097	0.0179	<b>0.0075</b>	<b>0.0169</b>
Putamen	0.0346	0.0452	0.0148	0.0192	<b>0.0112</b>	<b>0.0191</b>
Brainstem	0.0193	0.0402	0.0150	<b>0.0120</b>	<b>0.0131</b>	0.0217
Cerebellar WM	0.0461	0.0901	0.0444	<b>0.0209</b>	<b>0.0438</b>	0.0244
Cerebral WM	0.0289	0.0623	0.0221	0.0175	<b>0.0164</b>	<b>0.0147</b>

Table 4: The CV mean and standard deviation for DSC and volume (Vol) of each brain region computed using the segmentation result on the unharmonized, harmonized using HACA3, and harmonized using HACA3<sup>+</sup> T<sub>1</sub>-w images for the 5 subjects in Cohort II of the MASiVar dataset. The lowest DSC and Vol CV for each region are marked in **bold**.

Region	Unharmonized		HACA3		HACA3 <sup>+</sup>	
	DSC <sub>CV</sub> ↓	Vol <sub>CV</sub> ↓	DSC <sub>CV</sub> ↓	Vol <sub>CV</sub> ↓	DSC <sub>CV</sub> ↓	Vol <sub>CV</sub> ↓
Ventricles	0.0357±0.0096	0.0443±0.0122	<b>0.0140±0.0092</b>	<b>0.0267±0.0148</b>	0.0199±0.0167	0.0359±0.0175
Cerebellar GM	0.0540±0.0159	0.0343±0.0089	<b>0.0274±0.0128</b>	0.0194±0.0046	0.0284±0.0143	<b>0.0190±0.0077</b>
Cerebrum GM	0.0462±0.0080	0.0683±0.0114	0.0248±0.0108	<b>0.0147±0.0068</b>	<b>0.0231±0.0096</b>	0.0202±0.0053
Caudate	0.0255±0.0030	0.0592±0.0055	0.0232±0.0065	0.0359±0.0084	<b>0.0189±0.0140</b>	<b>0.0325±0.0095</b>
Thalamus	0.0233±0.0070	0.0420±0.0092	<b>0.0056±0.0021</b>	<b>0.0164±0.0021</b>	0.0096±0.0024	0.0243±0.0022
Putamen	0.0253±0.0119	0.0391±0.0012	<b>0.0139±0.0074</b>	0.0211±0.0112	0.0159±0.0124	<b>0.0197±0.0059</b>
Brainstem	0.0316±0.0048	0.0301±0.0103	<b>0.0143±0.0049</b>	<b>0.0144±0.0051</b>	0.0180±0.0059	0.0229±0.0028
Cerebellar WM	0.0827±0.0284	0.0704±0.0252	<b>0.0313±0.0182</b>	<b>0.0222±0.0080</b>	0.0399±0.0202	0.0346±0.0041
Cerebral WM	0.0388±0.0080	0.0585±0.0203	<b>0.0168±0.0083</b>	<b>0.0255±0.0042</b>	0.0170±0.0073	0.0351±0.0103

enhancement involved retraining the artifact encoder using more data and varying levels of simulated artifacts. The level of simulated artifacts was explicitly used as the margin in the triplet loss to develop a continuous space of low-to-high levels of artifacts. The second enhancement introduced a spatially-aware attention mechanism to distinguish foreground and background, thereby improving local adaptation within each 2D slice of the source images. The third enhancement included training on a substantially larger amount of data, encompassing images from 132 scanners—to our knowledge, a scale unmatched by previous structural MR harmonization efforts.

On our private traveling subjects dataset (14 subjects within 9 sites), we found that both HACA3 and HACA3<sup>+</sup> exhibited strong, statistically indistinguishable performance in terms of PSNR and SSIM, showing that both models generalize well to in-domain, multi-site data.

On the FTHP dataset involving a single subject who traveled to 116 different sites, harmonization with either model led to significant reductions in inter-scanner variability, measured as the CV for both DSC and regional brain volume using SLANT labels. Notably, HACA3<sup>+</sup> achieved the lowest CV across most brain regions and minor DSC improvements in the ventricles and deep gray matter structures such as the caudate and thalamus. Although direct comparisons between HACA3 and HACA3<sup>+</sup> were not statistically significant. We note that ground truth segmentation were not available, so comparisons relied on the reference segmentation from the initial scan.

Analysis on the MASiVar dataset involving five subjects scanned at multiple sites further confirmed these findings. Harmonization led to reduced CV for both DSC and regional brain volumes with each subject across the scanners, although the limited sample size ( $N = 5$ ) precluded statistical significance between methods.

Our study demonstrates that HACA3<sup>+</sup>’s methodological improvements do not compromise performance on standard full FOV datasets, maintaining strong harmonization reliability. Critically, the attention mechanism enhancement of HACA3<sup>+</sup> benefits in scenarios involving limited FOV or regional dropout, an area where HACA3 had limitations. Importantly, the model is conservative in its predictions. Regions absent in all source images are not imputed, which mitigates the risk of hallucination.

The third enhancement did not directly lead to statistically significant improvement in HACA3. We trained HACA3<sup>+</sup> on  $6\times$  more scanners than HACA3 and more than  $4\times$  the number of subjects. Although, we also did not modify the HACA3 architecture. Due to the substantial increase in training data, it is possible that the model architecture needs to be adapted to handle more complexity, which could be a reason why we did not observe a statistically significant improvement. Our experiments were geared towards validation of harmonization on traveling subject datasets and not on methodological development of a harmonization algorithm. This is an area for future exploration.

Our validation focuses primarily on healthy subjects and individuals with MS, reflecting the TREAT-MS pragmatic clinical trial dataset that contributed the bulk of our training data. The models’ performance in other populations has not been tested. Furthermore, harmonization is currently limited by the artifact encoder’s ability to detect and mitigate poor quality images. Extension to diverse disease populations and architecture development are promising directions of future work.

In summary, HACA3+ establishes a new benchmark in multi-site harmonization, trained on a large and diverse dataset, and rigorously validated on multiple traveling subject datasets. The model’s enhancements improve adaptability to real-world clinical variability without sacrificing in-domain reliability. The code and pre-trained weights are made publicly available to facilitate future research and clinical translation.

## Acknowledgments

This research is partially supported by the Johns Hopkins University Percy Pierre Fellowship (Hays) and the National Science Foundation Graduate Research Fellowship under Grant No. DGE-2139757 (Hays). Development is partially supported by FG-2008-36966 (Dewey), CDMRP W81XWH2010912 (Prince), NIH R01EB036013 (Prince), NIH R01 CA253923 (Landman), NIH R01 CA275015 (Landman), the National MS Society grant RG-1507-05243 (Pham) and Patient-Centered Outcomes Research Institute (PCORI) grant MS-1610-37115 (Newsome and Mowry). The statements in this publication are solely the responsibility of the authors and do not necessarily represent the views of the Patient-Centered Outcomes Research Institute (PCORI), its Board of Governors or Methodology Committee.

This research was supported in part by the Intramural Research Program of the National Institutes of Health (NIH). The contributions of the NIH author(s) were made as part of their official duties as NIH federal employees, are in compliance with agency policy requirements, and are considered Works of the United States Government. However, the findings and conclusions presented in this paper are those of the author(s) and do not necessarily reflect the views of the NIH or the U.S. Department of Health and Human Services.

Data were provided [in part] by the Human Connectome Project, WU-Minn Consortium (Principal Investigators: David Van Essen and Kamil Ugurbil; 1U54MH091657) funded by the 16 NIH Institutes and Centers that support the NIH Blueprint for Neuroscience Research; and by the McDonnell Center for Systems Neuroscience at Washington University.

Data were provided [in part] by OASIS-3: Longitudinal Multimodal Neuroimaging: Principal Investigators: T. Benzinger, D. Marcus, J. Morris; NIH P30 AG066444, P50 AG00561, P30 NS09857781, P01 AG026276, P01 AG003991, R01 AG043434, UL1 TR000448, R01 EB009352.

## References

- Brian B Avants et al. Advanced normalization tools (ANTS). *Insight J*, 2(365):1–35, 2009.
- Joanne C Beer et al. Longitudinal ComBat: A method for harmonizing longitudinal multi-scanner imaging data. *NeuroImage*, 220:117129, 2020.
- Leon Y. Cai et al. Masivar: Multisite, multiscanner, and multisubject acquisitions for studying variability in diffusion weighted mri. *Magnetic Resonance in Medicine*, 86(6): 3304–3320, 2021.
- Blake E Dewey et al. DeepHarmony: a deep learning approach to contrast harmonization across scanner changes. *Magnetic Resonance Imaging*, 64:160–170, 2019.
- Blake E Dewey et al. A disentangled latent space for cross-site MRI harmonization. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 720–729, 2020.
- Jean-Philippe Fortin et al. Harmonization of cortical thickness measurements across scanners and sites. *NeuroImage*, 167:104–120, 2018. ISSN 1053-8119.



- Savannah P Hays et al. Rescuing incomplete mr data: Anatomy imputation of restricted field of view images using multi-contrast mr images. In *Proceedings of the International Society for Magnetic Resonance in Medicine. ISMRM Annual Meeting*, volume 1045, 2025a.
- Savannah P Hays et al. An Unsupervised Approach for Artifact Severity Scoring in Multi-Contrast MR Images. In *Medical Imaging with Deep Learning*, 2025b.
- Hannah Horng et al. Generalized ComBat harmonization methods for radiomic features with multi-modal distributions and multiple batch effects. *Scientific Reports*, 12(1):4493, 2022.
- Xun Huang et al. Multimodal unsupervised image-to-image translation. In *Proceedings of the European Conference on Computer Vision*, pages 172–189, 2018.
- Yuankai Huo et al. 3D whole brain segmentation using spatially localized atlas network tiles. *NeuroImage*, 194:105–119, 2019.
- Pamela J. LaMontagne et al. OASIS-3: Longitudinal neuroimaging, clinical, and cognitive dataset for normal aging and Alzheimer disease. *medRxiv*, 2019.
- Mengting Liu et al. Style transfer using generative adversarial networks for multi-site MRI harmonization. In *24<sup>th</sup> International Conference on Medical Image Computing and Computer Assisted Intervention (MICCAI 2021)*, volume 12903, pages 313–322, 2021.
- Ming-Yu Liu et al. Unsupervised Image-to-image Translation Networks. In *Advances in Neural Information Processing Systems*, pages 700–708, 2017.
- Siyuan Liu and Pew-Thian Yap. Learning multi-site harmonization of magnetic resonance images without traveling human phantoms. *Communications Engineering*, 3(1):6, 2024.
- Yuan-Chiao Lu et al. An evaluation of image-based and statistical techniques for harmonizing brain volume measurements. *Imaging Neuroscience*, 3:IMAG.a.73, 2025.
- Gourav Modanwal et al. MRI image harmonization using cycle-consistent generative adversarial network. In *Medical Imaging 2020: Computer-Aided Diagnosis*, volume 11314, pages 259–264. SPIE, 2020.
- Ellen M Mowry et al. The TRaditional versus Early Aggressive Therapy for MS (TREAT-MS) trial: Design and participant characteristics at enrollment. *Contemporary Clinical Trials*, 159:108117, 2025.
- Roland Opfer et al. Automatic segmentation of the thalamus using a massively trained 3d convolutional neural network: higher sensitivity for the detection of reduced thalamus volume by improved inter-scanner stability. *European Radiology*, 33(3):1852–1861, 2023.
- Fernando Pérez-García et al. 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, 2021.

- Samuel W Remedios et al. Self-supervised super-resolution for anisotropic MR images with and without slice gap. In *International Workshop on Simulation and Synthesis in Medical Imaging*, pages 118–128. Springer, 2023.
- Susan M Resnick et al. One-year age changes in MRI brain volumes in older adults. *Cerebral Cortex*, 10(5):464–472, 2000.
- Nicholas J Tustison et al. N4ITK: Improved N3 bias correction. *IEEE Transactions on Medical Imaging*, 29(6):1310–1320, 2010.
- Lianrui Zuo et al. HACA3: A unified approach for multi-site MR image harmonization. *Computerized Medical Imaging and Graphics*, 109(102285), 2023.