# Block Coordinate Regularization by Denoising

Yu Sun\*, Student Member, IEEE, Jiaming Liu\*, Student Member, IEEE, and Ulugbek S. Kamilov, Member, IEEE

Abstract—We consider the problem of reconstructing an image from its noisy measurements using a prior specified only with an image denoiser. Recent work on plug-and-play priors (PnP) and regularization by denoising (RED) has shown the state-of-the-art performance of image reconstruction algorithms under such priors in a range of imaging problems. In this work, we develop a new block coordinate RED algorithm that decomposes a large-scale estimation problem into a sequence of updates over a small subset of the unknown variables. We theoretically analyze the convergence of the algorithm and discuss its relationship to the traditional proximal optimization. Our analysis complements and extends recent theoretical results for RED-based estimation methods. We numerically validate our method using several denoising priors, including those based on deep neural nets.

Index Terms—Regularized image reconstruction, plug-andplay priors, regularization by denoising, proximal optimization

#### I. Introduction

The problem of reconstructing an unknown image  $x \in \mathbb{R}^n$  from a set of noisy measurements  $y \in \mathbb{R}^m$  is common in computational imaging. Consider the scenario, where an image  $x \sim p_x$  is acquired via an imaging system characterized by its likelihood function  $p_{y|x}$  to produce the measurements y. When the inverse problem is ill-posed, it becomes essential to include the prior  $p_x$  in the reconstruction process. However, in high-dimensional settings, it is difficult to directly obtain the true prior  $p_x$  for natural images and one is restricted to various indirect sources of prior information. This paper considers the cases where the prior on x is specified only through an image denoiser,  $D: \mathbb{R}^n \to \mathbb{R}^n$ , designed for the removal of additive white Gaussian noise (AWGN).

There has been considerable recent interest in leveraging denoisers as priors for image recovery. One popular strategy, known as *plug-and-play priors* (*PnP*) [1], extends traditional proximal optimization [2] by replacing the proximal operator with a general off-the-shelf denoiser. It has been shown that the combination of proximal algorithms with advanced denoisers, such as BM3D [3] or DnCNN [4], leads to the state-of-the-art performance for various imaging problems [5]–[15]. A similar strategy has also been adopted in the context of a related class of

This material is based upon work supported by NSF award CCF-1813910. (Corresponding author: Ulugbek S. Kamilov.)

The material in this paper was presented in part at the 2019 32th Annual Conference on Neural Information Processing Systems (NeurIPS).

- Y. Sun is with the Department of Computer Science & Engineering, Washington University in St. Louis, MO 63130, USA.
- J. Liu is with the Department of Electrical & Systems Engineering, Washington University in St. Louis, MO 63130, USA.
- U. S. Kamilov (email: kamilov@wustl.edu) is with the Department of Computer Science & Engineering and the Department of Electrical & Systems Engineering, Washington University in St. Louis, MO 63130, USA.

algorithms known as approximate message passing (AMP) [16]–[20]. Regularization-by-denoising (RED) [21], and the closely related deep mean-shift priors [22], represent an alternative, where the denoiser specifies an explicit regularizer that has a simple gradient. Recent work has clarified the existence of explicit RED regularizers [23], demonstrated its excellent performance on phase retrieval [24], and further boosted its performance in combination with a deep image prior [25]. In short, the use of advanced denoisers has proven to be essential for achieving the state-of-the-art results in many contexts. However, solving the corresponding estimation problem is still a significant computational challenge, especially in the context of high-dimensional vectors  $\boldsymbol{x}$ , typical in imaging applications.

In this work, we extend the current family of RED algorithms by introducing a new block coordinate RED (BC-RED) algorithm. The algorithm relies on partial updates on x, which makes it scalable to images that would otherwise be prohibitively large for direct processing. Additionally, as we shall see, the overall computational complexity of BC-RED can sometimes be lower than corresponding methods operating on the full image. This behavior is consistent with the traditional coordinate descent methods that can outperform their full gradient counterparts by being able to better reuse local updates and take larger steps [26]–[30]. We present two theoretical results related to BC-RED. We first theoretically characterize the convergence of the algorithm under a set of transparent assumptions on the data-fidelity and the denoiser. Our analysis complements the recent theoretical analysis of the traditional RED algorithms in [23] by considering block-coordinate updates and establishing the explicit worst-case convergence rate. Our second result establishes backward compatibility of BC-RED with traditional proximal optimization. We show that when the denoiser corresponds to a proximal operator, BC-RED can be interpreted as an approximate MAP estimator, whose approximation error can be made arbitrarily small. To the best of our knowledge, this explicit link with proximal optimization is missing in the current literature on RED. BC-RED thus provides a flexible, scalable, and theoretically sound algorithm applicable to a wide variety of large-scale imaging problems. We demonstrate BC-RED on image recovery from linear measurements using several denoising priors, including those based on deep neural nets.

The outline for the rest of the paper is as follows. In Section II, we review some relevant background material on image reconstruction. In Section III, we introduce BC-RED and present its fixed point interpretation. In Section IV, we analyze the convergence of BC-RED under several transparent assumptions. In Section V, we provide numerical experiments that illustrate key properties of our method. Section VI concludes the paper. A preliminary version of this work was

<sup>\*</sup> Authors contributed equally to this work.

presented in [31]. The current paper contains all the proofs and additional simulations.

#### II. BACKGROUND

It is common to formulate image reconstruction as an optimization problem

$$\widehat{\boldsymbol{x}} = \mathop{\arg\min}_{\boldsymbol{x} \in \mathbb{R}^n} f(\boldsymbol{x}) \quad \text{with} \quad f(\boldsymbol{x}) = g(\boldsymbol{x}) + h(\boldsymbol{x}), \quad \ (1)$$

where g is the data-fidelity term and h is the regularizer. For example, the maximum a posteriori probability (MAP) estimator is obtained by setting

$$g(\mathbf{x}) = -\log(p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}))$$
 and  $h(\mathbf{x}) = -\log(p_{\mathbf{x}}(\mathbf{x})),$ 

where  $p_{\boldsymbol{y}|\boldsymbol{x}}$  is the likelihood that depends on  $\boldsymbol{y}$  and  $p_{\boldsymbol{x}}$  is the prior. One of the most popular data-fidelity terms is least-squares  $g(\boldsymbol{x}) = \frac{1}{2} \|\boldsymbol{y} - \boldsymbol{A}\boldsymbol{x}\|_2^2$ , which assumes a linear measurement model under AWGN. Similarly, one of the most popular regularizers is based on a sparsity-promoting penalty  $h(\boldsymbol{x}) = \tau \|\boldsymbol{D}\boldsymbol{x}\|_1$ , where  $\boldsymbol{D}$  is a linear transform and  $\tau > 0$  is the regularization parameter [32]-[35].

Many widely used regularizers, including the ones based on the  $\ell_1$ -norm, are nondifferentiable. Proximal algorithms [2], such as the proximal-gradient method (PGM) [36]–[39] and alternating direction method of multipliers (ADMM) [40]–[43], are a class of optimization methods that avoid differentiating nonsmooth regularizers by using the proximal operator

$$\operatorname{prox}_{\mu h}({m z}) \coloneqq \operatorname*{arg\,min}_{{m x} \in \mathbb{R}^n} \left\{ rac{1}{2} \|{m x} - {m z}\|_2^2 + \mu h({m x}) 
ight\}, \quad \ \ (2)$$

where  $\mu>0$  is a parameter. The observation that the proximal operator can be interpreted as the MAP denoiser for AWGN has prompted the development of PnP [1], where the proximal operator  $\text{prox}_{\mu h}(\cdot)$ , within ADMM or PGM, is replaced with a more general image denoiser  $D(\cdot)$ .

Consider the following alternative to PnP [21], [22]

$$\boldsymbol{x}^t \leftarrow \boldsymbol{x}^{t-1} - \gamma \mathsf{G}(\boldsymbol{x}^{t-1}), \quad \gamma > 0,$$
 (3a)

where the update direction also relies on a denoising function

$$G(x) := \nabla g(x) + \tau(x - D(x)), \quad \tau > 0.$$
 (3b)

Under some conditions on the denoiser, it is possible to relate  $H(x) := \tau(x - D(x))$  in (3a) to some explicit regularization function h. For example, when the denoiser is locally homogeneous and has a symmetric Jacobian [21], [23], the operator  $H(\cdot)$  corresponds to the gradient of the following function

$$h(\boldsymbol{x}) = \frac{\tau}{2} \boldsymbol{x}^{\mathsf{T}} (\boldsymbol{x} - \mathsf{D}(\boldsymbol{x})). \tag{4}$$

On the other hand, when the denoiser corresponds to the minimum mean squared error (MMSE) estimator  $D(z) = \mathbb{E}[x|z]$  for the AWGN denoising problem [22], [23], z = x + e, with  $x \sim p_x(x)$  and  $e \sim \mathcal{N}(\mathbf{0}, \sigma^2 I)$ , the operator  $H(\cdot)$  corresponds to the gradient of

$$h(\mathbf{x}) = -\tau \sigma^2 \log(p_{\mathbf{z}}(\mathbf{x})), \tag{5}$$

where

$$p_{\boldsymbol{z}}(\boldsymbol{x}) = (p_{\boldsymbol{x}} * p_{\boldsymbol{e}})(\boldsymbol{x}) = \int_{\mathbb{R}^n} p_{\boldsymbol{x}}(\boldsymbol{z}) \phi_{\sigma}(\boldsymbol{x} - \boldsymbol{z}) \, d\boldsymbol{z},$$

Algorithm 1 Block Coordinate Regularization by Denoising

- 1: **input:** initial value  $x^0 \in \mathbb{R}^n$ , parameter  $\tau > 0$ , and step-size  $\gamma > 0$ .
- 2: **for**  $k = 1, 2, 3, \dots$  **do**
- 3: Choose an index  $i_k \in \{1, \dots, b\}$
- 4:  $\mathbf{x}^k \leftarrow \mathbf{x}^{k-1} \gamma \mathsf{U}_{i_k} \mathsf{G}_{i_k} (\mathbf{x}^{k-1})$ where  $\mathsf{G}_i(\mathbf{x}) \coloneqq \mathsf{U}_i^\mathsf{T} \mathsf{G}(\mathbf{x})$ with  $\mathsf{G}(\mathbf{x}) \coloneqq \nabla q(\mathbf{x}) + \tau (\mathbf{x} - \mathsf{D}(\mathbf{x}))$ .
- 5: end for

where  $\phi_{\sigma}$  is the Gaussian probability density function of variance  $\sigma^2$  and \* denotes convolution. In this paper, we will use the term RED to denote the method in (3a). The key benefits of the RED methods [21]–[25] are their explicit separation of the forward model from the prior, their ability to accommodate powerful denoisers (such as the ones based on deep neural nets) without differentiating them, and their state-of-the-art performance on a number of imaging tasks. The next section further extends the scalability of RED by designing a new block coordinate RED algorithm.

#### III. BLOCK COORDINATE RED

All the current RED algorithms operate on vectors in  $\mathbb{R}^n$ . We propose BC-RED, shown in Algorithm 1, to allow for partial randomized updates on x. Consider the decomposition of  $\mathbb{R}^n$  into b > 1 subspaces

$$\mathbb{R}^n = \mathbb{R}^{n_1} \times \mathbb{R}^{n_2} \times \cdots \times \mathbb{R}^{n_b}$$
 with  $n = n_1 + n_2 + \cdots + n_b$ .

For each  $i \in \{1, \dots, b\}$ , we define the matrix  $U_i : \mathbb{R}^{n_i} \to \mathbb{R}^n$  that injects a vector in  $\mathbb{R}^{n_i}$  into  $\mathbb{R}^n$  and its transpose  $U_i^\mathsf{T}$  that extracts the ith block from a vector in  $\mathbb{R}^n$ . Then, for any  $\boldsymbol{x} = (\boldsymbol{x}_1, \dots, \boldsymbol{x}_b) \in \mathbb{R}^n$ 

$$\boldsymbol{x} = \sum_{i=1}^{b} \mathsf{U}_{i} \boldsymbol{x}_{i} \quad \text{with} \quad \boldsymbol{x}_{i} = \mathsf{U}_{i}^{\mathsf{T}} \boldsymbol{x} \in \mathbb{R}^{n_{i}}, \ i = 1, \dots, b \quad (6)$$

which is equivalent to  $\sum_{i=1}^b \mathsf{U}_i \mathsf{U}_i^\mathsf{T} = \mathsf{I}$ . Note that (6) directly implies the norm preservation  $\|\boldsymbol{x}\|_2^2 = \|\boldsymbol{x}_1\|_2^2 + \dots + \|\boldsymbol{x}_b\|_2^2$  for any  $\boldsymbol{x} \in \mathbb{R}^n$ . We are interested in a block-coordinate algorithm that uses only a subset of operator outputs corresponding to coordinates in some block  $i \in \{1,\dots,b\}$ . Hence, for an operator  $\mathsf{G}: \mathbb{R}^n \to \mathbb{R}^n$ , we define the block-coordinate operator  $\mathsf{G}_i: \mathbb{R}^n \to \mathbb{R}^{n_i}$  as

$$\mathsf{G}_i(x) \coloneqq [\mathsf{G}(x)]_i = \mathsf{U}_i^\mathsf{T} \mathsf{G}(x) \in \mathbb{R}^{n_i}, \quad x \in \mathbb{R}^n.$$
 (7)

The proposed BC-RED algorithm is summarized in Algorithm 1. Note that when b=1, we have  $n=n_1$  and  $U_1=U_1^{\mathsf{T}}=\mathsf{I}$ . Hence, the theoretical analysis in this paper is also applicable to the full-gradient RED algorithm in (3a).

As with traditional coordinate descent methods (see [29] for a review), BC-RED can be implemented using different block selection strategies. The strategy adopted for our theoretical analysis selects block indices  $i_k$  as independent and identically distributed (i.i.d.) random variables distributed uniformly over  $\{1, \ldots, b\}$ . An alternative is to proceed in epochs of b consecutive iterations, where at the start of each epoch the set  $\{1, \ldots, b\}$  is reshuffled, and  $i_k$  is then selected consecutively

from this ordered set. We numerically compare the convergence of both BC-RED variants in Section V.

BC-RED updates its iterates one randomly picked block at a time using the output of G. When the algorithm converges, it converges to the vectors in the zero set of G

$$G(\boldsymbol{x}^*) = \nabla g(\boldsymbol{x}^*) + \tau(\boldsymbol{x}^* - \mathsf{D}(\boldsymbol{x}^*)) = \mathbf{0}$$
  

$$\Leftrightarrow \quad \boldsymbol{x}^* \in \mathsf{zer}(\mathsf{G}) := \{ \boldsymbol{x} \in \mathbb{R}^n : \mathsf{G}(\boldsymbol{x}) = \mathbf{0} \}. \tag{8}$$

Consider the following two sets

$$zer(\nabla g) := \{ \boldsymbol{x} \in \mathbb{R}^n : \nabla g(\boldsymbol{x}) = \boldsymbol{0} \}$$
and 
$$fix(D) := \{ \boldsymbol{x} \in \mathbb{R}^n : \boldsymbol{x} = D(\boldsymbol{x}) \}, \tag{9}$$

where  $zer(\nabla g)$  is the set of all critical points of the datafidelity and fix(D) is the set of all fixed points of the denoiser. Intuitively, the fixed points of D correspond to all the vectors that are not denoised, and therefore can be interpreted as vectors that are *noise-free* according to the denoiser.

Note that if  $x^* \in \operatorname{zer}(\nabla g) \cap \operatorname{fix}(\mathsf{D})$ , then  $\mathsf{G}(x^*) = \mathbf{0}$  and  $x^*$  is one of the solutions of BC-RED. Hence, any vector that is consistent with the data for a convex g and noiseless according to  $\mathsf{D}$  is in the solution set. When  $\operatorname{zer}(\nabla g) \cap \operatorname{fix}(\mathsf{D}) = \varnothing$ ,  $x^* \in \operatorname{zer}(\mathsf{G})$  is an equilibrium point balancing the direction towards higher data-fit  $\nabla g(x)$  by the direction towards higher regularity  $(x - \mathsf{D}(x))$ , explicitly weighted by  $\tau > 0$  (see Fig. 3 in the appendix for an illustration). This explicit control is one of the key differences between RED and PnP.

BC-RED benefits from considerable *flexibility* compared to the full-gradient RED. Since each update is restricted to only one block of x, the algorithm is suitable for parallel implementations and can deal with problems where the vector x is distributed in space and in time. However, the maximal benefit of BC-RED is achieved when  $G_i$  is efficient to evaluate. Fortunately, it was systematically shown in [44] that many operators—common in machine learning, image processing, and compressive sensing—admit *coordinate friendly* updates.

For a specific example, consider the least-squares data-fidelity g and a patch-wise denoiser D. Define the residual vector  $r(x) \coloneqq Ax - y$  and consider a single iteration of BC-RED that produces  $x^+$  by updating the ith block of x. Then, the update direction and the residual update are computed as

$$G_i(\boldsymbol{x}) = \boldsymbol{A}_i^{\mathsf{T}} r(\boldsymbol{x}) + \tau(\boldsymbol{x}_i - \mathsf{D}(\boldsymbol{x}_i))$$
  
and  $r(\boldsymbol{x}^+) = r(\boldsymbol{x}) - \gamma \boldsymbol{A}_i \mathsf{G}_i(\boldsymbol{x}),$  (10)

where  $A_i \in \mathbb{R}^{m \times n_i}$  is a submatrix of A consisting of the columns corresponding to the ith block. In many problems of practical interest [44], the complexity of working with  $A_i$  is roughly b times lower than with A. Also, many advanced denoisers can be effectively applied on image patches rather than on the full image [45]–[47]. Therefore, in such settings, the speed of b iterations of BC-RED is expected to be comparable to a single iteration of the full-gradient RED (see also the discussion in Appendix D).

## IV. CONVERGENCE ANALYSIS AND COMPATIBILITY WITH PROXIMAL OPTIMIZATION

In this section, we present two theoretical results related to BC-RED. We first establish its convergence to an element of zer(G) and then discuss its compatibility with the theory of proximal optimization.

#### A. Fixed Point Convergence of BC-RED

Our analysis requires several assumptions that together serve as sufficient conditions for convergence.

**Assumption 1.** The operator G is such that  $zer(G) \neq \emptyset$ . There is a finite number  $R_0$  such that the distance of the initial  $x^0 \in \mathbb{R}^n$  to the farthest element of zer(G) is bounded, that is

$$\max_{\boldsymbol{x}^* \in \operatorname{zer}(\mathsf{G})} \|\boldsymbol{x}^0 - \boldsymbol{x}^*\|_2 \leq R_0.$$

This assumption is related to the existence of minimizers in the literature on traditional coordinate minimization [26]–[29]. The assumption that  $\operatorname{zer}(\mathsf{G}) \neq \varnothing$  is analogous to assuming that  $\operatorname{zer}(\nabla f) \neq \varnothing$  when minimizing a convex and smooth function f. In this setting,  $\operatorname{zer}(\nabla f)$  fully characterizes the minimizers of f. Thus,  $\operatorname{zer}(\nabla f) = \varnothing$ , implies that f does not have minimizers. Similarly,  $\operatorname{zer}(\mathsf{G}) = \varnothing$  implies that (8) has no solutions, which this assumption precludes.

The next two assumptions rely on Lipschitz constants along directions specified by specific blocks. We say that  $G_i$  is *block Lipschitz continuous* with constant  $\lambda_i > 0$  if

$$\|\mathsf{G}_i(\boldsymbol{x}) - \mathsf{G}_i(\boldsymbol{y})\|_2 \le \lambda_i \|\boldsymbol{h}_i\|_2,\tag{11}$$

where  $x = y + \bigcup_i h_i$ ,  $y \in \mathbb{R}^n$ , and  $h_i \in \mathbb{R}^{n_i}$ . When  $\lambda_i = 1$ , we say that  $G_i$  is *block nonexpansive*. Note that if an operator G is globally  $\lambda$ -Lipschitz continuous, then it is straightforward to see that each  $G_i = \bigcup_i^T G$  is also block  $\lambda$ -Lipschitz continuous.

**Assumption 2.** The function g is continuously differentiable and convex. Additionally, the block gradient  $\nabla_i g$  is block Lipschitz continuous with constant  $L_i > 0$  for each  $i \in \{1, \ldots, b\}$ . We define the largest block Lipschitz constant as  $L_{\max} := \max\{L_1, \ldots, L_b\}$ .

Let L>0 denote the global Lipschitz constant of  $\nabla g$ . We always have  $L_{\max} \leq L$  and, for some g, it may even happen that  $L_{\max} = L/b$  [29]. As we shall see, the largest possible step-size  $\gamma$  of BC-RED depends on  $L_{\max}$ , while that of the traditional full-gradient RED on L. Hence, one natural advantage of BC-RED is that it can often take more aggressive steps compared to the full-gradient RED.

**Assumption 3.** The denoiser D is such that each block denoiser  $D_i$  is block nonexpansive.

Since the proximal operator is nonexpansive [2], it automatically satisfies this assumption. We revisit this scenario in a greater depth in Section IV-B. We can now establish the following result for BC-RED.

**Theorem 1.** Run BC-RED for  $t \ge 1$  iterations with random i.i.d. block selection under Assumptions 1-3 using a fixed step-size  $0 < \gamma \le 1/(L_{\sf max} + 2\tau)$ . Then, we have

$$\mathbb{E}\left[\frac{1}{t}\sum_{k=1}^{t}\|\mathsf{G}(x^{k-1})\|_{2}^{2}\right] \leq \frac{b(L_{\mathsf{max}}+2\tau)}{\gamma t}R_{0}^{2}.$$
 (12)

A proof of the theorem is provided in Appendix A. Theorem 1 implies that  $\mathbb{E}[\|\mathsf{G}(\boldsymbol{x}^k)\|_2^2]$  is summable and  $\mathbb{E}[\|\mathsf{G}(\boldsymbol{x}^k)\|_2^2] \to 0$ , which establishes the fixed-point convergence of BC-RED in expectation to  $\mathsf{zer}(\mathsf{G})$  with O(1/t) rate, thus matching the rate of the traditional gradient-based methods [48]. The proof relies on the monotone operator theory [49], [50], widely used in the context of convex optimization [2], including in the unified analysis of various traditional coordinate descent algorithms [51], [52]. The theorem does *not* assume the existence of any regularizer h, which makes it applicable to denoisers beyond those characterized with explicit functions in (4) and (5).

Since  $L_{\text{max}} \leq L$ , one important implication of Theorem 1, is that the upper-bound on the convergence rate (in expectation) of b iterations of BC-RED is better than that of a single iteration of the full-gradient RED (to see this, note that the full-gradient rate is obtained by setting b=1,  $L_{\text{max}}=L$ , and removing the expectation in (12)). This implies that in *coordinate friendly settings* (as discussed at the end of Section III), the overall computational complexity of BC-RED can be lower than that of the full-gradient RED. This gain is primarily due to two factors: (a) possibility to pick a larger step-size  $\gamma=1/(L_{\text{max}}+2\tau)$ ; (b) immediate reuse of each local block-update when computing the next iterate (the full-gradient RED updates the full vector before computing the next iterate).

In the special case of  $D(x) = x - (1/\tau)\nabla h(x)$ , for some convex h, BC-RED reduces to the traditional coordinate descent method applied to (1). Hence, under the assumptions of Theorem 1, one can rely on the analysis of traditional randomized coordinate descent methods in [29] to obtain

$$\mathbb{E}\left[f(\boldsymbol{x}^t)\right] - f^* \le \frac{2b}{\gamma t} R_0^2 \tag{13}$$

where  $f^*$  is the minimum value in (1). A proof of (13) can be found in [29]. Therefore, such denoisers lead to explicit convex RED regularizers and O(1/t) convergence of BC-RED in terms of the objective. However, as discussed in Section IV-B, when the denoiser is a proximal operator of some convex h, BC-RED is *not* directly solving (1), but rather its approximation.

Note that Theorem 1 only provides *sufficient conditions* for the convergence of BC-RED. As corroborated by our numerical studies in Section V, the actual convergence of BC-RED is more general and often holds beyond nonexpansive denoisers. One plausible explanation for this is that such denoisers are *locally nonexpansive* over the set of input vectors used in testing. On the other hand, the recent techniques for spectral-normalization of deep neural nets [53]–[55] provide a convenient tool for building *globally nonexpansive* neural denoisers that result in provable convergence of BC-RED.

#### B. Convergence for Proximal Operators

One of the limitations of the current RED theory is in its limited backward compatibility with the theory of proximal optimization. For example, as discussed in [21] (see section "Can we mimic any prior?"), the popular total variation (TV) denoiser [32] cannot be justified with the original RED regularization function (4). In this section, we show that BC-RED (and hence also the full-gradient RED) can be used

to solve (1) for any convex, closed, and proper function h. We do this by establishing a formal link between RED and the concept of Moreau smoothing, widely used in nonsmooth optimization [56]–[58]. In particular, we consider the following proximal-operator denoiser

$$\begin{split} \mathsf{D}(\boldsymbol{z}) &= \mathsf{prox}_{(1/\tau)h}(\boldsymbol{z}) \\ &= \underset{\boldsymbol{x} \in \mathbb{R}^n}{\min} \left\{ \frac{1}{2} \|\boldsymbol{x} - \boldsymbol{z}\|_2^2 + (1/\tau)h(\boldsymbol{x}) \right\}, \end{split} \tag{14}$$

where  $\tau > 0$ ,  $z \in \mathbb{R}^n$ , and h is a closed, proper, and convex function [2]. Since the proximal operator is nonexpansive, Assumption 3 is automatically satisfied. Our analysis, however, requires an additional assumption using the constant  $R_0$  defined in Assumption 1.

**Assumption 4.** There is a finite number  $G_0$  that bounds the largest subgradient of h, that is

$$\max\{\|\boldsymbol{\xi}(\boldsymbol{x})\|_2: \boldsymbol{\xi}(\boldsymbol{x}) \in \partial h(\boldsymbol{x}), \boldsymbol{x} \in \mathcal{B}(\boldsymbol{x}^0, R_0)\} \leq G_0,$$

where  $\mathcal{B}(\mathbf{x}^0, R_0) := \{\mathbf{x} \in \mathbb{R}^n : ||\mathbf{x} - \mathbf{x}^0||_2 \le R_0\}$  denotes a ball of radius  $R_0$ , centered at  $\mathbf{x}^0$ .

This assumption on boundedness of the subgradients holds for a large number of regularizers used in practice, including both TV and the  $\ell_1$ -norm penalties. We can now establish the following result.

**Theorem 2.** Run BC-RED for  $t \ge 1$  iterations with random i.i.d. block selection and the denoiser (14) under Assumptions 1-4 using a fixed step-size  $0 < \gamma \le 1/(L_{\text{max}} + 2\tau)$ . Then

$$\mathbb{E}\left[f(\boldsymbol{x}^t)\right] - f^* \le \frac{2b}{\gamma t} R_0^2 + \frac{G_0^2}{2\tau},\tag{15}$$

where the function f is defined in (1) and  $f^*$  is its minimum.

The theorem is proved in Appendix B. It establishes that BC-RED in expectation *approximates* the solution of (1) with an error bounded by  $(G_0^2/(2\tau))$ . For example, by setting  $\tau = \sqrt{t}$  and  $\gamma = 1/(L_{\text{max}} + 2\sqrt{t})$ , one obtains the following bound

$$\mathbb{E}\left[f(\boldsymbol{x}^{t})\right] - f^* \le \frac{1}{\sqrt{t}} \left[2b(L_{\mathsf{max}} + 2)R_0^2 + G_0^2\right]. \tag{16}$$

When  $h(x) = -\log(p_x(x))$ , the proximal operator corresponds to the MAP denoiser, and the solution of BC-RED corresponds to an *approximate* MAP estimator. This approximation can be made as precise as desired by considering larger values for the parameter  $\tau > 0$ . Note that this further justifies the RED framework by establishing that it can be used to compute a minimizer of any proper, closed, and convex (but not necessarily differentiable) h. Therefore, our analysis strengthens RED by showing that it can accommodate a much larger class of explicit regularization functions, beyond those characterized in (4) and (5).

## V. NUMERICAL VALIDATION

There is a considerable recent interest in using advanced priors in the context of image recovery from underdetermined (m < n) and noisy measurements. Recent work [21]–[25] suggests significant performance improvements due to advanced

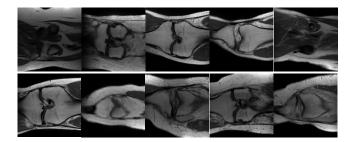


Fig. 1. Ten randomly selected test images from the fastMRI knee dataset [59].

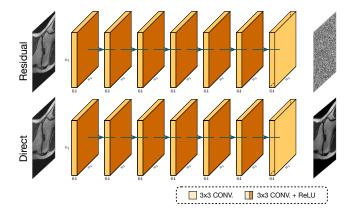


Fig. 2. The architecture of two variants of DnCNN\* used in our simulations. Each neural net is trained to remove AWGN from noisy input images. **Residual** denoiser is trained to predict the noise from the input. The final desired denoiser D is obtained by simply subtracting the predicted noise from the input  $D(z) = z - DnCNN^*(z)$ . **Direct** denoiser is trained to directly output a clean image from a noisy input  $D(z) = DnCNN^*(z)$ . In some experiments, we further constrain the Lipschitz constant (LC) of the direct denoiser to LC = 1 and of the residual denoiser to LC = 2 using spectral normalization [54]. LC = 1 implies a nonexpansive denoiser. A residual R = I - D with LC = 2 provides a necessary (but not sufficient) condition for a nonexpansive denoiser.

denoisers (such as BM3D [3] or DnCNN [4]) over traditional sparsity-driven priors (such as TV [32]). Our goal is to complement these studies with simulations validating our theoretical analysis and providing additional insights<sup>1</sup>. The simulations in this paper were performed on a machine equipped with an Intel Xeon E5-2620 v4 that has 8 cores of 2.1 GHz and 264 GBs of DDR memory. We trained all neural nets using NVIDIA RTX 2080 GPUs (see Appendix E for the details on training).

We consider inverse problems of form y = Ax + e, where  $e \in \mathbb{R}^m$  is an AWGN vector and  $A \in \mathbb{R}^{m \times n}$  is a matrix corresponding to either a sparse-view Radon transform, i.i.d. zero-mean Gaussian random matrix of variance 1/m, or radially subsampled two-dimensional Fourier transform. Such matrices are commonly used in the context of computerized tomography (CT) [60], compressive sensing [34], [35], and magnetic resonance imaging (MRI) [61], respectively. In all simulations, we set the measurement ratio to be approximately m/n = 0.5 with AWGN corresponding to input signal-to-noise ratio (SNR) of 30 dB and 40 dB. Throughout this paper, we

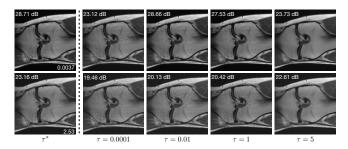


Fig. 3. Evolution of the images reconstructed by BC-RED using the DnCNN\* denoiser for different values of  $\tau$ . The first row corresponds to Fourier matrix with 30 dB noise, while the second row corresponds to the Radon matrix with 40 dB noise. Each reconstructed image is marked with its SNR value with respect to the ground truth image. The optimal parameters  $\tau^*$  for the two problems are 0.0037 and 2.35, respectively. The denoiser used in this simulation is the residual DnCNN\* with a Lipschitz constant LC = 2. This figure illustrates how  $\tau$  enables an explicit tradeoff between the data-fit and the regularization.

define SNR using the follows equation

$$SNR(\hat{\boldsymbol{y}}, \boldsymbol{y}) \triangleq 20 \log_{10} \left( \frac{\|\boldsymbol{y}\|_2}{\|\boldsymbol{y} - \hat{\boldsymbol{y}}\|_2} \right)$$

where  $\hat{y}$  represents the noisy vector and y denotes the ground truth. Fig. 1 shows the images we used, which correspond to 10 images randomly selected from the NYU fastMRI dataset [59], resized to be  $160 \times 160$  pixels. BC-RED is set to work with 16 blocks, each of size  $40 \times 40$  pixels. The reconstruction quality is quantified using SNR averaged over all ten test images.

In addition to well-studied denoisers, such as TV and BM3D, we design our own deep neural net denoiser denoted DnCNN\*, which is a simplified version of the popular DnCNN denoiser (see Fig. 2 for illustration and Sec. E in Appendix for details). This simplification reduces the computational complexity of denoising, which is important when running many iterations of BC-RED. Additionally, it makes it easier to control the global Lipschitz constant (LC) of the neural net via spectral-normalization [54]. We train DnCNN\* for the removal of AWGN at four noise levels corresponding to  $\sigma \in \{5, 10, 15, 20\}$ . For each experiment, we select the denoiser achieving the highest SNR value. Note that the  $\sigma$  parameter of BM3D is also fine-tuned for each experiment from the same set  $\{5, 10, 15, 20\}$ .

Theorem 1 establishes the convergence of BC-RED in expectation to an element of zer(G). This is illustrated in Fig. 4 (left) for the Radon matrix with 30 dB noise and a nonexpansive DnCNN\* denoiser. The average value of  $\|\mathsf{G}(\boldsymbol{x}^k)\|_2^2/\|\mathsf{G}(\boldsymbol{x}^0)\|_2^2$ is plotted against the iteration number for the full-gradient RED and BC-RED, with b updates of BC-RED (each modifying a single block) represented as one iteration. We numerically tested two block selection rules for BC-RED (i.i.d. and epoch) and observed that processing in randomized epochs leads to a faster convergence. For reference, the figure also plots the normalized squared norm of the gradient mapping vectors produced by the traditional PGM with TV [62]. The shaded areas indicate the range of values taken over 10 runs corresponding to each test image. The results highlight the potential of BC-RED to enjoy a better convergence rate compared to the full-gradient RED, with BC-RED (epoch) achieving the accuracy of  $10^{-10}$  in

<sup>&</sup>lt;sup>1</sup>Our code is available on GitHub https://github.com/wustl-cig/bcred.

TABLE I

AVERAGE SIGNAL-TO-NOISE RATIOS (SNRs) COMPUTED OVER 10 TEST IMAGES FOR DIFFERENT INVERSE PROBLEMS AND NOISE LEVELS. THE BEST SNR
FOR EACH EXPRIMENT IS HIGHLIGHTED IN BOLD-ITALIC, WHILE THE BEST DENOISER PRIOR IS IN LIGHT-GREEN.

Forward models		FISTA (TV)	UNet	RED			BC-RED		
				TV	BM3D	DnCNN*	TV	BM3D	DnCNN*
Radon	30 dB 40 dB	20.66 24.40	<b>21.90</b> 21.72	20.79 24.46	21.55 25.24	20.89 24.38	20.78 24.42	21.56 25.16	20.88 24.42
Random	30 dB 40 dB	26.07 28.42	16.37 16.40	25.64 28.30	26.46 27.82	26.53 28.05	25.70 28.39	26.50 27.88	<b>26.60</b> 28.12
Fourier	30 dB 40 dB	28.74 29.99	22.11 22.11	28.67 29.97	28.89 29.79	29.33 30.32	28.71 29.99	28.85 29.80	29.40 30.39
$\ G(\mathbf{x}^0)\ _2^2$	don (30 dl	В)	RED (Dn	M (TV)   CNN*)   (i.i.d.)			1		100

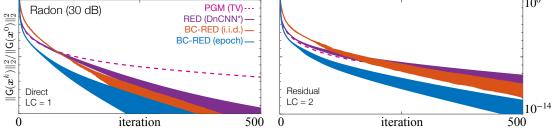


Fig. 4. Illustration of the convergence of BC-RED under two DnCNN\* priors. The left plot correspond to the **direct DnCNN**\* with the LC = 1, while the right plot correspond to the **residual DnCNN**\* with LC = 2. The average normalized distance to zer(G) is plotted against the iteration number for the Radon matrix with the shaded areas representing the range of values attained over all test images. Note that LC = 1 implies a nonexpansive denoiser, and LC = 2 provides a necessary (but not sufficient) condition for a nonexpansive denoiser.

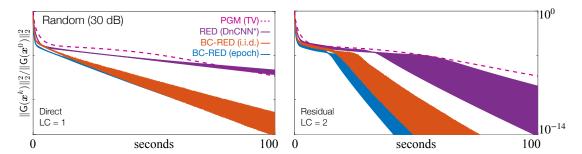


Fig. 5. Illustration of the run-time convergence of BC-RED under two  $DnCNN^*$  priors. The left plot correspond to the **direct DnCNN^\* with the LC = 1**, while the right plot correspond to the **residual DnCNN^\* with LC = 2**. The average normalized distance to zer(G) is plotted against the run-time for the Random matrix with the shaded areas representing the range of values attained over all test images. The run-time convergence of PGM with TV is also plotted for reference.

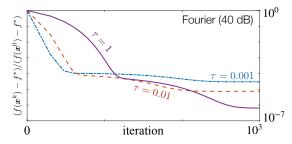


Fig. 6. Illustration of the influence of the parameter  $\tau>0$  for solving TV regularized least-squares problem using BC-RED. As  $\tau$  increases, BC-RED provides an increasingly accurate approximation to the TV optimization problem.

104 iterations, while the full-gradient RED achieves the same accuracy in 190 iterations. Additionally, the faster convergence in time of BC-RED, compared to the full-gradient RED, is

also empirically highlighted in Fig. 5 for the Random matrix with 30 dB noise and the nonexpansive DnCNN\* denoiser. Specifically, BC-RED (epoch) achieves the accuracy of  $10^{-13}$  in 100 seconds, while the full-gradient RED achieves the accuracy of  $10^{-9}$  in the same amount of time. This speedup shows the potential of of BC-RED to lead to faster convergence when applied to coordinate friendly reconstruction problems.

Theorem 2 establishes that for proximal-operator denoisers, BC-RED computes an approximate solution to (1) with an accuracy controlled by the parameter  $\tau$ . This is illustrated in Fig. 6 for the Fourier matrix with 40 dB noise and the TV regularized least-squares problem. The average value of  $(f(\boldsymbol{x}^k) - f^*)/(f(\boldsymbol{x}^0) - f^*)$  is plotted against the iteration number for BC-RED with  $\tau \in \{0.01, 0.1, 1\}$ . The optimal value  $f^*$  is obtained by running the traditional PGM until convergence. As before, the figure groups b updates of BC-

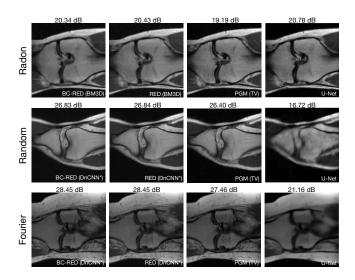


Fig. 7. Visual comparison between BC-RED and RED against PGM (TV) and U-Net for all three matrices with 30 dB noise. For BC-RED and RED, we selected the denoiser resulting in the best reconstruction performance. Every image is marked by its SNR value with respect to the ground truth. We highlight the excellent agreement between BC-RED and RED in all experiments. Note the strong degradation in the image quality for U-Net, due to the mismatch between the training and testing.

RED as a single iteration. The results are consistent with our theoretical analysis and show that as  $\tau$  increases BC-RED provides an increasingly accurate solution to TV. Since the range of possible values for the step-size  $\gamma$  depends on  $\tau$ , the speed of convergence to  $f^*$  is also influenced by  $\tau$ .

In BC-RED, the parameter  $\tau$  controls the tradeoff between  $\operatorname{zer}(\nabla g)$  and  $\operatorname{fix}(\mathsf{D})$ . Fig. 3 illustrates evolution of images reconstructed by BC-RED for different  $\tau$ . The first row corresponds to the reconstruction from the Fourier measurements with 30 dB noise, while the second row corresponds to the Radon measurements with 40 dB noise. The figure clearly shows how  $\tau$  explicitly adjusts the balance between the data-fit and the denoiser. In particular, small  $\tau$ , corresponding to weak denoising, results in unwanted artifacts in the reconstructed images, while large  $\tau$  promotes denoising strength but smooths out desired features and details. The leftmost images in Fig. 3 shows the optimal balance introduced by  $\tau^*$ .

The benefits of the full-gradient RED algorithms have been well discussed in prior work [21]-[25]. Table I summarizes the average SNR performance of BC-RED in comparison to the full-gradient RED for all three matrix types and several priors. Corresponding visual results are illustrated in Fig. 7. Unlike the full-gradient RED, BC-RED is implemented using block-wise denoisers that work on image patches rather than the full images. We empirically found that 40 pixel padding on the denoiser input is sufficient for BC-RED to match the performance of the full-gradient RED (see Appendix G for additional details). The table also includes the results for the traditional PGM with TV [62] and the widely-used end-to-end U-Net approach [63], [64]. The latter first backprojects the measurements into the image domain and then denoises the result using U-Net [65]. The model was specifically trained end-to-end for the Radon matrix with 30 dB noise and applied as such to other measurement settings. All the algorithms

were run until convergence with hyperparameters optimized for SNR. The DnCNN\* denoiser in the table corresponds to the residual network with LC = 2. The overall best SNR in the table is highlighted in bold-italic, while the best RED prior is highlighted in light-green. First, note the excellent agreement between BC-RED and the full-gradient RED. This close agreement between two methods is encouraging as BC-RED relies on block-wise denoising and our analysis does not establish uniqueness of the solution, yet, in practice, both methods seem to yield solutions of nearly identical quality. Second, note that BC-RED and RED provide excellent approximations to PGM-TV solutions. Third, note how (unlike U-Net) BC-RED and RED with DnCNN\* generalize to different measurement models. Finally, no prior seems to be universally good on all measurement settings, which indicates to the potential benefit of tailoring specific priors to specific measurement models.

Fig. 7 visually compares the images recovered by BC-RED and RED and two baseline methods. First, the images visually illustrate the excellent agreement between BC-RED and RED. Second, leveraging advanced denoisers in BC-RED largely improves the reconstruction quality over PGM with the traditional TV prior. For instance, BC-RED under DnCNN\* outperforms PGM under TV by 1 dB for Fourier matrix. Finally, we note the stability of BC-RED using the deep neural net denoiser versus the deteriorating performance of U-Net, which is trained end-to-end for Radon matrix with 30 dB noise. This fact highlights one key merit of the RED framework, that the denoiser, only trained once, can be directly applied in different scenarios for different tasks with no degradation.

#### VI. CONCLUSION

Coordinate descent methods have become increasingly important in optimization for solving large-scale problems arising in data analysis. We have introduced BC-RED as a coordinate descent extension to the current family of RED algorithms and theoretically analyzed its convergence. Our analysis provides two complementary interpretations for BC-RED: (i) equilibrium interpretation where the algorithm balances the direction towards higher data-fit by the direction towards higher regularity; (ii) minimization interpretation where the algorithm can perform traditional regularized inversion for arbitrary convex regularizers. Preliminary experiments suggest that BC-RED can be an effective tool in large-scale estimation problems arising in image recovery. More experiments are certainly needed to better asses the promise of this approach in various estimation tasks. The method and analysis presented in this paper can be extended in several complementary ways. One interesting direction would be to allow the algorithm to consider overlapping blocks, as is typically done in patch-based image denoising [3], [46]. Another interesting direction would be to consider improving the efficiency of the algorithm by designing data-adaptive block selection strategies [66]. Finally, one might consider going beyond gradient-based algorithms by designing block-coordinate variants of ADMM operating on image patches.

#### APPENDIX

#### A. Proof of Theorem 1

A fixed-point convergence of averaged operators is well-known under the name of Krasnosel'skii-Mann theorem (see Section 5.2 in [49]) and was recently applied to the analysis of PnP [13] and several full-gradient RED algorithms in [23]. We extend these results to the block-coordinate setting and provides explicit worst-case convergence rates for BC-RED. Our analysis of BC-RED relates to the analysis of block-coordinate optimization algorithms by Tseng [26], Nesterov [27], Beck and Tetruashvili [28], and Wright [29]. The key difference of our analysis from those prior works is that it does not require the prior to be expressible in the form of a regularization function, enabling BC-RED to exploit most effective image denoisers, such as those based on deep neural nets.

We consider the following operators  $G_i = \nabla_i g + H_i$  with  $H_i = \tau U_i^T (I - D)$ . and proceed in several steps.

- (a) Since  $\nabla_i g$  is block  $L_i$ -Lipschitz continuous, it is also block  $L_{\text{max}}$ -Lipschitz continuous. Hence, we know from Proposition 7 that it is block  $(1/L_{\text{max}})$ -cocoercive. Then from Proposition 4, we know that the operator  $(U_i^T (2/L_{\text{max}})\nabla_i g)$  is block nonexpansive.
- (b) From the definition of  $H_i$  and the fact that  $D_i$  is block nonexpansive, we know that  $(U_i^T (1/\tau)H_i) = D_i$  is block nonexpansive.
- (c) From Proposition 1, we know that a convex combination of block nonexpansive operators is also block nonexpansive, hence we conclude that

$$\begin{split} & \mathbf{U}_i^\mathsf{T} - \frac{2}{L_{\mathsf{max}} + 2\tau} \mathbf{G}_i = + \left( \frac{2}{L_{\mathsf{max}} + 2\tau} \cdot \frac{2\tau}{2} \right) \left[ \mathbf{U}_i^\mathsf{T} - \frac{1}{\tau} \mathbf{H}_i \right] \\ & = \left( \frac{2}{L_{\mathsf{max}} + 2\tau} \cdot \frac{L_{\mathsf{max}}}{2} \right) \left[ \mathbf{U}_i^\mathsf{T} - \frac{2}{L_{\mathsf{max}}} \nabla_i g \right], \end{split}$$

is block nonexpansive. Then from Proposition 4, we know that  $G_i$  is block  $1/(L_{\text{max}} + 2\tau)$ -cocoercive.

(d) Consider any  $x^* \in \text{zer}(\mathsf{G})$ , an index  $i \in \{1, \dots, b\}$  picked uniformly at random, and a single iteration of BC-RED  $x^+ = x - \gamma \mathsf{U}_i \mathsf{G}_i x$ . Define a vector  $h_i := \mathsf{U}_i^\mathsf{T}(x - x^*) \in \mathbb{R}^{n_i}$ . We then have

$$\|\boldsymbol{x}^{+} - \boldsymbol{x}^{*}\|^{2}$$

$$= \|\boldsymbol{x} - \boldsymbol{x}^{*} - \gamma \mathsf{U}_{i}\mathsf{G}_{i}\boldsymbol{x}\|^{2}$$

$$= \|\boldsymbol{x} - \boldsymbol{x}^{*}\|^{2} - 2\gamma(\mathsf{G}_{i}\boldsymbol{x} - \mathsf{G}_{i}\boldsymbol{x}^{*})^{\mathsf{T}}\boldsymbol{h}_{i} + \gamma^{2}\|\mathsf{G}_{i}\boldsymbol{x}\|^{2}$$

$$\leq \|\boldsymbol{x} - \boldsymbol{x}^{*}\|^{2} - \frac{2\gamma - (L_{\max} + 2\tau)\gamma^{2}}{L_{\max} + 2\tau}\|\mathsf{G}_{i}\boldsymbol{x}\|^{2}$$

$$\leq \|\boldsymbol{x} - \boldsymbol{x}^{*}\|^{2} - \frac{\gamma}{L_{\max} + 2\tau}\|\mathsf{G}_{i}\boldsymbol{x}\|^{2}, \tag{17}$$

where we used  $G_i x^* = U_i^T G x^* = 0$ , the block cocoercivity of  $G_i$ , and the fact that  $0 < \gamma \le 1/(L_{\text{max}} + 2\tau)$ .

(e) By taking a conditional expectation on both sides and rearranging the terms, we obtain

$$\begin{split} &\frac{\gamma}{L_{\mathsf{max}} + 2\tau} \mathbb{E}\left[\|\mathsf{G}_i \boldsymbol{x}\|^2 | \boldsymbol{x}\right] = \frac{\gamma}{b(L_{\mathsf{max}} + 2\tau)} \sum_{i=1}^b \|\mathsf{G}_i \boldsymbol{x}\|^2 \\ &\leq \mathbb{E}\left[\|\boldsymbol{x} - \boldsymbol{x}^*\|^2 - \|\boldsymbol{x}^+ - \boldsymbol{x}^*\|^2 | \boldsymbol{x}\right] \end{split}$$

(f) Hence by averaging over  $t \ge 1$  iterations and taking the total expectation

$$\mathbb{E}\left[\frac{1}{t}\sum_{k=1}^{t}\|\mathsf{G}\boldsymbol{x}^{k-1}\|^{2}\right] \leq \frac{1}{t}\left[\frac{b(L_{\mathsf{max}}+2\tau)}{\gamma}R_{0}^{2}\right]. \quad (18)$$

The last inequality directly leads to the result.

**Remark**. Eq. (17) implies that, under Assumptions 1-3, the iterates of BC-RED satisfy

$$\|\boldsymbol{x}^{t} - \boldsymbol{x}^{*}\| \le \|\boldsymbol{x}^{t-1} - \boldsymbol{x}^{*}\| \le \dots \le \|\boldsymbol{x}^{0} - \boldsymbol{x}^{*}\| \le R_{0},$$
 (19)

which means that the distance of the iterates of BC-RED to zer(G) is nonincreasing.

**Remark**. Our analysis in this section can be significantly strengthened if one adopts an additional assumption that g is strongly convex. This would imply that the algorithm corresponds to repeated applications of a *contractive operator* [49], which would establish the existence of a *unique* fixed point  $x^* \in \text{zer}(\mathsf{G})$  and the linear convergence of the algorithm. Our focus on the weaker form of convexity of g comes from its broader applicability in computational imaging.

#### B. Proof of Theorem 2

The concept of Moreau smoothing is well-known and has been extensively used in other contexts (see for example [58]). Our contribution is to formally connect the concept to RED-based algorithms, which leads to its novel justification as an approximate MAP estimator. The basic review of relevant concepts from proximal optimization is given in Appendix C.

For  $\tau > 0$ , we consider the Moreau envelope of h

$$h_{(1/\tau)}(\boldsymbol{x}) \, \coloneqq \, \min_{\boldsymbol{z} \in \mathbb{R}^n} \left\{ \frac{1}{2} \| \boldsymbol{z} - \boldsymbol{x} \|^2 + (1/\tau) h(\boldsymbol{z}) \right\}.$$

From Proposition 9, we know that

$$0 \le h(\boldsymbol{x}) - \tau h_{(1/\tau)}(\boldsymbol{x}) \le \frac{G_0}{2\tau}$$
 (20)

and from Proposition 8, we know that

$$\tau \nabla h_{(1/\tau)}(\boldsymbol{x}) = \tau(\boldsymbol{x} - \operatorname{prox}_{(1/\tau)h}(\boldsymbol{x})). \tag{21}$$

Hence, we can express the function f as follows

$$f(\mathbf{x}) = g(\mathbf{x}) + h(\mathbf{x})$$
  
=  $(g(\mathbf{x}) + \tau h_{(1/\tau)}(\mathbf{x})) + (h(\mathbf{x}) - \tau h_{(1/\tau)}(\mathbf{x}))$   
=  $f_{(1/\tau)}(\mathbf{x}) + (h(\mathbf{x}) - \tau h_{(1/\tau)}(\mathbf{x})),$ 

where  $f_{(1/\tau)} := g + \tau h_{(1/\tau)}$ . From eq. (21), we conclude that a single iteration of BC-RED

$$x^+ = x - \gamma \mathsf{U}_i \mathsf{G}_i x$$
 with  $\mathsf{G}_i = \mathsf{U}_i^\mathsf{T} (\nabla q(x) + \tau \nabla h_{(1/\tau)}(x))$ 

is performing a block-coordinate descent on the function  $f_{(1/\tau)}$ . From eq. (20) and the convexity of the Moreau envelope, we have

$$f_{(1/\tau)}^* = f_{(1/\tau)}(\boldsymbol{x}^*) \le f_{(1/\tau)}(\boldsymbol{x}) \le f(\boldsymbol{x}),$$

where  $x \in \mathbb{R}^n$ ,  $x^* \in \text{zer}(\mathsf{G})$ . Hence, there exists a finite  $f^*$ such that  $f(x) \geq f^*$  with  $f^*_{(1/\tau)} \leq f^*$ . Consider the iteration  $t \ge 1$  of BC-RED, then we have that

$$\begin{split} & \mathbb{E}[f(\boldsymbol{x}^t)] - f^* \leq \mathbb{E}[f(\boldsymbol{x}^t)] - f^*_{(1/\tau)} \\ & = (\mathbb{E}[f_{(1/\tau)}(\boldsymbol{x}^t)] - f^*_{(1/\tau)}) + \mathbb{E}[(h(\boldsymbol{x}^t) - \tau h_{(1/\tau)}(\boldsymbol{x}^t)]) \\ & \leq \frac{2b}{\gamma t} R_0^2 + \frac{G_0^2}{2\tau}, \end{split}$$

where we applied (13).

The proof of eq. (16) is directly obtained by setting  $\tau = \sqrt{t}$ ,  $\gamma = L_{\mathsf{max}} + 2\sqrt{t}$ , and noting that  $t \geq \sqrt{t}$ , for all  $t \geq 1$ .

#### C. Background Material

The results in this section are well-known in the optimization literature and can be found in different forms in standard textbooks [48], [49], [57], [67]. For completeness, we summarize the key results useful for our analysis by restating them in a block-coordinate form.

## Properties of Block-Coordinate Operators

Most of the concepts in this part come from the traditional monotone operator theory [49], [50] adapted for blockcoordinate operators.

**Definition 1.** We define the block-coordinate operator  $T_i$ :  $\mathbb{R}^n \to \mathbb{R}^{n_i}$  of  $\mathsf{T}: \mathbb{R}^n \to \mathbb{R}^n$  as

$$\mathsf{T}_i oldsymbol{x} \coloneqq [\mathsf{T} oldsymbol{x}]_i = \mathsf{U}_i^\mathsf{T} \mathsf{T} oldsymbol{x} \in \mathbb{R}^{n_i}, \quad oldsymbol{x} \in \mathbb{R}^n.$$

The operator  $T_i$  applies T to its input vector and then extracts the subset of outputs corresponding to the coordinates in the *block*  $i \in \{1, ..., b\}$ .

**Remark.** When b = 1, we have that  $n = n_1$  and  $U_1 =$  $U_1^T = I$ . Then, all the properties in this section reduce to their standard counterparts from the monotone operator theory in  $\mathbb{R}^n$ . In such settings, we simply drop the word *block* from the name of the property.

**Definition 2.**  $T_i$  is block Lipschitz continuous with constant  $\lambda_i > 0$  if

$$\|\mathsf{T}_i x - \mathsf{T}_i y\| \le \lambda_i \|h_i\|, \quad x = y + \mathsf{U}_i h_i, \quad y \in \mathbb{R}^n, h_i \in \mathbb{R}^{n_i}.$$

When  $\lambda_i = 1$ , we say that  $T_i$  is block nonexpansive.

**Definition 3.** An operator  $T_i$  is block cocoercive with constant  $\beta_i > 0$  if

$$(\mathsf{T}_i \boldsymbol{x} - \mathsf{T}_i \boldsymbol{y})^\mathsf{T} \boldsymbol{h}_i \ge \beta_i \|\mathsf{T}_i \boldsymbol{x} - \mathsf{T}_i \boldsymbol{y}\|^2,$$

where  $x = y + \bigcup_i h_i$ ,  $y \in \mathbb{R}^n$ ,  $h_i \in \mathbb{R}^{n_i}$ . When  $\beta_i = 1$ , we say that  $T_i$  is block firmly nonexpansive.

The following propositions are conclusions derived from the definition of above.

**Proposition 1.** Let  $T_{ij}: \mathbb{R}^n \to \mathbb{R}^{n_i}$  for  $j \in J$  be a set of block nonexpansive operators. Then, their convex combination

$$\mathsf{T}_i := \sum_{j \in J} \theta_j \mathsf{T}_{ij}, \quad \textit{with} \quad \theta_j > 0 \ \textit{and} \ \sum_{j \in J} \theta_j = 1,$$

is nonexpansive.

*Proof.* By using the triangular inequality and the definition of block nonexpansiveness, we obtain

$$\|\mathsf{T}_i oldsymbol{x} - \mathsf{T}_i oldsymbol{y}\| \le \sum_{j \in J} heta_j \|\mathsf{T}_{ij} oldsymbol{x} - \mathsf{T}_{ij} oldsymbol{y}\|$$
 $\le \left(\sum_{j \in J} heta_j\right) \|oldsymbol{h}_i\| = \|oldsymbol{h}_i\|,$ 

for all  $y \in \mathbb{R}^n$  and  $h_i \in \mathbb{R}^{n_i}$  where  $x = y + \bigcup_i h_i$ . 

**Proposition 2.** Consider  $R_i = U_i^T - T_i$  where  $T_i : \mathbb{R}^n \to \mathbb{R}^{n_i}$ .

 $T_i$  is block nonexpansive  $\Leftrightarrow R_i$  is (1/2)-block cocoercive.

*Proof.* First suppose that  $R_i$  is 1/2 block cocoercive. Let x = $y + \mathsf{U}_i h_i$  for all  $y \in \mathbb{R}^n$  and  $h_i \in \mathbb{R}^{n_i}$ . We then have

$$\frac{1}{2}\|\mathsf{R}_i\boldsymbol{x}-\mathsf{R}_i\boldsymbol{y}\|^2 \leq (\mathsf{R}_i\boldsymbol{x}-\mathsf{R}_i\boldsymbol{y})^\mathsf{T}\boldsymbol{h}_i = \|\boldsymbol{h}_i\|^2 - (\mathsf{T}_i\boldsymbol{x}-\mathsf{T}_i\boldsymbol{y})^\mathsf{T}\boldsymbol{h}_i.$$

We also have that

$$\frac{1}{2}\|\mathsf{R}_i\boldsymbol{x} - \mathsf{R}_i\boldsymbol{y}\|^2 = \frac{1}{2}\|\boldsymbol{h}_i\|^2 - (\mathsf{T}_i\boldsymbol{x} - \mathsf{T}_i\boldsymbol{y})^\mathsf{T}\boldsymbol{h}_i + \frac{1}{2}\|\mathsf{T}_i\boldsymbol{x} - \mathsf{T}_i\boldsymbol{y}\|^2.$$

By combining these two and simplifying the expression, we obtain that

$$\|\mathsf{T}_i \boldsymbol{x} - \mathsf{T}_i \boldsymbol{y}\| \leq \|\boldsymbol{h}_i\|.$$

The converse can be proved by following this logic in reverse.

**Block Averaged Operators** 

It is well known that the iteration of a nonexpansive operator does not necessarily converge. To see this consider a nonexpansive operator T = -I, where I is identity. However, it is also well known that the convergence can be established for averaged operators.

**Definition 4.** For a constant  $\alpha \in (0,1)$ , we say that the operator T is  $\alpha$ -averaged, if there exists a nonexpansive operator N such that  $T = (1 - \alpha)I + \alpha N$ .

**Definition 5.** For a constant  $\alpha \in (0,1)$ , we say that  $T_i : \mathbb{R}^n \to$  $\mathbb{R}^{n_i}$  is block  $\alpha$ -averaged, if there exists a block nonexpansive operator  $N_i$  such that  $T_i = (1 - \alpha)U_i^T + \alpha N_i$ .

**Remark.** It is clear that if T is  $\alpha$ -averaged, then  $T_i = U_i^T T$ is block  $\alpha$ -averaged.

The following characterization is often convenient.

**Proposition 3.** For a block nonexpansive operator  $T_i$ , a constant  $\alpha \in (0,1)$ , and the operator  $R_i := U_i^T - T_i$ , the following are equivalent

- (a)  $T_i$  is block  $\alpha$ -averaged
- (b)  $(1-1/\alpha)\mathsf{U}_i^\mathsf{T} + (1/\alpha)\mathsf{T}_i$  is block nonexpansive (c)  $\|\mathsf{T}_i \boldsymbol{x} \mathsf{T}_i \boldsymbol{y}\|^2 \le \|\boldsymbol{h}_i\|^2 \left(\frac{1-\alpha}{\alpha}\right)\|\mathsf{R}_i \boldsymbol{x} \mathsf{R}_i \boldsymbol{y}\|^2$ ,  $\boldsymbol{x} = \boldsymbol{y} + \mathsf{U}_i \boldsymbol{h}_i$ ,  $\boldsymbol{y} \in \mathbb{R}^n$ ,  $\boldsymbol{h}_i \in \mathbb{R}^{n_i}$

Proof. The equivalence of (a) and (b) is clear from the definition. To establish the equivalence with (c), consider an operator  $N_i$  and  $T_i = (1 - \alpha)U_i^T + \alpha N_i$ . Note that

$$R_i = U_i^\mathsf{T} - \mathsf{T}_i = \alpha(U_i^\mathsf{T} - \mathsf{N}_i).$$

Then, for all  $y \in \mathbb{R}^n$  and  $h_i \in \mathbb{R}^{n_i}$ , with  $x = y + \bigcup_i h_i$ , we have that

$$\|\mathbf{T}_{i}\boldsymbol{x} - \mathbf{T}_{i}\boldsymbol{y}\|^{2} = \|(1 - \alpha)\boldsymbol{h}_{i} + \alpha(\mathbf{N}_{i}\boldsymbol{x} - \mathbf{N}_{i}\boldsymbol{y})\|^{2}$$

$$= (1 - \alpha)\|\boldsymbol{h}_{i}\|^{2} + \alpha\|\mathbf{N}_{i}\boldsymbol{x} - \mathbf{N}_{i}\boldsymbol{y}\|^{2} -$$

$$\alpha(1 - \alpha)\|\boldsymbol{h}_{i} - (\mathbf{N}_{i}\boldsymbol{x} - \mathbf{N}_{i}\boldsymbol{y})\|^{2}$$

$$= (1 - \alpha)\|\boldsymbol{h}_{i}\|^{2} + \alpha\|\mathbf{N}_{i}\boldsymbol{x} - \mathbf{N}_{i}\boldsymbol{y}\|^{2}$$

$$- \left(\frac{1 - \alpha}{\alpha}\right)\|\mathbf{R}_{i}\boldsymbol{x} - \mathbf{R}_{i}\boldsymbol{y}\|^{2}, \tag{22}$$

where we used the fact that

$$\|(1-\alpha)\boldsymbol{x} + \alpha \boldsymbol{y}\|^2 = (1-\alpha)\|\boldsymbol{x}\|^2 + \alpha\|\boldsymbol{y}\|^2 - \alpha(1-\alpha)\|\boldsymbol{x} - \boldsymbol{y}\|^2,$$
 where  $\theta \in \mathbb{R}$  and  $\boldsymbol{x}, \boldsymbol{y} \in \mathbb{R}^n$ . Consider also

$$\|\boldsymbol{h}_i\|^2 - \left(\frac{1-\alpha}{\alpha}\right) \|\mathsf{R}_i \boldsymbol{x} - \mathsf{R}_i \boldsymbol{y}\|^2$$

$$= (1-\alpha)\|\boldsymbol{h}_i\|^2 + \alpha\|\boldsymbol{h}_i\|^2 - \left(\frac{1-\alpha}{\alpha}\right) \|\mathsf{R}_i \boldsymbol{x} - \mathsf{R}_i \boldsymbol{y}\|^2.$$
(23)

It is clear that we have

(22) 
$$\leq$$
 (23)  $\Leftrightarrow$   $\mathsf{N}_i$  is block nonexpansive  $\Leftrightarrow$   $\mathsf{T}_i$  is block  $\alpha$ -averaged, (24)

where we used the definition of block averagedness.

**Proposition 4.** Consider a block-coordinate operator  $\mathsf{T}_i = \mathsf{U}_i^\mathsf{T} \mathsf{T}$  with  $\mathsf{T} : \mathbb{R}^n \to \mathbb{R}^n$ . Let  $x = y + \mathsf{U}_i h$  with  $x \in \mathbb{R}^n$ ,  $h_i \in \mathbb{R}^{n_i}$  and consider  $\beta_i > 0$ . Then, the following are equivalent

- (a)  $T_i$  is block  $\beta_i$ -cocoercive
- (b)  $\beta_i \mathsf{T}_i$  is block firmly nonexpansive
- (c)  $\bigcup_{i=1}^{T} -\beta_{i} T_{i}$  is block firmly nonexpansive.
- (d)  $\beta_i \mathsf{T}_i$  is block (1/2)-averaged.
- (e)  $U_i^{\mathsf{T}} 2\beta_i \mathsf{T}_i$  is block nonexpansive.

*Proof.* The equivalence between (a) and (b) is readily observed by defining  $P_i := \beta_i T_i$  and noting that

$$(\mathsf{P}_{i}\boldsymbol{x} - \mathsf{P}_{i}\boldsymbol{y})^{\mathsf{T}}\boldsymbol{h}_{i} = \beta_{i}(\mathsf{T}_{i}\boldsymbol{x} - \mathsf{T}_{i}\boldsymbol{y})^{\mathsf{T}}\boldsymbol{h}_{i}$$
  
and  $\|\mathsf{P}_{i}\boldsymbol{x} - \mathsf{P}_{i}\boldsymbol{y}\|^{2} = \beta_{i}^{2}\|\mathsf{T}_{i}\boldsymbol{x} - \mathsf{T}_{i}\boldsymbol{y}\|.$  (25)

Define  $R_i := U_i^T - P_i$  and suppose (b) is true, then

$$(\mathsf{R}_{i}\boldsymbol{x} - \mathsf{R}_{i}\boldsymbol{y})^{\mathsf{T}}\boldsymbol{h}_{i} = \|\boldsymbol{h}_{i}\|^{2} - (\mathsf{P}_{i}\boldsymbol{x} - \mathsf{P}_{i}\boldsymbol{y})^{\mathsf{T}}\boldsymbol{h}_{i}$$

$$= \|\mathsf{R}_{i}\boldsymbol{x} - \mathsf{R}_{i}\boldsymbol{y}\|^{2} + (\mathsf{P}_{i}\boldsymbol{x} - \mathsf{P}_{i}\boldsymbol{y})^{\mathsf{T}}\boldsymbol{h}_{i} - \|\mathsf{P}_{i}\boldsymbol{x} - \mathsf{P}_{i}\boldsymbol{y}\|^{2}$$

$$\geq \|\mathsf{R}_{i}\boldsymbol{x} - \mathsf{R}_{i}\boldsymbol{y}\|^{2}.$$

By repeating the same argument for  $P_i = U_i^T - R_i$ , we establish the full equivalence between (b) and (c).

The full equivalence of (b) and (d) can be established by observing that

$$2\|\mathsf{P}_{i}\boldsymbol{x} - \mathsf{P}_{i}\boldsymbol{y}\|^{2} \leq 2(\mathsf{P}_{i}\boldsymbol{x} - \mathsf{P}_{i}\boldsymbol{y})^{\mathsf{T}}\boldsymbol{h}_{i}$$

$$\Leftrightarrow \|\mathsf{P}_{i}\boldsymbol{x} - \mathsf{P}_{i}\boldsymbol{y}\|^{2} \leq 2(\mathsf{P}_{i}\boldsymbol{x} - \mathsf{P}_{i}\boldsymbol{y})^{\mathsf{T}}\boldsymbol{h}_{i} - \|\mathsf{P}_{i}\boldsymbol{x} - \mathsf{P}_{i}\boldsymbol{y}\|^{2}$$

$$= \|\boldsymbol{h}_{i}\|^{2} - (\|\boldsymbol{h}_{i}\|^{2} - 2(\mathsf{P}_{i}\boldsymbol{x} - \mathsf{P}_{i}\boldsymbol{y})^{\mathsf{T}}\boldsymbol{h}_{i} + \|\mathsf{P}_{i}\boldsymbol{x} - \mathsf{P}_{i}\boldsymbol{y}\|^{2})$$

$$= \|\boldsymbol{h}_{i}\|^{2} - \|\mathsf{R}_{i}\boldsymbol{x} - \mathsf{R}_{i}\boldsymbol{y}\|^{2}.$$

To show the equivalence with (e), first suppose that  $N_i := U_i^\mathsf{T} - 2\mathsf{P}_i$  is block nonexpansive, then  $\mathsf{P}_i = \frac{1}{2}(\mathsf{U}_i^\mathsf{T} + (-\mathsf{N}_i))$  is block 1/2-averaged, which means that it is block firmly nonexpansive. On the other hand, if  $\mathsf{P}_i$  is block firmly nonexpansive, then it is block 1/2-averaged, which means that from Proposition 3(b) we have that  $(1-2)\mathsf{U}_i^\mathsf{T} + 2\mathsf{P}_i = 2\mathsf{P}_i - \mathsf{U}_i^\mathsf{T} = -\mathsf{N}_i$  is block nonexpansive. This directly means that  $\mathsf{N}_i$  is block nonexpansive.

Operator Properties for Convex Function

It is convenient to link properties of a function  $f: \mathbb{R}^n \to \mathbb{R}$ ,  $x \mapsto y = f(x)$ , to the properties of operators derived from it. The key properties for our analysis are related to continuity and convexity.

**Proposition 5.** Let f be continuously differentiable function with  $\nabla_i f$  that is block  $L_i$ -Lipschitz continuous. Then,

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\mathsf{T} (\mathbf{y} - \mathbf{x}) + \frac{L_i}{2} \|\mathbf{y} - \mathbf{x}\|^2$$
  
=  $f(\mathbf{x}) + \nabla_i f(\mathbf{x})^\mathsf{T} \mathbf{h}_i + \frac{L_i}{2} \|\mathbf{h}_i\|^2$ 

for all  $x \in \mathbb{R}^n$  and  $h_i \in \mathbb{R}^{n_i}$ , where  $y = x + \bigcup_i h_i$ .

*Proof.* The proof is a minor variation of the one presented in Section 2.1 of [48].

**Proposition 6.** Consider a continuously differentiable f such that  $\nabla_i f$  is block  $L_i$ -Lipschitz continuous. Let  $\mathbf{x}^* \in \mathbb{R}^n$  denote the global minimizer of f. Then, we have that

$$\frac{1}{2L_i} \|\nabla_i f(\boldsymbol{x})\|^2 \le (f(\boldsymbol{x}) - f(\boldsymbol{x}^*)) \le \frac{L_i}{2} \|\boldsymbol{x} - \boldsymbol{x}^*\|^2, \quad (26)$$

where  $x = x^* + \bigcup_i h_i$ ,  $x \in \mathbb{R}^n$ ,  $h_i \in \mathbb{R}^{n_i}$ .

*Proof.* The proof is a minor variation of the discussion in Section 9.1.2 of [67].  $\Box$ 

**Proposition 7.** For a convex and continuously differentiable function f, we have

$$\nabla_i f$$
 is block  $L_i$ -Lipschitz continuous  $\Leftrightarrow \nabla_i f$  is block  $(1/L_i)$ -cocoercive.

*Proof.* The proof is a minor variation of the one presented as Theorem 2.1.5 in Section 2.1 of [48].  $\Box$ 

Moreau smoothing and proximal operators

In this section, we consider a class of functions that are proper, closed, and convex, but are not necessarily differentiable. The proximal operator is a widely-used concept in such nonsmooth optimization problems [56], [57].

**Definition 6.** Consider a proper, closed, and convex h and a constant  $\mu > 0$ . We define the proximal operator

$$\mathrm{prox}_{\mu h}(\boldsymbol{x}) \, \coloneqq \, \operatorname*{arg\,min}_{\boldsymbol{z} \in \mathbb{R}^n} \left\{ \frac{1}{2} \|\boldsymbol{z} - \boldsymbol{x}\|^2 + \mu h(\boldsymbol{z}) \right\}$$

and the Moreau envelope

$$h_{\mu}(oldsymbol{x}) \coloneqq \min_{oldsymbol{z} \in \mathbb{R}^n} \left\{ rac{1}{2} \|oldsymbol{z} - oldsymbol{x}\|^2 + \mu h(oldsymbol{z}) 
ight\}.$$

**Proposition 8.** The function  $h_{\mu}$  is convex and continuously differentiable with a 1-Lipschitz gradient

$$\nabla h_{\mu}(\boldsymbol{x}) = \boldsymbol{x} - \operatorname{prox}_{uh}(\boldsymbol{x}), \quad \boldsymbol{x} \in \mathbb{R}^{n}.$$

*Proof.* We first show that  $h_{\mu}$  is convex. Consider

$$q(\boldsymbol{x}, \boldsymbol{z}) \coloneqq \frac{1}{2} \|\boldsymbol{z} - \boldsymbol{x}\|^2 + \mu h(\boldsymbol{z}),$$

which is convex (x, z). Then, for any  $0 \le \theta \le 1$  and  $(x_1, z_1), (x_2, z_2) \in \mathbb{R}^{2n}$ , we have

$$h_{\mu}(\theta x_{1} + (1 - \theta)x_{2}) \leq q(\theta x_{1} + (1 - \theta)x_{2}, \theta z_{1} + (1 - \theta)z_{2})$$

$$\leq \theta q(x_{1}, z_{1}) + (1 - \theta)q(x_{2}, z_{2}),$$
(27)

where we used the convexity of q. Since this inequality holds everywhere, we have

$$h_{\mu}(\theta x_1 + (1 - \theta)x_2) \le \theta h_{\mu}(x_1) + (1 - \theta)h_{\mu}(x_2),$$

with

$$h_{\mu}(\boldsymbol{x}_1) = \min_{\boldsymbol{z}_1} q(\boldsymbol{x}_1, \boldsymbol{z}_1) \quad \text{and} \quad h_{\mu}(\boldsymbol{x}_2) = \min_{\boldsymbol{z}_2} q(\boldsymbol{x}_2, \boldsymbol{z}_2).$$

To show the differentiability, note that

$$\begin{split} h_{\mu}(\boldsymbol{x}) &= \frac{1}{2} \|\boldsymbol{x}\|^2 - \max_{\boldsymbol{z} \in \mathbb{R}^n} \left\{ \boldsymbol{x}^\mathsf{T} \boldsymbol{z} - \mu h(\boldsymbol{z}) - \frac{1}{2} \|\boldsymbol{z}\|^2 \right\} \\ &= \frac{1}{2} \|\boldsymbol{x}\|^2 - \phi^{\star}(\boldsymbol{x}) \quad \text{with} \quad \phi(\boldsymbol{z}) \coloneqq \frac{1}{2} \|\boldsymbol{z}\|^2 + \mu h(\boldsymbol{z}), \end{split}$$

where  $\phi^*$  denotes the conjugate of  $\phi$ . The function  $\phi$  is closed and 1-strongly convex. Hence, we know that  $\phi^*$  is defined for all  $\mathbf{x} \in \mathbb{R}^n$  and is differentiable with gradient [67]

$$\nabla \phi^{\star}(\boldsymbol{x}) = \operatorname*{arg\,max}_{\boldsymbol{z} \in \mathbb{R}^n} \left\{ \boldsymbol{x}^{\mathsf{T}} \boldsymbol{z} - \mu h(\boldsymbol{z}) - \frac{1}{2} \|\boldsymbol{z}\|^2 \right\} = \mathrm{prox}_{\mu h}(\boldsymbol{x}).$$

Hence, we conclude that

$$\nabla h_{\mu}(\boldsymbol{x}) = \boldsymbol{x} - \nabla \phi^{\star}(\boldsymbol{x}) = \boldsymbol{x} - \text{prox}_{\mu h}(\boldsymbol{x}).$$

Note that since the proximal operator is firmly nonexpansive,  $\nabla h_{\mu}$  is also firmly nonexpansive, which means that it is 1-Lipschitz.

The next result shows that the Moreau envelope can serve as a smooth approximation to a nonsmooth function.

**Proposition 9.** Consider  $h \in \mathbb{R}^n$  and its Moreau envelope  $h_{\mu}(\mathbf{x})$  for  $\mu > 0$ . Then,

$$0 \le h(\boldsymbol{x}) - \frac{1}{\mu} h_{\mu}(\boldsymbol{x}) \le \frac{\mu}{2} G_{\boldsymbol{x}}^2$$
with  $G_{\boldsymbol{x}}^2 := \min_{\boldsymbol{\xi} \in \partial h(\boldsymbol{x})} \|\boldsymbol{\xi}\|^2$ ,  $\boldsymbol{x} \in \mathbb{R}^n$ .

Proof. First note that

$$rac{1}{\mu}h_{\mu}(oldsymbol{x}) = \min_{oldsymbol{z} \in \mathbb{R}^n} \left\{ rac{1}{2\mu} \|oldsymbol{z} - oldsymbol{x}\|^2 + h(oldsymbol{z}) 
ight\} \leq h(oldsymbol{x}), \quad oldsymbol{x} \in \mathbb{R}^n,$$

**Algorithm 2** BC-RED for a quadratic g and block-wise D

```
1: input: x^{0} \in \mathbb{R}^{n}, \tau > 0, and \gamma > 0.

2: initialize: r^{0} \leftarrow Ax^{0} - y

3: for k = 1, 2, 3, ... do

4: Choose an index i_{k} \in \{1, ..., b\}

5: x^{k} \leftarrow x^{k-1} - \gamma \bigcup_{i_{k}} \mathsf{G}_{i_{k}}(x^{k-1})

with \mathsf{G}_{i_{k}}(x^{k-1}) = A_{i_{k}}^{\mathsf{T}} r^{k-1} + \tau(x_{i_{k}} - \mathsf{D}(x_{i_{k}})).

6: r^{k} \leftarrow r^{k-1} - \gamma A_{i_{k}} \mathsf{G}_{i_{k}}(x^{k-1})

7: end for
```

which is due to the fact that z = x is potentially suboptimal. We additionally have for any  $\xi \in \partial h(x)$ 

$$\begin{split} h_{\mu}(\boldsymbol{x}) - \mu h(\boldsymbol{x}) &= \min_{\boldsymbol{z} \in \mathbb{R}^n} \left\{ \mu h(\boldsymbol{z}) - \mu h(\boldsymbol{x}) + \frac{1}{2} \|\boldsymbol{z} - \boldsymbol{x}\|^2 \right\} \\ &\geq \min_{\boldsymbol{z} \in \mathbb{R}^n} \left\{ \mu \boldsymbol{\xi}^{\mathsf{T}} (\boldsymbol{z} - \boldsymbol{x}) + \frac{1}{2} \|\boldsymbol{z} - \boldsymbol{x}\|^2 \right\} \\ &= \min_{\boldsymbol{z} \in \mathbb{R}^n} \left\{ \frac{1}{2} \|\boldsymbol{z} - (\boldsymbol{x} - \mu \boldsymbol{\xi})\|^2 - \frac{\mu^2}{2} \|\boldsymbol{\xi}\|^2 \right\} = -\frac{\mu^2}{2} \|\boldsymbol{\xi}\|^2. \end{split}$$

This directly leads to the conclusion.

#### D. Coordinate-Friendly Implementations

Theoretical analysis in Section IV of the main paper suggests that, if b updates of BC-RED (each modifying a single block) are counted as a single iteration, the worst-case convergence rate of BC-RED is expected to be better than that of the full-gradient RED. This fact was empirically validated in Section V, where we showed that in practice BC-RED needs much fewer iterations to converge. However, the overall computational complexity of two methods depends on their per-iteration complexities. In particular, the overall complexity of BC-RED is favorable when its total number of iterations required for convergence offsets the cost of solving the problem in a block-coordinate fashion. As for traditional coordinate descent methods [44], [68], in many problems of interest, the computational complexity of a single update of BC-RED will be roughly b times lower than that of the full-gradient method.

The computational complexity of each block-update will depend on the specifics of the data-fidelity term g and the denoiser D used in the estimation problem. For example, consider the problem where  $g(\boldsymbol{x}) = \frac{1}{2} \|\boldsymbol{A}\boldsymbol{x} - \boldsymbol{y}\|_2^2$ . Additionally, suppose that  $\boldsymbol{x}$  is such that it is sufficient represent its prior with a block-wise denoiser on each  $\boldsymbol{x}_i$ , rather than on the full  $\boldsymbol{x}$ . This situation is very common in image processing, where many popular denoisers are applied block-wise [47]. Then, one can obtain a very efficient implementation of BC-RED, illustrated in Algorithm 2.

The worst-case complexity of applying  $A_i$  and  $A_i^{\mathsf{T}}$  is  $O(mn_i)$ , which means that the cost of b updates such updates for  $i \in \{1, \ldots, b\}$  is O(mn). Additionally, if the complexity of b block-wise denoising operations is equivalent or less than the complexity of denoising the full vector (which is generally true for advanced denoisers), then the complexity of b updates of BC-RED will be equivalent or better than a single iteration of the full-gradient RED.

TABLE II

AVERAGE SNR ACHIEVED BY BC-RED FOR TWO VARIANTS OF  $DNCNN^*$  AT DIFFERENT LIPSCHITZ CONSTANT (LC) VALUES.

Variants of DnCNN*		Radon		Ran	dom	Fourier	
		30 dB	40 dB	30 dB	40 dB	30 dB	40 dB
Direct	Unconstrained	21.67	24.74	Diverges	Diverges	29.40	30.35
	LC = 1	19.33	22.98	19.89	20.26	25.06	25.40
Residual	Unconstrained	20.88	24.68	26.49	27.60	29.39	30.31
	LC = 2	20.88	24.42	26.60	28.12	29.40	30.39

Some of our simulations were conducted using denoisers applied on the full-image and others using patch-wise denoisers. In particular, the convergence simulations in Fig. 4 relied on the full-image denoisers, in order to use identical denoisers for both RED and BC-RED and be fully compatible with the theoretical analysis. On the other hand, the SNR results in Table I, Table II, Fig. 7, and Fig. 3 rely on block-wise denoisers, where the denoiser input includes an additional 40 pixel padding around the block and the output has the exact size of the block. The padding size was determined empirically in order to have a close match between BC-RED and RED. We have observed that having even larger paddings does not influence the results of BC-RED.

#### E. Architecture and Training of DnCNN\*

We designed DnCNN\* fully based on DnCNN architecture. The network contains three parts. The first part is a composite convolutional layer, consisting of a normal convolutional layer and a rectified linear units (ReLU) layer. It convolves the  $n_1 \times n_2$  input to  $n_1 \times n_2 \times 64$  features maps by using 64 filters of size  $3 \times 3$ . The second part is a sequence of 5 composite convolutional layers, each having 64 filters of size  $3 \times 3 \times 64$ . Those composite layers further processes the feature maps generated by the first part. The third part of the network, a single convolutional layer, generates the final output image by convolving the feature maps with a  $3 \times 3 \times 64$  filter. Every convolution is performed with a stride = 1, so that the intermediate feature maps share the same spatial size of the input image. Fig. 2 visualizes the architectural details. We generated 52000 training examples by adding AWGN to 13000 images ( $320 \times 320$ ) from the NYU fastMRI dataset [59] and cropping them into 4 sub-images of size  $160 \times 160$  pixels. We trained DnCNN\* to optimize the mean squared error by using the Adam optimizer.

#### F. Influence of the Lipschitz Constant on Performance

Theorem 1 assumes that the denoiser D is nonexpansive. It is relatively straightforward to control the global Lipschiz constants of deep neural nets via spectral normalization [53]–[55] and we have empirically tested the influence of nonexpansiveness to the quality of final image recovery.

Table II summarizes the SNR performance of BC-RED for two common variants of DnCNN\*. The first variant is trained to learn the *direct* mapping from a noisy input to a clean image, while the second variant relies on *residual learning* to map its input to noise (shown in Fig. 2). To gain insight

into the influence of the Lipschitz constant (LC) of a denoiser to its performance as a prior, we trained denoisers with both globally constrained and nonconstrained LCs via the spectralnormalization technique from [54]. For the direct network, we trained  $DnCNN^*$  with LC = 1, which corresponds to a nonexpansive denoiser. For the residual network, we considered LC = 2, which is a necessary (but not sufficient) condition for the nonexpansiveness. In our simulations, BC-RED converged for all the variants of DnCNN\*, except for the direct and unconstrained DnCNN\*, which confirms that our theoretical analysis provides only sufficient conditions for convergence. Nonetheless, our simulations reveal the performance loss of the algorithm for the direct and nonexpansive (LC = 1) DnCNN\*. On the other hand, the performance of the residual  $DnCNN^*$  with LC = 2 nearly matches the performance of fully unconstrained networks in all experiments.

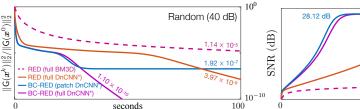
#### G. Influence of padding in patch-wise denoising

The procedure of block-wise image processing of BC-RED enables one to further reduce the overall computational complexity by using the patch-wise denoisers, where the denoising is performed only on the desired patch instead of the full image. The left table in Fig. 8 summarizes the averaged SNR values for the patch-wise residual DnCNN\* corresponding to the paddings of size {0, 5, 10, 20, 40} pixels. The lower SNR for 0 px suggests the non-separability of DnCNN\*; yet, a small 5 px padding is sufficient for matching the performance of the full-image DnCNN\*. Since the patchwise denoiser only approximates the full-image denoiser, the final accuracy of BC-RED under the patch-wise DnCNN\* to zer(G) is  $1.92 \times 10^{-7}$ , while the accuracy for the full-image DnCNN\* is  $1.10 \times 10^{-10}$ . However, the patch-wise DnCNN\* still matches the SNR performance of the full-image DnCNN\* and does it faster due to its reduced denoising complexity. The slight SNR improvement for 40 px patch-wise denoising over the full-image denoising is due the fact that  $\tau$  parameter of BC-RED was optimized for that case and reused in the rest of experiments. Note also the slow convergence of RED using the full-image BM3D, due to the lower convergence rate and the high complexity of denoising.

#### REFERENCES

- S. V. Venkatakrishnan, C. A. Bouman, and B. Wohlberg, "Plug-andplay priors for model based reconstruction," in *Proc. IEEE Global Conf.* Signal Process. and INf. Process. (GlobalSIP), 2013.
- [2] N. Parikh and S. Boyd, "Proximal algorithms," Foundations and Trends in Optimization, vol. 1, no. 3, pp. 123–231, 2014.

Padding	SNR			
0 px	$27.64~\mathrm{dB}$			
5 px	28.11  dB			
10  px	28.12  dB			
20 px	28.12  dB			
40  px	$28.12~\mathrm{dB}$			



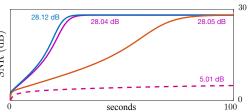


Fig. 8. *Left:* The averaged SNR values obtained by BC-RED for the Random matrix with 40 dB noise and *patch-wise* residual DnCNN\*, where the denoiser input includes an additional padding around the patch, while the output has the size of the patch. *Center and Right:* The convergence speed of BC-RED under *patch-wise* residual DnCNN\* with 40 px padding and the *full-image* residual DnCNN\*. Distance to zer(G) – corresponding to the *full-image* denoiser – and SNR are plotted against time. As a reference, we provide the convergence of RED using the full-image DnCNN\* and BM3D denoisers.

- [3] K. Dabov, A. Foi, V. Katkovnik, and K. Egiazarian, "Image denoising by sparse 3-D transform-domain collaborative filtering," *IEEE Trans. Image Process.*, vol. 16, no. 16, pp. 2080–2095, August 2007.
- [4] K. Zhang, W. Zuo, Y. Chen, D. Meng, and L. Zhang, "Beyond a Gaussian denoiser: Residual learning of deep CNN for image denoising," *IEEE Trans. Image Process.*, vol. 26, no. 7, pp. 3142–3155, July 2017.
- [5] A. Danielyan, V. Katkovnik, and K. Egiazarian, "BM3D frames and variational image deblurring," *IEEE Trans. Image Process.*, vol. 21, no. 4, pp. 1715–1728, April 2012.
- [6] S. H. Chan, X. Wang, and O. A. Elgendy, "Plug-and-play ADMM for image restoration: Fixed-point convergence and applications," *IEEE Trans. Comp. Imag.*, vol. 3, no. 1, pp. 84–98, March 2017.
- [7] S. et al., "Plug-and-play priors for bright field electron tomography and sparse interpolation," *IEEE Trans. Comp. Imag.*, vol. 2, no. 4, pp. 408–423, December 2016.
- [8] S. Ono, "Primal-dual plug-and-play image restoration," *IEEE Signal Process. Lett.*, vol. 24, no. 8, pp. 1108–1112, 2017.
- [9] U. S. Kamilov, H. Mansour, and B. Wohlberg, "A plug-and-play priors approach for solving nonlinear imaging inverse problems," *IEEE Signal Process. Lett.*, vol. 24, no. 12, pp. 1872–1876, December 2017.
- [10] T. Meinhardt, M. Moeller, C. Hazirbas, and D. Cremers, "Learning proximal operators: Using denoising networks for regularizing inverse imaging problems," in *Proc. IEEE Int. Conf. Comp. Vis. (ICCV)*, 2017.
- [11] K. Zhang, W. Zuo, S. Gu, and L. Zhang, "Learning deep CNN denoiser prior for image restoration," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [12] G. T. Buzzard, S. H. Chan, S. Sreehari, and C. A. Bouman, "Plugand-play unplugged: Optimization free reconstruction using consensus equilibrium," SIAM J. Imaging Sci., vol. 11, no. 3, pp. 2001–2020, 2018.
- [13] Y. Sun, B. Wohlberg, and U. S. Kamilov, "An online plug-and-play algorithm for regularized image reconstruction," *IEEE Trans. Comput. Imaging*, 2019.
- [14] A. M. Teodoro, J. M. Bioucas-Dias, and M. Figueiredo, "A convergent image fusion algorithm using scene-adapted Gaussian-mixture-based denoising," *IEEE Trans. Image Process.*, vol. 28, no. 1, pp. 451–463, Jan. 2019.
- [15] 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. 36th Int. Conf. Machine Learning (ICML), 2019.
- [16] J. Tan, Y. Ma, and D. Baron, "Compressive imaging via approximate message passing with image denoising," *IEEE Trans. Signal Process.*, vol. 63, no. 8, pp. 2085–2092, Apr. 2015.
- [17] C. A. Metzler, A. Maleki, and R. G. Baraniuk, "BM3D-AMP: A new image recovery algorithm based on BM3D denoising," in *Proc. IEEE Int. Conf. Image Proc. (ICIP 2015)*, 2015, pp. 3116–3120.
- [18] —, "From denoising to compressed sensing," *IEEE Trans. Inf. Theory*, vol. 62, no. 9, pp. 5117–5144, September 2016.
- [19] C. A. Metzler, A. Maleki, and R. Baraniuk, "BM3D-PRGAMP: Compressive phase retrieval based on BM3D denoising," in *Proc. IEEE Int. Conf. Image Proc.*, 2016.
- [20] A. Fletcher, S. Rangan, S. Sarkar, and P. Schniter, "Plug-in estimation in high-dimensional linear inverse problems: A rigorous analysis," in *Proc.* Advances in Neural Information Processing Systems 32, 2018.
- [21] Y. Romano, M. Elad, and P. Milanfar, "The little engine that could: Regularization by denoising (RED)," SIAM J. Imaging Sci., vol. 10, no. 4, pp. 1804–1844, 2017.
- [22] S. A. Bigdeli, M. Jin, P. Favaro, and M. Zwicker, "Deep mean-shift priors for image restoration," in *Proc. Advances in Neural Information Processing Systems* 31, 2017.

- [23] E. T. Reehorst and P. Schniter, "Regularization by denoising: Clarifications and new interpretations," *IEEE Trans. Comput. Imag.*, vol. 5, no. 1, pp. 52–67, Mar. 2019.
- [24] C. A. Metzler, P. Schniter, A. Veeraraghavan, and R. G. Baraniuk, "prDeep: Robust phase retrieval with a flexible deep network," in *Proc.* 35th Int. Conf. Machine Learning (ICML), 2018.
- [25] G. Mataev, M. Elad, and P. Milanfar, "DeepRED: Deep image prior powered by RED," in *Proc. IEEE Int. Conf. Comp. Vis. Workshops* (ICCVW), 2019.
- [26] P. Tseng, "Convergence of a block coordinate descent method for nondifferentiable minimization," *J. Optimiz. Theory App.*, vol. 109, no. 3, pp. 475–494, June 2001.
- [27] Y. Nesterov, "Efficiency of coordinate descent methods on huge-scale optimization problems," SIAM J. Optim., vol. 22, no. 2, pp. 341–362, 2012.
- [28] A. Beck and L. Tetruashvili, "On the convergence of block coordinate descent type methods," SIAM J. Optim., vol. 23, no. 4, pp. 2037–2060, Oct. 2013.
- [29] S. J. Wright, "Coordinate descent algorithms," Math. Program., vol. 151, no. 1, pp. 3–34, Jun. 2015.
- [30] O. Fercoq and A. Gramfort, "Coordinate descent methods," Lecture notes Optimization for Data Science, École polytechnique, 2018.
- [31] Y. Sun, J. Liu, and U. S. Kamilov, "Block coordinate regularization by denoising," in *Proc. Advances in Neural Information Processing Systems* 33, Vancouver, BC, Canada, December 8-14, 2019.
- [32] L. I. Rudin, S. Osher, and E. Fatemi, "Nonlinear total variation based noise removal algorithms," *Physica D*, vol. 60, no. 1–4, pp. 259–268, November 1992.
- [33] R. Tibshirani, "Regression and selection via the lasso," *J. R. Stat. Soc. Series B (Methodological)*, vol. 58, no. 1, pp. 267–288, 1996.
- [34] E. J. Candès, J. Romberg, and T. Tao, "Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information," *IEEE Trans. Inf. Theory*, vol. 52, no. 2, pp. 489–509, February 2006.
- [35] D. L. Donoho, "Compressed sensing," *IEEE Trans. Inf. Theory*, vol. 52, no. 4, pp. 1289–1306, April 2006.
- [36] M. A. T. Figueiredo and R. D. Nowak, "An EM algorithm for wavelet-based image restoration," *IEEE Trans. Image Process.*, vol. 12, no. 8, pp. 906–916, August 2003.
- [37] I. Daubechies, M. Defrise, and C. D. Mol, "An iterative thresholding algorithm for linear inverse problems with a sparsity constraint," *Commun. Pure Appl. Math.*, vol. 57, no. 11, pp. 1413–1457, November 2004.
- [38] J. Bect, L. Blanc-Feraud, G. Aubert, and A. Chambolle, "A ℓ₁-unified variational framework for image restoration," in *Proc. ECCV*, Springer, Ed., vol. 3024, New York, 2004, pp. 1–13.
- [39] A. Beck and M. Teboulle, "A fast iterative shrinkage-thresholding algorithm for linear inverse problems," SIAM J. Imaging Sciences, vol. 2, no. 1, pp. 183–202, 2009.
- [40] J. Eckstein and D. P. Bertsekas, "On the Douglas-Rachford splitting method and the proximal point algorithm for maximal monotone operators," *Mathematical Programming*, vol. 55, pp. 293–318, 1992.
- [41] M. V. Afonso, J. M.Bioucas-Dias, and M. A. T. Figueiredo, "Fast image recovery using variable splitting and constrained optimization," *IEEE Trans. Image Process.*, vol. 19, no. 9, pp. 2345–2356, September 2010.
- [42] M. K. Ng, P. Weiss, and X. Yuan, "Solving constrained total-variation image restoration and reconstruction problems via alternating direction methods," SIAM J. Sci. Comput., vol. 32, no. 5, pp. 2710–2736, August 2010.
- [43] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, "Distributed optimization and statistical learning via the alternating direction method

- of multipliers," *Foundations and Trends in Machine Learning*, vol. 3, no. 1, pp. 1–122, 2011.
- [44] Z. Peng, T. Wu, Y. Xu, M. Yan, and W. Yin, "Coordinate-friendly structures, algorithms and applications," *Adv. Math. Sci. Appl.*, vol. 1, no. 1, pp. 57–119, Apr. 2016.
- [45] M. Elad and M. Aharon, "Image denoising via sparse and redundant representations over learned dictionaries," *IEEE Trans. Image Process.*, vol. 15, no. 12, pp. 3736–3745, December 2006.
- [46] A. Buades, B. Coll, and J. M. Morel, "Image denoising methods. A new nonlocal principle," SIAM Rev, vol. 52, no. 1, pp. 113–147, 2010.
- [47] D. Zoran and Y. Weiss, "From learning models of natural image patches to whole image restoration," in *Proc. IEEE Int. Conf. Comp. Vis. (ICCV)*, 2011.
- [48] Y. Nesterov, Introductory Lectures on Convex Optimization: A Basic Course. Kluwer Academic Publishers, 2004.
- [49] H. H. Bauschke and P. L. Combettes, *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*, 2nd ed. Springer, 2017.
- [50] E. K. Ryu and S. Boyd, "A primer on monotone operator methods," Appl. Comput. Math., vol. 15, no. 1, pp. 3–43, 2016.
- [51] Z. Peng, Y. Xu, M. Yan, and W. Yin, "ARock: An algorithmic framework for asynchronous parallel coordinate updates," *SIAM J. Sci. Comput.*, vol. 38, no. 5, pp. A2851–A2879, 2016.
- [52] Y. T. Chow, T. Wu, and W. Yin, "Cyclic coordinate-update algorithms for fixed-point problems: Analysis and applications," SIAM J. Sci. Comput., vol. 39, no. 4, pp. A1280–A1300, 2017.
- [53] T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida, "Spectral normalization for generative adversarial networks," in *International Conference* on Learning Representations (ICLR), 2018.
- [54] H. Sedghi, V. Gupta, and P. M. Long, "The singular values of convolutional layers," in *International Conference on Learning Representations (ICLR)*, 2019.
- [55] H. Gouk, E. Frank, B. Pfahringer, and M. Cree, "Regularisation of neural networks by enforcing Lipschitz continuity," 2018, arXiv:1804.04368.
- [56] J. J. Moreau, "Proximité et dualité dans un espace hilbertien," Bull. Soc. Math. France, vol. 93, pp. 273–299, 1965.
- [57] R. T. Rockafellar and R. Wets, Variational Analysis. Springer, 1998.
- [58] Y.-L. Yu, "Better approximation and faster algorithm using the proximal average," in *Proc. Advances in Neural Information Processing Systems* 26, 2013.
- [59] Zbontar et al., "fastMRI: An open dataset and benchmarks for accelerated MRI," 2018, arXiv:1811.08839. [Online]. Available: http://arxiv.org/abs/1811.08839
- [60] A. C. Kak and M. Slaney, Principles of Computerized Tomographic Imaging. IEEE, 1988.
- [61] F. Knoll, K. Brendies, T. Pock, and R. Stollberger, "Second order total generalized variation (TGV) for MRI," *Magn. Reson. Med.*, vol. 65, no. 2, pp. 480–491, February 2011.
- [62] A. Beck and M. Teboulle, "Fast gradient-based algorithm for constrained total variation image denoising and deblurring problems," *IEEE Trans. Image Process.*, vol. 18, no. 11, pp. 2419–2434, November 2009.
- [63] K. H. Jin, M. T. McCann, E. Froustey, and M. Unser, "Deep convolutional neural network for inverse problems in imaging," *IEEE Trans. Image Process.*, vol. 26, no. 9, pp. 4509–4522, Sep. 2017.
- [64] Y. S. Han, J. Yoo, and J. C. Ye, "Deep learning with domain adaptation for accelerated projection reconstruction MR," *Magn. Reson. Med.*, vol. 80, no. 3, pp. 1189–1205, Sep. 2017.
- [65] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional networks for biomedical image segmentation," in *Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, 2015.
- [66] E. Ndiaye, O. Fercoq, Alex, re Gramfort, and J. Salmon, "Gap safe screening rules for sparsity enforcing penalties," *Journal of Machine Learning Research*, vol. 18, no. 128, pp. 1–33, 2017.
- [67] S. Boyd and L. Vandenberghe, Convex Optimization. Cambridge Univ. Press, 2004.
- [68] F. Niu, B. Recht, C. Ré, and S. J. Wright, "Hogwild!: A lock-free approach to parallelizing stochastic gradient descent," in *Proc. Advances* in Neural Information Processing Systems 24, 2011.



Yu Sun (S'18) received the B.Eng. degree in electronics and information from Sichuan University, Chengdu, China, in 2015, and the M.S. degree in data analytics and statistics, in 2017, from Washington University in St. Louis, St. Louis, MO, USA, where he is currently working toward the Ph.D. degree with the Computational Imaging Group. His research interests include computational imaging, machine learning, deep learning, and optimization.



**Jiaming Liu** (S'19) received the B.Sc. degree in electronic and information engineering from University of Electronic Science and Technology of China, Chengdu, China, in 2017, and the M.S. degree in electrical and engineering, in 2018, from Washington University in St. Louis, St. Louis, MO, USA, where he is currently pursuing the Ph.D. degree with the Computational Imaging Group. His research interests include machine learning, image processing, and optimization.



Ulugbek S. Kamilov (S'11–M'15) is an Assistant Professor and Director of Computational Imaging Group at Washington University in St. Louis. From 2015 to 2017, he was a Research Scientist at Mitsubishi Electric Research Laboratories (MERL), Cambridge, MA, USA. He received the BSc and MSc degrees in Communication Systems, and the PhD degree in Electrical Engineering from EPFL, Switzerland, in 2008, 2011, and 2015, respectively. His main research area is computational imaging with an emphasis on the mathematical and computational

aspects of image reconstruction.

He is a recipient of the IEEE Signal Processing Society's 2017 Best Paper Award (with V. Goyal and S. Rangan). His PhD thesis was selected as a finalist for EPFL's Doctorate Award in 2016. He is an Associate Editor for IEEE Transactions on Computational Imaging Technical Committee of IEEE Signal Processing Society since 2016.