Efficient Approximation of Jacobian Matrices Involving a Non-Uniform Fast Fourier Transform (NUFFT)

Guanhua Wang and Jeffrey A. Fessler

Abstract—There is growing interest in learning Fourier domain sampling strategies (particularly for magnetic resonance imaging, MRI) using optimization approaches. For non-Cartesian sampling, the system models typically involve non-uniform fast Fourier transform (NUFFT) operations. Commonly used NUFFT algorithms contain frequency domain interpolation, which is not differentiable with respect to the sampling pattern, complicating the use of gradient methods. This paper describes an efficient and accurate approach for computing approximate gradients involving NUFFTs. Multiple numerical experiments validate the improved accuracy and efficiency of the proposed approximation. As an application to computational imaging, the NUFFT Jacobians were used to optimize non-Cartesian MRI sampling trajectories via data-driven stochastic optimization. Specifically, the sampling patterns were learned with respect to various modelbased image reconstruction (MBIR) algorithms. The proposed approach enables sampling optimization for image sizes that are infeasible with standard auto-differentiation methods due to memory limits. The synergistic acquisition and reconstruction design leads to remarkably improved image quality. In fact, we show that model-based image reconstruction methods with suitably optimized imaging parameters can perform nearly as well as CNN-based methods.

Index Terms—NUFFT, auto-differentiation, MRI k-space trajectory, accelerated MRI, data-driven optimization, machine learning

I. INTRODUCTION

There are several computational imaging modalities where the raw measurements can be modeled as samples of the imaged object's spectrum, where those samples need not lie on the Cartesian grid, including radar [1], diffraction ultrasound tomography [2], parallel-beam tomography [3], and MRI [4], [5]. Image reconstruction methods for such modalities may use non-uniform fast Fourier transform (NUFFT) operations to accelerate computation [6], [7]. The quality of the reconstructed image depends both on the image reconstruction method and on the characteristics of the frequency domain sampling pattern.

MRI has particular flexibility in designing frequency domain sampling patterns. Many MR sampling patterns are discrete subsets of the Cartesian grid, and the corresponding optimization/learning strategies include greedy algorithms [8], [9], [10], reparameterization [11], [12], [13], [14], [15], Bayesian

This work supported in part by NIH Grants R01 EB023618 and U01 EB026977 and NSF Grant IIS 1838179.

optimization [16], [17], and system matrix analysis [18], [19], [20], [21]. The other type is non-Cartesian sampling, which uses a collection of continuous functions in k-space. Several studies applied gradient methods to optimize non-Cartesian sampling trajectories [22], [23], [24], [25], and it is also possible to use derivative-free optimization algorithms in certain applications [26]. This paper develops efficient tools for applying gradient methods to non-Cartesian sampling pattern optimization.

1

Some data-driven optimization methods for non-Cartesian sampling solve an optimization problem involving both forward system models and image reconstruction methods [22], [23], [24]. The forward models and reconstruction methods both depend on NUFFT operations. In principle, the Fourier transform operation is a continuous function of the k-space sample locations and thus should be applicable to gradient-based optimization methods. In practice, the NUFFT $(\mathcal{O}(N \log N))$ operations) is an approximation to the nonuniform discrete Fourier transform (NUDFT, $\mathcal{O}(N^2)$) operations) and that approximation often is implemented using nondifferentiable lookup table operations or other interpolation techniques [27], [28]. Such approximations are sufficient for image reconstruction (forward mode), but have problematic efficiency and accuracy if one attempts to use standard autodifferentiation tools for gradient-based optimization. Standard auto-differentiation methods using subgradients can lead to incorrect NUFFT Jacobians. They also require prohibitively large amounts of memory for back-propagation through certain algorithm stages such as conjugate gradient (CG) steps that involve NUFFT operations.

This paper proposes an efficient approach that replaces memory-intensive and inaccurate auto-differentiation steps with fast Jacobian approximations that are themselves based on NUFFT operations. The proposed approach requires substantially less memory for iterative updates like CG steps.

As a direct application, we used the proposed Jacobian to learn MRI sampling trajectories via stochastic optimization. By applying the forward system model and subsequent reconstruction, reconstructed images were simulated from reference images in the training set. The similarity between simulated and reference images was the metric for updating the sampling trajectory. We used model-based reconstruction methods, such as regularized least-squares and compressed sensing. In comparison with previous works using reconstruction neural networks (NN) [22], [24], such model-based reconstruction methods can be more robust and require less training data.

G. Wang is with the Department of Biomedical Engineering, University of Michigan, Ann Arbor, MI 48109 USA (e-mail:guanhuaw@umich.edu).

J. A. Fessler is with the Department of EECS, University of Michigan, Ann Arbor, MI 48109 USA (e-mail:fessler@umich.edu).

In addition to simple NUFFT-based sensing matrices, we also considered several scenarios in MR sampling and reconstruction, including the multi-coil (sensitivity-encoded) imaging [29] system models that account for B_0 field inhomogeneity [30]. The derivation also includes fast Jacobian approximations for Gram and "data consistency" operations commonly used in iterative reconstruction methods.

Jacobians with respect to the non-Cartesian sampling pattern are also relevant for tomographic image reconstruction problems with unknown view angles (like cryo-EM) where the view angles must be estimated [31].

The remainder of this paper is organized as follows. Section II derives the efficient Jacobian approximations. Section III details how to optimize MRI sampling patterns using learning-based methods. Section IV provides empirical validation of the approach, showing the efficacy and accuracy of the proposed approach. The appendix includes an error analysis of the proposed method.

The methods in this paper were used to assist the design of k-space sampling for a CNN-based reconstruction approach in our previous work [24]. This paper derives the theory in detail and considers k-space sampling optimization for general model-based reconstruction methods. Preliminary results were shown in an earlier short conference abstract [32].

II. JACOBIAN EXPRESSIONS

This section derives the key Jacobian expressions and their efficient approximations based on NUFFT operations. These approximations enable the applications that follow.

A. Lemmas

We denote matrices, vectors and scalars by A, a and a, respectively. A', A^T and A^* denote the Hermitian transpose, the transpose and the complex conjugate of A, respectively.

Consider a scalar function f(z), $z=x+yi\in\mathbb{C}$, $x,y\in\mathbb{R}$. Following the conventions in Wirtinger calculus [33, p. 67], the differential operators are defined as

$$\frac{\partial}{\partial z} = \frac{1}{2} \frac{\partial}{\partial x} - \frac{\imath}{2} \frac{\partial}{\partial y}, \ \frac{\partial}{\partial z^*} = \frac{1}{2} \frac{\partial}{\partial x} + \frac{\imath}{2} \frac{\partial}{\partial y}.$$

A function f is complex differentiable or holomorphic iff $\frac{\partial f}{\partial z^*} = 0$ (Cauchy–Riemann equation) [33, p. 66]. In the context of optimization, a cost function L (usually a real scalar) is not holomorphic w.r.t. complex variables. A common approach (as adopted by PyTorch and TensorFlow) regards the real and imaginary components of a complex variable as two real-valued variables, and updates them separately, similar to the real-valued calculus [34]. For example, the nth gradient descent step uses the update

$$\boldsymbol{z}_{n+1} = \boldsymbol{z}_n - \alpha \left(\frac{\partial L}{\partial \boldsymbol{x}} + \imath \frac{\partial L}{\partial \boldsymbol{y}} \right) = \boldsymbol{z}_n - 2\alpha \frac{\partial L}{\partial \boldsymbol{z}^*},$$

where $\alpha \in \mathbb{R}^+$ denotes the step size. The chain rule still applies to calculating $\frac{\partial L}{\partial z^*}$ [35] [33, p. 68]; for s = f(z):

$$\frac{\partial L}{\partial z^*} = \left(\frac{\partial L}{\partial s^*}\right)^* \frac{\partial s}{\partial z^*} + \frac{\partial L}{\partial s^*} \left(\frac{\partial s}{\partial z}\right)^*. \tag{1}$$

For Jacobian matrices, we follow the "numerator-layout" notation [36]. For example, the derivative of an m-element column vector \boldsymbol{y} w.r.t. an n-element vector \boldsymbol{x} is an $m \times n$ matrix:

$$\frac{\partial \boldsymbol{y}}{\partial \boldsymbol{x}} \triangleq \begin{bmatrix}
\frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \cdots & \frac{\partial y_1}{\partial x_n} \\
\frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \cdots & \frac{\partial y_2}{\partial x_n} \\
\vdots & \vdots & \ddots & \vdots \\
\frac{\partial y_m}{\partial x_1} & \frac{\partial y_m}{\partial x_2} & \cdots & \frac{\partial y_m}{\partial x_n}
\end{bmatrix} .$$
(2)

However, this convention does not handle scenarios such as the derivatives of the elements of one matrix w.r.t. the elements of another matrix. Thus, we adopt a natural extension by using the vec (vectorization) operation. Specifically, for a $M \times N$ matrix \boldsymbol{A} that is a function of a $P \times Q$ matrix \boldsymbol{B} , we write the derivative as a $MN \times PQ$ matrix by applying (2) to the vec of each matrix:

$$\mathcal{D}_{\mathbf{B}} \mathbf{A} = \mathcal{D}_{\mathbf{B}} \mathbf{A}(\mathbf{B}) \triangleq \frac{\partial \operatorname{vec}(\mathbf{A})}{\partial \operatorname{vec}(\mathbf{B})}.$$
 (3)

The following equalities are useful in our derivations. (Equalities involving products all assume the sizes are compatible.) For $A \in \mathbb{C}^{K \times L}$, $B \in \mathbb{C}^{L \times M}$, $C \in \mathbb{C}^{M \times N}$:

$$vec(\boldsymbol{A}\boldsymbol{B}\boldsymbol{C}) = (\boldsymbol{I}_N \otimes \boldsymbol{A}\boldsymbol{B}) \operatorname{vec}(\boldsymbol{C})$$
$$= (\boldsymbol{C}^T\boldsymbol{B}^T \otimes \boldsymbol{I}_K) \operatorname{vec}(\boldsymbol{A}). \tag{P1}$$

In general:

$$(A \otimes B)(C \otimes D) = (AC) \otimes (BD).$$
 (P2)

For $\mathbf{A} \in \mathbb{C}^{K \times L}$, $\mathbf{B} \in \mathbb{C}^{M \times N}$:

$$A \otimes B = (I_K \otimes B)(A \otimes I_N) = (A \otimes I_M)(I_L \otimes B).$$
 (P3)

For $\mathbf{A} \in \mathbb{C}^{M \times N}$, $\mathbf{x} \in \mathbb{C}^N$:

$$\mathcal{D}_{A}(Ax) = x^{T} \otimes I_{M}, \ \mathcal{D}_{A^{*}}(Ax) = 0.$$
 (P4)

For an invertible matrix A:

$$A \in \mathbb{C}^{N \times N} \implies \mathcal{D}_A A^{-1} = -(A^T)^{-1} \otimes A^{-1},$$

$$\mathcal{D}_{A^*} A^{-1} = 0. \tag{P5}$$

The chain rule still holds for the extended Jacobian formulation. Suppose $F: \mathbb{C}^{K \times L} \to \mathbb{C}^{M \times N}$ and $G: \mathbb{C}^{M \times N} \to \mathbb{C}^{P \times Q}$ are both holomorphic. For $\boldsymbol{X} \in \mathbb{C}^{K \times L}$, the Jacobian of the composite function is:

$$\underbrace{\mathcal{D}_{\boldsymbol{X}} G(F(\boldsymbol{X}))}_{PQ \times KL} = \underbrace{\mathcal{D}_{\boldsymbol{Y}} G(\boldsymbol{Y})|_{\boldsymbol{Y} = F(\boldsymbol{X})}}_{PQ \times MN} \underbrace{\mathcal{D}_{\boldsymbol{X}} F(\boldsymbol{X})}_{MN \times KL},$$

$$\mathcal{D}_{\boldsymbol{X}^*} G(F(\boldsymbol{X})) = \mathbf{0}. \tag{P6}$$

Equalities (P1)-(P3) are common matrix vectorization properties. See [37, Ch. 9] for (P4), [35] for (P5) and (P6).

B. System Model

Consider the (single-coil, initially) MRI measurement model for non-Cartesian sampling based on the NUDFT [5]:

$$y = Ax + \varepsilon$$
,

where $y \in \mathbb{C}^M$ denotes the measured k-space data, $x \in \mathbb{C}^N$ denotes the unknown image to be reconstructed, and $A \in$

 $\mathbb{C}^{M \times N}$ denotes the system matrix or encoding matrix, where $m{A} = m{A}(m{\omega})$ has elements

$$a_{ij} = e^{-i\vec{\omega}_i \cdot \vec{r}_j}, \quad i = 1, \dots, M, \quad j = 1, \dots, N$$
 (4)

for $\vec{\omega}_i \in \mathbb{R}^D$ and $\vec{r}_j \in \mathbb{R}^D$ where $D \in \{1,2,3\ldots\}$ denotes the image dimension, and where

$$\boldsymbol{\omega} = [\boldsymbol{\omega}^{[1]} \ \boldsymbol{\omega}^{[2]} \dots \boldsymbol{\omega}^{[D]}]$$

is the $M \times d$ matrix consisting of all the k-space sampling locations and $\boldsymbol{\omega}^{[d]} \in \mathbb{R}^M$ denotes its dth column. (For simplicity here, we ignore other physical effects like field inhomogeneity and relaxation that are sometimes included in the forward model in MRI [5].) The center locations of voxels $\{\vec{r_j}\}$ usually lie on a Cartesian grid, but the k-space sample locations $\boldsymbol{\omega}$ in principle can be arbitrary subject to the Nyquist constraint.

Typically *A* is approximated by a NUFFT [6]. Usually, the NUFFT operator involves frequency-domain interpolation operations that are often non-differentiable. One previous trajectory optimization approach that used auto-differentiation [22] replaced the non-differentiable lookup table with a bilinear interpolator. Bilinear interpolation is differentiable everywhere except at the sample locations. Auto-differentiation of bilinear interpolation involves differentiating some floor and ceiling operations and those derivatives are defined to be zero in popular deep learning frameworks such as PyTorch and TensorFlow, leading to suboptimal sub-gradient calculations. Nearest-neighbor interpolation has even worse properties for auto-differentiation because its derivative is zero almost everywhere, leading to a completely vanishing gradient.

In the following derivations, we investigate a different approach where we analyze the Jacobians w.r.t. ω and x using the NUDFT expression (4). Then for efficient implementation, we replace the NUDFT operations within the Jacobians with NUFFT approximations. This approach enables faster computation and requires substantially less memory.

C. Forward Operator

We first focus on the forward operation $A(\omega)x$ and determine Jacobian matrices with respect to x and ω . The $M \times N$ Jacobian matrix of the forward linear operation with respect to x is

$$\frac{\partial Ax}{\partial x} = A, \qquad \frac{\partial Ax}{\partial x^*} = 0.$$

For the dth column of the spectrum sampling pattern ω , the Jacobian has elements

$$\begin{split} \left[\frac{\partial \boldsymbol{A}\boldsymbol{x}}{\partial \boldsymbol{\omega}^{[d]}}\right]_{il} &= \frac{\partial [\boldsymbol{A}\boldsymbol{x}]_i}{\partial \omega_l^{[d]}} = \frac{\partial}{\partial \omega_l^{[d]}} \sum_{j=1}^N \mathrm{e}^{-\imath \vec{\omega}_i \cdot \vec{r}_j} x_j \\ &= \begin{cases} -\imath \sum_{j=1}^N \mathrm{e}^{-\imath \vec{\omega}_i \cdot \vec{r}_j} x_j r_j^{[d]}, & i = l \\ 0, & \text{otherwise,} \end{cases} \end{split}$$

for $i, l = 1, \ldots, M$. The above summation is the product of the ith row of $-i \mathbf{A}$ with $\mathbf{x} \odot \mathbf{r}^{[d]}$. Thus the $M \times M$ Jacobian matrix for the partial derivatives of $\mathbf{A}\mathbf{x}$ w.r.t. $\boldsymbol{\omega}^{[d]}$ is:

$$\frac{\partial Ax}{\partial \iota_{\iota_{\iota}}[d]} = -\iota \operatorname{diag} \left\{ A(x \odot r^{[d]}) \right\}. \tag{5}$$

Consequently, the Jacobian calculation should apply A to vector $x \odot r^{[d]}$ once. In the above derivation, A is a NUDFT operator. In the practical implementation, we use a NUFFT to approximate A, both for the forward model and for the Jacobian calculation.

D. Adjoint Operator

Derivations of the Jacobians for the adjoint operation $A'(\omega) y$ follow a similar approach. For y:

$$rac{\partial A'y}{\partial y}=A', \qquad rac{\partial A'y}{\partial y^*}=0.$$

For the dth column of ω , the $N \times M$ Jacobian matrix has elements:

$$\begin{bmatrix} \frac{\partial \mathbf{A}' \mathbf{y}}{\partial \boldsymbol{\omega}^{[d]}} \end{bmatrix}_{jl} = \frac{\partial [\mathbf{A}' \mathbf{y}]_{j}}{\partial \omega_{l}^{[d]}} = \frac{\partial \sum_{i=1}^{M} e^{i\vec{\omega}_{i} \cdot \vec{r}_{j}} y_{i}}{\partial \omega_{l}^{[d]}} \\
= i e^{i\vec{\omega}_{i} \cdot \vec{r}_{j}} y_{i} r_{j}^{[d]}.$$

Thus the Jacobian matrix is

$$\frac{\partial \mathbf{A}' \mathbf{y}}{\partial \boldsymbol{\omega}^{[d]}} = i \operatorname{diag} \left\{ \mathbf{r}^{[d]} \right\} \mathbf{A}' \operatorname{diag} \left\{ \mathbf{y} \right\}. \tag{6}$$

E. Gram Matrix

The product $A'(\omega) A(\omega) x$ of the Gram matrix of the NUDFT with a vector also arises in optimization steps and requires appropriate Jacobian matrices. For x:

$$rac{\partial A'Ax}{\partial x} = A'A, \qquad rac{\partial A'Ax}{\partial x^*} = 0.$$

The (k, j)th element of the $N \times N$ matrix containing the partial derivatives of the Gram matrix w.r.t. $\omega_i^{[d]}$ is

$$\left[\frac{\partial \mathbf{A}' \mathbf{A}}{\partial \omega_l^{[d]}}\right]_{k,j} = \frac{\partial}{\partial \omega_l^{[d]}} \sum_{i=1}^M e^{-\imath \vec{\omega}_i \cdot (\vec{r}_j - \vec{r}_k)}$$

$$= -\imath \left(r_j^{[d]} - r_k^{[d]}\right) e^{-\imath \vec{\omega}_l \cdot (\vec{r}_j - \vec{r}_k)}$$

$$= -\imath \left(r_j^{[d]} - r_k^{[d]}\right) a_{lk}^* a_{lj}. \tag{7}$$

In matrix form:

$$\frac{\partial \mathbf{A}' \mathbf{A}}{\partial \omega_l^{[d]}} = i \operatorname{diag} \left\{ \mathbf{r}^{[d]} \right\} \mathbf{A}' \mathbf{e}_l \mathbf{e}_l' \mathbf{A} - i \mathbf{A}' \mathbf{e}_l \mathbf{e}_l' \mathbf{A} \operatorname{diag} \left\{ \mathbf{r}^{[d]} \right\}.$$
(8)

When multiplying the Jacobian with a vector x:

$$\frac{\partial \mathbf{A}' \mathbf{A}}{\partial \omega_l^{[d]}} \mathbf{x} = i \operatorname{diag} \left\{ \mathbf{r}^{[d]} \right\} \mathbf{a}_l(\mathbf{a}'_l \mathbf{x}) - i \mathbf{a}_l \mathbf{a}'_l \operatorname{diag} \left\{ \mathbf{r}^{[d]} \right\} \mathbf{x}$$

$$= i \left(\mathbf{a}'_l \mathbf{x} \right) \left(\mathbf{r}^{[d]} \odot \mathbf{a}_l \right) - i \left(\mathbf{a}'_l (\mathbf{x} \odot \mathbf{r}^{[d]}) \right) \mathbf{a}_l, \tag{9}$$

where $a_l = A'e_l$ denotes the lth column of A'.

Consider the extended Jacobian expression:

$$\mathcal{D}_{\omega_l^{[d]}} \, oldsymbol{A}' oldsymbol{A} = ext{vec} \left(rac{\partial oldsymbol{A}' oldsymbol{A}}{\partial \omega_l^{[d]}}
ight).$$

,

Multiplying by x yields:

$$\begin{split} \frac{\partial \boldsymbol{A}' \boldsymbol{A}}{\partial \omega_{l}^{[d]}} \boldsymbol{x} &= \operatorname{vec} \left(\frac{\partial \boldsymbol{A}' \boldsymbol{A}}{\partial \omega_{l}^{[d]}} \boldsymbol{x} \right) \\ &= (\boldsymbol{x}^{T} \otimes \boldsymbol{I}_{N}) \operatorname{vec} \left(\frac{\partial \boldsymbol{A}' \boldsymbol{A}}{\partial \omega_{l}^{[d]}} \right) \qquad \text{(use P1)} \\ &= (\boldsymbol{x}^{T} \otimes \boldsymbol{I}_{N}) \left(\mathcal{D}_{\omega_{l}^{[d]}} \boldsymbol{A}' \boldsymbol{A} \right) \\ &= (\mathcal{D}_{\boldsymbol{A}' \boldsymbol{A}} \boldsymbol{A}' \boldsymbol{A} \boldsymbol{x}) \left(\mathcal{D}_{\omega_{l}^{[d]}} \boldsymbol{A}' \boldsymbol{A} \right) \qquad \text{(use P4)} \\ &= \mathcal{D}_{\omega_{l}^{[d]}} \boldsymbol{A}' \boldsymbol{A} \boldsymbol{x}. \qquad \qquad \text{(use P6)} \end{split}$$

Concatenating (9) by columns leads to the matrix

$$\left[\frac{\partial \mathbf{A}'\mathbf{A}}{\partial \omega_{1}^{[d]}} \cdots \frac{\partial \mathbf{A}'\mathbf{A}}{\partial \omega_{M}^{[d]}}\right] \mathbf{x} = -i \mathbf{A}' \operatorname{diag} \left\{ \mathbf{A}(\mathbf{x} \odot \mathbf{r}^{[d]}) \right\} + i \operatorname{diag} \left\{ \mathbf{r}^{[d]} \right\} \mathbf{A}' \operatorname{diag} \left\{ \mathbf{A} \mathbf{x} \right\}. \quad (10)$$

Alternatively, we can express the extended Jacobian as

$$\begin{bmatrix}
\frac{\partial \mathbf{A}' \mathbf{A}}{\partial \omega_{1}^{[d]}} & \dots & \frac{\partial \mathbf{A}' \mathbf{A}}{\partial \omega_{M}^{[d]}} \mathbf{x} \\
&= (\mathbf{x}^{T} \otimes \mathbf{I}_{n}) (\mathcal{D}_{\boldsymbol{\omega}^{[d]}} \mathbf{A}' \mathbf{A}) \\
&= (\mathcal{D}_{\mathbf{A}' \mathbf{A}} \mathbf{A}' \mathbf{A} \mathbf{x}) (\mathcal{D}_{\boldsymbol{\omega}^{[d]}} \mathbf{A}' \mathbf{A}) \\
&= \mathcal{D}_{\boldsymbol{\omega}^{[d]}} \mathbf{A}' \mathbf{A} \mathbf{x}.
\end{cases} (11)$$

Again we use NUFFT operations to efficiently approximate (11).

F. Inverse of Positive Semidefinite (PSD) Matrix

Image reconstruction methods based on algorithms like the augmented Lagrangian approach [38] use "data consistency" steps [39], [40], [41] that often involve least-squares problems with solutions in the following form:

$$(\boldsymbol{A}'\boldsymbol{A} + \lambda \boldsymbol{I})^{-1}\boldsymbol{x},$$

for some vector $\boldsymbol{x} \in \boldsymbol{C}^N$, or

$$(\mathbf{A}'\mathbf{A} + \lambda \mathbf{T}'\mathbf{T})^{-1}\mathbf{x},\tag{12}$$

where T denotes a linear regularization operator that is independent of ω . In both cases, $\lambda>0$ and the null spaces of T and A are disjoint, so the Hessian matrix is invertible. A few iterations of a CG method usually suffices to efficiently compute the approximate product of such a matrix inverse with a vector. The direct inverse is impractical for large-scale problems like MRI. Following [39], we treat CG as solving the above equations accurately, so that we can derive efficient approximations as follows. Otherwise, attempting to auto-differentiate through a finite number of CG iterations would require large amounts of memory. Here we derive the corresponding Jacobian matrices for the exact inverse to (12) and then apply fast approximations. For x, the $N \times N$ Jacobian is

$$\frac{\partial (\mathbf{A}'\mathbf{A} + \lambda \mathbf{T}'\mathbf{T})^{-1}\mathbf{x}}{\partial \mathbf{x}} = (\mathbf{A}'\mathbf{A} + \lambda \mathbf{T}'\mathbf{T})^{-1},$$
$$\frac{\partial (\mathbf{A}'\mathbf{A} + \lambda \mathbf{T}'\mathbf{T})^{-1}\mathbf{x}}{\partial \mathbf{x}^*} = 0.$$

We can still use CG (with NUFFT) to efficiently multiply this Jacobian by a vector, albeit approximately.

To consider the Jacobian w.r.t. the sampling pattern $\omega^{[d]}$, define $z = (A'A + \lambda T'T)^{-1}x$ and $F = A'A + \lambda T'T$. We assume that A and T have disjoint null spaces, so that F is positive definite and hence invertible. Applying equalities derived above leads to the following expression for the $M \times N$ Jacobian:

$$\mathcal{D}_{\boldsymbol{\omega}^{[d]}} \boldsymbol{F}^{-1} \boldsymbol{x}$$

$$= \left(\mathcal{D}_{\boldsymbol{F}} \boldsymbol{F}^{-1} \boldsymbol{x} \right) \left(\mathcal{D}_{\boldsymbol{\omega}^{[d]}} \boldsymbol{F} \right) \qquad \text{use P6}$$

$$= -(\boldsymbol{x}^T \otimes \boldsymbol{I}) ((\boldsymbol{F}^T)^{-1} \otimes \boldsymbol{F}^{-1}) \left(\mathcal{D}_{\boldsymbol{\omega}^{[d]}} \boldsymbol{F} \right) \qquad \text{use P5}$$

$$= -\left((\boldsymbol{x}^T (\boldsymbol{F}^T)^{-1}) \otimes \boldsymbol{F}^{-1} \right) \left(\mathcal{D}_{\boldsymbol{\omega}^{[d]}} \boldsymbol{F} \right) \qquad \text{use P2}$$

$$= -\boldsymbol{F}^{-1} (\boldsymbol{x}^T (\boldsymbol{F}^T)^{-1} \otimes \boldsymbol{I}) \left(\mathcal{D}_{\boldsymbol{\omega}^{[d]}} \boldsymbol{F} \right) \qquad \text{use P3}$$

$$= -\boldsymbol{F}^{-1} \left(\mathcal{D}_{\boldsymbol{\omega}^{[d]}} \boldsymbol{F} \boldsymbol{z} \right) \qquad \text{use P4}$$

$$= -(\boldsymbol{A}' \boldsymbol{A} + \lambda \boldsymbol{T}' \boldsymbol{T})^{-1} \left(-\imath \boldsymbol{A}' \operatorname{diag} \left\{ \boldsymbol{A} (\boldsymbol{z} \odot \boldsymbol{r}^{[d]}) \right\} \right)$$

$$+ \imath \operatorname{diag} \left\{ \boldsymbol{r}^{[d]} \right\} \boldsymbol{A}' \operatorname{diag} \left\{ \boldsymbol{A} \boldsymbol{z} \right\} \right) \qquad \text{use (11)}.$$

$$(13)$$

We apply this Jacobian to a vector by using four NUFFT operations followed by running CG to approximate the product of \mathbf{F}^{-1} times a vector. Notably, the memory cost of (13) is constant w.r.t the number of iterations, whereas the standard auto-differentiation approach has linear memory cost. Using the proposed method, one may apply enough iterations to ensure convergence to a desired tolerance. This new fast and low-memory Jacobian approximation is particularly important for the MRI applications shown in the following sections. Without this approximation, memory cost can be prohibitively large.

G. Sensitivity Maps

In multi-coil (parallel) acquisition, the MRI system model contains another linear operator

$$oldsymbol{S} = egin{bmatrix} oldsymbol{S}_1 \ dots \ oldsymbol{S}_{N_2} \end{bmatrix},$$

where $S_i = \operatorname{diag}\{s_i\}$ denotes a diagonal matrix containing the receiver coil sensitivity map [29]. The total number of receiver channels is N_c . The system matrix (E) for MRI in this case becomes ($I_{N_c} \otimes A$)S. Because of the special block-diagonal structure of S, all the Jacobian matrices in previous sections still hold by simply replacing A with E.

The Jacobian derivations are as follows. For the forward operator (§II-C), one can show

$$egin{aligned} rac{\partial oldsymbol{E} oldsymbol{x}}{\partial oldsymbol{\omega}^{[d]}} &= egin{aligned} egin{aligned} oldsymbol{A} oldsymbol{S}_{N_c} oldsymbol{x} \ & dots \ -\imath \operatorname{diag} ig\{ oldsymbol{A} oldsymbol{s}_{N_c} \odot oldsymbol{x} \odot oldsymbol{r}^{[d]} ig) ig\} \ &= \imath \operatorname{diag} ig\{ oldsymbol{E} oldsymbol{x} \odot oldsymbol{A} oldsymbol{S} oldsymbol{x} \odot oldsymbol{r}^{[d]} ig) ig\} \ &= \imath \operatorname{diag} ig\{ oldsymbol{E} oldsymbol{x} \odot oldsymbol{r}^{[d]} ig) ig\} \,. \end{aligned}$$

The adjoint operator (§II-D) follows the same proof and produces

$$rac{\partial oldsymbol{E}'oldsymbol{y}}{\partial \omega^{[d]}} = \imath \operatorname{diag} \Bigl\{ oldsymbol{r}^{[d]} \Bigr\} oldsymbol{E}' \mathrm{diag} \{oldsymbol{y}\} \,.$$

For the gram operator (§II-E) we have

$$\frac{\partial \mathbf{E}' \mathbf{E} \mathbf{x}}{\partial \boldsymbol{\omega}^{[d]}} = \sum_{i} \frac{\partial \mathbf{S}'_{i} \mathbf{A}' \mathbf{A} \mathbf{S}_{i} \mathbf{x}}{\partial \boldsymbol{\omega}^{[d]}} = \sum_{i} \mathbf{S}'_{i} \frac{\partial \mathbf{A}' \mathbf{A} \mathbf{S}_{i} \mathbf{x}}{\partial \boldsymbol{\omega}^{[d]}}$$

$$= \sum_{i} -i \mathbf{S}'_{i} \mathbf{A}' \operatorname{diag} \left\{ \mathbf{A} (\mathbf{S} \mathbf{x} \odot \mathbf{r}^{[d]}) \right\}$$

$$+ i \mathbf{S}'_{i} \operatorname{diag} \left\{ \mathbf{r}^{[d]} \right\} \mathbf{A}' \operatorname{diag} \left\{ \mathbf{A} \mathbf{S}_{i} \mathbf{x} \right\}$$

$$= \sum_{i} -i \mathbf{S}'_{i} \mathbf{A}' \operatorname{diag} \left\{ \mathbf{A} (\mathbf{S} \mathbf{x} \odot \mathbf{r}^{[d]}) \right\}$$

$$+ i \operatorname{diag} \left\{ \mathbf{r}^{[d]} \right\} \mathbf{S}'_{i} \mathbf{A}' \operatorname{diag} \left\{ \mathbf{A} \mathbf{S}_{i} \mathbf{x} \right\}$$

$$= -i \mathbf{E}' \operatorname{diag} \left\{ \mathbf{E} (\mathbf{x} \odot \mathbf{r}^{[d]}) \right\}$$

$$+ i \operatorname{diag} \left\{ \mathbf{r}^{[d]} \right\} \mathbf{E}' \operatorname{diag} \left\{ \mathbf{E} \mathbf{x} \right\}. \tag{14}$$

For the inverse of the PSD matrix (§II-F), let $G = E'E + \lambda T'T$ and $z = G^{-1}x$ (in the usual case where the regularizer matrix T is designed such that G is invertible). Combining (13) and (14) produces:

$$\begin{split} \frac{\partial \left(\boldsymbol{E}' \boldsymbol{E} + \lambda \boldsymbol{T}' \boldsymbol{T} \right)^{-1} \boldsymbol{x}}{\partial \boldsymbol{\omega}^{[d]}} \\ &= -\boldsymbol{G}^{-1} \left(\boldsymbol{x}^T (\boldsymbol{G}^T)^{-1} \otimes \boldsymbol{I} \right) \mathcal{D}_{\boldsymbol{\omega}^{[d]}} \boldsymbol{G} \\ &= - (\boldsymbol{E}' \boldsymbol{E} + \lambda \boldsymbol{T}' \boldsymbol{T})^{-1} \Big(- \imath \, \boldsymbol{E}' \mathrm{diag} \Big\{ \boldsymbol{E} (\boldsymbol{z} \odot \boldsymbol{r}^{[d]}) \Big\} \\ &+ \imath \, \mathrm{diag} \Big\{ \boldsymbol{r}^{[d]} \Big\} \, \boldsymbol{E}' \mathrm{diag} \big\{ \boldsymbol{E} \boldsymbol{z} \big\} \Big). \end{split}$$

Again, we apply this Jacobian matrix to a vector by combining NUFFTs and CG.

H. Field Inhomogeneity

For MRI scans with long readouts, one should also consider the effects of off-resonance (e.g., B_0 field inhomogeneity), in which case the system matrix elements are given by [5]

$$a_{ij} = e^{-i\vec{\omega}_i \cdot \vec{r}_j} e^{-iw_j t_i},$$

where w_j denotes the field map value at the jth voxel and t_i is the time of the ith readout sample.

This form is no longer a Fourier transform operation, but there are fast and accurate approximations [42] that enable the use of $O(N\log N)$ NUFFT steps and avoid the very slow $O(N^2)$ matrix-vector multiplication. Such approximations of system matrix \boldsymbol{E} usually have the form:

$$oldsymbol{E}_{\mathrm{f}} pprox \sum_{l=1}^{L} \mathrm{diag}\{b_{il}\} \, oldsymbol{A}(oldsymbol{\omega}) \, \mathrm{diag}\{c_{lj}\} \, ,$$

where A denotes the usual (possibly non-uniform) DFT that is usually approximated by a NUFFT, $b_{il} \in \mathbb{C}^M$, and $b_{il} \in \mathbb{C}^N$. It is relatively straightforward to generalize the Jacobian expressions in this paper to handle the case of field inhomogeneity, by simply replacing A with E_f , similar to the sensitivity map case.

III. OPTIMIZING SAMPLING PATTERNS

For modern MRI systems, the sampling trajectory ω is a programmable parameter. Traditionally ω is a geometrical curve controlled by few parameters (such as radial spokes or spiral leaves), and its tuning relies on derivative-free optimizers such as grid-search. In this paper, we optimize ω by minimizing a training loss function from image reconstruction, where the descent direction of ω is the negative gradient of that loss [22], [24]. We adopt such a "reconstruction loss" because the terminal goal of sampling pattern optimization is to improve image quality. To learn from large datasets, the optimization uses stochastic gradient descent (SGD)-like algorithms. Additionally, the loss function may include other terms, such as a penalty on the maximum gradient strength and slew rate [22], [24] or peripheral nerve stimulation effects [43].

For image reconstruction, consider a convex and smooth regularizer $R(\cdot)$ for simplicity. Since the noise statistics are Gaussian, a typical regularized cost function used for model-based image reconstruction is [5]

$$\Psi(x) = \frac{1}{2} ||Ax - y||_2^2 + R(x).$$
 (15)

During training, the observation y can be retrospectively simulated using $y = A(\omega)x^{\text{true}}$. For illustration, consider applying the kth step of gradient descent (GD) to that cost function:

$$x_{k+1} = x_k - \alpha \nabla \Psi(x_k)$$

= $x_k - \alpha A(\omega)' (A(\omega) x_k - y) - \alpha \nabla R(x_k),$

where $x_k \in \mathbb{C}^N$, $\alpha \in \mathbb{R}^+$ is the step size. After K iterations, we have a reconstructed image (batch) $x_K = x_K(\omega) = f(\omega, y)$, where the reconstruction method $f(\omega, y)$ is a function of both the data y and the sampling pattern ω . To learn/update the sampling pattern ω , consider a simple loss function for a single training example:

$$L(\boldsymbol{\omega}) = \left\| \boldsymbol{x}_K(\boldsymbol{\omega}) - \boldsymbol{x}^{\text{true}} \right\|_2^2$$
 (16)

where $\boldsymbol{x}^{\text{true}}$ is the reference fully-sampled image (batch). Learning $\boldsymbol{\omega}$ via backpropagation (or chain-rule) requires differentiating L w.r.t. the sampling pattern $\boldsymbol{\omega}$, which in turn involves Jacobians of quantities like $\boldsymbol{A}(\boldsymbol{\omega})$ that we derived above.

Here we use the forward operator (§II-C) as an example to illustrate one step in propagation. As needed in a backpropagation step (Jacobian-vector product, JVP), the Jacobian (5) is multiplied with a gradient vector $\boldsymbol{v} = \frac{\partial L}{\partial (\boldsymbol{A}\boldsymbol{x})^*} \in \mathbb{C}^M$ calculated in the prior step. Using (1), the corresponding JVP is

$$\frac{\partial L}{\partial w} = \text{real}\left\{ (-\imath \ \mathbf{A}(\mathbf{x} \odot \mathbf{r}^{[d]}))' \odot \mathbf{v} \right\}. \tag{17}$$

Efficiently computation can simply apply a NUFFT operation to $\boldsymbol{x} \odot \boldsymbol{r}^{[d]}$, followed by a point-wise multiplication with \boldsymbol{v} . The Gram and PSD inverse ("data consistency") term in §II-E and §II-F follow a similar pattern during backpropagation. See our open-source codes for implementation details.

¹https://github.com/guanhuaw/Bjork

Although we illustrate the GD algorithm with a simple smooth regularizer, more generally, the reconstruction method $f(\omega, y)$ can involve more sophisticated regularizers such as neural networks [22], [24] or non-smooth sparsity models [44] used in compressed sensing. In such cases, backpropagation uses sub-gradients, instead of gradients, as is common in stochastic optimization. The loss JVPs are backpropagated through iterative reconstruction steps to compute a gradient w.r.t. ω .

The proposed approach is applicable only to non-Cartesian MRI, because Cartesian sampling pattern design is usually a discrete optimization problem, incompatible with gradient-based methods. However, one could optimize phase-encoding locations continuously (in 2D or 3D) with the frequency-encoding direction being fully sampled, which is a hybrid Cartesian / non-Cartesian approach [23].

IV. EXPERIMENTS

This section validates the accuracy and efficiency of the proposed methods. It also showcases the application to MRI sampling trajectory optimization.

A. Accuracy and Efficiency

The appendix discusses error bounds for the Jacobian approximations of §II-C and §II-D.

We performed numerical experiments to examine the following test cases:

$$\frac{\partial \left\| f(\boldsymbol{x}) \right\|_2^2}{\partial \boldsymbol{\omega}^{[d]}} \text{ and } \frac{\partial \left\| f(\boldsymbol{x}) \right\|_2^2}{\partial \boldsymbol{x}^*},$$

where $f(\cdot)$ denotes multiplication by A, by the Gram matrix A'A, or by the 'inverse of PSD matrix' (§II-F) of sensitivityinformed NUFFTs (§II-G). The Gram and inverse experiments implicitly test the adjoint operator's approximations. The xadopted is a 40×40 patch cropped from the center of a Shepp-Logan phantom with random additional phases uniformly distributed in $[-\pi, \pi]$. S is a simulated 8-channel sensitivity map, and ω is one radial spoke crossing the k-space center. The Jacobian calculation methods are: (1) auto-differentiation of NUFFT; the lookup table operation [27] is replaced by bilinear interpolation to enable autodifferentiation, similar to [22], (2) our approximation described above, (3) auto-differentiation of exact non-uniform discrete Fourier transform (NUDFT), implemented with single precision. We regard method 3 (NUDFT) as the ground truth. Since NUDFT (in its simplest form) involves only one exponential function, multiplication and addition for each element, its backpropagation introduces minimal numerical errors. For the PSD inverse, we applied 20 CG iterations for all three methods, which was sufficiently close to convergence based on the residual norm ||r||/||b|| (the definition follows [45, (45)]).

Fig. 1 and Fig. 2 illustrate representative profiles of the gradients w.r.t. \boldsymbol{x} and $\boldsymbol{\omega}$. For $\boldsymbol{\omega}$, the auto-differentiation (method 1) approach has larger deviations from method 3 (NUDFT) because of the non-differentiability of interpolation operations w.r.t. coordinates. For the gradient w.r.t. \boldsymbol{x} , both method 1 and method 2 generate accurate results for forward

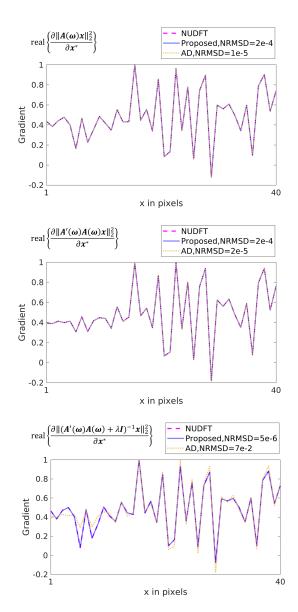


Fig. 1. Examples of gradients w.r.t. \boldsymbol{x}^* (the real part is plotted). Plots show one representative row of a 40×40 matrix (rescaled to [-1,1]). The rows are the forward, Gram, and PSD inverse operator cases. The horizontal axis is the pixel index. The legend reports the normalized root-mean-square difference (NRMSD) compared with the reference NUDFT calculation.

and Gram operators. The reason is that in method 1 (auto-diff), the interpolation operation w.r.t \boldsymbol{x} is linear, hence accurately differentiable. For the PSD inverse, method 1 led to a slightly inaccurate gradient, stemming from the accumulated errors of backpropagating CG iterations.

Table I and Table II compare the time and memory cost of methods 1 (auto-diff) and 2 (proposed). The CPU is Intel(R) Xeon(R) Gold 6138 CPU @ 2.00GHz and the GPU is an Nvidia(R) RTX2080Ti. We used PyTorch 1.9.1 and torchkbnufft 1.1.0. The memory usage was tracked by torch.cuda.max_memory_allocated on Nvidia GPUs. We implemented the numerical experiments with torchkbnufft² [46] and MIRTorch³ toolboxes.

²https://github.com/mmuckley/torchkbnufft

³https://github.com/guanhuaw/MIRTorch

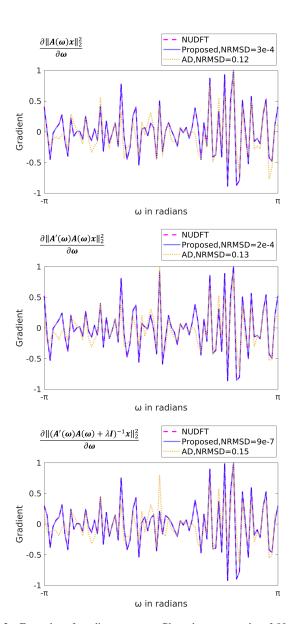


Fig. 2. Examples of gradients w.r.t. ω . Plots show one spoke of 80 points (rescaled to [-1,1]). The rows are the forward, Gram, and PSD inverse operator cases. The proposed approximation better matches the gradient of the NUDFT. The legend reports the normalized root-mean-square difference (NRMSD) compared with the reference NUDFT calculation. The proposed approach has at least $400\times$ smaller NRMSD for this nonlinear case.

Our method is much faster than auto-differentiation on both GPUs and CPUs, and uses less memory. Importantly, the PSD inverse Jacobian is impractical for the 3D case, whereas the proposed approach fit comfortably in GPU's onboard memory.

B. MRI Trajectory Optimization

This experiment optimized the MRI sampling trajectory using the proposed Jacobian approximations and stochastic optimization. The reconstruction methods (15) here consider two types of algorithms, namely smooth (regularized) least-squares reconstruction and sparsity-based reconstruction.

The smooth reconstruction method uses the cost function

$$\Psi(x) = \frac{1}{2} \|E(\omega)x - y\|_2^2 + \frac{\lambda}{2} \|Tx\|_2^2,$$

TABLE I COMPUTATION TIME OF THE TEST CASE.

	Gram		Inverse		
	auto-diff	proposed	auto-diff	proposed	
Large image - GPU	0.3s	0.2s	4.3s	2.5s	
Small image - GPU	0.1s	0.1s	1.3s	0.9s	
Large image - CPU	5.2s	1.7s	276.2s	48.5s	
Small image - CPU	0.8s	0.5s	27.4s	6.9s	

Large size: 400×400 ; small size: 40×40

20 CG iterations were applied in the PSD inverse cases.

TABLE II
MEMORY USE OF THE TEST CASE.

	Gram		Inverse		
	standard	proposed	standard	proposed	
Small	3.1MB	2.8MB	145.7MB	2.9MB	
Large	375.9MB	267.5MB	5673.2MB	272.0MB	
3D	N/A	10.1GB	N/A	10.8GB	

Large size: 400×400 ; small size: 40×40 ;

3D size: $200 \times 200 \times 100$.

N/A: the memory usage was too large for a single GPU.

20 CG iterations were applied in the PSD inverse cases.

where T is a finite-difference operator encouraging smoothness. Correspondingly, the reconstructed image is:

$$\boldsymbol{x}_K = (\boldsymbol{E}'\boldsymbol{E} + \lambda \boldsymbol{T}'\boldsymbol{T})^{-1}\boldsymbol{E}'\boldsymbol{y},$$

which we solved using CG. The following sections refer to this method as quadratic penalized least-squares (QPLS). We also implemented a simpler case, where T = I, which is referred as CG-SENSE [47]. In both scenarios, we set λ to 10^{-3} empirically and still applied 20 CG iterations. The initialization of CG used the density compensated reconstruction [48].

The sparsity-based compressed (CS) sensing algorithm adopts a wavelets-based sparsity penalty, and has the following objective function

$$\Psi(x) = \frac{1}{2} ||E(\omega)x - y||_2^2 + \lambda ||Wx||_1,$$

where W is an orthogonal DWT matrix and we set $\lambda = 10^{-5}$ empirically. We used 40 iterations of the proximal optimized gradient method (POGM) [49], [50] to solve this non-smooth optimization problem.

For the purpose of comparing trajectories and image quality, we also applied the proposed approximations to an unrolled neural network (UNN) reconstruction method that followed the definition of [39] ([24] extends it to non-Cartesian cases.) We used the same network configuration as in [24].

To optimize the k-space trajectory for each of these reconstruction methods, the training loss (16) is:

$$L(\boldsymbol{\omega}) = \|\boldsymbol{x}_K(\boldsymbol{\omega}) - \boldsymbol{x}^{\text{true}}\|_2^2 + \mu_1 \boldsymbol{\phi}_{\gamma \Delta t g_{\text{max}}}(|\boldsymbol{D}_1 \boldsymbol{\omega}|) + \mu_2 \boldsymbol{\phi}_{\gamma \Delta t^2 s_{\text{max}}}(|\boldsymbol{D}_2 \boldsymbol{\omega}|),$$

where $\boldsymbol{x}^{\mathrm{true}}$ is the conjugate phase reconstruction of fully sampled Cartesian data [51]. The second and third terms applied a soft constraint on gradient strength and slew rate according to [24, Eqn. 2], where $\phi_{\lambda}(|\boldsymbol{x}|) = \sum \max(|\boldsymbol{x}| - \lambda, 0)$. The maximum gradient strength (g_{\max}) was 5 Gauss/cm and the maximum slew rate (s_{\max}) was 15 Gauss/cm/ms. $\mu_1 = \mu_2 = 10$. We estimated sensitivity maps in \boldsymbol{E} using ESPIRiT

[52], and simulated noiseless raw signals $y = E(\omega)x^{\rm true}$ retrospectively w.r.t. candidate trajectories. The training used the fastMRI brain dataset [53] containing 15902 T1w slices, 16020 T2w slices, and 3311 FLAIR slices cropped to size 320×320 . The number of coils ranges from 2 to 28. We used the Adam optimizer [54], with step size 10^{-4} and minibatch size 12. We used 6 epochs for training model-based methods (CG-SENSE, QPLS and CS) and 60 epochs for the UNN training. The initialization of learned trajectories was an under-sampled radial trajectory in all experiments. The initialization had 16 "spokes" and each spoke was 5ms long with 1280 sampling points. We also adopted the k-space parameterization trick detailed in [24, Eqn. 3] to avoid suboptimal local minima, and parameterized each shot with 40 quadratic spline kernels.

Fig. 3 showcases the trajectories optimized for each of the reconstruction methods. The centers of trajectories optimized with quadratic regularizers (CG-SENSE and QPLS) are not aligned with the k-space origin. We hypothesize that regularizers (and corresponding iterative algorithms) handle image phases differently, resulting in distinct trajectory centers.

Table III reports the average image reconstruction quality (PSNR and SSIM [55], fully sampled image as the ground truth) on 500 test slices. It also showcases the image quality of these learned trajectories with reconstruction methods different from the training phase. All learned trajectories led to improved reconstruction quality compared to the initial radial trajectory (unopt.), even with different reconstruction methods. Importantly, the same reconstruction algorithm across training and test led to the greatest improvement (the bold diagonal entries). Fig. 4 shows reconstruction examples.

TABLE III
AVERAGE RECONSTRUCTION QUALITY ON TEST SET WITH TRAJECTORIES
LEARNED FOR DIFFERENT RECONSTRUCTION METHODS.

SSIM							
Test	QPLS	SENSE	CS	UNN	unopt.		
QPLS	0.963	0.963	0.962	0.961	0.947		
SENSE	0.964	0.964	0.963	0.961	0.946		
CS	0.962	0.963	0.966	0.964	0.946		
UNN	0.960	0.960	0.958	0.964	0.950		
PSNR (in dB) Learn Test	QPLS	SENSE	CS	UNN	unopt.		
OPLS	35.1	35.1	34.9	35.0	33.1		
SENSE	35.2	35.2	34.9	35.1	33.1		
CS	34.8	34.9	35.4	35.2	33.0		
UNN	34.6	34.6	34.5	35.0	33.5		
Learn: the reconstruction method the trajectory is jointly trained with							

Learn: the reconstruction method the trajectory is jointly trained with. Test: the reconstruction method in the test phase.

unopt.: the undersampled radial trajectory (unoptimized initialization).

C. Accelerated Learning with Low-Accuracy NUFFT

The major computation cost of trajectory learning is proportional to NUFFTs and their Jacobian calculations. An empirical acceleration method is to use faster NUFFT approximations (low over-sampling factors and/or small interpolation neighborhoods) in training. Later, when the learned trajectory is deployed on test data or prospectively acquired data, one

could use default NUFFT accuracy. We investigated learning trajectories with two different NUFFT accuracies: (1) gridding size = $1.25 \times$ image size, interpolation kernel size = 5 and (2) gridding size = $2 \times$ image size, interpolation kernel size = 6 which is a commonly used setting. On our GPUs, the loweraccuracy setting was $1.4 \times$ faster than the higher-accuracy one. We used the CS-based reconstruction and corresponding training strategy described in the previous subsection. Fig. 5 shows the trajectory optimized for the two NUFFT accuracy levels. To compare the trajectory optimized by these two settings, we used the reconstruction image quality as the evaluation metric. We simulated and reconstructed images using the two trajectories on the test data (same as the previous experiment). The trajectories optimized with the "low accuracy" and "high accuracy" NUFFT had mean PSNR values of 35.4±4.6 dB and 35.4 ± 4.7 dB.

V. DISCUSSION

This paper presents a model-based approximation of Jacobian matrices involving NUFFTs. Compared to direct autodifferentiation, the proposed method is faster, needs less memory, and better approximates the reference NUDFT results. As discussed in II-B, the error of auto-differentiation is not a software limitation, but rather a problem that stems from the non-differentiability of interpolation or lookup table operations. NUFFT alternatives such as re-gridding or filtered back-projection also suffer from similar non-differentiability issues and are less effective than (NUFFT-based) iterative reconstruction. Our previous studies [24, Fig. 14] compared the trajectory optimization results of the proposed method and standard auto-differentiation. The trajectory optimized by the proposed approximation generated superior image quality, and conformed better to the empirical criteria [56], [57]: sampling points should not overlap or be too distant from each other.

Sampling patterns learned with different reconstruction methods showed distinct characteristics in §IV-B. This phenomenon was also observed in previous literature [8], [10]. The differences in sampling patterns may stem from different regularizers, as well as different iterative algorithms. Importantly, as shown in Table III, synergistic sampling and reconstruction led to the best image quality. Several previous studies [22], [24], [23] only used NN-based reconstruction methods. while the stability and generalizability of NN-based reconstruction are still being investigated. In comparison, using our method delineated in §III, one may optimize trajectories for model-based reconstruction methods that may be more robust. Our results show that with a suitably tailored sampling pattern, traditional model-based reconstruction can compete with NNbased reconstruction, reinforcing related observations in recent studies [58]. Additionally, sampling optimization for modelbased reconstruction requires less training data than for NNbased reconstruction. This property is beneficial for medical imaging where the data availability is often limited.

The training used discrete-space image datasets, whereas the actual objects in practice are continuous. Ideally, using an accurate continuous image model could better approximate the actual situation. This implicit bias is common for learningbased methods, and may lead to suboptimal results, such as

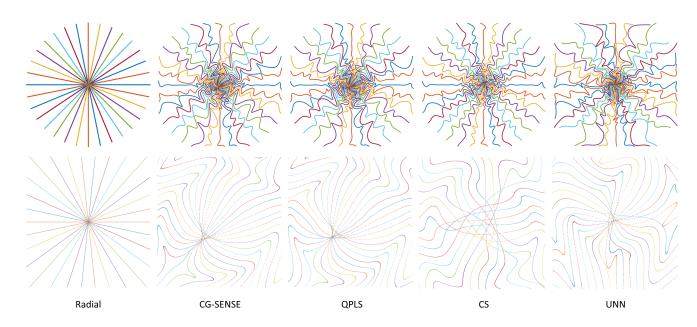


Fig. 3. Optimized sampling trajectories for several iterative reconstruction methods. The left column shows the uniform radial initialization. The second row shows the $8\times$ zoomed-in central k-space.

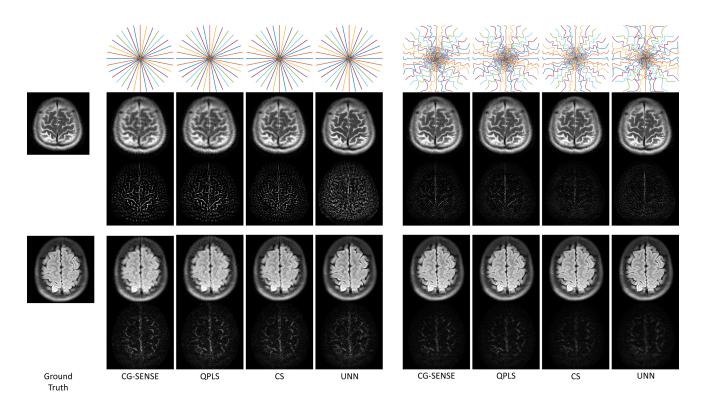


Fig. 4. Examples of the reconstructed images with unoptimized (left) and optimized trajectories (right). Rows 3 and 5 show corresponding error maps.

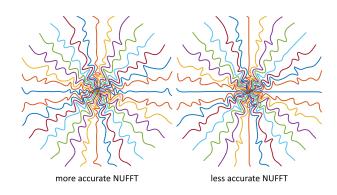


Fig. 5. Learned trajectories with different NUFFT accuracies.

the backtracking in the edge/corner of k-space (Fig. 3). The training also ignored physical processes such as relaxation and magnetization transfer. Future studies may consider these processes in the forward system model. The mismatch or domain shift from training to prospective scans may influence the results. For example, there exist differences in protocols (RF pulses, FOVs, and resolutions), hardware (field strengths and Tx/Rx coils), system imperfections (eddy currents, gradient non-linearity, and inhomogeneity), demography, and pathology. Our previous studies [24] tested the optimized trajectory in a prospective in-vivo experiment, and discussed practical issues, including eddy currents, and contrast/SNR mismatch between the training set and prospective protocols. Subsequent studies should evaluate the robustness of learned sampling trajectories in more scenarios.

ACKNOWLEDGMENT

The authors gratefully thank Dr. Douglas Noll, Dr. Tianrui Luo, Naveen Murthy, Yuran Zhu, and the anonymous reviewers for helpful advice.

REFERENCES

- [1] D. C. Munson and J. L. Sanz, "Image reconstruction from frequency-offset Fourier data," *Proc. IEEE*, vol. 72, no. 6, pp. 661–9, Jun. 1984.
- [2] M. M. Bronstein, A. M. Bronstein, M. Zibulevsky, and H. Azhari, "Reconstruction in diffraction ultrasound tomography using nonuniform FFT," *IEEE Trans. Med. Imag.*, vol. 21, no. 11, pp. 1395–1401, Nov. 2002.
- [3] S. Matej, J. A. Fessler, and I. G. Kazantsev, "Iterative tomographic image reconstruction using Fourier-based forward and back- projectors," *IEEE Trans. Med. Imag.*, vol. 23, no. 4, pp. 401–12, Apr. 2004.
- [4] K. L. Wright, J. I. Hamilton, M. A. Griswold, V. Gulani, and N. Seiberlich, "Non-Cartesian parallel imaging reconstruction," *J. Magn. Reson. Imag.*, vol. 40, no. 5, pp. 1022–1040, 2014.
- [5] J. A. Fessler, "Model-based image reconstruction for MRI," *IEEE Sig. Proc. Mag.*, vol. 27, no. 4, pp. 81–9, Jul. 2010.
- [6] J. A. Fessler and B. P. Sutton, "Nonuniform fast Fourier transforms using min-fmax interpolation," *IEEE Trans. Sig. Proc.*, vol. 51, no. 2, pp. 560–74, Feb. 2003.
- [7] Z. Yang and M. Jacob, "Mean square optimal NUFFT approximation for efficient non-Cartesian MRI reconstruction," *J. Mag. Res.*, vol. 242, pp. 126–35, May 2014.
- [8] M. V. W. Zibetti, G. T. Herman, and R. R. Regatte, "Fast data-driven learning of parallel MRI sampling patterns for large scale problems," *Sci Rep*, vol. 11, no. 1, p. 19312, Sep. 2021.
- [9] T. Sanchez et al., "Scalable learning-based sampling optimization for compressive dynamic MRI," in 2020 IEEE Intl. Conf. on Acous., Speech and Sig. Proc. (ICASSP), May 2020, pp. 8584–8588.

- [10] B. Gözcü, T. Sanchez, and V. Cevher, "Rethinking Sampling in Parallel MRI: A Data-Driven Approach," in 2019 27th Euro. Sig. Proc. Conf. (EUSIPCO), Sep. 2019, pp. 1–5.
- [11] Y. Cao and D. N. Levin, "Feature-recognizing MRI," Magn. Reson. Med., vol. 30, no. 3, pp. 305–317, 1993.
- [12] F. Knoll, C. Clason, C. Diwoky, and R. Stollberger, "Adapted random sampling patterns for accelerated MRI," Mag. Res. Mat. Phys. Bio. Med., vol. 24, no. 1, pp. 43–50, Feb. 2011.
- [13] C. D. Bahadir, A. Q. Wang, A. V. Dalca, and M. R. Sabuncu, "Deep-learning-based optimization of the under-sampling pattern in MRI," *IEEE Trans. Comput. Imag.*, vol. 6, pp. 1139–1152, 2020.
- [14] F. Sherry, M. Benning, J. C. D. . Reyes, M. J. Graves, G. Maierhofer, G. Williams, C.-B. Schonlieb, and M. J. Ehrhardt, "Learning the sampling pattern for MRI," *IEEE Trans. Med. Imag.*, vol. 39, no. 12, pp. 4310–21, Dec. 2020.
- [15] I. A. M. Huijben, B. S. Veeling, and R. J. G. van Sloun, "Learning sampling and model-based signal recovery for compressed sensing MRI," in 2020 IEEE Intl. Conf. on Acous., Speech and Sig. Proc. (ICASSP), May 2020, pp. 8906–8910.
- [16] M. Seeger, H. Nickisch, R. Pohmann, and B. Schölkopf, "Optimization of k-space trajectories for compressed sensing by Bayesian experimental design," *Magn. Reson. Med.*, vol. 63, no. 1, pp. 116–126, 2010.
- [17] J. P. Haldar and D. Kim, "OEDIPUS: An Experiment Design Framework for Sparsity-Constrained MRI," *IEEE Trans. Med. Imaging*, vol. 38, no. 7, pp. 1545–1558, Jul. 2019.
- [18] M. von Kienlin and R. Mejia, "Spectral localization with optimal pointspread function," *J. Magn. Reson.* (1969), vol. 94, no. 2, pp. 268– 287, Sep. 1991.
- [19] Y. Gao and S. Reeves, "Optimal k-space sampling in MRSI for images with a limited region of support," *IEEE Trans. Med. Imag.*, vol. 19, no. 12, pp. 1168–1178, Dec. 2000.
- [20] D. Xu, M. Jacob, and Z. Liang, "Optimal sampling of k-space with cartesian grids for parallel MR imaging," in *Proc. Intl. Soc. Mag. Res.* Med., 2005, p. 2450.
- [21] E. Levine and B. Hargreaves, "On-the-fly adaptive k-Space sampling for linear MRI reconstruction using moment-based spectral analysis," *IEEE Trans. Med. Imag.*, vol. 37, no. 2, pp. 557–567, 2017.
- [22] T. Weiss, O. Senouf, S. Vedula, O. Michailovich, M. Zibulevsky, and A. Bronstein, "PILOT: Physics-informed learned optimized trajectories for accelerated MRI," MELBA, pp. 1–23, 2021.
- [23] H. K. Aggarwal and M. Jacob, "J-MoDL: Joint model-based deep learning for optimized sampling and reconstruction," *IEEE J. Sel. Top. Sig. Proc.*, vol. 14, no. 6, pp. 1151–62, Oct. 2020.
- [24] G. Wang, T. Luo, J.-F. Nielsen, D. C. Noll, and J. A. Fessler, "B-Spline Parameterized Joint Optimization of Reconstruction and K-Space Trajectories (BJORK) for Accelerated 2D MRI," *IEEE Trans. Med. Imag.*, vol. 41, no. 9, pp. 2318–2330, Sep. 2022.
- [25] E. Scope Crafts, H. Lu, H. Ye, L. L. Wald, and B. Zhao, "An efficient approach to optimal experimental design for magnetic resonance fingerprinting with B-splines," *Magn. Reson. Med.*, vol. 88, no. 1, pp. 239–253, 2022.
- [26] S. P. Jordan, S. Hu, I. Rozada, D. F. McGivney, R. Boyaciouglu, D. C. Jacob, S. Huang, M. Beverland, H. G. Katzgraber, M. Troyer et al., "Automated design of pulse sequences for magnetic resonance fingerprinting using physics-inspired optimization," Proc. Natl. Acad. Sci., vol. 118, no. 40, p. e2020516118, 2021.
- [27] B. Dale, M. Wendt, and J. L. Duerk, "A rapid look-up table method for reconstructing MR images from arbitrary K-space trajectories," *IEEE Trans. Med. Imag.*, vol. 20, no. 3, pp. 207–17, Mar. 2001.
- [28] P. J. Beatty, D. G. Nishimura, and J. M. Pauly, "Rapid gridding reconstruction with a minimal oversampling ratio," *IEEE Trans. Med. Imag.*, vol. 24, no. 6, pp. 799–808, Jun. 2005.
- [29] K. P. Pruessmann, M. Weiger, M. B. Scheidegger, and P. Boesiger, "SENSE: sensitivity encoding for fast MRI," *Magn. Reson. Med.*, vol. 42, no. 5, pp. 952–962, 1999.
- [30] B. Sutton, D. Noll, and J. Fessler, "Fast, iterative image reconstruction for MRI in the presence of field inhomogeneities," *IEEE Trans. Med. Imaging*, vol. 22, no. 2, pp. 178–188, Feb. 2003.
- [31] M. Zehni, L. Donati, E. Soubies, Z. Zhao, and M. Unser, "Joint angular refinement and reconstruction for single-particle cryo-EM," *IEEE Trans. Im. Proc.*, vol. 29, pp. 6151–63, 2020.
- [32] G. Wang, D. C. Noll, and J. A. Fessler, "Reconstruction may benefit from tailored sampling trajectories: optimizing non-Cartesian trajectories for model-based reconstruction," in *Proc. Intl. Soc. Mag. Res. Med.*, 2022, p. 5011. [Online]. Available: https://submissions. mirasmart.com/ISMRM2022/Itinerary/Files/PDFFiles/5011.html

- [33] R. Remmert, Theory of complex functions. Springer Science & Business Media, 1991, vol. 122.
- [34] K. Kreutz-Delgado, "The complex gradient operator and the CR-calculus," arXiv preprint arXiv:0906.4835, 2009.
- [35] A. Hjorungnes and D. Gesbert, "Complex-valued matrix differentiation: Techniques and key results," *IEEE Trans. Sig. Proc.*, vol. 55, no. 6, pp. 2740–2746, 2007.
- [36] Wikipedia contributors, "Numerator layout notation," 2022, [Accessed 2022-11-08]. [Online]. Available: https://en.wikipedia.org/wiki/Matrix_ calculus#Numerator-layout_notation
- [37] J. R. Magnus and H. Neudecker, Matrix differential calculus with applications in statistics and econometrics. John Wiley & Sons, 2019.
- [38] M. R. Hestenes, "Multiplier and gradient methods," J. Optim. Theory Appl., vol. 4, no. 5, pp. 303–320, Nov. 1969.
- [39] H. K. Aggarwal, M. P. Mani, and M. Jacob, "MoDL: model-based deep learning architecture for inverse problems," *IEEE Trans. Med. Imag.*, vol. 38, no. 2, pp. 394–405, Feb. 2019.
- [40] S. Ramani and J. A. Fessler, "Parallel MR image reconstruction using augmented Lagrangian methods," *IEEE Trans. Med. Imag.*, vol. 30, no. 3, pp. 694–706, Mar. 2011.
- [41] S. H. Chan, X. Wang, and O. A. Elgendy, "Plug-and-play ADMM for image restoration: fixed-point convergence and applications," *IEEE Trans. Computational Imaging*, vol. 3, no. 1, pp. 84–98, Mar. 2017.
- [42] J. A. Fessler, S. Lee, V. T. Olafsson, H. R. Shi, and D. C. Noll, "Toeplitz-based iterative image reconstruction for MRI with correction for magnetic field inhomogeneity," *IEEE Trans. Sig. Proc.*, vol. 53, no. 9, pp. 3393–402, Sep. 2005.
- [43] G. Wang, J.-F. Nielsen, J. A. Fessler, and D. C. Noll, "Stochastic optimization of 3D non-Cartesian sampling trajectory (SNOPY)," 2022. [Online]. Available: http://arxiv.org/abs/2209.11030
- [44] M. Lustig, D. L. Donoho, J. M. Santos, and J. M. Pauly, "Compressed sensing MRI," *IEEE Signal Process. Mag.*, vol. 25, no. 2, pp. 72–82, Mar. 2008.
- [45] J. R. Shewchuk et al., "An introduction to the conjugate gradient method without the agonizing pain," 1994. [Online]. Available: https://www.cs.cmu.edu/~quake-papers/painless-conjugate-gradient.pdf
- [46] M. J. Muckley, R. Stern, T. Murrell, and F. Knoll, "TorchKbNufft: A high-level, hardware-agnostic non-uniform fast fourier transform," in ISMRM Workshop on Data Sampling & Image Reconstruction, 2020.
- [47] O. Maier et al., "CG-SENSE revisited: Results from the first ISMRM reproducibility challenge," Magn. Reson. Med., vol. 85, no. 4, pp. 1821– 1839, 2021.
- [48] R. D. Hoge, R. K. Kwan, and G. Bruce Pike, "Density compensation functions for spiral MRI," *Magn. Reson. Med.*, vol. 38, no. 1, pp. 117– 128, 1997.
- [49] D. Kim and J. A. Fessler, "Adaptive restart of the optimized gradient method for convex optimization," *J Optim Theory Appl*, vol. 178, no. 1, pp. 240–263, Jul. 2018.
- [50] J. A. Fessler, "Optimization methods for MR image reconstruction," IEEE Sig. Proc. Mag., vol. 37, no. 1, pp. 33–40, Jan. 2020.
- [51] D. C. Noll, J. A. Fessler, and B. P. Sutton, "Conjugate phase MRI reconstruction with spatially variant sample density correction," *IEEE Trans. Med. Imag.*, vol. 24, no. 3, pp. 325–336, 2005.
- [52] M. Uecker et al., "ESPIRiT an eigenvalue approach to autocalibrating parallel MRI: where SENSE meets GRAPPA," Mag. Reson. Med., vol. 71, no. 3, pp. 990–1001, Mar. 2014.
- [53] J. Zbontar et al., "fastMRI: An open dataset and benchmarks for accelerated MRI," 2018. [Online]. Available: http://arxiv.org/abs/1811. 08839
- [54] D. P. Kingma and J. Ba, "Adam: a method for stochastic optimization," 2017. [Online]. Available: http://arxiv.org/abs/1412.6980
- [55] A. Horé and D. Ziou, "Image quality metrics: PSNR vs. SSIM," in Intl. Conf. on Patn. Recog. (ICPR), Aug. 2010, pp. 2366–2369.
- [56] C. Lazarus et al., "SPARKLING: variable-density k-space filling curves for accelerated T2*-weighted MRI," Mag. Res. Med., vol. 81, no. 6, pp. 3643–61, Jun. 2019.
- [57] C. Boyer, N. Chauffert, P. Ciuciu, J. Kahn, and P. Weiss, "On the generation of sampling schemes for magnetic resonance imaging," SIAM J. Imaging Sci., vol. 9, no. 4, pp. 2039–2072, 2016.
- [58] H. Gu, B. Yaman, S. Moeller, J. Ellermann, K. Ugurbil, and M. Ak-cakaya, "Revisiting 11-wavelet compressed-sensing MRI in the era of deep learning," *Proc. Natl. Acad. Sci.*, vol. 119, no. 33, p. e2201062119, Aug. 2022.

APPENDIX

This appendix analyzes the error of approximations based on (5) and (6), by comparing Jacobians computed when A is an exact NUDFT to those for an NUFFT, denoted \tilde{A} . For simplicity, the analysis is 1D, though the conclusion extends easily to multi-dimensional NUFFTs.

The system matrix $\boldsymbol{A} \in \mathbb{C}^{M \times N}$ has elements

$$a_{mn} = e^{-i\omega_m n}, \quad m = 1, \dots, M, \quad n = 1, \dots, N.$$

Typically, an NUFFT involves three steps. The first step applies scaling factors s_n to the signal x_n . The second step applies a K-point FFT to the scaled signal, where $K \geq N$ via zero-padding. The third step interpolates K frequency locations into M sampling locations of ω . For efficiency, the interpolator usually has finite support, denoted J>0. The NUFFT \tilde{A} has elements as follows:

$$\tilde{a}_{mn} = \sum_{j=1}^{J} u_j^*(\omega_m) s_n e^{-i\gamma (k_m+j) n},$$

where u denotes interpolation coefficients, k_m is an elementwise offset, and $\gamma = 2\pi/K$ [6].

Define the NUFFT error matrix as $E = \tilde{A} - A$. The worst-case NUFFT error has a bound that can be written as

$$\|\boldsymbol{E}\boldsymbol{x}\|_{\infty} \leq \varepsilon_p \|\boldsymbol{x}\|_2$$
,

where ε_p is tabulated numerically for various choices of interpolation parameters p, e.g., in [6, Fig. 12].

The Jacobian of the forward operator (5) is

$$oldsymbol{J} = rac{\partial oldsymbol{A} oldsymbol{x}}{\partial oldsymbol{\omega}} = -\imath \operatorname{diag} \{oldsymbol{A} (oldsymbol{x} \odot oldsymbol{r}) \} \,.$$

Let \tilde{J} denote the case where an NUFFT is applied. Since the backpropagation uses Jacobians in the JVP calculation, here we analyze the error of JVPs using J and \tilde{J} . We define the worst-case relative error for a JVP with a (gradient) vector v as follows:

$$\begin{split} E_{1}(\boldsymbol{\omega}, \boldsymbol{x}, p) &\triangleq \max_{\|\boldsymbol{v}\|_{\infty} = 1} \|\tilde{\boldsymbol{J}}\boldsymbol{v} - \boldsymbol{J}\boldsymbol{v}\|_{\infty} / \|\boldsymbol{x}\|_{2} \\ &= \max_{\|\boldsymbol{v}\|_{\infty} = 1} \|(\boldsymbol{E}(\boldsymbol{x} \odot \boldsymbol{r})) \odot \boldsymbol{v}\|_{\infty} / \|\boldsymbol{x}\|_{2} \\ &= \|\boldsymbol{E}(\boldsymbol{x} \odot \boldsymbol{r})\|_{\infty} / \|\boldsymbol{x}\|_{2} \\ &\leq \varepsilon_{p} \|\boldsymbol{x} \odot \boldsymbol{r}\|_{2} / \|\boldsymbol{x}\|_{2} \leq \varepsilon_{p} \|\boldsymbol{r}\|_{\infty}. \end{split}$$

Similarly, the worst-case relative error of a JVP with (6) is bounded by

$$\begin{split} E_{2}(\boldsymbol{\omega}, \boldsymbol{x}, p) &\triangleq \max_{\|\boldsymbol{v}\|_{\infty} = 1} \|\operatorname{diag}\{\boldsymbol{r}\} \, \boldsymbol{E}' \operatorname{diag}\{\boldsymbol{y}\} \, \boldsymbol{v}\|_{\infty} \, / \, \|\boldsymbol{y}\|_{2} \\ &\leq \max_{\|\boldsymbol{v}\|_{\infty} = 1} \|\boldsymbol{r}\|_{\infty} \, \|\boldsymbol{E}'(\boldsymbol{y} \odot \boldsymbol{v})\|_{\infty} \, / \, \|\boldsymbol{y}\|_{2} \\ &\leq \varepsilon_{p} \, \|\boldsymbol{r}\|_{\infty} \, \max_{\|\boldsymbol{v}\|_{\infty} = 1} \|\boldsymbol{y} \odot \boldsymbol{v}\|_{2} \, / \, \|\boldsymbol{y}\|_{2} \\ &\leq \varepsilon_{p} \, \|\boldsymbol{r}\|_{\infty} \, \|\boldsymbol{y}\|_{2} \, / \, \|\boldsymbol{y}\|_{2} \leq \varepsilon_{p} \, \|\boldsymbol{r}\|_{\infty} \, . \end{split}$$

In both cases, the worst-case error of the NUFFT approximation for a JVP is bounded by the usual NUFFT error multiplied by a constant $\|r\|_{\infty}$ that is usually half of the field of view (FOV) in imaging applications. This constant is expected from unit analysis. If the sampling grid r_i has a unit

in cm, then the sample locations ω have units in radians/cm. Corresponding, the Jacobian matrices in (5) and (6) have units in cm, because A is unitless. The NUFFT error ε_p is unitless, so there is an r-related factor in the JVP error E. In other words, the error bounds above depend on the choice of units. One could express the FOV in voxels to get the unitless error bound $\varepsilon_p N/2$. However, the accuracy of JVPs does not necessarily deteriorate with larger N. Above we normalized the error by $\|x\|_2$ or $\|y\|_2$, whereas the Jacobians are "scaled" with $\|x\odot r\|_2$ or $\|y\|_2\|r\|_2$. A relative error could better describe the effect on optimization.

An alternate definition uses the worst-case in the numerator relative to an average case in the denominator, considering the stochastic gradient descent-like optimizers. For example, this relative error for the JVP of Jacobian (5) is

$$\begin{split} \epsilon &\triangleq \frac{\max_{\|\boldsymbol{x}\|_{2}=1} \|\tilde{\boldsymbol{J}} - \boldsymbol{J}\|_{\mathrm{F}}}{\sqrt{\mathsf{E}_{p(\boldsymbol{x})}[\|\boldsymbol{J}\|_{\mathrm{F}}^{2}]}} = \frac{\max_{\|\boldsymbol{x}\|_{2}=1} \|\boldsymbol{E}\left(\boldsymbol{x}\odot\boldsymbol{r}\right)\|_{2}}{\sqrt{\mathsf{E}_{p(\boldsymbol{x})}[\|\boldsymbol{A}\left(\boldsymbol{x}\odot\boldsymbol{r}\right)\|_{2}^{2}]}} \\ &\leq \frac{\max_{\|\boldsymbol{x}\|_{2}=1} \sqrt{M} \|\boldsymbol{E}\left(\boldsymbol{x}\odot\boldsymbol{r}\right)\|_{\infty}}{\sqrt{\mathsf{E}_{p(\boldsymbol{x})}[\|\boldsymbol{A}\left(\boldsymbol{x}\odot\boldsymbol{r}\right)\|_{2}^{2}]}} \leq \frac{\sqrt{M}\varepsilon_{p} \|\boldsymbol{r}\|_{\infty}}{\sqrt{\mathsf{E}_{p(\boldsymbol{x})}[\|\boldsymbol{A}\left(\boldsymbol{x}\odot\boldsymbol{r}\right)\|_{2}^{2}]}}, \end{split}$$

where $\mathsf{E}_{p(\boldsymbol{x})}[\cdot]$ denotes expectation w.r.t. a certain distribution $p(\boldsymbol{x})$. For parity with the unit sphere constraint in the numerator, we consider the case where $p(\cdot)$ is the random distribution on the unit N-sphere. Use the cyclic property of the trace:

$$\begin{split} \left\| oldsymbol{A} \left(oldsymbol{x} \odot oldsymbol{r}
ight)
ight\|_2^2 &= oldsymbol{x}' \mathrm{diag} \{ oldsymbol{r} \} \, oldsymbol{A}' oldsymbol{A} \mathrm{diag} \{ oldsymbol{r} \} \, oldsymbol{x} oldsymbol{x}' \} \, . \end{split}$$

Since the covariance of random points on the N-sphere is $(1/N)\mathbf{I}$, the denominator's expectation is

$$\begin{split} \mathsf{E}_{p(\boldsymbol{x})}[\|\boldsymbol{A}\left(\boldsymbol{x}\odot\boldsymbol{r}\right)\|_{2}^{2}] &= \mathrm{Tr}\big\{\mathrm{diag}\{\boldsymbol{r}\}\,\boldsymbol{A}'\boldsymbol{A}\mathrm{diag}\{\boldsymbol{r}\}\,\mathsf{E}_{p(\boldsymbol{x})}[\boldsymbol{x}\boldsymbol{x}']\big\} \\ &= \frac{1}{N}\mathrm{Tr}\big\{\mathrm{diag}\{\boldsymbol{r}\}\,\boldsymbol{A}'\boldsymbol{A}\mathrm{diag}\{\boldsymbol{r}\}\big\} \\ &= \frac{1}{N}\sum_{j}r_{j}^{2}[\boldsymbol{A}'\boldsymbol{A}]_{jj} = \frac{M}{N}\sum_{j}r_{j}^{2} = \frac{M}{N}\left\|\boldsymbol{r}\right\|_{2}^{2}. \end{split}$$

Thus we have the following bound for the relative error:

$$\epsilon \leq rac{\sqrt{M}arepsilon_p \|m{r}\|_{\infty}}{\sqrt{M/N} \|m{r}\|_2} = arepsilon_p \sqrt{N} rac{\|m{r}\|_{\infty}}{\|m{r}\|_2} \leq arepsilon_p \sqrt{N}.$$

Note that the bound can be tighter when considering specific formulations of r. Similarly, for the Jacobian operator (6), the alternate error of the JVP is

$$\epsilon \le \varepsilon_p \sqrt{M}$$
.