# CONVOLUTIONAL ANALYSIS OPERATOR LEARNING BY END-TO-END TRAINING OF ITERATIVE NEURAL NETWORKS

*Andreas Kofler[1], Christian Wald[2], Tobias Schaeffter[1,3,4], Markus Haltmeier[5], Christoph Kolbitsch[1,3]*

[1] Physikalisch-Technische Bundesanstalt, Berlin and Braunschweig, Germany
[2] Department of Radiology, Charité - Universitätsmedizin Berlin, Berlin, Germany
[3] School of Imaging Sciences and Biomedical Engineering, King's College London, London, UK
[4] Department of Biomedical Engineering, Technical University of Berlin, Berlin, Germany
[5] Department of Mathematics, University of Innsbruck, Innsbruck, Austria

## ABSTRACT

The concept of sparsity has been extensively applied for regularization in image reconstruction. Typically, sparsifying transforms are either pre-trained on ground-truth images or adaptively trained during the reconstruction. Thereby, learning algorithms are designed to minimize some target function which encodes the desired properties of the transform. However, this procedure ignores the subsequently employed reconstruction algorithm as well as the physical model which is responsible for the image formation process. Iterative neural networks - which contain the physical model - can overcome these issues. In this work, we demonstrate how convolutional sparsifying filters can be efficiently learned by end-to-end training of iterative neural networks. We evaluated our approach on a non-Cartesian 2D cardiac cine MRI example and show that the obtained filters are better suitable for the corresponding reconstruction algorithm than the ones obtained by decoupled pre-training.

***Index Terms***— Iterative Neural Networks, Sparsity, Analysis Operator, Compressed Sensing, Cardiac Cine MRI

## 1. INTRODUCTION

Recently, iterative convolutional neural networks (CNNs) have been successfully applied to image reconstruction problems and seem to define the state-of-the-art across many imaging modalities, see e.g. [1], [2], [3], [4]. Iterative CNNs resemble iterative reconstruction schemes of finite length in which the regularizer is parametrized by convolutional operations and can be learned in a supervised manner by end-to-end training of the network. Their success seems to be attributable to i) the fact that the physical model is inherently present in the learning process - which has been reported to lower the expected maximum error-bound [5] - and ii) because the regularizers are trained in conjunction with the reconstruction algorithm that is used to reconstruct the images.

Regardless of their success, neural networks have also been reported to possibly suffer from instabilities [6] and still operate as black-boxes. This is an issue especially for a field such as medical imaging where the image content directly impacts diagnosis and treatment planning or decisions. In contrast, more classical learning-based regularization approaches typically come with solid mathematical theory, see e.g. [7], [8]. However, in these algorithms - in contrast to iterative neural networks - the physical model is not integrated in the learning process and typically, training refers to minimizing some object function which reflects the desired properties of the regularizer rather than being optimal for the purpose they have to serve in a subsequent reconstruction process.

In this work, we combine the best of the two worlds by using iterative neural networks to train a classical data-driven method based on learned sparsifying transforms given as convolutional filters, similar as in [8]. In contrast to iterative NNs using many convolutional layers, the role of the learned regularizer is more transparent. Further, unlike in [8], where the filters are pre-trained on a set of ground-truth images, in our network the filters are learned to be optimal with respect to the reconstruction algorithm and the number of iterations that the network uses to reconstruct the images and are adapted to the operator of the inverse problem. Our work also differs from [2] which stems from the field-of-experts model [9] and uses a Landweber iteration. We instead use a splitting approach, and because of the used formulation, the required non-linear activation function is given by the soft-thresholding operator. In addition, the presented approach differs from the work in [10], where the filters are trained in a greedy fashion (i.e. layer-by-layer), and in the application.

## 2. METHODS

We consider the general type of inverse problem of the form

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

where $\mathbf{A}$ denotes the forward model, $\mathbf{x}$ the (unknown) image, $\mathbf{e}$ random noise and $\mathbf{y}$ the measured data. Problem (1)

---

Our implementation of the network is available under `www.github.com/koflera/ConvSparsityNNs`.

can be ill-posed for different reasons. For example, if $\mathbf{y}$ has less entries than $\mathbf{x}$, the problem is underdetermined and there exists an infinite number of solutions. For properly designed overdetermined systems, a solution can be obtained by solving the normal equations, but the stability of the inversion process depends on the condition of the operator $\mathbf{A}^\mathsf{T}\mathbf{A}$. In this work, we investigate a regularization method given by the assumption that the image $\mathbf{x}$ is sparse with respect to convolutional filters. Assuming a *fixed* set of $K$ sparsifying filters $\{h_k\}_k$, one can formulate the reconstruction problem as a minimization problem

$$\min_{\mathbf{x}} \frac{1}{2}\|\mathbf{A}\mathbf{x}-\mathbf{y}\|_2^2 + \alpha \sum_{k=1}^{K}\|h_k * \mathbf{x}\|_1 \tag{2}$$

over the image $\mathbf{x}$ with $\alpha > 0$. Because $\mathbf{x}$ is coupled to the operator $\mathbf{A}$ as well as to the filters $h_k$ which appear in the $L_1$-norm, directly solving problem (2) is challenging. A possible solution strategy is to introduce $K$ auxiliary variables $\mathbf{s}_k$ to transfer $h_k * \mathbf{x}$ out of the $L_1$-norm and to relax the equality constraint by including it in a quadratic penalty term, i.e.

$$\min_{\mathbf{x},\{\mathbf{s}_k\}_k} \frac{1}{2}\|\mathbf{A}\mathbf{x}-\mathbf{y}\|_2^2 + \frac{\lambda}{2}\sum_{k=1}^{K}\|h_k * \mathbf{x}-\mathbf{s}_k\|_2^2 + \alpha \sum_{k=1}^{K}\|\mathbf{s}_k\|_1, \tag{3}$$

where $\lambda > 0$. A possible approach for minimizing (3) uses alternating minimization of (3) with respect to $\mathbf{x}$ and $\{\mathbf{s}_k\}_k$ in an iterative manner [11]. For fixed $\mathbf{x}$, problem (3) is separable with respect to $k$ and thus, the solution for $\mathbf{s}_k$ is given by applying the soft-thresholding operator to $h_k * \mathbf{x}$ for all $k$. For fixed $\{\mathbf{s}_k\}_k$, the minimization with respect to $\mathbf{x}$ corresponds to solving a linear system $\mathbf{H}\mathbf{x} = \mathbf{b}_j$ with

$$\mathbf{H} = \mathbf{A}^\mathsf{H}\mathbf{A} + \lambda \sum_{k=1}^{K} h_k^\mathsf{T} * h_k \tag{4}$$

$$\mathbf{b}_j = \mathbf{A}^\mathsf{H}\mathbf{y} + \lambda \sum_{k=1}^{K} h_k^\mathsf{T} * \mathbf{s}_k, \tag{5}$$

where we see that the operator $\mathbf{H}$ depends on the filters. Since we aim to train the set of filters $\{h_k\}_k$ by training an iterative network in an end-to-end fashion, this alternating-minimization scheme can be compuationally demanding for realistic large-scale applications, e.g. for the later discussed dynamic cardiac MRI problem. Therefore, motivated by the backward-backward splitting method [12], similar to previous works [8], [10], we approach the minimization of (3) by

$$\mathbf{z}_j = \sum_{k=1}^{K} h_k^\mathsf{T} * \mathcal{S}_{\alpha/\lambda}(h_k * \mathbf{x}_j) \tag{6}$$

$$\mathbf{x}_{j+1} = \arg\min_{\mathbf{x}} \frac{1}{2}\|\mathbf{A}\mathbf{x}-\mathbf{y}\|_2^2 + \frac{\lambda}{2}\|\mathbf{x}-\mathbf{z}_j\|_2^2, \tag{7}$$

for $0 \le j \le T$, with $\mathbf{x}_0 := \mathbf{A}^\sharp\mathbf{y}$, where $\mathbf{A}^\sharp$ denotes some pseudo-inverse of $\mathbf{A}$. In (6), $\mathcal{S}_{\alpha/\lambda}$ denotes the soft-thresholding operator with threshold $\alpha/\lambda$ and $h_k^\mathsf{T}$ denotes the adjoint of $h_k$. Under an orthonormal basis assumption, the
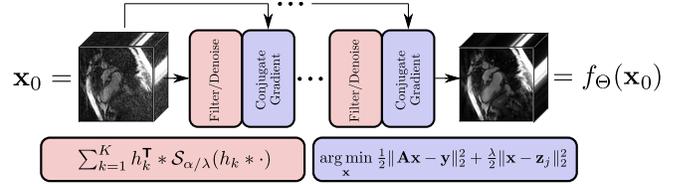


**Fig. 1**. Proposed Network structure. The image is first filtered, soft-thresholded and filtered with the transposed filters. Then, the denoised image is used as regularizing prior in a regularized functional. The filters $\{h_k\}_k$ as well as the regularization parameters $\lambda$ and $\alpha$ are obtained by end-to-end training of the entire network.

sequence defined by (6) and (7) reduces to the backward-backward splitting algorithm for (2), known to converge to a minimizer of (3) [12]. The minimizer of (7) can be obtained by solving a linear system $\mathbf{H}\mathbf{x} = \mathbf{b}_j$ with

$$\mathbf{H} = \mathbf{A}^\mathsf{H}\mathbf{A} + \lambda\,\mathbf{I} \tag{8}$$

$$\mathbf{b}_j = \mathbf{A}^\mathsf{H}\mathbf{y} + \lambda\,\mathbf{z}_j, \tag{9}$$

where the $\mathbf{H}$ does not depend on the filters and is thus computationally favorable.

**Proposed Reconstruction Network**: We propose to train the filters $\{h_k\}_k$ by constructing a network $f_\Theta$ which corresponds to a sequence of alternating steps which implement the operations in (6) and (7). In the network, the filters are treated as trainable parameters, i.e. $\Theta = \cup_k\{h_k\}$, and can therefore be learned by back-propagation in a supervised manner on a set of $M$ data-pairs $\mathcal{D} = \{(\mathbf{x}_0^i, \mathbf{x}_f^i)_{i=1}^M\}$, where $\mathbf{x}_f$ denotes a ground-truth image. Further, we can learn the optimal regularization parameters $\lambda$ and $\alpha$ as well. In order to constrain the regularization parameters to be strictly positive, we apply a Soft-Plus activation to $\alpha$ and $\lambda$.

Because the soft-thresholding operator $\mathcal{S}_{\alpha/\lambda}$ is not differentiable with respect to its threshold $\alpha/\lambda$, following [13], we smoothly approximate it by

$$\tilde{\mathcal{S}}_t(z) = z + \frac{1}{2}\left(\sqrt{(z-t)^2+b} - \sqrt{(z+t)^2+b}\right) \tag{10}$$

to be able to learn the optimal threshold by back-propagation, where $b > 0$ is a parameter which we fixed to $b = 0.001$. In the network, the complex-valued images are treated as two-channeled real-valued images and the the real and the imaginary parts of the images share the same filters. For the convolutional layers, we employ circular padding. Figure 1 illustrates the proposed network architecture.

## 3. EXPERIMENTS

In the following, we tested our proposed method on an accelerated radial cardiac cine MR image reconstruction problem.
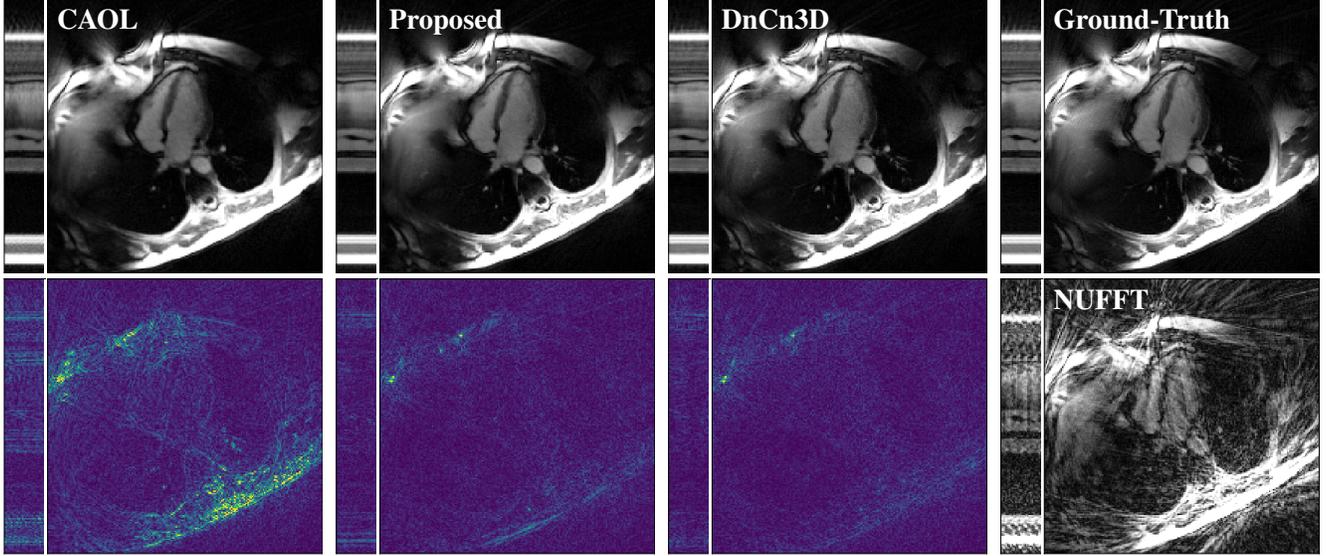
**Fig. 2**. An example of reconstructions and corresponding point-wise error-images of the test set for the proposed reconstruction method using CAOL filters [8] for $K = 16, k_f = 3$ and the ones obtained by our proposed end-to-end training-approach for $K = 16$ and $k_f = 7$ as well as for the deep CNN-cascade DnCn3D [3]. Although DnCn3D yields a slightly lower point-wise error, our proposed approach shows competitive performance with the advantage that the role of the regularizing kernels is fully interpretable. All results are shown for the best combination of hyper-parameters (chosen on the validation set) for each respective method.

Similar as in [14], the operator $\mathbf{A}$ in (1) is given by

$$\mathbf{A} := (\mathbf{I}_{N_c} \otimes \mathbf{E})\mathbf{C}, \qquad (11)$$

for a complex valued image $\mathbf{x} = [\mathbf{x}_1, \ldots, \mathbf{x}_{N_t}]^{\mathsf{T}} \in \mathbb{C}^N$ with $N = N_x \times N_y \times N_t$. The operator $\mathbf{I}_{N_c}$ denotes an identity operator and $\mathbf{C}$ contains the $N_c$ coil-sensitivity maps which are multiplied to the cine MR image, i.e. $\mathbf{C} = [\mathbf{C}_1, \ldots, \mathbf{C}_{N_c}]^{\mathsf{T}}$, with $\mathbf{C}_j = \mathrm{diag}(\mathbf{c}_j, \mathbf{c}_j, \ldots, \mathbf{c}_j) \in \mathbb{C}^{N \times N}$ and $\mathbf{c}_j \in \mathbb{C}^{N_x \times N_y}$. The operator $\mathbf{E} = \mathrm{diag}(\mathbf{E}_1, \ldots, \mathbf{E}_{N_t})$ consists of different 2D non-uniform (NUFFT) Fourier-encoding operators $\mathbf{E}_t$ which for each point $t \in \{1, \ldots, N_t\}$ sample a 2D image $\mathbf{x}_t \in \mathbb{C}^{N_x \times N_y}$ along radial lines in Fourier-space. To accelerate the acquisition process, we only acquire a subset of the $k$-space coefficients which are needed to sample a 2D image $\mathbf{x}_t$ at Nyquist limit, which we denote by $I \subset J = \{1, \ldots, N_{\mathrm{rad}}\}$. Finally, by $\mathbf{A}_I$, we denote the undersampled 2D radial encoding operator which samples all $k$-space coefficients in the set $I = I_1 \cup \ldots \cup I_{N_t}$ with $I_t \subset J$ for all $t = 1, \ldots, N_t$ according to a golden-angle radial pattern [15]. The operator $\mathbf{A}_I$ was implemented using `TorchKBNufft` [16].

As often done for non-Cartesian sampling schemes, in the data-consistency term in (2), the $k$-space data is multiplied by a diagonal operator $\mathbf{W}^{1/2}$ which contains the entries of the density-compensation function and is used to pre-condition the problem. By doing so, the operator $\mathbf{A}_I^{\sharp}$ takes the form $\mathbf{A}_I^{\sharp} := \mathbf{A}_I^{\mathsf{H}} \mathbf{W}^{1/2}$. Accordingly, in Section 2, the operators $\mathbf{A}_I^{\mathsf{H}} \mathbf{A}_I$ in (8) and $\mathbf{A}_I^{\mathsf{H}}$ in (9) must be replaced by $\mathbf{A}_I^{\sharp} \mathbf{A}_I$ and $\mathbf{A}_I^{\sharp}$, respectively.

**Dataset**: We used a set of 15 healthy volunteers and four patients which amounted to 216 cine MR images of shape $320 \times 320 \times 30$. We split the data into 12/3/4 subjects (144/36/36 dynamic images) for training, validation and testing where the test set consisted of the four patients. The initial $k$-space data was retrospectively simulated using an acceleration factor of approximately $R \approx 18$ and $N_c = 12$ coil-sensitivity maps and was further corrupted by Gaussian noise with a standard deviation of $\sigma = 0.02$.

**Methods of Comparison and Evaluation**: Since our proposed method is a method for training sparsifying convolutional filters, the first method of comparison is the one in [8], which we denote by CAOL. After having trained the filters with CAOL, we fixed them in our reconstruction network and only trained the regularization parameters. We also compared our method to a deep cascade of convolutional neural networks [3], which we abbreviate by DnCn3D. For DnCn3D, we used $T = 4$ and each block has two convolutional layers with 16 filters, amounting to a total number of 31.232 trainable parameters. Note that the original work in [3] was presented for a single-coil Cartesian acquisition scheme. For our comparison, we extended the method to be applicable to non-Cartesian multi-coil data-acquisitions by replacing the data-consistency layer in [3] by a CG module. For CAOL, at test time, the length of the network was increased to $T = 24$ as it further decreased the NRMSE . All results were evaluated in terms of PSNR, NRMSE, structural similarity index
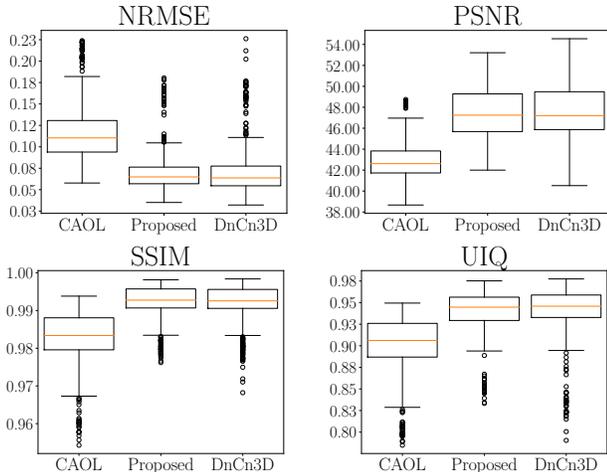
**Fig. 4**. Training- and validation-error (solid/dashed) during the optimization of the convolutional filters with our proposed reconstruction network for $K = 16$ and different $k_f$ (only shown for $K = 16$ for presentation purposes). For CAOL, only $\lambda$ and $\alpha$ were trained.

**Fig. 3**. Box-plots of the quantitative results obtained for CAOL [8] for $K = 16, k_f = 3$, our proposed method for $K = 16, k_f = 7$ and DnCn3D [3]. Our proposed method yields similar results as [3] while only having $110 < K/2 \cdot k_f^3 + 2 < 4.118$ trainable parameters (i.e. the filters and the regularization parameters $\alpha$ and $\lambda$) compared to 31.233 for DnCn3D.

measure [17] (SSIM) and universal image quality index [18] (UIQ) which were calculated over a central squared ROI of $160 \times 160$ pixels for all cardiac phases.

**Network Training**: We trained different networks with $K = 8, 16, 24$ for different 3D kernels of shape $k_f \times k_f \times k_f$ for $k_f = 3, 5, 7$ to minimize the squared $L_2$-error between the estimated output and the target-images. Because the application of the NUFFT-operator is computationally expensive and problem (3) is separable with respect to the time points, for our method, we reduced the number of cardiac phases to $N_t = 8$ during training. We set $T = 4$ and used $n_{CG} = 4$ iterations to solve (7). All methods were trained using the ADAM optimizer with an initial learning rate of $10^{-4}$. Our network was trained for 75 epochs ($\approx$ 9 hours), while DnCn3D was trained for 500 epochs ($\approx$ 4 days).

## 4. RESULTS AND DISCUSSION

In Figure 2, we see an example of a reconstruction for our method using the CAOL-filters with $K = 16$ and $k_f = 3$ and the ones obtained by end-to-end training, which yield a visibly smaller point-wise error and better preserve image details. These results are also supported in terms of the reported quantitative measures, as can be seen in the box-plots in Figure 3. Our proposed method and DnCn3D clearly surpass CAOL. Further, the proposed reconstruction method yields comparable results to DnCn3D and further seems to be slightly more stable, as can be seen from the outliers in the box-plots. This can most probably be attributed to the fact that it contains
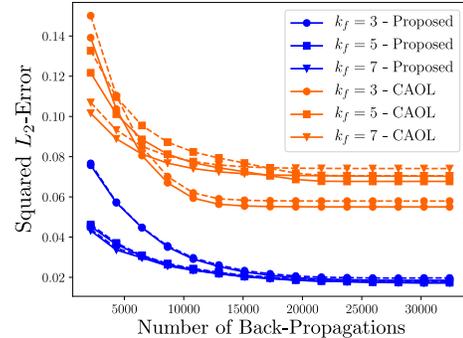
significantly fewer trainable parameters. Note that the results for CAOL and our proposed method are shown for the best configuration of $K$ and $k_f$ based on the validation set. This can be seen from Figure 4 which shows the training and validation errors for our network and for CAOL for $K = 16$. Increasing the filter size $k_f$ slightly reduces the achievable validation error for our method. Interestingly, we found that CAOL performs better with smaller kernel-sizes. Although this might seem somewhat counter-intuitive, this aspect shows that choosing the optimal hyper-parameters for decoupled methods is challenging. In contrast, using iterative networks to train the convolutional filters, larger filter-sizes tend to lead to smaller reconstruction errors and the filters are optimally adjusted to be used with the employed reconstruction algorithm regardless of the chosen hyper-parameters.

## 5. CONCLUSION

In this work, we have shown that end-to-end trained iterative neural networks can be used to learn classical sparsity-based regularization methods in a task-driven and physics-informed manner. The obtained sparsifying transforms are better tailored to the employed reconstruction algorithm compared to the ones obtained by the corresponding decoupled method. Further, we have evaluated our method on a realistic large-scale dynamic cardiac MR problem and found that the proposed method yields results which are on par with the ones obtained by a state-of-the-art method employing a deep cascade of neural networks. In addition, in our method, the exact role of the regularizer is fully explainable and allows for a theoretical analysis of the reconstruction algorithm which we leave for future work. Although we have presented the approach for a dynamic non-Cartesian MR image reconstruction example, we point out that the method may be applicable to other imaging modalities as well.

# 6. ACKNOWLEDGMENTS

# 7. COMPLIANCE WITH ETHICAL STANDARDS

All subjects gave written informed consent before participation, in accordance with the ethical committee of the responsible institution.

# 8. REFERENCES

[1] Jonas Adler and Ozan Öktem, "Solving ill-posed inverse problems using iterative deep neural networks," *Inverse Problems*, vol. 33, no. 12, pp. 124007, 2017.

[2] Kerstin Hammernik, Teresa Klatzer, Erich Kobler, Michael P Recht, Daniel K Sodickson, Thomas Pock, and Florian Knoll, "Learning a variational network for reconstruction of accelerated MRI data," *Magnetic Resonance in Medicine*, vol. 79, no. 6, pp. 3055–3071, 2018.

[3] Jo Schlemper, Jose Caballero, Joseph V Hajnal, Anthony N Price, and Daniel Rueckert, "A deep cascade of convolutional neural networks for dynamic MR image reconstruction," *IEEE Transactions on Medical Imaging*, vol. 37, no. 2, pp. 491–503, 2018.

[4] Andreas Hauptmann, Felix Lucka, Marta Betcke, Nam Huynh, Jonas Adler, Ben Cox, Paul Beard, Sebastien Ourselin, and Simon Arridge, "Model-based learning for accelerated, limited-view 3-d photoacoustic tomography," *IEEE transactions on medical imaging*, vol. 37, no. 6, pp. 1382–1393, 2018.

[5] Andreas K Maier, Christopher Syben, Bernhard Stimpel, Tobias Würfl, Mathis Hoffmann, Frank Schebesch, Weilin Fu, Leonid Mill, Lasse Kling, and Silke Christiansen, "Learning with known operators reduces maximum error bounds," *Nature machine intelligence*, vol. 1, no. 8, pp. 373–380, 2019.

[6] Vegard Antun, Francesco Renna, Clarice Poon, Ben Adcock, and Anders C Hansen, "On instabilities of deep learning in image reconstruction and the potential costs of ai," *Proceedings of the National Academy of Sciences*, vol. 117, no. 48, pp. 30088–30095, 2020.

[7] Il Yong Chun and Jeffrey A Fessler, "Convolutional dictionary learning: Acceleration and convergence," *IEEE Transactions on Image Processing*, vol. 27, no. 4, pp. 1697–1712, 2017.

[8] Il Yong Chun and J. Fessler, "Convolutional analysis operator learning: Acceleration and convergence," *IEEE Transactions on Image Processing*, vol. 29, pp. 2108–2122, 2020.

[9] Stefan Roth and Michael J Black, "Fields of experts," *International Journal of Computer Vision*, vol. 82, no. 2, pp. 205–229, 2009.

[10] Il Yong Chun, Zhengyu Huang, Hongki Lim, and Jeff Fessler, "Momentum-net: Fast and convergent iterative neural network for inverse problems," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.

[11] Yilun Wang, Junfeng Yang, Wotao Yin, and Yin Zhang, "A new alternating minimization algorithm for total variation image reconstruction," *SIAM Journal on Imaging Sciences*, vol. 1, no. 3, pp. 248–272, 2008.

[12] Patrick L Combettes and Jean-Christophe Pesquet, "Proximal splitting methods in signal processing," in *Fixed-point algorithms for inverse problems in science and engineering*, pp. 185–212. Springer, 2011.

[13] Xiao-Ping Zhang, "Thresholding neural network for adaptive noise reduction," *IEEE Transactions on Neural Networks*, vol. 12, no. 3, pp. 567–584, 2001.

[14] Andreas Kofler, Markus Haltmeier, Tobias Schaeffter, and Christoph Kolbitsch, "An end-to-end-trainable iterative network architecture for accelerated radial multi-coil 2d cine MR image reconstruction," *Medical Physics*, vol. 48, no. 5, pp. 2412–2425, 2021.

[15] Stefanie Winkelmann, Tobias Schaeffter, Thomas Koehler, Holger Eggers, and Olaf Doessel, "An optimal radial profile order based on the golden ratio for time-resolved MRI," *IEEE Transactions on Medical Imaging*, vol. 26, no. 1, pp. 68–76, 2006.

[16] Matthew J Muckley, Ruben Stern, Tullie Murrell, and Florian Knoll, "Torchkbnufft: A high-level, hardware-agnostic non-uniform fast fourier transform," in *ISMRM Workshop on Data Sampling & Image Reconstruction*, 2020.

[17] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE Transactions on Image Processing*, vol. 13, no. 4, pp. 600–612, 2004.

[18] Zhou Wang and Alan C Bovik, "A universal image quality index," *IEEE signal processing letters*, vol. 9, no. 3, pp. 81–84, 2002.