



Deep Unfolding Network with Spatial Alignment for multi-modal MRI reconstruction

Hao Zhang^a, Qi Wang^a, Jun Shi^b, Shihui Ying^{a,*}, Zhijie Wen^a

^aDepartment of Mathematics, School of Science, Shanghai University, Shanghai 200444, China

^bSchool of Communication and Information Engineering, Shanghai University, Shanghai 200444, China

ARTICLE INFO

Article history:

Keywords: Multi-modal MRI reconstruction, deep unfolding network, spatial alignment, denoising and inter-modality prior

ABSTRACT

Multi-modal Magnetic Resonance Imaging (MRI) offers complementary diagnostic information, but some modalities are limited by the long scanning time. To accelerate the whole acquisition process, MRI reconstruction of one modality from highly undersampled k-space data with another fully-sampled reference modality is an efficient solution. However, the misalignment between modalities, which is common in clinic practice, can negatively affect reconstruction quality. Existing deep learning-based methods that account for inter-modality misalignment perform better, but still share two main common limitations: (1) The spatial alignment task is not adaptively integrated with the reconstruction process, resulting in insufficient complementarity between the two tasks; (2) the entire framework has weak interpretability. In this paper, we construct a novel Deep Unfolding Network with Spatial Alignment, termed DUN-SA, to appropriately embed the spatial alignment task into the reconstruction process. Concretely, we derive a novel joint alignment-reconstruction model with a specially designed cross-modal spatial alignment term. By relaxing the model into cross-modal spatial alignment and multi-modal reconstruction tasks, we propose an effective algorithm to solve this model alternatively. Then, we unfold the iterative steps of the proposed algorithm and design corresponding network modules to build DUN-SA with interpretability. Through end-to-end training, we effectively compensate for spatial misalignment using only reconstruction loss, and utilize the progressively aligned reference modality to provide inter-modality prior to improve the reconstruction of the target modality. Comprehensive experiments on three real datasets demonstrate that our method exhibits superior reconstruction performance compared to state-of-the-art methods.

© 2023 Elsevier B. V. All rights reserved.

1. Introduction

Magnetic Resonance Imaging (MRI), due to its non-invasive, high resolution, and significant soft tissue contrast, has become a widely used medical imaging technique. However, the MR scan is relatively slow due to the repetitive acquisition of MR signal spatial encoding and hardware limitations. This time-consuming process can lead to discomfort for patients, caus-

ing them to move and introducing motion artifacts in the images, which may negatively affect subsequent disease diagnosis. Therefore, accelerating MRI acquisition is of great importance in clinic practice. To accelerate MRI acquisition, a feasible strategy is to reduce the amount of k-space data collected and then reconstruct fully-sampled images by undersampled data.

The Compressed Sensing MRI (CS-MRI) methods enable the accurate reconstruction from undersampled data at sampling rates significantly below those required by the Nyquist sampling theorem (Lustig et al., 2007). They aim at designing some hand-crafted regularizers based on different priors (e.g., struc-

*Corresponding author:
e-mail: shying@shu.edu.cn (Shihui Ying)

tered sparsity (Lai et al., 2016; Yang et al., 2015), non-local sparsity (Qu et al., 2014; Eksioğlu, 2016)) of MR images and formulating them into the algorithm optimization to constrain the solution space. Regardless of its theoretical guarantees, it is challenging to handcraft an optimal regularizer. Alternatively, deep learning-based methods have attracted widespread attention in MRI reconstruction due to their accuracy and speed (Wang et al., 2016; Sriram et al., 2020; Zhu et al., 2018). By learning robust feature representations, the deep learning-based methods have achieved impressive reconstruction performance. However, most deep learning-based methods are a black-box process and lack interpretability, which is required in clinical practice. To alleviate this black-box issue, deep unfolding networks (Yang et al., 2020a; Xin et al., 2022; Huang et al., 2023; Zhang et al., 2022; Jiang et al., 2023b,a) have been proposed to incorporate the imaging model and domain knowledge into the network. They unfold the iterations of an optimization algorithm into deep neural networks, thus making the learning process interpretable.

Despite the potential of deep unfolding networks in MRI reconstruction, most of them focus on utilizing information from a single modality. However, in clinical practice, it is common to acquire MR images of different contrasts because each modality reveals distinct tissue and organ characteristics, and the complementary information among modalities contributes to more accurate diagnoses. While different modalities of MR images display different signal types, they are spatially corresponding and depict the same anatomical structures. Researches (Xiang et al., 2018; Feng et al., 2023; Bian et al., 2022; Lei et al., 2023) have shown that reconstruction of one modality (target modality) can be improved by utilizing the information from another modality (reference modality). But, these multi-modal MRI reconstruction methods are all based on the assumption that the images are perfectly aligned, which is rare in practice. The misalignment may negatively affect the reconstruction performance due to insufficiently exploring the correlation of different modalities.

To mitigate the misalignment between modalities, Lai et al. (2017) alternating iteration of the registration and the reconstruction to align the reference modality with the intermediate reconstruction result of the under-sampled target modality. However, the use of conventional iterative optimization is relatively time-consuming. Deep learning-based methods that take inter-modality misalignment into consideration perform better in terms of reconstruction and spatial alignment accuracy, but still share two main common limitations: (1) The spatial alignment is not adaptively integrated with the reconstruction process, resulting in insufficient complementarity between the two tasks; (2) the entire framework has weak interpretability. For example, before reconstruction, Xuan et al. (2022) and Liu et al. (2021) simply align the images of the reference modality with the undersampled images of the target modality, neglecting the negative effects that artifacts in the undersampled images can have on spatial alignment.

To alleviate the aforementioned limitations, we propose a Deep Unfolding Network with Spatial Alignment (DUN-SA) framework in this paper. Specifically, we first derive a novel joint alignment-reconstruction model, in which a cross-modal

spatial alignment term is developed to compensate for misalignment between modalities. Subsequently, we relax it into cross-modal spatial alignment and multi-modal reconstruction tasks, and propose an optimization algorithm for these two tasks. Concretely, we use the gradient-based algorithm for spatial alignment and the Half-Quadratic Splitting (HQS) algorithm for reconstruction. By alternately optimizing the two tasks, we solve this model. Further, we unfold the iterative steps of the proposed algorithm and design corresponding network modules. Finally, we propose an end-to-end, deep unfolding network with interpretability. The main contributions of this paper are as follows:

- We propose a novel joint alignment-reconstruction model for multi-modal MRI reconstruction, in which a cross-modal spatial alignment term is developed to compensate for misalignment between modalities. Utilizing gradient-based and HQS techniques, we design an optimization algorithm to alternately solve this model.
- By unfolding the iterative steps of the proposed algorithm and integrating them with specially designed network modules, we construct a deep unfolding network, termed DUN-SA, which exhibits clear interpretability.
- We design the Spatial Alignment Augmented Multi-Modal Learning Block (SAAMMLB) to learn inter-modality prior through aligned images from reference modality, and utilize the Denoising Block (DB) to fully exploit intra-modality prior.
- Through extensive experiments on the fastMRI brain, the IXI dataset, and an in-house dataset, we demonstrate that the proposed DUN-SA outperforms existing state-of-the-art methods in terms of quantitative and qualitative reconstruction results.

2. Related Work

2.1. Single-Modal CS-MRI

The single-modal MRI reconstruction problem is to reconstruct the original image from its partial acquisition. Traditionally, model-based methods usually utilize image prior to improve reconstruction performance. Lai et al. (2016); Yang et al. (2015) use wavelet transform to sparsely represent magnetic resonance images in iterative image reconstructions, capitalizing on sparse prior. Qu et al. (2014); Eksioğlu (2016) further sparsify magnetic resonance images by exploiting the similarity of image patches, emphasizing non-local regularization. Ravishanker and Bresler (2010); Zhan et al. (2015) propose dictionary learning for adaptively learning the sparsifying transform (dictionary), and reconstructing the image simultaneously from highly undersampled data.

Deep learning methods, with their powerful ability to represent features, can achieve superior reconstruction quality and acceleration compared to non-deep learning approaches. Wang et al. (2016) utilize CNN architectures to learn the mapping relationship between the MR images obtained from zero-filled

and fully-sampled k-space data, significantly improving reconstruction speed while maintaining image quality. Yang et al. (2017); Defazio et al. (2020) not only reconstruct images but also retain more image texture details and reduce more image artifacts through adversarial learning. Contrary to models that reconstruct under-sampled inputs in the image domain, Sriram et al. (2020); Zhu et al. (2018) operate on under-sampled k-space, achieving better reconstruction performance. Shaul et al. (2020); Ran et al. (2020); Zhang et al. (2019); Eo et al. (2018) utilize cross-domain networks for image reconstruction, leveraging information in both image domain and k-space, outperforming single-domain methods. Furthermore, Wang et al. (2020) incorporate the wavelet domain, using information across three domains, further improving reconstruction accuracy.

2.2. Multi-Modal CS-MRI

In multi-modal MRI reconstruction, a reference modality can guide the reconstruction of the target modality. Model-based methods explore the relationship between modalities based on prior knowledge. Ehrhardt and Betcke (2016) introduce two types of total variation based on position and direction, taking into account the structural prior of MR images. Weizman et al. (2016) acknowledge the differences between modalities and proposes an iterative weighted reconstruction approach. Song et al. (2020) present a method based on coupled dictionary learning for multi-contrast MRI reconstruction, effectively utilizing the structural dependencies between different contrasts. Lai et al. (2017) address the misalignment issue between modalities, proposing a novel MRI image reconstruction method that learns prior from multi-contrast images through a graph wavelet representation, modeling it as a bi-level optimization problem to allow misalignment between these images.

Deep learning-based methods (Sun et al., 2019; Xiang et al., 2019; Dar et al., 2020) combine T1-weighted and T2-weighted images as dual-channel inputs to a deep learning model, improving reconstruction quality by leveraging inter-modality complementary information. (Zhou and Zhou, 2020) integrate multi-modal and dual-domain information, further enhancing reconstruction accuracy. Feng et al. (2023) employ a Transformer structure, using multi-head attention mechanism to deeply capture multi-modal information, offering more global information compared to existing CNN-based methods. Methods that take misalignment between modalities into consideration perform better. Liu et al. (2021) utilize CNN to learn rigid transformation for compensating misalignment between unregistered paired multi-modal MR images, thereby making more efficient use of the reference modality. Xuan et al. (2022) employ a spatial alignment network to compensate misalignment between modalities and incorporate a novel loss function that efficiently trains both the spatial alignment network and the reconstruction network simultaneously.

2.3. Deep Unfolding Network

Deep Unfolding Networks have achieved impressive results in many medical applications (e.g. dynamic MR imaging (Huang et al., 2021), metal artifact reduction (Wang et al., 2023,

2021), MRI Super-resolution (Yang et al., 2023, 2022; Lei et al., 2023) and CS-MRI (Yang et al., 2020a,b; Lei et al., 2023)), thereby attracting widespread attention recently. By unfolding certain optimization algorithms with network modules, They integrate model-based and learning-based methods well.

In single-modal CS-MRI, for example, by unfolding different optimization algorithms such as Alternating Direction Method of Multipliers (ADMM) (Yang et al., 2020a; Huang et al., 2023; Jiang et al., 2023b), Alternating Iterative Shrinkage-thresholding Algorithm (ISTA) (Zhang et al., 2022), Half-Quadratic Splitting (HQS) (Xin et al., 2022; Jiang et al., 2023a) and combining them with different deep neural networks (e.g., CNNs, S-Nets), deep unfolding networks achieve sparsity in the denoising perspective. This approach not only improves reconstruction quality but also offers interpretability. Additionally, research has shown that incorporating multi-modal information can improve the quality of reconstruction, prompting the development of deep unfolding networks for multi-modal MRI reconstruction. Yang et al. (2020b) capitalize on cross-modal prior and combine channel and spatial attention mechanisms to construct proximal operators, significantly enhancing image reconstruction accuracy. Lei et al. (2023) employ Convolutional Sparse Coding (CSC) techniques for multi-modal MRI modeling. It transfers the common texture information from the reference modality images to the target modality images while avoiding the interference of inconsistent information and integrates deep network modules, achieving superior reconstruction results.

3. Method

In Section 3, we propose a joint alignment-reconstruction model for multi-modal MRI reconstruction and design an optimization algorithm in Section 3.1. Then, we unfold the iterative steps of the proposed algorithm with corresponding network modules and build DUN-SA in Section 3.2. Finally, we introduce the learnable parameters in DUN-SA and the loss function in Section 3.3.

3.1. Joint alignment-reconstruction model

3.1.1. The objective functions

Let $x \in \mathbb{R}^N$ represent the unknown fully-sampled MRI image, and $k \in \mathbb{C}^M$ denote the measured fully-sampled frequency-domain signal, also known as k-space. In the classical compressed sensing problem for MRI, $\tilde{k} \in \mathbb{C}^m$ is used to represent the undersampled signal, derived by applying a Fourier transform and a masking operator to the fully-sampled image. Consequently, the relationship between the fully-sampled MRI image and the undersampled k-space signal can be expressed as:

$$\tilde{k} = F_m x + n, \quad (1)$$

where $F_m = MF$ represents masked Fourier transform, F is the Fourier transform operator, M is the masking operator, and n is the measurement noise generally considered as Gaussian noise. Additionally, $m \ll M$, and $\frac{m}{M}$ indicates the compression ratio. m is significantly smaller than M , which leads to an underdetermined inverse problem.

The classical compressed sensing method achieves reconstruction by optimizing the following energy function:

$$\min_x \frac{1}{2} \|F_m x - \tilde{k}\|_2^2 + \eta R(x), \quad (2)$$

where the first term is the fidelity term, ensuring that the reconstructed image x is consistent with the undersampled signal \tilde{k} in k-space; the second term $R(x)$ is the regularization term, enforcing prior constraint like sparsity, and η is the balancing factor.

In multi-modal MRI reconstruction, an auxiliary term measures the correlation between x and x_{ref} , where the reference modality provides extra inter-modality prior to assist the reconstruction of the target modality. The energy function becomes:

$$\min_x \frac{1}{2} \|F_m x - \tilde{k}\|_2^2 + \lambda \Psi(x, x_{\text{ref}}) + \eta R(x), \quad (3)$$

where x_{ref} represents the fully-sampled image from the reference modality, $\Psi(x, x_{\text{ref}})$ models the correlation between x and x_{ref} , and λ is the balancing coefficient.

Furthermore, to compensate for the spatial misalignment between x and x_{ref} , we transform the auxiliary term into a cross-modal spatial alignment term, integrating the spatial alignment task into the reconstruction process, yielding the final energy function:

$$\min_{x, \phi} \frac{1}{2} \|F_m x - \tilde{k}\|_2^2 + \lambda \Psi(x, \mathcal{T}(x_{\text{ref}}, \phi)) + \eta R(x), \quad (4)$$

where ϕ represents the displacement field, \mathcal{T} denotes the warp operation, and $\mathcal{T}(x_{\text{ref}}, \phi)$ represents the aligned reference modality image. We refer to $\Psi(x, \mathcal{T}(x_{\text{ref}}, \phi))$ as the cross-modal spatial alignment term, which models the correlation between target modality and aligned reference modality.

We set Eq. (4) as our final objective function. By eliminating misalignment between modalities, we can utilize the information between corresponding points in aligned reference modality and target modality images to learn inter-modality prior more efficiently, and further improve the reconstruction results.

3.1.2. The optimization algorithm

To solve this model effectively, we relax Eq. (4) into cross-modal spatial alignment and multi-modal reconstruction tasks and optimize them alternatively:

$$\hat{\phi} = \arg \min_{\phi} \lambda \Psi(x, \mathcal{T}(x_{\text{ref}}, \phi)), \quad (5a)$$

$$\hat{x} = \arg \min_x \frac{1}{2} \|F_m x - \tilde{k}\|_2^2 + \lambda \Psi(x, \mathcal{T}(x_{\text{ref}}, \phi)) + \eta R(x). \quad (5b)$$

Update ϕ : Eq. (5a) describes the cross-modal spatial alignment task. We optimize this by employing the gradient-based algorithm, and the $(t+1)$ -th optimization step can be expressed as:

$$\phi^{t+1} = \phi^t - \alpha \nabla_{\phi} \Psi(x^t, \mathcal{T}(x_{\text{ref}}, \phi^t)), \quad (6)$$

where $\alpha = \rho \lambda$, ρ is the step size and $\nabla_{\phi} \Psi(x^t, \mathcal{T}(x_{\text{ref}}, \phi^t))$ is the gradient of the cross-modal spatial alignment term with respect to the displacement field.

Eq. (5b) represents the image reconstruction task. At this point, we regard the cross-modal spatial alignment term as a regularization term. Given that the HQS method has been proven effective for image inverse problems (Geman and Reynolds, 1992; He et al., 2013), by introducing two auxiliary variables z and s , we further transform Eq. (5b) to an unconstrained optimization problem:

$$\min_{x, s, z} \frac{1}{2} \|F_m x - \tilde{k}\|_2^2 + \lambda \Psi(z, \mathcal{T}(x_{\text{ref}}, \phi)) + \eta R(s) + \frac{\beta_1}{2} \|x - z\|_2^2 + \frac{\beta_2}{2} \|x - s\|_2^2, \quad (7)$$

where β_1 and β_2 are penalty parameters. As β_1 and β_2 approach infinity, the result of minimizing Eq. (7) will converge to the result of minimizing equation Eq. (5b). The optimization for the $(t+1)$ -th iteration can be represented as:

$$z^{t+1} = \arg \min_z \frac{\beta_1}{2} \|x^t - z\|_2^2 + \lambda \Psi(z, \mathcal{T}(x_{\text{ref}}, \phi^{t+1})), \quad (8a)$$

$$s^{t+1} = \arg \min_s \frac{\beta_2}{2} \|x^t - s\|_2^2 + \eta R(s), \quad (8b)$$

$$x^{t+1} = \arg \min_x \frac{1}{2} \|F_m x - \tilde{k}\|_2^2 + \frac{\beta_1}{2} \|x - z^{t+1}\|_2^2 + \frac{\beta_2}{2} \|x - s^{t+1}\|_2^2. \quad (8c)$$

Update z : Given target modality image x^t and aligned reference modality image $\mathcal{T}(x_{\text{ref}}, \phi^{t+1})$. We define the proximal operator $\text{prox}_{\frac{\beta_1}{2} \Psi(\cdot, \mathcal{T}(x_{\text{ref}}, \phi^{t+1}))}(\cdot)$ such that $\text{prox}_{\frac{\beta_1}{2} \Psi(\cdot, \mathcal{T}(x_{\text{ref}}, \phi^{t+1}))}(x) = \arg \min_z \frac{\beta_1}{2} \|x - z\|_2^2 + \lambda \Psi(z, \mathcal{T}(x_{\text{ref}}, \phi^{t+1}))$. Eq. (8a) is then solved by the following equation:

$$z^{t+1} = \text{prox}_{\frac{\beta_1}{2} \Psi(\cdot, \mathcal{T}(x_{\text{ref}}, \phi^{t+1}))}(x^t). \quad (9)$$

Update s : Given the target modality image x^t . We also define the proximal operator $\text{prox}_{\frac{\beta_2}{2} R}(\cdot)$ such that $\text{prox}_{\frac{\beta_2}{2} R}(x) = \arg \min_s \frac{\beta_2}{2} \|x - s\|_2^2 + \eta R(s)$. The solution to Eq. (8b) can then be expressed as:

$$s^{t+1} = \text{prox}_{\frac{\beta_2}{2} R}(x^t). \quad (10)$$

Update x : Eq. (8c) represents a quadratic regularized least squares problem and has a closed-form solution:

$$\begin{aligned} x^{t+1} &= (F_m^H F_m + (\beta_1 + \beta_2) I)^{-1} (F_m^H \tilde{k} + \beta_1 z^{t+1} + \beta_2 s^{t+1}) \\ &= F_m^H \Lambda^{-1} (M^H \tilde{k} + \beta_1 F z^{t+1} + \beta_2 F s^{t+1}), \end{aligned} \quad (11)$$

where, $\Lambda = \text{diag}(\beta_1 + \beta_2)$, F_m^H is the Hermitian transpose of F_m , and I is the identity matrix.

Through the iterative updates of ϕ , z , s , x , we can address the joint alignment-reconstruction model. We unfold this optimization algorithm into a deep unfolding network and name it Deep Unfolding Network with Spatial Alignment (DUN-SA), as shown in Fig. 1. The details are presented in Section 3.2.

3.2. Deep Unfolding Network with Spatial Alignment

In Section 3.1, we propose an optimization algorithm for joint alignment-reconstruction model. However, this algorithm

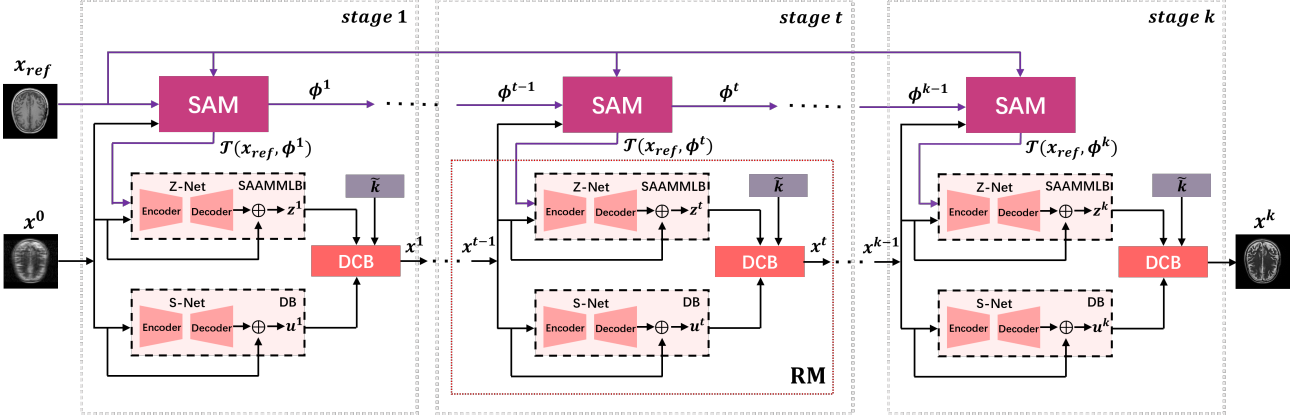


Fig. 1. The overall structure of the proposed Deep Unfolding Network with Spatial Alignment (DUN-SA) consists of SAM (Spatial Alignment Module) and RM (Reconstruction Module). The RM is composed of SAAMMLB (Spatial Alignment Augmented Multi-Modal Learning Block), DB (Denoising Block), and DCB (Data Consistency Block). SAM is used to solve spatial alignment task, while RM is for reconstruction task. Specifically, SAAMMLB is used to learn spatial alignment augmented inter-modality prior, DB is used to learn denoising prior, and DCB is used to enforce data consistency constraint.

has some limitations. Firstly, it requires incorporating the spatial alignment task into the reconstruction process and effectively solving both alignment and reconstruction tasks, making the design of the cross-modal spatial alignment term complicated. Secondly, similar to traditional optimization algorithms, this method requires dozens of iterations to converge, thus increasing time and computational costs significantly. To alleviate these problems, inspired by deep unfolding networks, we unfold the optimization algorithm with the network modules.

3.2.1. Model Overview

The core concept of the Deep Unfolding Network with Spatial Alignment is to solve the joint alignment-reconstruction model by unfolding the proposed algorithm with network modules and learning all parameters in an end-to-end manner.

We regard Eq. (6) as the cross-modal spatial alignment process and design the Spatial Alignment Module (SAM) accordingly; we view Eq. (7) as the multi-modal reconstruction task, and propose corresponding Reconstruction Module (RM). Here are detailed introductions to both modules:

Within SAM, we utilize the gradient of the cross-modal spatial alignment term to update the displacement field. However, the gradient is difficult to represent accurately due to the non-linearity and high complexity of the cross-modal spatial alignment term. Therefore, we use a deep network to learn the gradient information and refine the displacement field.

In RM, we iteratively update z , s , x to achieve high-quality reconstruction results. First, based on Eq. (9), we design the SAAMMLB. By fully leveraging the relationship between corresponding points in aligned reference and target modality images, the SAAMMLB is able to effectively learn inter-modality prior. Further, we consider Eq. (10) as a denoising process and consequently design the DB to learn intra-modality denoising priors. Finally, we design the DCB according to the closed-form solution in Eq. (11). Subsequently, we will provide a detailed explanation of each module.

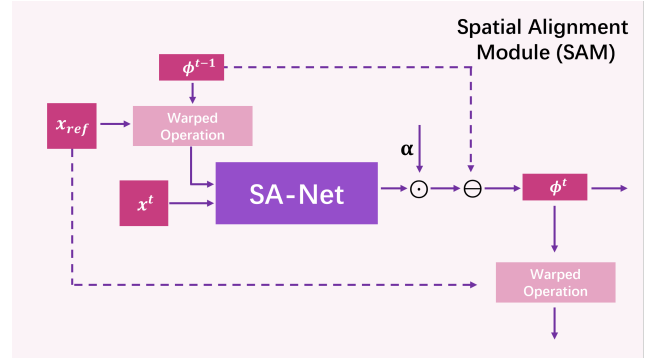


Fig. 2. Architecture of Spatial Alignment Module (SAM).

3.2.2. Spatial Alignment Module (SAM)

In SAM, in order to maintain the form of Eq. (6), considering the complexity of the cross-modal spatial alignment term, we employ a CNN to learn the gradient of the displacement field. This design refines the spatial alignment result by aligning the reference modality with the intermediate reconstruction result of target modality, thus reducing the impact of artifacts in the undersampled image on spatial alignment process.

Throughout the optimization process, the CNN takes the aligned reference modality image from the previous step and the image generated by the reconstruction module, denoted as $(x^t, \mathcal{T}(x_{ref}, \phi^{t+1}))$, as input and outputs the corresponding gradient to update ϕ as in Eq. (6). It is worth noting that x^0 is initialized as zero-filled reconstruction, Furthermore, to effectively capture the multi-scale information between images, we employ a CNN architecture known as U-Net (Ronneberger et al., 2015) and name it SA-Net (spatial alignment network). Ultimately, Eq. (6) can be expressed as:

$$\phi^{t+1} = \phi^t - \alpha \text{SA-Net}(x^t, \mathcal{T}(x_{ref}, \phi^t)). \quad (12)$$

The learnable parameters include the scaling factor α , which is initialized to be 1, and the parameters within the SA-Net, which is initialized to be 0, shown in Fig. 2.

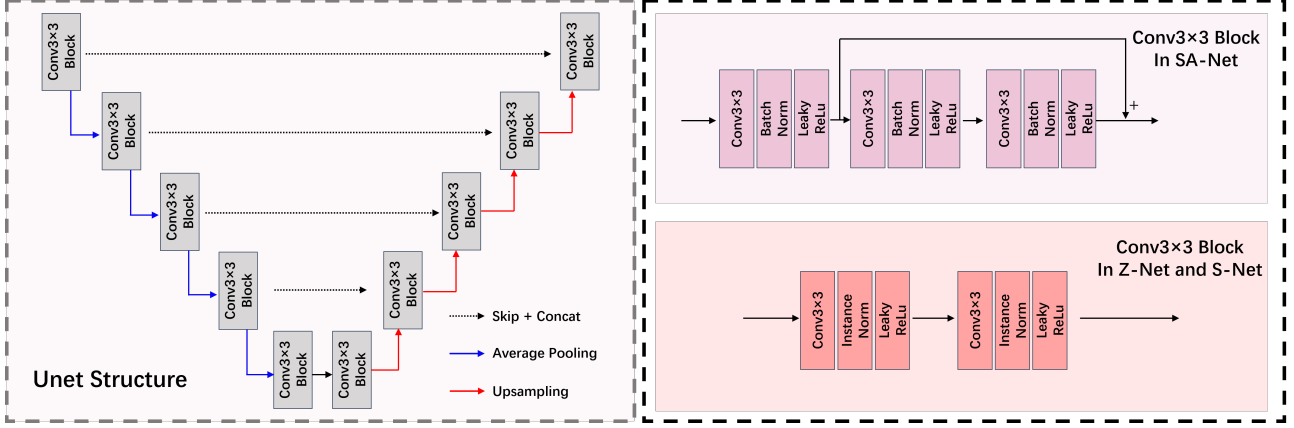


Fig. 3. Detailed configurations of SA-Net, Z-Net and S-Net.

3.2.3. Spatial Alignment Augmented Multi-Modal Learning Block (SAAMMLB)

In SAAMMLB, instead of designing a hand-crafted regularization to learn inter-modality prior, we replace the proximal operator in Eq. (9) by a learnable network. Previous works (Xiang et al., 2018; Sun et al., 2019; Xiang et al., 2019; Dar et al., 2020) have demonstrated the effectiveness of U-Net in multi-modal feature extraction. By concatenating aligned images of two modalities along the channel dimension and feeding them into the network, the forward convolution layers will fuse information and features from corresponding positions between channels. In this paper, we adopted a variant of U-Net as the backbone of SAAMMLB, naming it Z-Net (network used to update z). By feeding the target modality image and the aligned reference modality image, denoted as $(x^t, \mathcal{T}(x_{\text{ref}}, \phi^{t+1}))$, into Z-Net, we obtain an image refined with inter-modality prior. Thus, Eq. (9) can be transformed to:

$$z^{t+1} = \text{Z-Net}(x^t, \mathcal{T}(x_{\text{ref}}, \phi^{t+1})). \quad (13)$$

However, directly obtaining the refined image might lead to instability in the reconstruction process. Inspired by ResNet structure (He et al., 2016), We carry out this process based on the reconstructed image from the previous step, and further integrating the information learned from the inter-modality prior to stabilize the reconstruction process, resulting in the equation:

$$z^{t+1} = x^t + \text{Z-Net}(x^t, \mathcal{T}(x_{\text{ref}}, \phi^{t+1})). \quad (14)$$

The structure of Z-Net is illustrated in the Fig. 3. The learnable parameters include those within the Z-Net.

3.2.4. Denoising Block (DB)

In DB, we consider Eq. (10) as a denoising process driven by a denoising prior (Aggarwal et al., 2019; Xin et al., 2022; Yang et al., 2020a). In general, any existing image denoising network can be applied here. Given the simplicity and effectiveness of U-Net in learning the intra-modality prior (Sriram et al., 2020; Xin et al., 2022), we employ U-Net to replace the proximal mapping, naming it S-Net (network used to update s), leading to a transformation of Eq. (10) into:

$$u^{t+1} = \text{S-Net}(x^t) \quad (15)$$

To stabilize both the training and reconstruction processes, we also incorporate the architecture of ResNet, resulting in the following equation:

$$u^{t+1} = x^t + \text{S-Net}(x^t) \quad (16)$$

The structure of S-Net is shown in Fig. 3. The learnable parameters include those within the S-Net.

3.2.5. Data Consistency Block (DCB)

In DCB, to maintain the structure of the algorithm, the reconstruction process can be directly expressed by Eq. (11). The closed-form solution can be obtained under the given (z^{t+1}, s^{t+1}) from SAAMMLB and DB. The learnable parameters include the scaling factors β_1 and β_2 , which are both initialized to be 1.

3.3. Network Parameters and Loss Function

Throughout the training of all model parameters, including α and SA-Net in SAM; Z-Net in SAAMMLB; S-Net in DB; and β_1, β_2 in DCB, we optimize them by minimizing the loss function. The training loss for each training pair is defined as the distance between the reconstructed images and the ground truth images. Structural Similarity Loss (SSIM) focuses on image structure, luminance, and contrast, making it well-suited for our reconstruction task. Specifically, the loss function is defined as:

$$\mathcal{L} = \sum_{i=1}^N \text{SSIM}(x_{\text{rec},i}, x_{\text{gt},i}) \quad (17)$$

where N represents the number of training samples, $x_{\text{rec},i}$ denotes the i^{th} image reconstructed by the network, and $x_{\text{gt},i}$ signifies the i^{th} ground-truth image.

4. Experiments settings

4.1. Datasets

We evaluate our method using three datasets, namely the fastMRI dataset, the IXI dataset and an in-house dataset.

FastMRI Dataset¹: To maintain consistency in the experiments, we follow the setups described in (Xuan *et al.*, 2022) and select 340 pairs of T1-weighted and T2-weighted axial brain MRIs. We use 170 volumes (2720 pairs of slices) as the training set, 68 volumes (1088 pairs of slices) as the validation set, and 102 volumes (1632 pairs of slices) as the test set. The in-plane size of all T1-weighted images and T2-weighted images is 320×320 , with the resolution of $0.68\text{mm} \times 0.68\text{mm}$ and the slice spacing of 5mm.

IXI Dataset²: The IXI dataset contains 576 paired multi-modal 3D brain MRIs. We select 570 pairs of PD-weighted and T2-weighted axial brain MRIs for our experiments. Specifically, we use 285 volumes as the training set, 115 volumes as the validation set, and 170 volumes as the test set. The in-plane size of all PD-weighted images and T2-weighted images is 256×256 . Consistent with (Lei *et al.*, 2023), we select the middle 100 slices from each volume for our experiments.

In-house Dataset: The in-house dataset contains 3D brain MRIs of 34 subjects, including paired T1-weighted images and T2-weighted images of their whole brains. Among them, 24 pairs of volumes (1797 pairs of slices) are employed as the training set. 4 pairs of volumes (300 pairs of slices) refer to the validation set, and 6 pairs of volumes (450 pairs of slices) are used for testing. The in-plane size of all T1-weighted images and T2-weighted images is 320×320 , with the resolution of $0.68\text{mm} \times 0.68\text{mm}$ and the slice spacing of 5mm.

4.2. Compared Methods

To evaluate the performance of our model (DUN-SA), we conduct comparisons with state of the art methods, including two single-modal MRI reconstruction methods: E2E-Varnet (Sriram *et al.*, 2020), HQS-Unet (Xin *et al.*, 2022); and four multi-modal MRI reconstruction methods: MD-DUN (Yang *et al.*, 2020b), MM-E2E-Varnet (a variant of E2E-VarNet for multi-modal input), SAN (Xuan *et al.*, 2022), and MC-CDic (Lei *et al.*, 2023). The details are presented as follows.

E2E-Varnet employs variational techniques to solve the MRI reconstruction problem, unfolding it into a network structure combined with U-Net for end-to-end learning.

HQS-Unet utilizes the HQS algorithm, integrating U-Net with residual and buffering designs to learn a denoising prior in an end-to-end manner.

MD-DUN builds upon the unfolded HQS algorithm and incorporates spatial and channel attention mechanisms within the denoising module, enabling more effective integration of reference modality information.

MM-E2E-Varnet modifies the original E2E-Varnet by replacing its architecture with a variant of U-Net consistent with multi-modal input, and utilizes reference modality information.

SAN performs image alignment between the target and reference modalities in advance, more effectively utilizing the information of the reference modality.

MC-CDic leverages the convolutional dictionary learning-based model and proximal gradient algorithm to propose a

corresponding multi-scale convolutional dictionary network for multi-modal MRI reconstruction.

4.3. Performance Evaluation Metric

To accurately and comprehensively evaluate the performance of image reconstruction, we select three core metrics: Peak Signal-to-Noise Ratio (PSNR), Structural Similarity Index Metric (SSIM), and Mean Absolute Error (MAE). These metrics are employed to quantify the similarity between the reconstructed images and the corresponding ground truth. It is noteworthy that higher values of PSNR and SSIM, along with lower MAE values, indicate superior reconstruction performance of the model.

4.4. Implementation Details

In our experiments, we follow the paradigm set by the fastMRI challenge, where under-sampled MRIs are obtained by masking the corresponding fully-sampled MRIs in k-space using a Cartesian sampling pattern. Specifically, we employ two types of sampling patterns: random and equispaced. These are applied at sampling ratios of 25% (4x acceleration) and 12.5% (8x acceleration), respectively. Acknowledging that low-frequency signals contain the majority of energy in k-space, we allocate 32% of the sampling to these low frequencies in both patterns. The remaining portion of the sampling is distributed either randomly or in an equispaced manner.

Concerning the selection of hyperparameters for the network, we choose a model size with stage=12 for our experiments.

All network models used in these experiments are implemented using the PyTorch framework and are run on a computer equipped with four NVIDIA GeForce GTX 3090 GPUs. During the model training phase, we set the batch size to be 2, adopt an end-to-end training strategy, and use the SSIM loss function with a coefficient of 1. All parameters are optimized using the Adam optimizer with a learning rate of 1×10^{-4} . To prevent overfitting, we used a validation set to select the optimal parameters during the training process. The final experimental results are based on these optimal parameters, as selected and validated through the validation set, and applied on the test set.

5. Experimental results

5.1. Results on the fastMRI dataset

We first run our proposed method on the fastMRI dataset. In this part of the experiments, we use the fully-sampled T1-weighted image as the reference modality to assist the k-space undersampled data of the T2 modality in reconstruction.

Table 1 and Table 2 display the quantitative evaluations on the fastMRI dataset for 1D equispaced and random masks under 4x and 8x acceleration. We observe that images reconstructed by multi-modal methods show significant improvements in the PSNR, SSIM, and MAE metrics across all experimental setups compared to those reconstructed by single-modal methods. This enhancement is particularly pronounced for 8x acceleration. The introduction of spatial alignment in the multi-modal reconstruction further improves the PSNR, SSIM, and

¹<https://fastMRI.med.nyu.edu/>.

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

Table 1. Quantitative evaluation of DUN-SA vs. other methods on the fastMRI dataset for 4x and 8x acceleration under equispaced 1D subsampling masks, where T1-weighted images are used as reference modality to assist the reconstruction of T2-weighted images. Best results are emphasized in bold, and the second best are emphasized with an underline.

Methods	Equispaced 4x acceleration			Equispaced 8x acceleration		
	PSNR	SSIM	MAE	PSNR	SSIM	MAE
Zero-filling	26.92±1.02	0.7321±0.0279	0.0238±0.0034	24.36±1.11	0.6428±0.0350	0.0328±0.0052
E2E-Varnet	38.81±1.61	0.9761±0.0065	0.0058±0.0012	36.66±1.65	0.9658±0.0087	0.0074±0.0015
HQS-Unet	39.31±1.69	0.9773±0.0063	0.0056±0.0011	37.12±1.68	0.9683±0.0085	0.0071±0.0017
MD-DUN	40.31±1.84	0.9721±0.0069	0.0053±0.0012	38.45±1.78	0.9618±0.0096	0.0061±0.0014
MM-E2E-Varnet	40.28±1.80	0.9806±0.0067	0.0051±0.0011	38.35±1.89	0.9728±0.0091	0.0063±0.0015
SAN	40.48±1.77	0.9819±0.0065	0.0048±0.0011	38.93±1.87	<u>0.9764±0.0082</u>	0.0059±0.0014
MC-CDic	40.72±1.80	0.9814±0.0067	<u>0.0047±0.0011</u>	38.25±1.85	0.9724±0.0100	0.0062±0.0014
DUN-SA	41.48±1.84	0.9838±0.0060	0.0045±0.0010	40.23±1.88	0.9802±0.0072	0.0052±0.0012

Table 2. Quantitative evaluation of DUN-SA vs. other methods on the fastMRI dataset for 4x and 8x acceleration under random 1D subsampling masks, where T1-weighted images are used as reference modality to assist the reconstruction of T2-weighted images. Best results are emphasized in bold, and the second best are emphasized with an underline.

Methods	Random 4x acceleration			Random 8x acceleration		
	PSNR	SSIM	MAE	PSNR	SSIM	MAE
Zero-filling	27.12±1.03	0.7373±0.0282	0.0233±0.0032	23.88±1.70	0.5818±0.0431	0.0342±0.0067
E2E-Varnet	42.60±1.71	0.9868±0.0044	0.0041±0.0009	35.82±1.63	0.9612±0.0094	0.0080±0.0017
HQS-Unet	43.05±1.87	0.9873±0.0042	0.0040±0.0008	36.22±1.73	0.9637±0.0093	0.0077±0.0016
MD-DUN	43.45±1.80	0.9832±0.0068	0.0037±0.0008	37.28±1.86	0.9541±0.0099	0.0069±0.0017
MM-E2E-Varnet	43.47±1.82	0.9879±0.0045	0.0038±0.0009	37.40±1.89	0.9705±0.0093	0.0071±0.0015
SAN	43.87±1.82	0.9886±0.0043	0.0038±0.0009	37.84±1.94	0.9724±0.0091	0.0069±0.0015
MC-CDic	44.01±1.82	0.9885±0.0038	0.0036±0.0008	37.60±1.82	0.9705±0.0100	0.0067±0.0015
DUN-SA	45.06±1.85	0.9907±0.0029	0.0032±0.0007	39.22±1.84	0.9770±0.0075	0.0057±0.0012

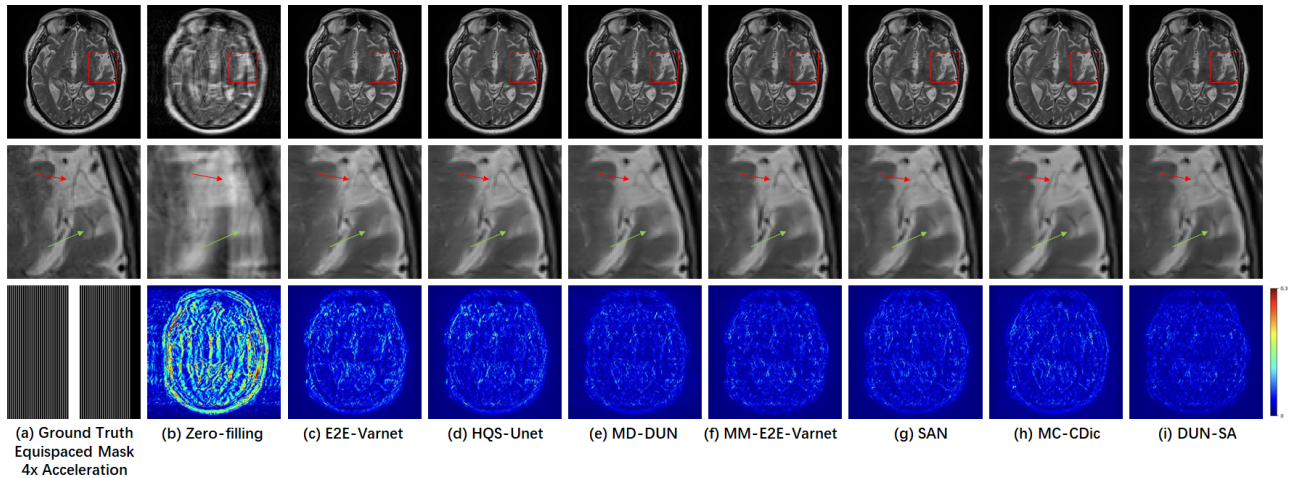


Fig. 4. Visual comparison with representative methods for 4x acceleration under 1D equispaced subsampling mask on fast MRI dataset. First row: Reconstructed images by different methods; second row: Zoomed-in region of interest; third row: Equispaced mask of 4x acceleration and error maps of different methods.

MAE values. It is evident that DUN-SA outperforms all others in every metric across all settings. In comparison with the second-best values, for 8x acceleration under random 1D subsampling masks, the highest differences in PSNR, SSIM, and MAE are 1.38db, 0.0046, and 0.001, respectively. In contrast, for 4x acceleration under equispaced 1D subsampling masks, the smallest differences are 0.76db, 0.0019, and 0.0002. These results suggest that DUN-SA can effectively utilize the reference modality to aid the MRI reconstruction of the target modality, leading to a noticeable improvement in reconstruction quality.

Fig. 4 illustrates a qualitative comparison of the reconstructed images and their corresponding error maps using various methods on the fastMRI dataset for 4x acceleration under 1D equispaced mask. The first row displays the reconstruction results from different methods, the second row showcases zoomed-in regions of interest, and the third row demonstrates the error maps. The patterns in the error maps represent reconstruction errors; a less intense color on the map (closer to the cool end of the spectrum) signifies superior reconstruction quality. Clearly, the Zero-filling image exhibits aliasing artifacts and lacks anatomical details. Significantly, compared to

Table 3. Quantitative evaluation of DUN-SA vs. other methods on the IXI dataset for 4x and 8x acceleration under equispaced 1D subsampling masks, where PD-weighted images are used as reference modality to assist the reconstruction of T2-weighted images. Best results are emphasized in bold, and the second best are emphasized with an underline.

Methods	Equispaced 4x acceleration			Equispaced 8x acceleration		
	PSNR	SSIM	MAE	PSNR	SSIM	MAE
Zero-filling	26.56±2.04	0.6337±0.0449	0.0287±0.0064	24.08±2.01	0.5517±0.0535	0.0369±0.0084
E2E-Varnet	41.87±2.29	0.9859±0.0044	0.0042±0.0011	34.51±2.10	0.9521±0.0109	0.0095±0.0022
HQS-Unet	42.73±2.33	0.9874±0.0036	0.0040±0.0009	35.41±2.12	0.9575±0.0090	0.0086±0.0020
MD-DUN	45.17±2.39	0.9874±0.0038	0.0034±0.0009	40.52±2.21	0.9729±0.0065	0.0058±0.0014
MM-E2E-Varnet	45.96±2.40	0.9921±0.0030	0.0030±0.0008	41.19±2.25	0.9832±0.0056	0.0049±0.0012
SAN	46.17±2.44	0.9923±0.0030	0.0029±0.0007	41.21±2.27	0.9833±0.0057	0.0049±0.0012
MC-CDic	45.79±2.41	0.9916±0.0032	0.0031±0.0006	40.82±2.23	0.9823±0.0060	0.0048±0.0012
DUN-SA	47.25±2.53	0.9936±0.0026	0.0025±0.0007	41.84±2.33	0.9850±0.0053	0.0046±0.0011

Table 4. Quantitative evaluation of DUN-SA vs. other methods on the IXI dataset for 4x and 8x acceleration under random 1D subsampling masks, where PD-weighted images are used as reference modality to assist the reconstruction of T2-weighted images. Best results are emphasized in bold, and the second best are emphasized with an underline.

Methods	Random 4x acceleration			Random 8x acceleration		
	PSNR	SSIM	MAE	PSNR	SSIM	MAE
Zero-filling	26.32±2.04	0.6177±0.0465	0.0296±0.0068	23.91±2.01	0.5371±0.0538	0.0380±0.0086
E2E-Varnet	40.10±2.20	0.9808±0.0055	0.0054±0.0013	33.33±2.10	0.9417±0.0126	0.0107±0.0025
HQS-Unet	41.06±2.23	0.9838±0.0044	0.0049±0.0011	34.23±2.05	0.9486±0.0116	0.0098±0.0022
MD-DUN	43.82±2.30	0.9844±0.0053	0.0038±0.0010	39.81±2.18	0.9699±0.0067	0.0060±0.0015
MM-E2E-Varnet	44.77±2.41	0.9905±0.0034	0.0034±0.0009	40.89±2.16	0.9821±0.0061	0.0052±0.0012
SAN	44.75±2.46	0.9907±0.0034	0.0034±0.0009	41.09±2.26	0.9829±0.0058	0.0050±0.0012
MC-CDic	44.39±2.36	0.9877±0.0037	0.0037±0.0010	40.78±2.23	0.9792±0.0062	0.0052±0.0012
DUN-SA	45.38±2.51	0.9915±0.0032	0.0032±0.0008	41.62±2.29	0.9846±0.0054	0.0047±0.0012

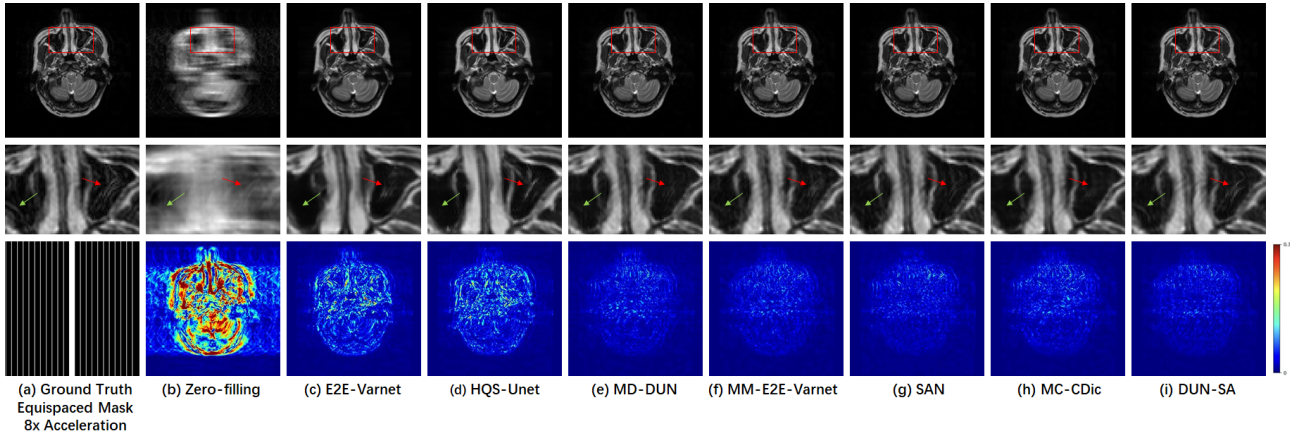


Fig. 5. Visual comparison with representative methods for 8x acceleration under 1D equispaced subsampling mask on the IXI dataset. First row: Reconstructed images by different methods; second row: Zoomed-in region of interest; third row: Equispaced mask of 8x acceleration and error maps of different methods.

other methods, our method reconstructs images with minimal visible artifacts and reconstruction discrepancies compared to the ground truth. Specifically, as depicted in the zoomed-in regions of interest in the second row, the detail indicated by the green arrow is only recovered by DUN-SA and MC-CDic. Moreover, at the location indicated by the red arrow, MC-CDic reconstruct a texture that originally does not exist, introducing a mistake. These findings verify that our method not only shows remarkable enhancements in all metrics but also reconstructs more textural details.

5.2. Results on the IXI dataset

We validate the effectiveness of our method on a larger dataset, IXI. In this part of the experiments, the fully sampled PD-weighted image is used as reference modality image to assist the reconstruction of the paired undersampled T2-weighted image.

Table 3 and Table 4 present the quantitative evaluations on the IXI dataset for 1D equispaced and random masks under 4× and 8× acceleration. From the metrics, it can be observed that the reference modality can significantly enhance the reconstruction performance, and the spatial alignment further improves all metrics, which is consistent with the experimental results on the fastMRI dataset. Our method once again brings noticeable

Table 5. Quantitative evaluation of DUN-SA vs. other methods on the in-house dataset for 4x and 8x acceleration under equispaced 1D subsampling masks, where T1-weighted images are used as reference modality to assist the reconstruction of T2-weighted images. Best results are emphasized in bold, and the second best are emphasized with an underline.

Methods	Equispaced 4x acceleration			Equispaced 8x acceleration		
	PSNR	SSIM	MAE	PSNR	SSIM	MAE
Zero-filling	26.63±0.67	0.6864±0.0153	0.0272±0.0026	24.32±0.66	0.5962±0.0152	0.0362±0.0033
E2E-Varnet	35.42±0.77	0.9643±0.0031	0.0083±0.0009	33.90±0.79	0.9526±0.0044	0.0100±0.0011
HQS-Unet	35.62±0.79	0.9651±0.0032	0.0082±0.0009	34.13±0.77	0.9535±0.0043	0.0098±0.0011
MD-DUN	36.19±1.24	0.9587±0.0059	0.0079±0.0011	34.68±1.33	0.9451±0.0090	0.0092±0.0017
MM-E2E-Varnet	36.17±1.24	0.9684±0.0062	0.0078±0.0012	34.88±1.38	0.9600±0.0092	0.0090±0.0017
SAN	36.58±1.25	0.9705±0.0055	0.0075±0.0012	35.17±1.71	0.9613±0.0112	0.0089±0.0019
MC-CDic	36.21±1.28	0.9682±0.0064	0.0078±0.0013	34.55±1.36	0.9574±0.0097	0.0094±0.0018
DUN-SA	37.41±1.08	0.9747±0.0039	0.0069±0.0009	36.31±1.41	0.9687±0.0070	0.0078±0.0014

Table 6. Quantitative evaluation of DUN-SA vs. other methods on the in-house dataset for 4x and 8x acceleration under random 1D subsampling masks, where T1-weighted images are used as reference modality to assist the reconstruction of T2-weighted images. Best results are emphasized in bold, and the second best are emphasized with an underline.

Methods	Random 4x acceleration			Random 8x acceleration		
	PSNR	SSIM	MAE	PSNR	SSIM	MAE
Zero-filling	26.82±0.66	0.7075±0.0144	0.0261±0.0024	24.27±0.66	0.5982±0.0164	0.0359±0.0033
E2E-Varnet	37.57±0.78	0.9743±0.0020	0.0068±0.0007	32.31±0.78	0.9441±0.0044	0.0124±0.0013
HQS-Unet	37.88±0.77	0.9757±0.0020	0.0067±0.0007	32.64±0.76	0.9457±0.0043	0.0120±0.0012
MD-DUN	38.44±1.19	0.9702±0.0035	0.0065±0.0009	33.41±1.43	0.9386±0.0123	0.107±0.0021
MM-E2E-Varnet	38.61±1.18	0.9789±0.0038	0.0064±0.0010	33.46±1.41	0.9491±0.0123	0.0107±0.0020
SAN	38.92±1.23	0.9799±0.0038	0.0061±0.0010	33.74±1.53	0.9511±0.0136	0.0102±0.0023
MC-CDic	38.84±1.20	0.9797±0.0037	0.0062±0.0009	33.20±1.43	0.9484±0.0128	0.0110±0.0022
DUN-SA	40.47±1.09	0.9845±0.0024	0.0052±0.0007	34.84±1.38	0.9596±0.0081	0.0091±0.0017

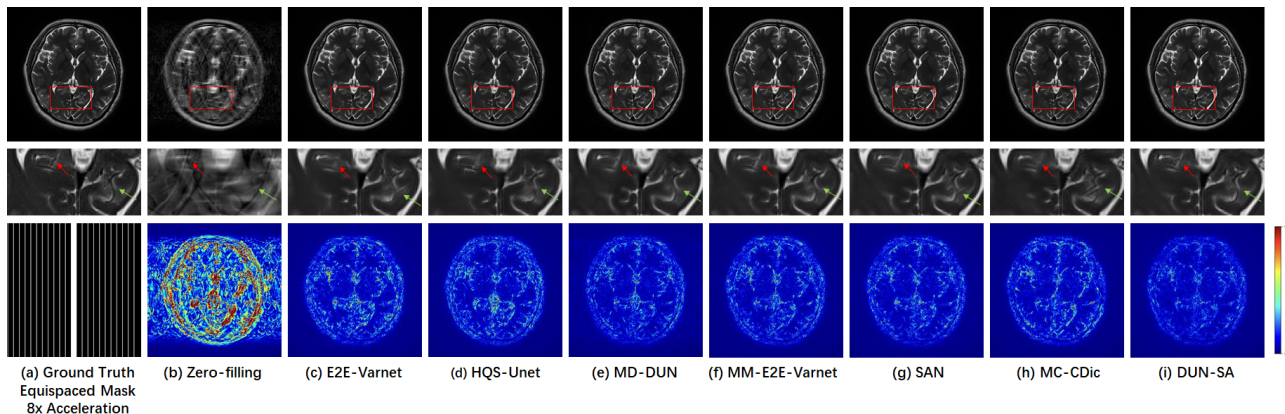


Fig. 6. Visual comparison with representative methods for 8x acceleration under 1D equispaced subsampling mask on the in-house dataset. First row: Reconstructed images by different methods; second row: Zoomed-in region of interest; third row: Equispaced mask of 8x acceleration and error maps of different methods.

improvements for every metric across all settings compared to other methods. For instance, compared to the second-best values, under equispaced 1D subsampling masks of 4x acceleration, DUN-SA improves the PSNR value from SAN's 46.17db to 47.25db. Meanwhile, under random 1D subsampling masks of 8x acceleration, DUN-SA enhances SAN's SSIM value from 0.9829 to 0.9846 and reduces the MAE value from 0.0050 to 0.0047.

Fig. 5 illustrates a qualitative comparison of the reconstructed images of different methods and their corresponding error maps on the IXI dataset for 8x acceleration under equispaced 1D subsampling mask. The content displayed is consistent with Fig. 4. Specifically, from the error maps, it's

evident that our method yields the minimal error. Observing the zoomed-in region of interest, in the area indicated by the green arrow, only DUN-SA successfully reconstructs the relevant texture details. In contrast, in the region pointed out by the red arrow, other methods either introduce excessive noise or fail to reconstruct the corresponding texture, but the proposed DUN-SA can reconstruct relatively clear details without being significantly affected by noise. These results suggest that our approach is capable of producing high-quality reconstruction results that present tissues with fewer artifacts and noise compared to other methods.

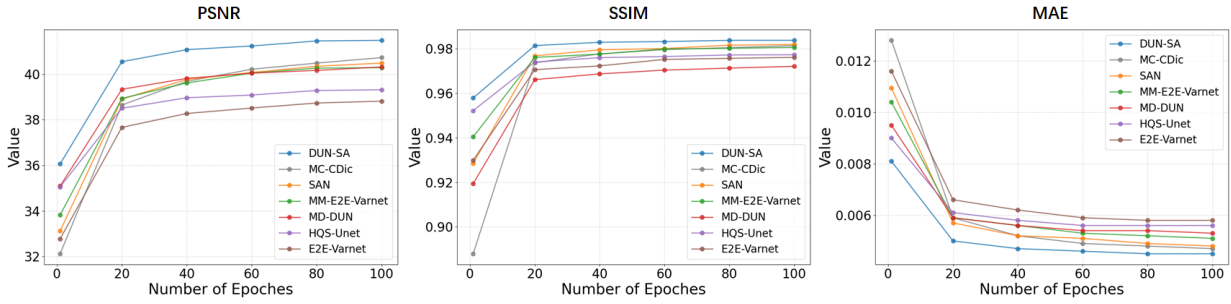


Fig. 7. Comparison on the learning trajectories of different models on the fastMRI dataset for 4x acceleration under equispaced mask.

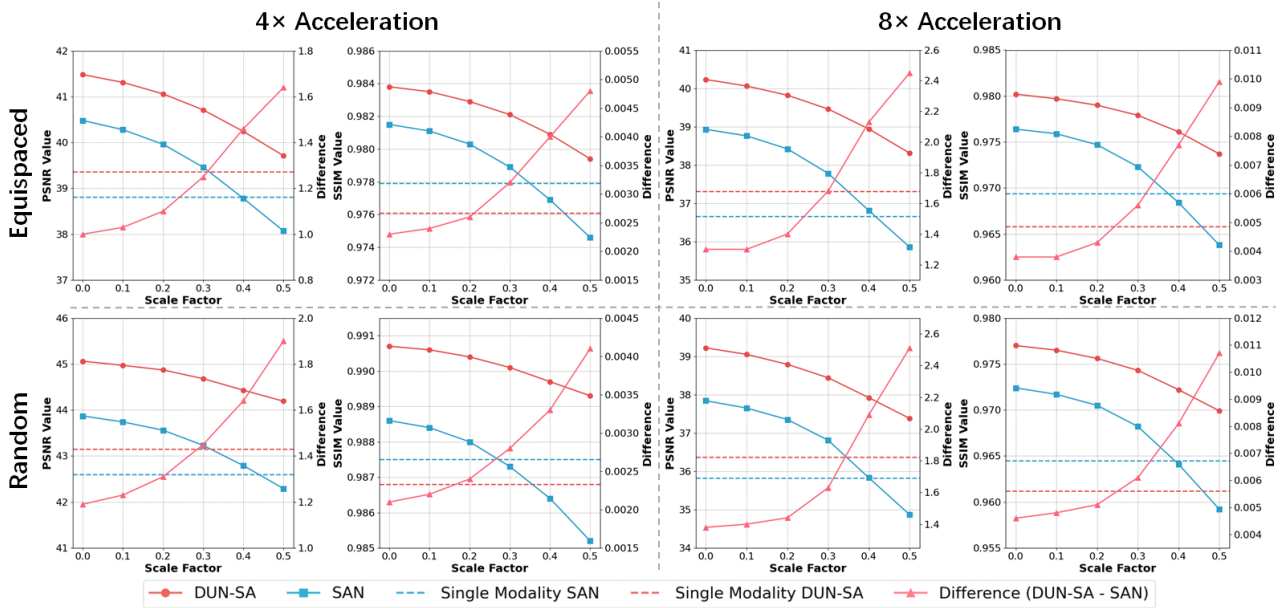


Fig. 8. Quantitative comparison of multi-modal MRI reconstruction on the fastMRI dataset with different scales of simulated spatial misalignment. Left y-axes are for reconstruction performances (“DUN-SA”, “SAN”, “Single Modality SAN” and “Single Modality DUN-SA”) while the right y-axes are for the “Difference” between “DUN-SA” and “SAN”.

5.3. Results on the in-house dataset

In this section, we conduct experiments on the in-house dataset, where the T1-weighted image is used as the reference modality to assist in the reconstruction of the T2-weighted image.

Quantitative results on the in-house dataset for 1D equispaced and random masks under 4x and 8x acceleration are presented in Table 5 and Table 6, showing that the proposed DUN-SA outperforms other methods in achieving the best results for every metric across all settings.

Qualitative results on the in-house dataset for 1D equispaced mask under 8x acceleration are depicted in Fig. 6. Similarly, by comparing the reconstructed images of different methods and their corresponding error maps, it can be observed that the proposed DUN-SA obtains the best results, reconstructing more details. Specifically, in zoomed-in regions of interest in the second row, the details indicated by the red arrow and green row are only reconstructed by the proposed DUN-SA. This is in line with the performance improvement observed in the aforementioned two datasets and indicates the promising generalizability of DUN-SA.

6. Discussion

In this section, we first compare the convergence performance of the proposed DUN-SA with other methods during the training process. Then, we conduct experiments to evaluate the performance under different scales of misalignment. We further compare its spatial alignment performance under different accelerations and conduct ablation studies to validate the benefits of the key components of the DUN-SA at last.

6.1. Convergence of DUN-SA vs. other methods

In this section, we compare the convergence performance of the proposed DUN-SA with other methods by recording and comparing the reconstruction performance of each method throughout the training process. Specifically, we record the PSNR, SSIM, and MAE metrics during the training process on the fastMRI dataset for 4x acceleration under equispaced mask. The curves in Fig. 7 illustrate that the convergence performance of DUN-SA is superior compared to the other methods because it consistently maintains higher PSNR and SSIM values and a lower MAE throughout the entire training process.

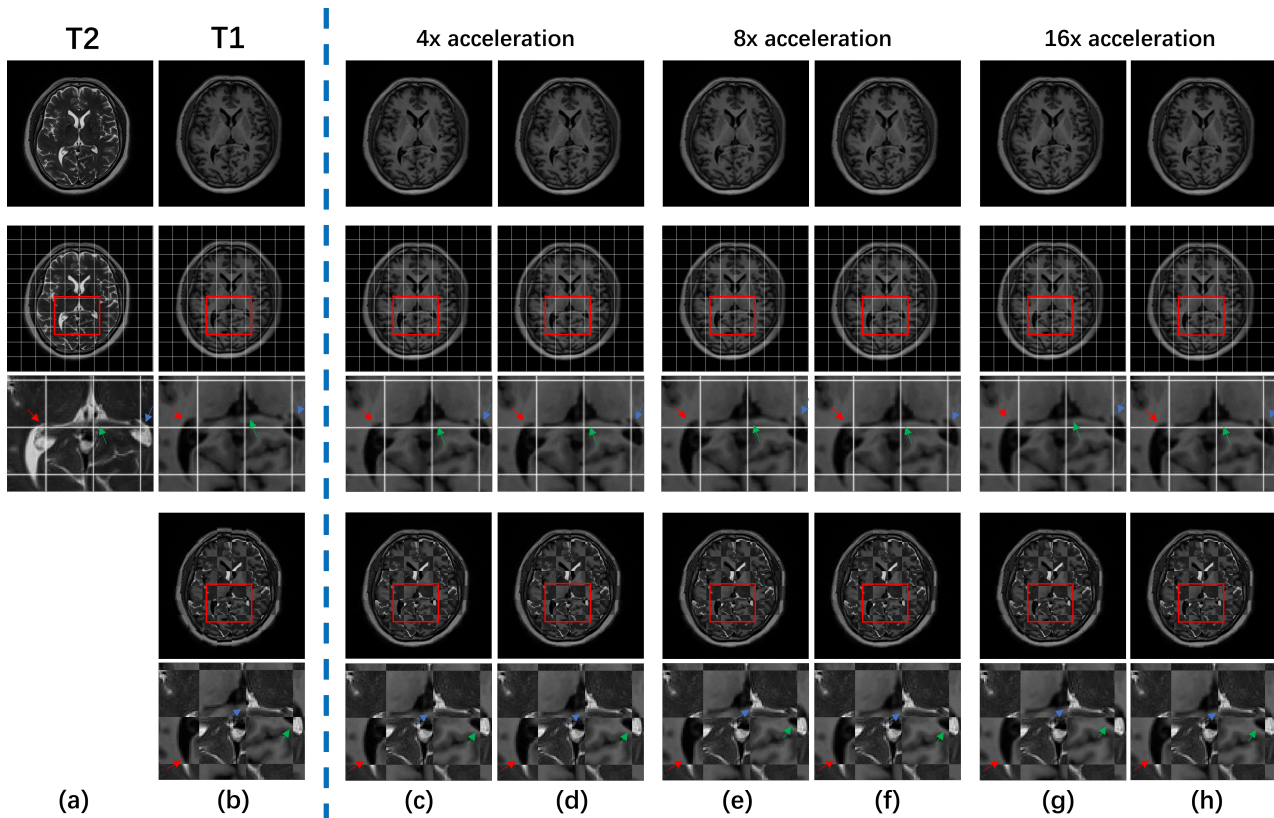


Fig. 9. The visualization of the effects of spatial alignment on the in-house dataset. (a), (b) represent fully-sampled T2-weighted image and fully-sampled T1-weighted image, respectively. (c), (e), and (g) depict T1-weighted images aligned using SAN under different acceleration factors, while (d), (f), and (h) display T1-weighted images aligned using DUN-SA under different acceleration factors. In the second row, a grid is used to facilitate observation of the spatial position of each aforementioned image, with zoomed-in views presented in the third row. In the fourth row, checkerboard visualizations are employed to demonstrate the misalignment between T2-weighted image and T1-weighted image/aligned T1-weighted image, and the last row magnifies the corresponding areas to display the details more clearly.

6.2. Performance evaluation under different scales of misalignment between modalities

To demonstrate the superior robustness of the proposed DUN-SA to misalignment between modalities, we compare it with SAN, another reconstruction method that considers misalignment. We perform different scales of spatial misalignment on the reference modality and then test using the parameters in Section 5.1. This experiment is conducted on the fastMRI dataset for 4x and 8x acceleration under both equispaced and random 1D subsampling masks. The method for simulating spatial misalignment is consistent with the approach in SAN; specifically, we employ random rotations within the range of $[-0.01\pi\sigma, 0.01\pi\sigma]$, translations between $[-0.05N\sigma, 0.05N\sigma]$, and displacement field bicubically interpolated from 9×9 control-points (with displacements uniformly sampled within $[-0.02N\sigma, 0.02N\sigma]$ in both directions. Here, N represents the size of the MR images and σ indicates the factor controlling the degree of spatial misalignment). Fig. 8 compares the reconstruction performance of the two methods under different degrees of spatial misalignment. It's noteworthy that as the scale of spatial misalignment increases, the performance of both methods declines. However, DUN-SA degrades more slowly, whereas the performance of SAN decreases relatively faster, leading to a gradually increasing difference between the two

methods. This highlights the greater robustness of the proposed DUN-SA to spatial misalignment.

6.3. Comparison of spatial alignment under different acceleration factors

In this section, we evaluate the spatial alignment performance of the proposed DUN-SA under different acceleration factors and compare it with SAN. In Fig. 9, (a) and (b) represent fully-sampled T2-weighted image and fully-sampled T1-weighted image, respectively. (c), (e), and (g) depict T1-weighted images aligned using SAN under different acceleration factors, while (d), (f), and (h) display T1-weighted images aligned using DUN-SA under different acceleration factors. We use a grid to facilitate the observation of spatial positions and employ a checkerboard to visualize misalignment. We find that both DUN-SA and SAN exhibit good spatial alignment performance at 4x acceleration. However, as the acceleration factor increases, SAN gradually fails to align T1 and T2 well, whereas DUN-SA is minimally affected. For instance, we can observe in the zoomed-in views of the third and fifth rows, at the locations indicated by the red, blue, and green arrows, misalignment still exists in (e) and (g), but this misalignment has been alleviated in (f) and (h). This demonstrates that DUN-SA, by iteratively solving the proposed model, appropriately integrates the spa-

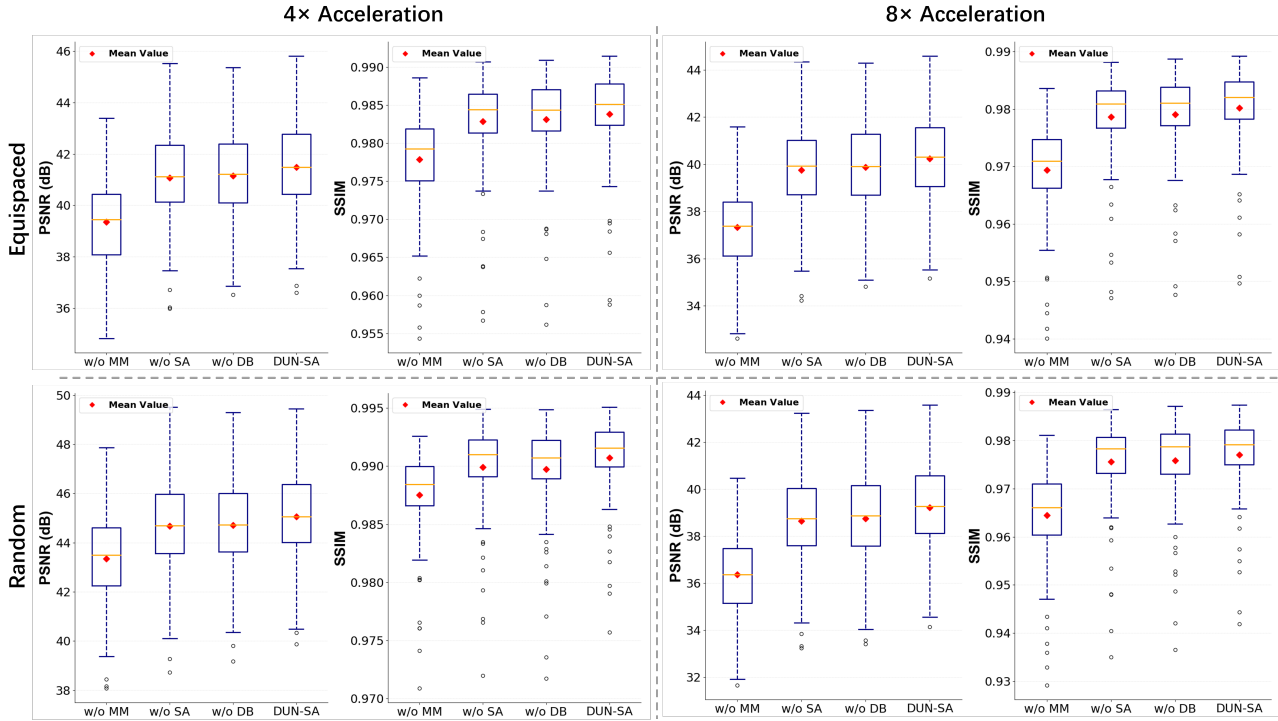


Fig. 10. Quantitative comparison of MRI reconstruction (in PSNR/SSIM) with proposed DUN-SA lacking of different key components on the fastMRI dataset for 4x and 8x acceleration under equispaced and random subsampling masks. In each box plot, the average value is marked with the red diamond symbol.

tial alignment task into the reconstruction process and achieves more precise spatial alignment results under high acceleration factors.

6.4. Ablation Study

To determine the optimal network architecture, we conduct two ablation studies. The first study focus on the number of stages, and the second emphasize each component within the network.

6.4.1. Effect of number of stages

To demonstrate how the number of stages k affects the reconstruction performance, we carry out a quantitative comparison of DUN-SA under different numbers of stages on the fastMRI dataset for 4x and 8x acceleration under equispaced and random subsampling masks. Fig. 11 shows the mean SSIM and PSNR values for stages ranging from 1 to 16. We observed that as the number of stages increases, the quality of reconstruction improves across all settings. However, when $k = 12$, further increasing the stages leads to only subtle improvements in reconstruction performance. Taking both performance and model complexity into consideration, we chose the model with $k = 12$.

6.4.2. Effect of key components

Spatial Alignment Module (SAM): The ablation study for SAM involves removing the SAM module from the network. In this scenario, the spatial alignment operator is ignored within the network. This configuration is termed the proposed DUN-SA without spatial alignment (w/o SA).

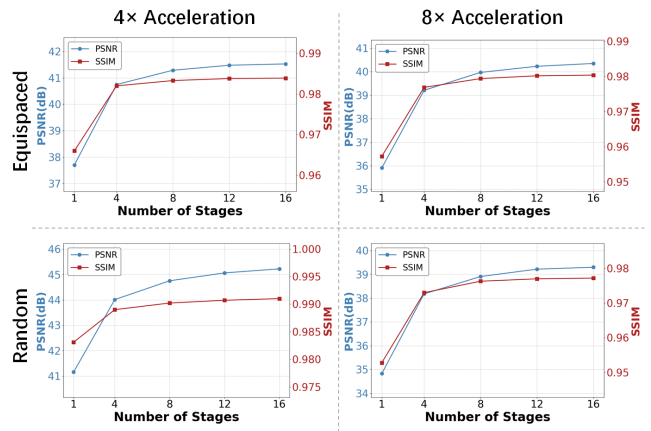


Fig. 11. The PSNR and SSIM curves on the fastMRI dataset with different numbers of stages k

Spatial Alignment Augmented Multi-Modal Learning Block (SAAMMLB): The ablation study for SAAMMLB means removing the SAAMMLB module from the network. It's worth noting that, when the SAAMMLB module is removed, the assistance from the reference modality will also be omitted, and the spatial alignment operator will be simultaneously disregarded. This configuration is termed the proposed DUN-SA without multi-modality (w/o MM).

Denosing Block (DB): The ablation study for DB refers to the removal of the DB module from the network. This configuration is termed the proposed DUN-SA without denoising block (w/o DB).

Table 7. Effect of each component on the performance of DUN-SA on the fastMRI dataset for 4x and 8x acceleration under equispaced and random subsampling masks, measured in PSNR and SSIM.

Methods	Equispaced 4x	Equispaced 8x	Random 4x	Random 8x
Zero-filling	26.92/0.7321	24.36/0.6428	27.12/0.7373	23.88/0.5818
w/o MM	39.36/0.9779	37.32/0.9694	43.35/0.9875	36.37/0.9645
w/o SA	41.07/0.9829	39.75/0.9787	44.68/0.9899	38.65/0.9755
w/o DB	41.16/0.9831	39.86/0.9791	44.70/0.9900	38.75/0.9758
DUN-SA	41.48/0.9838	40.23/0.9802	45.06/0.9907	39.23/0.9770

Firstly, Fig. 10 and Table 7 quantitatively demonstrate the superiority of our model. In Fig. 10, we use a box plot to depict the performance of each method. The orange horizontal line represents the median, and the red diamond represents the mean. From left to right, they are w/o MM, w/o SA, w/o DB and our proposed model (DUN-SA). It can be seen from the Fig. 10 and Table 7 that the reconstruction results on the fastMRI dataset using a single modality (w/o MM) are the poorest across all settings; the reconstruction performance decreases without the spatial alignment, indicating that after aligning the reference modality image, it can provide richer inter-modality prior to assist in the reconstruction task. Similarly, the absence of the denoising block (w/o DB) results in a noticeable decline in performance. Overall, each key component can help the model enhance reconstruction performance, and the optimal model can be obtained by combining all the modules.

7. Conclusion

In this paper, we propose a novel joint alignment and reconstruction model for multi-modal MRI reconstruction. By developing a cross-modal spatial alignment term, we integrate the spatial alignment task into the reconstruction process. We design an optimization algorithm for solving it and then unfold each iterative step into the corresponding network module. As a result, we have constructed a deep unfolding network with interpretability, termed DUN-SA. Through end-to-end training, we fully leverage both intra-modality and inter-modality priors. Comprehensive experiments conducted on three real datasets have demonstrated that the proposed DUN-SA outperforms current state-of-the-art methods in both quantitative and qualitative assessments. Additionally, we have verified that DUN-SA is relatively robust to misalignment, with minimal impact on spatial alignment even as acceleration factors increase.

Acknowledgments

This research is supported by the National Key R & D Program of China (No. 2021YFA1003004) and the National Natural Science Foundation of China (No. 11971296).

References

Aggarwal, H.K., Mani, M.P., Jacob, M., 2019. MoDL: Model-based deep learning architecture for inverse problems. *IEEE Transactions on Medical Imaging* 38, 394–405.

- Bian, W., Zhang, Q., Ye, X., Chen, Y., 2022. A learnable variational model for joint multimodal MRI reconstruction and synthesis, in: *International Conference on Medical Image Computing and Computer-Assisted Intervention*, Springer. pp. 354–364.
- Dar, S.U., Yurt, M., Shahdloo, M., Ildiz, M.E., Tinaz, B., Çukur, T., 2020. Prior-guided image reconstruction for accelerated multi-contrast MRI via generative adversarial networks. *IEEE Journal of Selected Topics in Signal Processing* 14, 1072–1087.
- Defazio, A., Murrell, T., Recht, M., 2020. MRI banding removal via adversarial training, in: *Advances in Neural Information Processing Systems*, pp. 7660–7670.
- Ehrhardt, M.J., Betcke, M.M., 2016. Multicontrast MRI reconstruction with structure-guided total variation. *SIAM Journal on Imaging Sciences* 9, 1084–1106.
- Eksioglu, E.M., 2016. Decoupled algorithm for MRI reconstruction using non-local block matching model: BM3D-MRI. *Journal of Mathematical Imaging and Vision* 56, 430–440.
- EO, T., Jun, Y., Kim, T., Jang, J., Lee, H.J., Hwang, D., 2018. KIKI-net: cross-domain convolutional neural networks for reconstructing undersampled magnetic resonance images. *Magnetic Resonance in Medicine* 80, 2188–2201.
- Feng, C.M., Yan, Y., Chen, G., Xu, Y., Hu, Y., Shao, L., Fu, H., 2023. Multimodal transformer for accelerated MR imaging. *IEEE Transactions on Medical Imaging* 42, 2804–2816.
- Geman, D., Reynolds, G., 1992. Constrained restoration and the recovery of discontinuities. *IEEE Transactions on Pattern Analysis & Machine Intelligence* 14, 367–383.
- He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning for image recognition, in: *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778.
- He, R., Zheng, W.S., Tan, T., Sun, Z., 2013. Half-quadratic-based iterative minimization for robust sparse representation. *IEEE transactions on pattern analysis and machine intelligence* 36, 261–275.
- Huang, P., Zhang, C., Zhang, X., Li, X., Dong, L., Ying, L., 2023. Self-supervised deep unrolled reconstruction using regularization by denoising. *IEEE Transactions on Medical Imaging*, 1–1.
- Huang, W., Ke, Z., Cui, Z.X., Cheng, J., Qiu, Z., Jia, S., Ying, L., Zhu, Y., Liang, D., 2021. Deep low-rank plus sparse network for dynamic MR imaging. *Medical Image Analysis* 73, 102190.
- Jiang, J., Chen, J., Xu, H., Feng, Y., Zheng, J., 2023a. GA-HQS: MRI reconstruction via a generically accelerated unfolding approach, in: *2023 IEEE International Conference on Multimedia and Expo (ICME)*, IEEE. pp. 186–191.
- Jiang, J., Feng, Y., Chen, J., Guo, D., Zheng, J., 2023b. Latent-space unfolding for MRI reconstruction, in: *Proceedings of the 31st ACM International Conference on Multimedia*, pp. 1294–1302.
- Lai, Z., Qu, X., Liu, Y., Guo, D., Ye, J., Zhan, Z., Chen, Z., 2016. Image reconstruction of compressed sensing MRI using graph-based redundant wavelet transform. *Medical Image Analysis* 27, 93–104.
- Lai, Z., Qu, X., Lu, H., Peng, X., Guo, D., Yang, Y., Guo, G., Chen, Z., 2017. Sparse MRI reconstruction using multi-contrast image guided graph representation. *Magnetic Resonance Imaging* 43, 95–104.
- Lei, P., Fang, F., Zhang, G., Xu, M., 2023. Deep unfolding convolutional dictionary model for multi-contrast MRI super-resolution and reconstruction, in: Elkind, E. (Ed.), *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence, IJCAI-23, International Joint Conferences on Artificial Intelligence Organization*. pp. 1008–1016. Main Track.
- Liu, X., Wang, J., Jin, J., Li, M., Tang, F., Crozier, S., Liu, F., 2021. Deep unregistered multi-contrast MRI reconstruction. *Magnetic Resonance Imaging* 81, 33–41.

- Lustig, M., Donoho, D., Pauly, J.M., 2007. Sparse MRI: The application of compressed sensing for rapid MR imaging. *Magnetic Resonance in Medicine: An Official Journal of the International Society for Magnetic Resonance in Medicine* 58, 1182–1195.
- Qu, X., Hou, Y., Lam, F., Guo, D., Zhong, J., Chen, Z., 2014. Magnetic resonance image reconstruction from undersampled measurements using a patch-based nonlocal operator. *Medical Image Analysis* 18, 843–856.
- Ran, M., Xia, W., Huang, Y., Lu, Z., Bao, P., Liu, Y., Sun, H., Zhou, J., Zhang, Y., 2020. MD-Recon-Net: A parallel dual-domain convolutional neural network for compressed sensing MRI. *IEEE Transactions on Radiation and Plasma Medical Sciences* 5, 120–135.
- Ravishanker, S., Bresler, Y., 2010. MR image reconstruction from highly undersampled k-space data by dictionary learning. *IEEE Transactions on Medical Imaging* 30, 1028–1041.
- Ronneberger, O., Fischer, P., Brox, T., 2015. U-Net: Convolutional networks for biomedical image segmentation, in: Navab, N., Hornegger, J., Wells, W.M., Frangi, A.F. (Eds.), *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*, Springer International Publishing, Cham. pp. 234–241.
- Shaul, R., David, I., Shitrit, O., Raviv, T.R., 2020. Subsampled brain MRI reconstruction by generative adversarial neural networks. *Medical Image Analysis* 65, 101747.
- Song, P., Weizman, L., Mota, J.F.C., Eldar, Y.C., Rodrigues, M.R.D., 2020. Coupled dictionary learning for multi-contrast MRI reconstruction. *IEEE Transactions on Medical Imaging* 39, 621–633.
- Sriram, A., Zbontar, J., Murrell, T., Defazio, A., Zitnick, C.L., Yakubova, N., Knoll, F., Johnson, P., 2020. End-to-end variational networks for accelerated MRI reconstruction, in: Martel, A.L., Abolmaesumi, P., Stoyanov, D., Mateus, D., Zuluaga, M.A., Zhou, S.K., Racoceanu, D., Joskowicz, L. (Eds.), *Medical Image Computing and Computer Assisted Intervention – MICCAI 2020*, Springer International Publishing, Cham. pp. 64–73.
- Sun, L., Fan, Z., Fu, X., Huang, Y., Ding, X., Paisley, J., 2019. A deep information sharing network for multi-contrast compressed sensing MRI reconstruction. *IEEE Transactions on Image Processing* 28, 6141–6153.
- Wang, H., Li, Y., Zhang, H., Chen, J., Ma, K., Meng, D., Zheng, Y., 2021. InDuDoNet: An interpretable dual domain network for CT metal artifact reduction, in: de Bruijne, M., Cattin, P.C., Cotin, S., Padoy, N., Speidel, S., Zheng, Y., Essert, C. (Eds.), *Medical Image Computing and Computer Assisted Intervention – MICCAI 2021*, Springer International Publishing, Cham. pp. 107–118.
- Wang, H., Li, Y., Zhang, H., Meng, D., Zheng, Y., 2023. InDuDoNet+: A deep unfolding dual domain network for metal artifact reduction in CT images. *Medical Image Analysis* 85, 102729.
- Wang, S., Su, Z., Ying, L., Peng, X., Zhu, S., Liang, F., Feng, D., Liang, D., 2016. Accelerating magnetic resonance imaging via deep learning, in: 2016 IEEE 13th international symposium on biomedical imaging (ISBI), IEEE. pp. 514–517.
- Wang, Z., Jiang, H., Du, H., Xu, J., Qiu, B., 2020. IKWI-net: A cross-domain convolutional neural network for undersampled magnetic resonance image reconstruction. *Magnetic Resonance Imaging* 73, 1–10.
- Weizman, L., Eldar, Y.C., Ben Bashat, D., 2016. Reference-based MRI. *Medical Physics* 43, 5357–5369.
- Xiang, L., Chen, Y., Chang, W., Zhan, Y., Lin, W., Wang, Q., Shen, D., 2018. Ultra-fast T2-Weighted MR reconstruction using complementary T1-Weighted information, in: Frangi, A.F., Schnabel, J.A., Davatzikos, C., Alberola-López, C., Fichtinger, G. (Eds.), *Medical Image Computing and Computer Assisted Intervention – MICCAI 2018*, Springer International Publishing, Cham. pp. 215–223.
- Xiang, L., Chen, Y., Chang, W., Zhan, Y., Lin, W., Wang, Q., Shen, D., 2019. Deep-learning-based multi-modal fusion for fast MR reconstruction. *IEEE Transactions on Biomedical Engineering* 66, 2105–2114.
- Xin, B., Phan, T., Axel, L., Metaxas, D., 2022. Learned half-quadratic splitting network for MR image reconstruction, in: Konukoglu, E., Menze, B., Venkataraman, A., Baumgartner, C., Dou, Q., Albarqouni, S. (Eds.), *Proceedings of The 5th International Conference on Medical Imaging with Deep Learning*, PMLR. pp. 1403–1412.
- Xuan, K., Xiang, L., Huang, X., Zhang, L., Liao, S., Shen, D., Wang, Q., 2022. Multimodal MRI reconstruction assisted with spatial alignment network. *IEEE Transactions on Medical Imaging* 41, 2499–2509.
- Yang, G., Yu, S., Dong, H., Slabaugh, G., Dragotti, P.L., Ye, X., Liu, F., Arridge, S., Keegan, J., Guo, Y., et al., 2017. DAGAN: Deep de-aliasing generative adversarial networks for fast compressed sensing MRI reconstruction. *IEEE Transactions on Medical Imaging* 37, 1310–1321.
- Yang, G., Zhang, L., Liu, A., Fu, X., Chen, X., Wang, R., 2023. MGDUN: An interpretable network for multi-contrast MRI image super-resolution reconstruction. *Computers in Biology and Medicine* 167, 107605.
- Yang, G., Zhang, L., Zhou, M., Liu, A., Chen, X., Xiong, Z., Wu, F., 2022. Model-guided multi-contrast deep unfolding network for MRI super-resolution reconstruction, in: *Proceedings of the 30th ACM International Conference on Multimedia*, Association for Computing Machinery. p. 3974–3982.
- Yang, Y., Liu, F., Xu, W., Crozier, S., 2015. Compressed sensing MRI via two-stage reconstruction. *IEEE Transactions on Biomedical Engineering* 62, 110–118.
- Yang, Y., Sun, J., Li, H., Xu, Z., 2020a. ADMM-CSNet: A deep learning approach for image compressive sensing. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 42, 521–538.
- Yang, Y., Wang, N., Yang, H., Sun, J., Xu, Z., 2020b. Model-driven deep attention network for ultra-fast compressive sensing MRI guided by cross-contrast MR image, in: Martel, A.L., Abolmaesumi, P., Stoyanov, D., Mateus, D., Zuluaga, M.A., Zhou, S.K., Racoceanu, D., Joskowicz, L. (Eds.), *Medical Image Computing and Computer Assisted Intervention – MICCAI 2020*, Springer International Publishing, Cham. pp. 188–198.
- Zhan, Z., Cai, J.F., Guo, D., Liu, Y., Chen, Z., Qu, X., 2015. Fast multiclass dictionaries learning with geometrical directions in MRI reconstruction. *IEEE Transactions on Biomedical Engineering* 63, 1850–1861.
- Zhang, J., Zhang, Z., Xie, J., Zhang, Y., 2022. High-throughput deep unfolding network for compressive sensing MRI. *IEEE Journal of Selected Topics in Signal Processing* 16, 750–761.
- Zhang, Y., Yap, P.T., Qu, L., Cheng, J.Z., Shen, D., 2019. Dual-domain convolutional neural networks for improving structural information in 3 T MRI. *Magnetic Resonance Imaging* 64, 90–100.
- Zhou, B., Zhou, S.K., 2020. DuDoRNet: Learning a dual-domain recurrent network for fast MRI reconstruction with deep t1 prior, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 4272–4281.
- Zhu, B., Liu, J.Z., Cauley, S.F., Rosen, B.R., Rosen, M.S., 2018. Image reconstruction by domain-transform manifold learning. *Nature* 555, 487–492.