# Self-Supervised Deep Equilibrium Models for Inverse Problems with Theoretical Guarantees

Weijie Gan[1], Chunwei Ying[3], Parna Eshraghi[2], Tongyao Wang[2], Cihat Eldeniz[3], Yuyang Hu[4], Jiaming Liu[4], Yasheng Chen[5], Hongyu An[2,3,4,5], Ulugbek S. Kamilov[1,4]

[1]Department of Computer Science & Engineering, Washington University in St. Louis, St. Louis

[2]Department of Biomedical Engineering, Washington University in St. Louis, St. Louis

[3]Mallinckrodt Institute of Radiology, Washington University in St. Louis, St. Louis

[4]Department of Electrical & System Engineering, Washington University in St. Louis, St. Louis

[5]Department of Neurology, Washington University in St. Louis, St. Louis

{weijie.gan,chunwei.ying,p.eshrag,tongyaow,cihat.eldeniz,h.yuyang,
jiaming.liu,yasheng.chen,hongyuan,kamilov}@wustl.edu

## Abstract

Deep equilibrium models (DEQ) have emerged as a powerful alternative to deep unfolding (DU) for image reconstruction. DEQ models—*implicit neural networks* with effectively infinite number of layers—were shown to achieve state-of-the-art image reconstruction without the memory complexity associated with DU. While the performance of DEQ has been widely investigated, the existing work has primarily focused on the settings where groundtruth data is available for training. We present *self-supervised deep equilibrium model (SelfDEQ)* as the first self-supervised reconstruction framework for training model-based implicit networks from undersampled and noisy MRI measurements. Our theoretical results show that SelfDEQ can compensate for unbalanced sampling across multiple acquisitions and match the performance of fully supervised DEQ. Our numerical results on *in-vivo* MRI data show that SelfDEQ leads to state-of-the-art performance using only undersampled and noisy training data.

## 1 Introduction

We consider an *inverse problem* where one seeks to recover an unknown image $x \in \mathbb{C}^n$ from its undersampled and noisy measurements $y \in \mathbb{C}^m$. Inverse problems are ubiquitous across medical imaging, bio-microscopy, and computational photography. In particular, *compressed sensing magnetic resonance imaging (CS-MRI)* is a well known inverse problem that aims to recover diagnostic quality images from undersampled and noisy $k$-space measurements [26]. *Deep learning (DL)* has recently gained popularity in inverse problems due to its state-of-the-art performance [25,29]. Traditional DL methods train *convolutional neural networks (CNNs)* to map acquired measurements to the desired images [17,43]. Recent work has shown that *deep unfolding (DU)* can perform better than generic CNNs by accounting for the physics of the imaging system [1,35]. DU models are often obtained from optimization methods by interpreting a *fixed number* of iterations as layers of a deep architecture and training it end-to-end. Despite the empirical success of DU in some applications, the high memory complexity of training DU models limits its use in large-scale imaging applications (e.g., 3D/4D MRI).

Recently, *neural ODEs* [7,19] and *deep equilibrium models (DEQ)* [4,9] have emerged as frameworks for training deep models with effectively infinite number of layers without the associated memory cost. The potential of DEQ to address imaging inverse problems was recently shown in [11]. Training a DEQ model for inverse problems is analogous to training an *infinite-depth* DU model with constant memory complexity. However, DEQ is traditionally trained using *supervised learning*, which limits its applicability to problems with no groundtruth training data. While there has been substantial interest in developing *self-supervised learning* methods that use undersampled and noisy measurements for training [2,36,41], the potential of self-supervised learning has never been explored in the context

of DEQ. This work bridges this gap by proposing *self-supervised deep equilibrium model (SelfDEQ)* as a framework for training *implicit neural networks* for MRI without groundtruth data. Our contributions are as follows:

- We introduce SelfDEQ as an image reconstruction framework for CS-MRI based on training a model-based implicit neural network directly on undersampled and noisy measurements. SelfDEQ extends the line of work based on *Noise2Noise (N2N)* [22] by introducing a model-based implicit architecture, a specialized loss function that accounts for unbalanced sampling, and a memory-efficient training method using *Jacobian-Free Backpropagation (JFB)* [9].

- We present new theoretical results showing that for certain measurement operators SelfDEQ computes updates that match those obtained by fully-supervised DEQ. In the context of CS-MRI, our results imply that under a set of explicitly specified assumptions, SelfDEQ can provably match the performance of DEQ trained using the groundtruth MRI images. It is worth highlighting that the theoretical guarantees provided by our analysis leverage the proposed correction for unbalanced sampling.

- We present new numerical results on experimentally-collected *in-vivo* brain MRI data. Our results show that SelfDEQ can (a) outperform recent self-supervised DU methods; (b) match the performance of fully-supervised DEQ, corroborating our theoretical analysis; and (c) enable highly-accelerated data-collection in parallel MRI.

## 2 Background

### 2.1 Imaging Inverse Problems

We consider inverse problems where the measurements $\boldsymbol{y}$ are specified by a linear system

$$\boldsymbol{y} = \boldsymbol{M}\boldsymbol{A}\boldsymbol{x} + \boldsymbol{e}\,, \tag{1}$$

where $\boldsymbol{x}$ is the unknown image, $\boldsymbol{e} \in \mathbb{C}^m$ is *additive white Gaussian noise (AWGN)*, $\boldsymbol{A} \in \mathbb{C}^{n \times n}$ is a measurement matrix, and $\boldsymbol{M} \in \{0, 1\}^{m \times n}$ is a diagonal sampling matrix. A well-known application of (1) is CS-MRI [26], where the measurements correspond to the noisy samples in the Fourier domain (referred to as *k*-space).

Inverse problems are generally ill-posed. Traditional methods recover $\boldsymbol{x}$ by solving a regularized optimization

$$\widehat{\boldsymbol{x}} = \arg\min_{\boldsymbol{x}} f(\boldsymbol{x}) \quad \text{with} \quad f(\boldsymbol{x}) = g(\boldsymbol{x}) + h(\boldsymbol{x}), \tag{2}$$

where $g$ is the data-fidelity term that quantifies the discrepancy between the measurements and the solution, and $h$ is a regularizer that imposes prior knowledge on the unknown image. Well-known examples in the context of imaging inverse problems are the *least-squares* and *total variation (TV)*

$$g(\boldsymbol{x}) = (1/2)\left\|\boldsymbol{y} - \boldsymbol{M}\boldsymbol{A}\boldsymbol{x}\right\|_2^2 \text{ and } h(x) = \tau\left\|\boldsymbol{D}\boldsymbol{x}\right\|_1\,, \tag{3}$$

where $\boldsymbol{D}$ is an image gradient and $\tau > 0$ is the regularization parameter.

### 2.2 Deep Learning

The focus in the area has recently moved to DL (see recent reviews in [25, 29]). A widely-used DL approach is to train a CNN to learn a mapping from the measurements to the corresponding groundtruth images [17, 43]. There is also a growing interest in *deep model-based architectures (DMBAs)* that can combine physical measurement models and learned image priors specified using CNNs. Well known examples of DMBAs are *plug-and-play priors (PnP)* [18, 39], *Regularized by Denoiser (RED)* [32], and *deep unfolding (DU)* [1, 14, 35]. In particular, DU has gained notoriety due to its ability to achieve the state-of-the-art performance, while providing robustness to changes in data acquisition. DU architectures are typically obtained by unfolding iterations of an image reconstruction algorithm as layers, representing the regularizer within image reconstruction as a CNN, and training the resulting network end-to-end. DU architectures, however, are usually limited to a small number of unfolded iterations due to the high memory complexity of training [35].
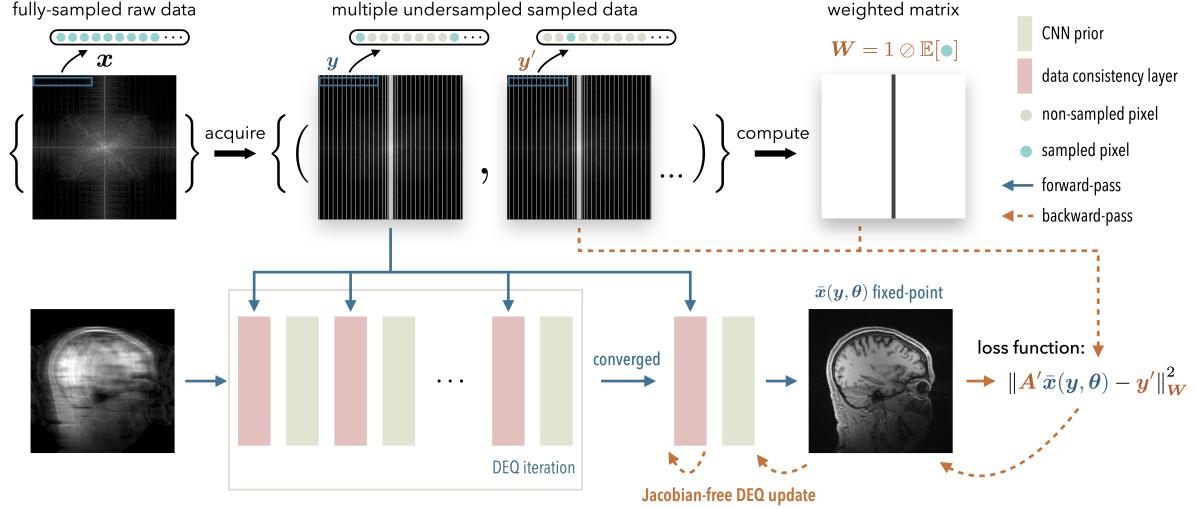
Figure 1: Illustration of SelfDEQ for CS-MRI. The forward pass of SelfDEQ computes a *fixed-point* of an operator consisting of data consistency layer and a CNN prior. The backward pass of SelfDEQ computes a descent direction using the Jacobian-free update that can be used to optimize the training parameters. SelfDEQ is trained using the proposed weighted loss that directly maps pairs of undersampled and noisy measurements of the same object to each other without fully-sampled groundtruth.

## 2.3 Deep Equilibrium Models

DEQ has emerged as a framework for training recursive networks that have *infinitely* many layers without storing intermediate latent variables [4, 9, 11, 24, 30, 31, 42]. It is implemented by running two consecutive steps in each training iteration, namely the *forward* pass and the *backward* pass. The forward pass computes a fixed point $\bar{x}$ of an operator $\mathsf{T}_{\boldsymbol{\theta}}$ parameterized by weights $\boldsymbol{\theta}$

$$\bar{\boldsymbol{x}} = \mathsf{T}_{\boldsymbol{\theta}}(\bar{\boldsymbol{x}}, \boldsymbol{y}) \, , \tag{4}$$

where $\boldsymbol{y}$ is the measurement vector. The fixed point $\bar{x}$ is often computed by running a fixed-point iteration with an acceleration algorithm (such as Anderson acceleration [3]). It is worth noting that, when $\mathsf{T}_{\boldsymbol{\theta}}$ denotes a step of DMBA, the DEQ forward pass is equivalent to DU with infinitely many unfolded layers. Given a loss function, the backward pass produces gradients with respect to $\boldsymbol{\theta}$ by implicitly differentiating through the fixed points without the knowledge of how they are estimated (see Sec. 4.2 in [11] for more details). DEQ does not require storing the intermediate variables for computing the gradient, which dramatically reduces the memory complexity of training. There have been several applications of DEQ in imaging, including applications to MRI [11, 30, 31], computed tomography (CT) [24] and video snapshot imaging [42].

## 2.4 Self-Supervised Deep Image Reconstruction

There is a growing interests in developing DL methods that reduce the dependence on the groundtruth training data (see recent reviews [2, 36, 41]). Some well-known strategies include *Noise2Noise (N2N)* [22], *Noise2Void (N2V)* [21], *deep image prior (DIP)* [38], *Compressive Sensing using Generative Models (CSGM)* [5,13], and equivariant imaging [6]. In particular, N2N is one of the most widely-used self-supervised DL frameworks for image restoration that directly uses noisy observations $\{\widehat{\boldsymbol{x}}_{i,j} = \boldsymbol{x}_i + \boldsymbol{e}_{i,j}\}$ of groundtruth images $\{\boldsymbol{x}_i\}$ for training. The N2N training can be formulated as

$$\arg\min_{\boldsymbol{\theta}} \sum_i \sum_{j \neq j'} \|\mathsf{f}_{\boldsymbol{\theta}}(\widehat{\boldsymbol{x}}_{i,j}) - \widehat{\boldsymbol{x}}_{i,j'}\|_2^2 \, , \tag{5}$$

where $\mathsf{f}_{\boldsymbol{\theta}}$ denotes the DL model with trainable parameters $\boldsymbol{\theta}$. There have been many extensions of N2N to different imaging problems, such as MRI [10, 27, 40], OCT [16], and CT [15]. In particular, SSDU [40] is a recent state-of-the-

art method based on training a DU model without groundtruth by dividing a single k-space MRI acquisition into two subsets that are used as training targets for each other. The work [27] has provided a theoretical justification for SSDU by extending Noisier2Noise [28] to variable-density subsampled MRI data.

## 2.5 Our Contributions

While DEQ has been shown to achieve the state-of-the-art imaging performance, the existing work has focused on settings where groundtruth data is available for training. Our work addresses this gap by enabling DEQ training on noisy and undersampled sensor measurements, which has not been investigated before. The proposed SelfDEQ framework consists of several synergistic elements: (a) a model-based implicit network that integrates measurement operators and CNN priors; (b) a self-supervised loss that accounts for sampling imbalances; (c) a Jacobian-free backward pass that leads to efficient training.

# 3 Method

## 3.1 Weighted Self-Supervised Loss

Consider the training set of measurement pairs $\{\boldsymbol{y}_i, \boldsymbol{y}_i'\}_{i=1}^N$ with each pair $\boldsymbol{y}_i, \boldsymbol{y}_i'$ corresponding to the same object $\boldsymbol{x}_i$

$$\boldsymbol{y}_i = \boldsymbol{M}_i \boldsymbol{A} \boldsymbol{x}_i + \boldsymbol{e}_i \text{ and } \boldsymbol{y}_i' = \boldsymbol{M}_i' \boldsymbol{A} \boldsymbol{x}_i + \boldsymbol{e}_i' . \tag{6}$$

Here, $N \geq 1$ denotes the number training pairs. One can obtain measurement pairs by physically conducting two acquisitions or splitting each acquisition into two subsets.

Existing algorithms based on N2N directly map the measurement pairs to each other during training. However, the measurements in the training dataset often have a significant overlap. For example, each acquisition may share the *auto calibration signal (ACS)* region [37], thus giving more weight to corresponding regions of the k-space (see SelfDEQ *(unweigted)* in Fig. 6). We introduce a diagonal weighted matrix $\overline{\boldsymbol{W}} = \mathsf{diag}(\overline{w_0}, \overline{w_1}, ..., \overline{w_n}) \in \mathbb{R}^{n \times n}$ that accounts for the oversampled regions in the loss function. We set the diagonal entries of $\overline{\boldsymbol{W}}$ as follows

$$\overline{w_k} = \begin{cases} \frac{1}{\sqrt{\mathbb{E}[\boldsymbol{M'}^\mathsf{T}\boldsymbol{M'}]_{k,k}}} & \sqrt{\mathbb{E}[\boldsymbol{M'}^\mathsf{T}\boldsymbol{M'}]_{k,k}} \neq 0 \\ 0 & \sqrt{\mathbb{E}[\boldsymbol{M'}^\mathsf{T}\boldsymbol{M'}]_{k,k}} = 0 \end{cases} , \tag{7}$$

where, in practice, the expectation over random sampling patterns can be replaced with an empirical average over the training set. We can then define the following self-supervised training loss function

$$\ell_{\mathsf{self}}(\boldsymbol{\theta}) = \mathbb{E} \left\| \boldsymbol{M'}\boldsymbol{A'}\bar{\boldsymbol{x}}(\boldsymbol{\theta}) - \boldsymbol{y'} \right\|_{\boldsymbol{W}}^2 , \tag{8}$$

where $\boldsymbol{W} = \boldsymbol{M'}\overline{\boldsymbol{W}}(\boldsymbol{M'}\overline{\boldsymbol{W}})^\mathsf{T} \in \mathbb{R}^{m \times m}$ denotes a subsampled variant of $\overline{\boldsymbol{W}}$ given $\boldsymbol{M'}$, and $\bar{\boldsymbol{x}} = \mathsf{T}_{\boldsymbol{\theta}}(\bar{\boldsymbol{x}}, \boldsymbol{y})$ denotes the fixed-point of $\mathsf{T}_{\boldsymbol{\theta}}$ for the input $\boldsymbol{y}$ and weights $\boldsymbol{\theta}$.

## 3.2 Forward and Backward Passes

The SelfDEQ forward pass is a fixed-point iteration

$$\boldsymbol{x}^k = \mathsf{T}_{\boldsymbol{\theta}}(\boldsymbol{x}^{k-1}, \boldsymbol{y}) , \tag{9}$$

where

$$\begin{aligned} \mathsf{T}_{\boldsymbol{\theta}}(\boldsymbol{x}) &= \alpha \mathsf{f}_{\boldsymbol{\theta}}(\boldsymbol{s}) + (1 - \alpha)\boldsymbol{s} \\ \text{with } \boldsymbol{s} &= \boldsymbol{x} - \gamma \nabla g(\boldsymbol{x}) \end{aligned} \tag{10}$$

The vector $\boldsymbol{x}^k$ denotes the image at the $k$th layer of the implicit network, $\gamma$ and $\alpha$ are two hyper-parameters, and $\mathsf{f}_{\boldsymbol{\theta}}$ is the CNN prior with trainable parameters $\boldsymbol{\theta}$. The implicit neural network is initialized using the pseudoinverse of
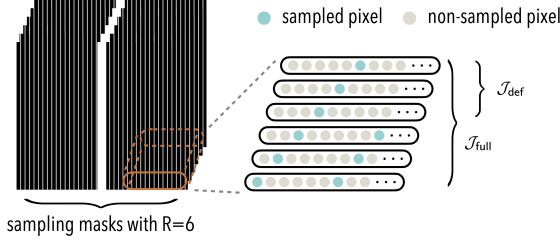
Figure 2: An illustration of sampling masks used for our numerical evelution for the acceleration factor $R = 6$. We consider two settings: full-rank ($\mathcal{J}_{\mathsf{full}}$) and rank-deficient ($\mathcal{J}_{\mathsf{def}}$). In the full-rank setting, the union of sampling masks across training data covers the whole space $\mathbb{C}^n$, thus satisfying Assumption 2. In the rank-deficient setting, some of the frequencies are never sampled in the training dataset. Note that each individual sampling mask is always undersampled.

the raw measurements $\boldsymbol{x}^0 = \boldsymbol{A}^\dagger \boldsymbol{y}$, which corresponds to the zero-filled solution in CS-MRI. Given $\boldsymbol{y}$, we run the forward-pass until convergence to a fixed-point $\bar{\boldsymbol{x}} = \mathsf{T}_{\boldsymbol{\theta}}(\bar{\boldsymbol{x}}, \boldsymbol{y})$. We use U-Net for $\mathsf{f}_{\boldsymbol{\theta}}$ [33].

Given the self-supervised loss in eq. (8), the SelfDEQ backward pass uses a Jacobian-free backward pass (JFB) to compute the DEQ update direction without computing the inverse-Jacobian. Specifically, the JFB update of $\ell_{\mathsf{self}}$ in term of $\boldsymbol{\theta}$ is given by

$$\mathsf{JFB}_{\ell_{\mathsf{self}}}(\boldsymbol{\theta}) = \mathsf{Real}\left(\left(\nabla_{\boldsymbol{\theta}} \mathsf{T}_{\boldsymbol{\theta}}(\bar{\boldsymbol{x}})\right)^{\mathsf{H}} \left[\frac{\partial \ell_{\mathsf{self}}}{\partial \bar{\boldsymbol{x}}}\right]^{\mathsf{T}}\right). \tag{11}$$

JFB was theoretically shown to provide valid descent directions for training implicit networks in [9].

## 3.3 Theoretical Analysis

We now present the theoretical analysis of SelfDEQ learning under two explicitly-specified assumptions.

**Assumption 1** The training samples correspond to the setting in (6) with $\boldsymbol{x} \sim p_{\boldsymbol{x}}$, $\boldsymbol{M} \sim p_{\boldsymbol{M}}$, $\boldsymbol{M}' \sim p_{\boldsymbol{M}}$, $\boldsymbol{e} \sim \mathcal{N}(0, \sigma^2)$ and $\boldsymbol{e}' \sim \mathcal{N}(0, \sigma^2)$ drawn i.i.d. from their respective distributions.

This mild assumption simply states that the sampling matrices, images, and noise are all sampled independently from each other, which is a reasonable assumption.

**Assumption 2** $\mathbb{E}_{\boldsymbol{M}}[\boldsymbol{M}^{\mathsf{T}}\boldsymbol{M}]$ has a full rank and $\boldsymbol{A}$ is an orthogonal matrix, where the expectation is taken over $p_{\boldsymbol{M}}$.

This assumption implies that *union* of all the sampling matrices $\{\boldsymbol{M}\}$ covers the complete measurement domain. Note that each individual $\boldsymbol{M}$ can still be undersampled.

**Theorem 1** *Under Assumptions 1-2, the JFB update of the weighted self-supervised loss* (8) *is equivalent to its supervised counterpart, namely we have that*

$$\mathsf{JFB}_{\ell_{\mathsf{self}}}(\boldsymbol{\theta}) = \mathsf{JFB}_{\ell_{\mathsf{sup}}}(\boldsymbol{\theta}). \tag{12}$$

*where*

$$\ell_{\mathsf{sup}}(\boldsymbol{\theta}) = \mathbb{E}\left[\frac{1}{2}\|\bar{\boldsymbol{x}}(\boldsymbol{\theta}) - \boldsymbol{x}\|_2^2\right]. \tag{13}$$

The proof is provided in the supplementary material. Theorem 1 states that the JFB updates from our weighted self-supervised loss theoretically match those obtained using conventional supervised learning on DEQ. These updates can be easily integrated into any DL optimization algorithm, such as SGD and Adam [20].

# 4 Numerical Validation

We now presents numerical results evaluating SelfDEQ on both simulated and *in-vivo* MRI data. The measurement matrix in parallel MRI can be expressed as $A_i = FS_i$ where $F$ is the Fourier transform operator, and $S_i$ denotes the sensitivity profiles of the $i$th receiver coil. We assume that $S$ is known and normalized to satisfy $\sum_i S_i^H S_i = I$. Since the Fourier transform operator is orthogonal, the matrix $A$ in MRI naturally satisfies Assumption 2. Note that, in order to estimate $S_i$ in practice, $M$ has a *fixed* ACS in the low-frequency region of $k$-space [37]. The random valuables in $p_M$ in our experiments are the randomly sampled non-ACS indices of the $k$-space.

We ran the forward-pass of SelfDEQ with a maximum number of iterations of 100 and the stopping criterion of the relative norm difference between iterations being less than $10^{-3}$. We added spectral normalization to all the layers of $f_\theta$ for stability [34]. We empirically determined the best values of $\alpha$ and $\gamma$ to be $\alpha = 0.5$ and $\gamma = 1$. We used Adam [20] as the optimizer with the learning rate $10^{-4}$. We set the mini-batch size to 8 and training epochs to 100 and 300 for real and simulated data, respectively. We performed all our experiments on a machine equipped with an AMD Ryzen Threadripper 3960X Processor and an NVIDIA GeForce RTX 3090 GPU. We used two widely used quantitative metrics, *peak signal-to-noise ratio (PSNR)* measured in dB and *structural similarity index (SSIM)*, to evaluate the quality of reconstructed images.
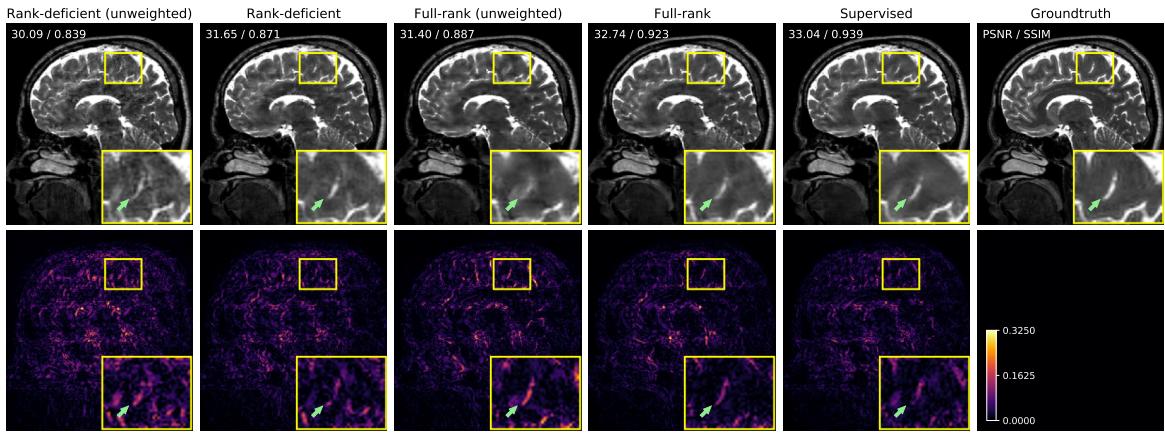


Figure 3: An illustration of reconstructed results obtained from several ablated variants of SelfDEQ using the acceleration factor $R = 6$ on simulated data. We highlight the visually significant differences using green arrows. This figure highlights that in the *Full-rank* setting with the weighting matrix in the loss function, SelfDEQ nearly matches the performance of its supervised counterpart. This figure also shows that the proposed weighting matrix within the self-supervised loss significantly improves the imaging quality, even when Assumption 2 is not satisfied (compare *Rank-deficient (unweighted)* and *Rank-deficient*).
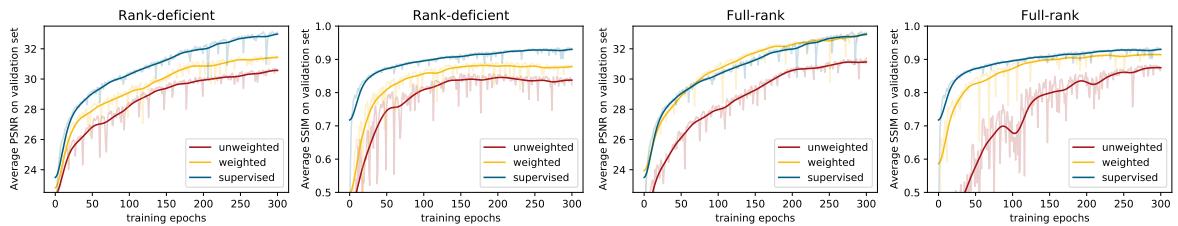


Figure 4: PSNR/SSIM plotted against the training epochs using the simulated validation dataset. This figure highlights that performance of SelfDEQ in the *Full-rank* setting with a weighting matrix closely tracks that of the *supervised* DEQ. The figure also shows the benefit of using the weighting matrix for self-supervised training in both *Rank-deficient* and *Full-rank* scenarios.

## 4.1 Ablation Study on Simulated Data

### 4.1.1 Dataset

We simulated multi-coil undersampled measurements from an open-access T2-weighted human brain MRI data[1], which was collected by [1]. This MRI dataset has 360 and 160 slices of fully-sampled $k$-space measurements for training and testing, respectively. We extracted 60 slices from the training set for validation. The image domain matrix size for each slice is $256 \times 232$. The number of receiver coils is 12. The coil sensitivity maps for each slice are also provided, which are pre-computed using ESPIRiT algorithm [37]. These fully-sampled data correspond to the groundtruth in the imaging system (*i.e.,* $x_i$ in (6)). We simulated a Cartesian sampling pattern that subsamples and fully samples along $k_y$ and $k_x$ dimensions, respectively. The simulated sampling mask in the $ky$ dimension has fixed ACS lines and equispaced non-ACS lines. Let $R$ be the acceleration factor. We set the size of ACS lines to $92/R$. We conducted experiments on three acceleration factors $R = 4, 6$, and 8, corresponding to 31%, 21% and 16% sampling rate, respectively. As illustrated in Fig. 2, given a acceleration factor $R$, one can simulate $R$ different sampling masks. Let $\mathcal{J}$ denote a *subset* of those simulated masks. When simulating the measurements, we sampled sampling masks uniformly at random from $\mathcal{J}$ (*i.e.,* $M_i$ and $M_i'$ in (6)). We set the standard deviation of the AWGN (*i.e.,* $e_i$ and $e_i'$ in (6)) to 0.01.

### 4.1.2 Results

We consider two different sampling settings: full-rank ($\mathcal{J}_{\mathsf{full}}$) and rank-deficient ($\mathcal{J}_{\mathsf{def}}$). In the full-rank setting the sampling masks across training data cover all possible frequencies, while in the rank-deficient setting the union of all sampling masks only covers *half* of the k-space. Fig. 2 visually illustrates both settings $\mathcal{J}_{\mathsf{full}}$ and $\mathcal{J}_{\mathsf{def}}$ for the acceleration factor $R = 6$. Note that the sampling masks selected from $\mathcal{J}_{\mathsf{full}}$ naturally satisfy Assumption 2, since $\mathbb{E}_{M \in \mathcal{J}_{\mathsf{full}}}[M^\mathsf{T} M]$ has full rank. On the other hand, $\mathbb{E}_{M \in \mathcal{J}_{\mathsf{def}}}[M^\mathsf{T} M]$ does not satisfy Assumption 2. Under these two sampling settings, we ran the following experiments using ablated variants of SelfDEQ: (a) *Rank-deficient* trains SelfDEQ on $\mathcal{J}_{\mathsf{def}}$; (b) *Rank-deficient (unweighted)* is similar to *Rank-deficient*, but uses the self-supervised loss without $W$; (c) *Full-rank* trains SelfDEQ on $\mathcal{J}_{\mathsf{full}}$; (d) *Full-rank (unweighted)* is similar to *Full-rank*, but uses the self-supervised loss without $W$; (e) *Supervised* is similar to *Full-rank* but uses the supervised loss in (13), corresponding to the oracle DEQ performance.

| Metrics | PSNR (dB) | | | SSIM | | |
|---|---|---|---|---|---|---|
| Acceleration rate | ×8 | ×6 | ×4 | ×8 | ×6 | ×4 |
| *Zero-Filled* | 22.27 | 23.40 | 26.59 | 0.771 | 0.798 | 0.861 |
| *Rank-deficient*[w/o] | 29.58 | 31.86 | 38.54 | 0.843 | 0.856 | 0.905 |
| *Rank-deficient* | 30.31 | 33.45 | 39.91 | 0.880 | 0.894 | 0.936 |
| *Full-rank*[w/o] | 29.95 | 32.73 | 38.61 | 0.861 | 0.892 | 0.922 |
| *Full-rank* | **31.28** | **34.14** | **41.56** | **0.903** | **0.931** | **0.951** |
| *Supervised* | 32.35 | 34.68 | 42.02 | 0.925 | 0.955 | 0.981 |

*Rank-deficient*[w/o]: *Rank-deficient (unweighted)*;
*Full-rank*[w/o]: *Full-rank (unweighted)*.

Table 1: Average PSNR and SSIM values from the ablation study on simulated data. *Full-rank* satisfies assumptions used for the theoretical analysis and is trained using weighted self-supervised loss. This table validates our theoretical analysis and highlights the importance of the weighting for the self-supervised loss function in both rank-deficient and full-rank settings.

---

[1]This dataset is available at https://drive.google.com/file/d/1qp-l9kJbRfQU1W5wCjOQZi7I3T6jwA37/view?usp=sharing
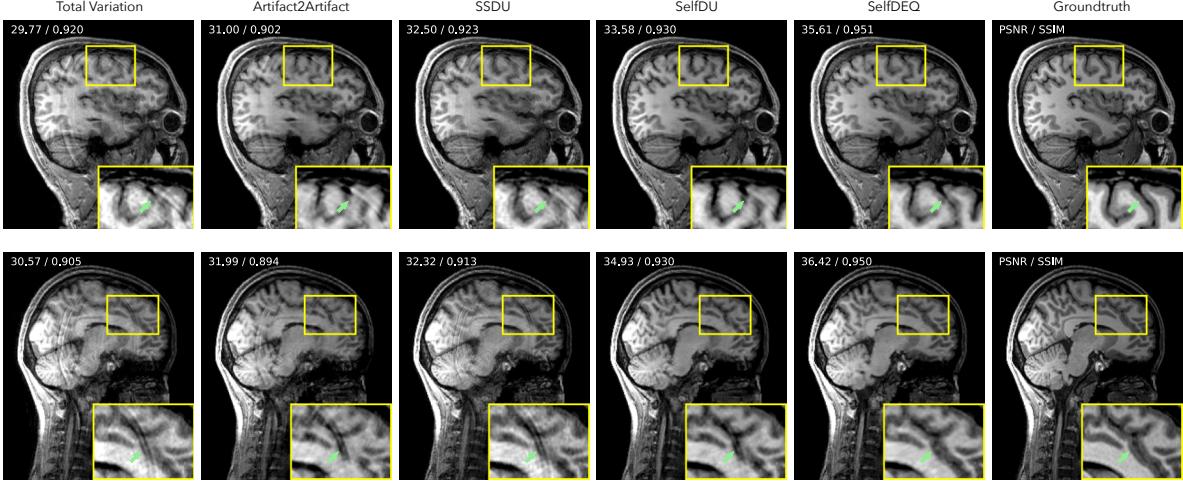
Figure 5: Reconstruction results on in-vivo data comparing several reconstruction algorithms on two different slices under the acceleration factor $R = 6$. Note that all the DL methods in the figure are based on self-supervised learning. We highlight visually significant differences with green arrows. Note the regions where SelfDEQ reconstructs images with fine details, but other methods result in ghosting artifacts. This figure highlights that SelfDEQ can achieve superior performance compared to recent self-supervised DU methods on experimentally-acquired parallel MRI data.

### 4.1.3 Discussion

Fig. 3 illustrates the reconstruction results of all the ablated methods on $R = 6$. Fig. 3 shows that when Assumption 2 is satisfied and the proposed weighting scheme is used, the performance of self-supervised learning nearly matches that of fully-supervised learning. Fig. 3 also highlights that using weighted self-supervised loss improves the imaging quality even when Assumption 2 is not satisfied (i.e., the union of all the sampling masks doesn't cover the full k-space). For instance, note how the brain tissue highlighted using a green arrow is blurry for *Full-rank (unweighted)*, while *Full-rank* can reconstruct it with fine details. Also note that while the settings $\mathcal{J}_{\mathsf{full}}$ provide better reconstruction performances compared to those of $\mathcal{J}_{\mathsf{def}}$ under the same losses, *Rank-deficient* outperforms *Full-rank (unweighted)*, highlighting the effectiveness of the weighted matrix $\boldsymbol{W}$ for self-supervised learning. Table 1 summarizes PSNR/SSIM values of ablation methods on the testing dataset, thus quantitatively corroborating the visual results.

Fig. 4 plots PSNR and SSIM values against training epochs on the validation set with $R = 8$. Fig. 4 shows that *Full-rank* with weighted matrix has approximately the same PSNR/SSIM curve as the *supervised* baseline. Fig. 4 also shows that using a weighting matrix $\boldsymbol{W}$ in the self-supervised loss can significantly improve imaging quality, even in the setting $\mathcal{J}_{\mathsf{def}}$ where Assumption 2 does not hold.

## 4.2 Experimentally Collected In-Vivo Data

### 4.2.1 Dataset

Data acquisition was performed on a Siemens 3T Prisma scanner (Siemens Healthcare, Erlangen, Germany) with 64-channel Head/Neck coils. We collected images using the Sagittal T1 magnetization-prepared rapid gradient-echo (MPRAGE) sequence. The acquisition parameters were as follows: repetition time (TR) = 2400 ms, echo time (TE) = 2.62 ms, inversion time (TI) = 1000 ms, flip angle (FA) = 8 degrees, FOV = 256×256 mm, voxel size = 1×1×1 mm, slices per slab = 176, slice and phase resolution = 100% and slice and phase partial Fourier off. A 2× oversampling was used in the frequency encoding direction, and the asymmetric echo was turned on to allow short TE. The sampling pattern is equispaced 1D Cartesian with ACS lines. Upon the approval of our Institutional Review Board, we used brain MRI data from 14, 1, and 5 participants in this study for training, validation, and testing, respectively. We acquired the training data with GRAPPA = 2 in phase encoding (PE) direction with 24 ACS lines, the total acquisition time was 5 minutes and 35 seconds, and the raw measurements correspond to approximately 65% sampling rate. The validation

and testing data were fully-sampled measurements acquired with GRAPPA turned off, and the total acquisition time was 10 minutes and 16 seconds. We considered groundtruth as the *root-sum-square (RSS)* reconstruction from the fully-sampled data.

Experiments used the acceleration factors of $R = 4$, $R = 6$, and $R = 8$, corresponding to the retrospectively sampling rate of $30\%$, $20\%$, and $16\%$, respectively. We obtained the two training measurements of the same subjects (*i.e.*, $\boldsymbol{y}_i'$ and $\boldsymbol{y}_i$ in (6)) by allocating acquired Cartesian lines into two bins. Note that no groundtruth data was used during training. We applied 1D Fourier transform on the $k_z$ dimension of the raw data and then reconstructed the images slice by slice. The raw measurement of each slice is of size $512 \times 256 \times 64$ with $512 \times 256$ being $k_x \times k_y$ dimension and $64$ being the numbers of receiver coils. Note that the high number of receiver coils makes DU methods impractical due to the increase in the computation and memory complexity.

| Metrics | PSNR (dB) | | | SSIM | | | Mem[1] |
|---|---|---|---|---|---|---|---|
| Acceleration rate | $\times 8$ | $\times 6$ | $\times 4$ | $\times 8$ | $\times 6$ | $\times 4$ | |
| *Zero-Filled* | 16.86 | 17.61 | 19.30 | 0.698 | 0.735 | 0.802 | N/A |
| *TV* | 26.57 | 30.53 | 38.54 | 0.862 | 0.913 | 0.971 | 1745 |
| *A2A* | 29.80 | 31.84 | 35.03 | 0.874 | 0.903 | 0.940 | 7859 |
| *SSDU* | 31.10 | 32.21 | 36.43 | 0.895 | 0.920 | 0.961 | 21981 |
| *SelfDU* | 32.19 | 34.44 | 37.65 | 0.906 | 0.931 | 0.961 | 21981 |
| *SelfDEQ*[w/o] | 33.80 | 36.05 | 39.22 | 0.928 | 0.949 | 0.973 | 9325 |
| *SelfDEQ* | **34.05** | **36.79** | **39.64** | **0.929** | **0.950** | **0.973** | 9325 |

[1] GPU memory demand for training (MB).
*SelfDEQ*[w/o]: SelfDEQ *(unweighted)*.

Table 2: Summary of the PSNR and SSIM values on the experimentally-collected data. This table highlights that SelfDEQ can outperform several self-supervised methods, including those based on DU.
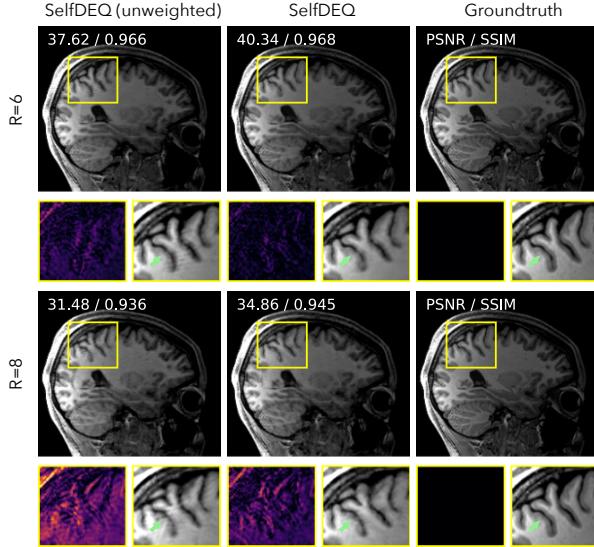


Figure 6: The reconstruction performance of SelfDEQ and its ablated variant that does not use the weighting matrix in the self-supervised loss. This figure shows results on acceleration factors $R = 6$ and $R = 8$, using the same image slice from the *experimentally-collected* data. Note the improvement due to the proposed weighting matrix $\boldsymbol{W}$.

#### 4.2.2 Comparison

We compared SelfDEQ against several standard baseline methods and state-of-the-art self-supervised MRI methods. (a) *Total Variation (TV)*: an optimization-based method using total variation regularizer in (3). We optimized the trade-off parameter $\tau$ using grid search. (b) *Artifact2Artifact (A2A)* [23]: trains U-Net by mapping corrupted MR images of the same subject to each other. (c) $SSDU^2$ [40]: a recent self-supervised method that trains a DU network by dividing each k-space MRI acquisition into two subsets and mapping them to each other. (d) *SelfDU*: a DU network trained on the same DL architecture and using the *weighted* loss function of SelfDEQ. We set DU iterations of *SSDU* and *SelfDU* to 7, which is the maximum number achivable under memory constraints of our machine. We also implemented SelfDEQ *(unweighted)* as an ablated method that trains SelfDEQ on the self-supervised loss *without* $\boldsymbol{W}$.

#### 4.2.3 Discussion

Fig. 5 illustrates the reconstruction results of all baseline methods on $R = 6$. *Total variation* suffers from detail loss due to the well-known "staircase effect." While *Artifact2Artifact* has better performance than *Total variation* by learning the prior from data, *SSDU* and *SelfDU* outperform it due to their model-based DU architectures. Overall, SelfDEQ achieves the best performance in artifact removal and sharpness. For instance, the reconstructed images obtained using SelfDEQ are sharper and have more fine details, especially in the brain tissues highlighted by green arrows. On the other hand, other methods show ghosting artifacts in their reconstructed images (see also zoom-in regions in Fig. 5). Table 2 summarizes the average PSNR and SSIM values of all the baseline methods on the testing dataset. The quantitative evaluations in Table 2 show the superior performance achieved by SelfDEQ. Table 2 also provides the GPU memory requirements of each method for training, highlighting that SelfDEQ can achieve better results with lower GPU memory demand than the DU-based methods.

Fig. 6 illustrates reconstruction results of SelfDEQ and SelfDEQ *(unweighted)* on different acceleration factors and on the same image slice. Fig. 6 shows that using the weighted matrix $\boldsymbol{W}$ in the self-supervised loss function can improve the imaging quality at different sampling rates.

# 5 Discussion and Conclusion

## 5.1 Applicability

Practically obtaining training data for SelfDEQ is straightforward. According to Assumption 2, it is sufficient to have a set of forward operators, where each operator subsamples the measurement domain, but their union over the training data covers the full space. For example, in MRI, one can implement a set of sampling masks illustrated in Fig. 2, then randomly pick one of the sampling masks from the set during scanning. In this example, individual undersampled measurements are still compatible with widely-used imaging techniques, such as GRAPPA [12] or ESPIRiT [37].

## 5.2 Future work

One benefit of SelfDEQ is its memory efficiency, which is well suited for large-scale imaging problems with high dimensional data. In our experiments on real MRI data, we have applied SelfDEQ on parallel MRI where the dimension of the receiver coils are high for conventional DU methods. In future work, we will apply SelfDEQ on other high-dimensional data such as 4D free-breathing MRI [8] or 2D+time cardiac MRI [31], where it is also challenging to obtain high-quality groundtruth training data.

## 5.3 Conclusion

This work presents SelfDEQ as a novel self-supervised learning framework for training model-based deep implicit neural networks for image reconstruction in accelerated MRI. The motivation behind SelfDEQ is to enable *efficient* and *effective* training of *implicit networks* directly on undersampled and noisy MRI measurements without any

---

[2]SSDU implementation is available on GitHub: https://github.com/byaman14/SSDU.

groundtruth. The SelfDEQ framework consists of several synergistic elements: (a) a model-based implicit network that integrates measurement operators and CNN priors; (b) a self-supervised loss that accounts for sampling imbalances; (c) a Jacobian-free backward pass that leads to efficient training. We theoretical analysed SelfDEQ showing that it can do as well as the supervised learning. We tested the SelfDEQ framework on real MRI data, showing that it (i) outperforms recent DU based self-supervised methods; (ii) matches the performance of fully-supervised DEQ; and (iii) enables highly-accelerated data-collection in parallel MRI.

## Acknowledgements

## References

[1] H. K. Aggarwal, M. P. Mani, and M. Jacob. MoDL: Model-Based Deep Learning Architecture for Inverse Problems. *IEEE Trans. Med. Imaging*, 38(2):394–405, February 2019.

[2] M. Akçakaya, B. Yaman, H. Chung, and J. C. Ye. Unsupervised deep learning methods for biological image reconstruction and enhancement: An overview from a signal processing perspective. *IEEE Signal Process. Mag.*, 39(2):28–44, 2022.

[3] D. G Anderson. Iterative procedures for nonlinear integral equations. *Journal of the ACM (JACM)*, 12(4):547–560, 1965.

[4] S. Bai, J Z. Kolter, and V. Koltun. Deep Equilibrium Models. In *Proc. Adv. Neural Inf. Process. Syst.*, volume 32, 2019.

[5] A. Bora, E. Price, and A. G Dimakis. Ambientgan: Generative models from lossy measurements. In *Int. Conf. Learn. Represent.*, 2018.

[6] D. Chen, J. Tachella, and M. E Davies. Equivariant imaging: Learning beyond the range space. In *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, pages 4379–4388, 2021.

[7] R. TQ Chen, Y. Rubanova, J. Bettencourt, and D. K Duvenaud. Neural ordinary differential equations. In *Proc. Adv. Neural Inf. Process. Syst.*, volume 31, 2018.

[8] C. Eldeniz, W. Gan, S. Chen, T. J. Fraum, D. R. Ludwig, Y. Yan, J. Liu, T. Vahle, U. Krishnamurthy, U. S. Kamilov, and H. An. Phase2Phase: Respiratory Motion-Resolved Reconstruction of Free-Breathing Magnetic Resonance Imaging Using Deep Learning Without a Ground Truth for Improved Liver Imaging. *Invest. Radiol.*, 56(12):809–819, May 2021.

[9] S. W. Fung, H. Heaton, Q. Li, D. McKenzie, S. Osher, and W. Yin. JFB: Jacobian-Free Backpropagation for Implicit Networks. In *arXiv:2103.12803*, December 2021.

[10] W. Gan, Y. Sun, C. Eldeniz, J. Liu, H. An, and U. S. Kamilov. Deformation-Compensated Learning for Image Reconstruction without Ground Truth. *IEEE Trans. Med. Imaging*, pages 1–1, 2022.

[11] D. Gilton, G. Ongie, and R. Willett. Deep Equilibrium Architectures for Inverse Problems in Imaging. *IEEE Trans. Comput. Imaging*, 7:1123–1133, June 2021.

[12] M. A. Griswold, P. M. Jakob, R. M. Heidemann, M. Nittka, V. Jellus, J. Wang, B. Kiefer, and A. Haase. Generalized autocalibrating partially parallel acquisitions (GRAPPA). *Magn. Reson. Med.*, 47(6):1202–1210, June 2002.

[13] H. Gupta, M. T McCann, L. Donati, and M. Unser. Cryogan: a new reconstruction paradigm for single-particle cryo-em via deep adversarial learning. *IEEE Trans. Comput. Imaging*, 7:759–774, 2021.

[14] K. Hammernik, T. Klatzer, E. Kobler, M. P Recht, D. K Sodickson, T. Pock, and F. Knoll. Learning a variational network for reconstruction of accelerated MRI data. *Magn. Reson. Med.*, 79(6):3055–3071, 2018.

[15] A. A. Hendriksen, D. M. Pelt, and K. J. Batenburg. Noise2Inverse: Self-Supervised Deep Convolutional Denoising for Tomography. *IEEE Trans. Comput. Imaging*, 6:1320–1335, 2020.

[16] Z. Jiang, Z. Huang, B. Qiu, X. Meng, Y. You, X. Liu, M. Geng, G. Liu, C. Zhou, K. Yang, A. Maier, Q. Ren, and Y. Lu. Weakly Supervised Deep Learning-Based Optical Coherence Tomography Angiography. *IEEE Trans. Med. Imaging*, 40(2):688–698, February 2021.

[17] K. H. Jin, M. T. McCann, E. Froustey, and M. Unser. Deep Convolutional Neural Network for Inverse Problems in Imaging. *IEEE Trans. on Image Process.*, 26(9):4509–4522, September 2017.

[18] U. S. Kamilov, C. A. Bouman, G. T. Buzzard, and B. Wohlberg. Plug-and-play methods for integrating physical and learned models in computational imaging. *IEEE Signal Process. Mag.*, 2022. arXiv:2203.17061.

[19] J. Kelly, J. Bettencourt, M. J Johnson, and D. K Duvenaud. Learning differential equations that are easy to solve. In *Proc. Adv. Neural Inf. Process. Syst.*, volume 33, pages 4370–4380, 2020.

[20] D. P. Kingma and J. Ba. Adam: A Method for Stochastic Optimization. *arXiv:1412.6980*, January 2017.

[21] A. Krull, T.-O. Buchholz, and F. Jug. Noise2Void - Learning Denoising From Single Noisy Images. In *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, pages 2124–2132, June 2019.

[22] J. Lehtinen, J. Munkberg, J. Hasselgren, S. Laine, T. Karras, M. Aittala, and T. Aila. Noise2Noise: Learning image restoration without clean data. In *Proc. Int. Conf. Mach. Learn.*, 2018.

[23] J. Liu, Y. Sun, C. Eldeniz, W. Gan, H. An, and U. S. Kamilov. RARE: Image Reconstruction Using Deep Priors Learned Without Groundtruth. *IEEE J. Sel. Top. Signal Process.*, 14(6):1088–1099, October 2020.

[24] J. Liu, X. Xu, W. Gan, S. Shoushtari, and U. S. Kamilov. Online Deep Equilibrium Learning for Regularization by Denoising. *arXiv:2205.13051*, May 2022.

[25] A. Lucas, M. Iliadis, R. Molina, and A. K Katsaggelos. Using deep neural networks for inverse problems in imaging: Beyond analytical methods. *IEEE Signal Process. Mag.*, 35(1):20–36, 2018.

[26] M. Lustig, D. Donoho, and J. M Pauly. Sparse MRI: The application of compressed sensing for rapid MR imaging. *Magn. Reson. Med.*, 58(6):1182–1195, 2007.

[27] C. Millard and M. Chiew. A framework for self-supervised MR image reconstruction using sub-sampling via Noisier2Noise. *arXiv:2205.10278*, June 2022.

[28] N. Moran, D. Schmidt, Y. Zhong, and P. Coady. Noisier2Noise: Learning to Denoise From Unpaired Noisy Data. In *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, pages 12061–12069, June 2020.

[29] G. Ongie, A. Jalal, C. A Metzler, R. G Baraniuk, A. G Dimakis, and R. Willett. Deep learning techniques for inverse problems in imaging. *IEEE J. Sel. Areas Inf. Theory*, 1(1):39–56, 2020.

[30] A. Pramanik and M. Jacob. Improved Model based Deep Learning using Monotone Operator Learning (MOL). In *Proc. Int. Symp. Biomedical Imaging*, 2022.

[31] A. Pramanik and M. Jacob. Stable and memory-efficient image recovery using monotone operator learning (MOL). *arXiv:2206.04797*, June 2022.

[32] Y. Romano, M. Elad, and P. Milanfar. The Little Engine That Could: Regularization by Denoising (RED). *SIAM J. Imaging Sci.*, 10(4):1804–1844, January 2017.

[33] O. Ronneberger, P. Fischer, and T. Brox. U-Net: Convolutional networks for biomedical image segmentation. In *Proc. Medical Image Computing and Computer-Assisted Intervention*, pages 234–241, 2015.

[34] E. K Ryu, J. Liu, S. Wang, X. Chen, Z. Wang, and W. Yin. Plug-and-Play Methods Provably Converge with Properly Trained Denoisers. In *Proc. Int. Conf. Mach. Learn.*, page 12, 2019.

[35] J. Schlemper, J. Caballero, J. V. Hajnal, A. N. Price, and D. Rueckert. A Deep Cascade of Convolutional Neural Networks for Dynamic MR Image Reconstruction. *IEEE Trans. Med. Imaging*, 37(2):491–503, February 2018.

[36] J. Tachella, D. Chen, and M. Davies. Sampling Theorems for Unsupervised Learning in Linear Inverse Problems. *arXiv:2203.12513*, March 2022.

[37] M. Uecker, P. Lai, M. J. Murphy, P. Virtue, M. Elad, J. M. Pauly, S. S. Vasanawala, and M. Lustig. ESPIRiT-an eigenvalue approach to autocalibrating parallel MRI: Where SENSE meets GRAPPA. *Magn. Reson. Med.*, 71(3):990–1001, March 2014.

[38] D. Ulyanov, A. Vedaldi, and V. Lempitsky. Deep image prior. In *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, pages 9446–9454, 2018.

[39] S. V. Venkatakrishnan, C. A. Bouman, and B. Wohlberg. Plug-and-Play priors for model based reconstruction. In *Proc. IEEE Glob. Conf. Signal Process. Inf. Process.*, pages 945–948, December 2013.

[40] B. Yaman, S. A. H. Hosseini, S. Moeller, J. Ellermann, K. Uğurbil, and M. Akçakaya. Self-supervised learning of physics-guided reconstruction neural networks without fully sampled reference data. *Magn. Reson. Med.*, 84(6):3172–3191, December 2020.

[41] G. Zeng, Y. Guo, J. Zhan, Z. Wang, Z. Lai, X. Du, X. Qu, and D. Guo. A review on deep learning mri reconstruction without fully sampled k-space. *BMC Medical Imaging*, 21(1):1–11, 2021.

[42] Y. Zhao, S. Zheng, and X. Yuan. Deep Equilibrium Models for Video Snapshot Compressive Imaging. *arXiv:2201.06931*, January 2022.

[43] B. Zhu, J. Z. Liu, S. F. Cauley, B. R. Rosen, and M. S. Rosen. Image reconstruction by domain-transform manifold learning. *Nature*, 555(7697):487–492, March 2018.

# Supplementary Materials

Our main manuscript presents *self-supervised deep equilibrium model (SelfDEQ)* as the first self-supervised MRI reconstruction framework for training model-based implicit neural networks from undersampled and noisy measurements. This supplementary document presents the details of our theoretical analysis.

We use the same notations as in the main manuscript. We consider reconstruction of an image $\boldsymbol{x} \in \mathbb{C}^n$ from its noisy and undersampled measurement

$$\boldsymbol{y} = \boldsymbol{M}\boldsymbol{A}\boldsymbol{x} + \boldsymbol{e}, \tag{14}$$

where $\boldsymbol{e} \in \mathbb{C}^m$ is a *additive white Gaussian noise (AWGN)* vector, $\boldsymbol{A} \in \mathbb{C}^{n \times n}$ is a measurement matrix, and $\boldsymbol{M} \in \{0,1\}^{m \times n}$ is a diagonal sampling matrix. We define the *average weighting matrix* $\overline{\boldsymbol{W}}$ to account for unbalanced sampling across multiple acquisitions in the training data

$$\overline{\boldsymbol{W}} = \mathsf{diag}(\overline{w_0}, \overline{w_1}, ..., \overline{w_n}) \in \mathbb{R}^{n \times n} \,, \tag{15}$$

where

$$\overline{w_k} = \begin{cases} \frac{1}{\sqrt{\mathbb{E}[\boldsymbol{M'}^\mathsf{T}\boldsymbol{M'}]_{k,k}}} & \sqrt{\mathbb{E}[\boldsymbol{M'}^\mathsf{T}\boldsymbol{M'}]_{k,k}} \neq 0 \\ 0 & \sqrt{\mathbb{E}[\boldsymbol{M'}^\mathsf{T}\boldsymbol{M'}]_{k,k}} = 0 \end{cases} \,.$$

We define the *weighting matrix* $\boldsymbol{W}$ as a subsampled variant of $\overline{\boldsymbol{W}}$ given an individual subsampling operator $\boldsymbol{M'}$, which is used in our self-supervised loss function

$$\boldsymbol{W} = \boldsymbol{M'}\overline{\boldsymbol{W}}(\boldsymbol{M'}\overline{\boldsymbol{W}})^\mathsf{T} \in \mathbb{R}^{m \times m} \,. \tag{16}$$

Consider the training set of measurement pairs $\{\boldsymbol{y}_i, \boldsymbol{y}_i'\}_{i=1}^N$ with each pair $\boldsymbol{y}_i, \boldsymbol{y}_i'$ corresponding to the same object $\boldsymbol{x}_i$

$$\boldsymbol{y}_i = \boldsymbol{M}_i\boldsymbol{A}\boldsymbol{x}_i + \boldsymbol{e}_i \text{ and } \boldsymbol{y}_i' = \boldsymbol{M}_i'\boldsymbol{A}\boldsymbol{x}_i + \boldsymbol{e}_i' \,. \tag{17}$$

The value $N \geq 1$ denotes the number training pairs.

The traditional DEQ gradient for the loss $\ell$ is given by

$$\nabla\ell(\boldsymbol{\theta}) = \mathsf{Real}\Big( \big(\nabla_{\boldsymbol{\theta}}\mathsf{T}_{\boldsymbol{\theta}}(\bar{\boldsymbol{x}})\big)^\mathsf{H}\boldsymbol{b} \Big)$$
$$\text{where } \boldsymbol{b} = \big(\boldsymbol{I} - \nabla_{\boldsymbol{x}}\mathsf{T}_{\boldsymbol{\theta}}(\bar{\boldsymbol{x}})\big)^{-\mathsf{T}}\Big[\frac{\partial\ell}{\partial\bar{\boldsymbol{x}}}\Big]^\mathsf{T} \,. \tag{18}$$

Jacobian-free backpropagation (JFB) approximates (18) as

$$\mathsf{JFB}_\ell(\boldsymbol{\theta}) = \mathsf{Real}\Big( \big(\nabla_{\boldsymbol{\theta}}\mathsf{T}_{\boldsymbol{\theta}}(\bar{\boldsymbol{x}})\big)^\mathsf{H}\Big[\frac{\partial\ell}{\partial\bar{\boldsymbol{x}}}\Big]^\mathsf{T} \Big) \,. \tag{19}$$

The analysis below shows that the JFB update using the self-supervised loss matches that using the supervised loss.

**Proposition 1** *When Assumption 2 is satisfied,*

$$\mathbb{E}\big[(\boldsymbol{M'}\boldsymbol{A})^\mathsf{H}\boldsymbol{W}\boldsymbol{M'}\boldsymbol{A}\big] = \boldsymbol{I} \,,$$

*where the expectation is with respect to $p_{\boldsymbol{M}}$.*

*Proof:* Since Assumption 2 implies that $\mathbb{E}[\boldsymbol{M'}]_{k,k} \neq 0$

$$\overline{w_k} = \frac{1}{\sqrt{\mathbb{E}[\boldsymbol{M'}^\mathsf{T}\boldsymbol{M'}]_{k,k}}} \,. \tag{20}$$

Since $M'^{\mathsf{T}} M' \in \{0,1\}^{n \times n}$ and $\overline{W} \in \mathbb{R}^{n \times n}$ are both diagonal matrices, we have

$$
\begin{aligned}
&\mathbb{E}[M'^{\mathsf{T}} W M'] \\
&= \mathbb{E}[M'^{\mathsf{T}} M' \overline{W}\, \overline{W}^{\mathsf{T}} M'^{\mathsf{T}} M'] \\
&= \overline{W}\, \overline{W}^{\mathsf{T}} \mathbb{E}[M'^{\mathsf{T}} M' M'^{\mathsf{T}} M'] \\
&= I \, .
\end{aligned}
\tag{21}
$$

Now we can establish the desired result

$$
\begin{aligned}
\mathbb{E}\big[(M'A)^{\mathsf{H}} W M' A\big] &= A^{\mathsf{H}} \mathbb{E}\big[M'^{\mathsf{T}} W M'\big] A \\
&= A^{\mathsf{H}} A = I \, .
\end{aligned}
\tag{22}
$$

where the second equation are due to (21), and the last equation is because $A$ is an orthogonal matrix.

**Theorem 2** *Under Assumptions 1-2, the JFB update of the weighted self-supervised loss ($\ell_{\mathsf{self}}$) is equivalent to its supervised counterpart ($\ell_{\mathsf{sup}}$), namely we have that*

$$
\mathsf{JFB}_{\ell_{\mathsf{self}}}(\boldsymbol{\theta}) = \mathsf{JFB}_{\ell_{\mathsf{sup}}}(\boldsymbol{\theta}) \, .
\tag{23}
$$

*where*

$$
\ell_{\mathsf{sup}} = \mathbb{E}\left[\frac{1}{2}\,\|\bar{\boldsymbol{x}} - \boldsymbol{x}\|_2^2\right]
\tag{24}
$$

*and*

$$
\ell_{\mathsf{self}} = \mathbb{E}\left[\frac{1}{2}\,\|M' A \bar{\boldsymbol{x}} - \boldsymbol{y}'\|_W^2\right] \, .
\tag{25}
$$

*The vector $\bar{\boldsymbol{x}} = \mathsf{T}_{\boldsymbol{\theta}}(\bar{\boldsymbol{x}}, \boldsymbol{y})$ is the fixed-point of $\mathsf{T}_{\boldsymbol{\theta}}$ for $\boldsymbol{y}$.*

*Proof:* In order to simplify the notations in the following analysis, we directly use complex valued quantities and assume that the real part is taken at the end.

The supervised update $\mathsf{JFB}_{\ell_{\mathsf{sup}}}(\boldsymbol{\theta})$ is given by:

$$
\begin{aligned}
\mathsf{JFB}_{\ell_{\mathsf{sup}}}(\boldsymbol{\theta}) &= \mathbb{E}\left[\left(\nabla_{\boldsymbol{\theta}} \mathsf{T}_{\boldsymbol{\theta}}(\bar{\boldsymbol{x}})\right)^{\mathsf{H}}\left[\frac{\partial \ell_{\mathsf{sup}}}{\bar{\boldsymbol{x}}}\right]^{\mathsf{T}}\right] \\
&= \mathbb{E}\left[\left(\nabla_{\boldsymbol{\theta}} \mathsf{T}_{\boldsymbol{\theta}}(\bar{\boldsymbol{x}})\right)^{\mathsf{H}}(\bar{\boldsymbol{x}} - \boldsymbol{x})\right] \, .
\end{aligned}
\tag{26}
$$

On the other hand, we can re-write the weighted self-supervised update as $\mathsf{JFB}_{\ell_{\mathsf{self}}}(\boldsymbol{\theta})$:

$$
\begin{aligned}
\mathsf{JFB}_{\ell_{\mathsf{self}}}(\boldsymbol{\theta}) &= \mathbb{E}\left[\left(\nabla_{\boldsymbol{\theta}} \mathsf{T}_{\boldsymbol{\theta}}(\bar{\boldsymbol{x}})\right)^{\mathsf{H}}\left[\frac{\partial \ell_{\mathsf{self}}}{\partial \bar{\boldsymbol{x}}}\right]^{\mathsf{T}}\right] \\
&= \mathbb{E}\left[\left(\nabla_{\boldsymbol{\theta}} \mathsf{T}_{\boldsymbol{\theta}}(\bar{\boldsymbol{x}})\right)^{\mathsf{H}}\left(H'^{\mathsf{H}}(H'\bar{\boldsymbol{x}} - \sqrt{W}\boldsymbol{y}')\right)\right] \\
&= \mathbb{E}\left[\left(\nabla_{\boldsymbol{\theta}} \mathsf{T}_{\boldsymbol{\theta}}(\bar{\boldsymbol{x}})\right)^{\mathsf{H}}\mathbb{E}\left[H'^{\mathsf{H}}(H'\bar{\boldsymbol{x}} - \sqrt{W}\boldsymbol{y}')\big|\boldsymbol{x}, M, e\right]\right] \, ,
\end{aligned}
\tag{27}
$$

where $H' := \sqrt{W} M' A$. The last equation is true since $\nabla_{\boldsymbol{\theta}} \mathsf{T}_{\boldsymbol{\theta}}(\bar{\boldsymbol{x}})$ is deterministic when conditioned on $\boldsymbol{x}$, $M$ and $e$. We also have that

$$
\begin{aligned}
&\mathbb{E}\big[H'^{\mathsf{H}}(H'\bar{\boldsymbol{x}} - \sqrt{W}\boldsymbol{y}')\big|\boldsymbol{x}, M, e\big] \\
&= \mathbb{E}\big[H'^{\mathsf{H}}\big(H'\bar{\boldsymbol{x}} - \sqrt{W}(M' A \boldsymbol{x} + e')\big)\big|\boldsymbol{x}, M, e\big] \\
&= \mathbb{E}\big[H'^{\mathsf{H}}\big(H'(\bar{\boldsymbol{x}} - \boldsymbol{x}) + \sqrt{W}e'\big)\big|\boldsymbol{x}, M, e\big] \\
&= \mathbb{E}\big[H'^{\mathsf{H}} H'(\bar{\boldsymbol{x}} - \boldsymbol{x})\big|\boldsymbol{x}, M, e\big] + \mathbb{E}\big[H'^{\mathsf{H}}\sqrt{W}e'\big|\boldsymbol{x}, M, e\big],
\end{aligned}
\tag{28}
$$

where in the second row we used $\boldsymbol{y} = \boldsymbol{M}'\boldsymbol{A}\boldsymbol{x} + \boldsymbol{e}'$. The first term in (28) can also be expressed as

$$
\begin{aligned}
&\mathbb{E}\big[\boldsymbol{H}'^{\mathsf{H}}\boldsymbol{H}'(\bar{\boldsymbol{x}} - \boldsymbol{x})\big|\boldsymbol{x}, \boldsymbol{M}, \boldsymbol{e}\big] \\
&= \mathbb{E}\big[\boldsymbol{H}'^{\mathsf{H}}\boldsymbol{H}'\big](\bar{\boldsymbol{x}} - \boldsymbol{x}) \\
&= \bar{\boldsymbol{x}} - \boldsymbol{x} \,,
\end{aligned}
\tag{29}
$$

where the second equation is due to independence of $\boldsymbol{M}'$ from $\boldsymbol{x}$, $\boldsymbol{M}$ and $\boldsymbol{e}$, and the last equation is due to Proposition 1. The second term of (28) can be expressed as

$$
\begin{aligned}
&\mathbb{E}\big[\boldsymbol{H}'^{\mathsf{H}}\sqrt{\boldsymbol{W}}\boldsymbol{e}'\big|\boldsymbol{x}, \boldsymbol{M}, \boldsymbol{e}\big] \\
&= \mathbb{E}\big[\boldsymbol{H}'^{\mathsf{H}}\sqrt{\boldsymbol{W}}\boldsymbol{e}'\big] = \mathbb{E}\big[\boldsymbol{H}'^{\mathsf{H}}\sqrt{\boldsymbol{W}}\big]\mathbb{E}\big[\boldsymbol{e}'\big] = \boldsymbol{0} \,,
\end{aligned}
\tag{30}
$$

where the first equation is due to the independence of $\boldsymbol{M}'$ in $\boldsymbol{H}'$ and $\boldsymbol{e}'$ from $\boldsymbol{x}$, $\boldsymbol{M}$ and $\boldsymbol{e}$, the second equation is due to the independence of $\boldsymbol{M}'$ from $\boldsymbol{e}'$, and the last equation is due to $\boldsymbol{e}' \sim \mathcal{N}(0, \sigma^2 \boldsymbol{I})$. Combining (27), (28), (29) and (30), we get

$$
\mathsf{JFB}_{\ell_{\text{self}}}(\boldsymbol{\theta}) = \mathbb{E}\Big[\big(\nabla_{\boldsymbol{\theta}}\mathsf{T}_{\boldsymbol{\theta}}(\bar{\boldsymbol{x}})\big)^{\mathsf{H}}(\bar{\boldsymbol{x}} - \boldsymbol{x})\Big] \,,
\tag{31}
$$

which establishes the desired results.