

1 **ACCELERATION OF PRIMAL-DUAL METHODS BY**
2 **PRECONDITIONING AND FIXED NUMBER OF INNER LOOPS***

3 YANLI LIU[†], YUNBEI XU[‡], AND WOTAO YIN[†]

4 **Abstract.** Primal-Dual Hybrid Gradient (PDHG) and Alternating Direction Method of Multipliers (ADMM) have been widely used due to their wide applicability and easy implementation.
5
6 However, they may suffer from slow tail convergence. In view of this, many techniques have been
7 proposed to accelerate them, including preconditioning and inexact solve of the subproblems. In this
8 paper, we integrate these two techniques to achieve a further acceleration. Specifically, we give a
9 criterion for choosing good preconditioners, and propose to solve one of the subproblems by only a
10 fixed number (usually very few) of inner loops of several common routines. Global convergence is
11 established for the proposed scheme. Since our method overcomes the previous restriction of choosing
12 only diagonal preconditioners, we obtain significant accelerations on several popular applications.

13 **Key words.** Primal-Dual Hybrid Gradient, Alternating Direction Method of Multipliers,
14 preconditioning, fixed number of inner loops, structured subproblem, suitable subproblem solver

15 **AMS subject classifications.** 49M29, 65K10, 65Y20, 90C25

16 **1. Introduction.** In this paper, we consider the following optimization problem:

17 (1.1)
$$\underset{x \in \mathbb{R}^n}{\text{minimize}} f(x) + g(Ax),$$

18

19 together with its dual problem:

20 (1.2)
$$\underset{z \in \mathbb{R}^m}{\text{minimize}} f^*(-A^T z) + g^*(z).$$

21

22 Here $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ and $g : \mathbb{R}^m \rightarrow \mathbb{R} \cup \{+\infty\}$ are closed proper convex, and
23 $A \in \mathbb{R}^{m \times n}$ is a matrix, f^* and g^* are the convex conjugates of f and g , respectively.

24 Many practical problems can be formulated in the form of (1.1) or (1.2), for
25 example, image restoration [39], magnetic resonance imaging [35], network optimization
26 [15], computer vision [30], and earth mover's distance [22].

27 Primal-Dual algorithms such as Primal-Dual Hybrid Gradient (PDHG) and Al-
28 ternating Direction Method of Multipliers (ADMM) can be applied to solve (1.1).
29 However, PDHG and ADMM suffer from slow (tail) convergence in practice. They may
30 take more than a few thousand iterations and still cannot reach four digits of accuracy.
31 In general, their performance is very sensitive to problem conditions. Therefore, efforts
32 have been made to accelerate them. In the next subsection, we review two common
33 acceleration techniques: preconditioning and inexact solve of subproblems.

34 **1.1. Background.** The convergence rate of PDHG depends on its step sizes,
35 which need an estimate of the operator norm of A . To accelerate PDHG and avoid
36 estimating the norm of A , diagonal preconditioning [29] was proposed and analyzed.
37 This technique improves the iteration complexity and adds only little computational
38 cost per iteration. However, non-diagonal preconditioners can further reduce the
39 iteration complexity significantly, but it remains open to apply such preconditioners

*Submitted to the editors XXX 2018.

Funding: This work is supported in part by NSF grant DMS-1720237 and ONR grant N000141712162.

[†]Department of Mathematics, University of California, Los Angeles, CA (yanli@math.ucla.edu, wotaoyin@math.ucla.edu).

[‡]Graduate School of Business, Columbia University, New York, NY (yunbei.xu@gsb.columbia.edu).

40 while still maintaining the computation cost per iteration. As another acceleration
 41 technique, *inexact PDHG* allows the PDHG subproblems to be solved approximately.
 42 To ensure convergence, [31] uses three types of errors to control the solution errors of
 43 the subproblems; all of them need to be summable over the iterations. Therefore, [31]
 44 requires the subproblems to be solved with increasing accuracies.

45 Unlike PDHG, a subproblem of ADMM minimizes the sum of $f(x)$ and a squared
 46 term involving Ax . In general, it may not have a closed form solution. Several versions
 47 of inexact ADMM have been studied. An *absolute error criterion* is introduced in
 48 [11], where the subproblem errors are controlled by sequences of error tolerances that
 49 are summable. To simplify the choice of the sequences, the *relative error criterion*
 50 was adopted in several works, where the subproblem errors are controlled by a single
 51 parameter and quantities generated naturally by the algorithm. In [26], the parameters
 52 need to be square summable. In [21], the parameters are constants but both objectives
 53 are required to be Lipschitz differentiable. In [12, 13], two possible outcomes of the
 54 algorithm are described: (i) infinite outer loops and finite inner loops, and (ii) finite
 55 out loops and the last inner loop is infinite, both of them guaranteeing convergence to
 56 a solution. On the other hand, it is unclear how to distinguish them, and since there
 57 is no bound on the number of inner loops in case (i), one may recognize it as case (ii)
 58 and stop the algorithm before it converges.

59 Certain types of preconditioning have been applied to accelerate ADMM. In [17],
 60 diagonal preconditioning is used with ADMM. After that, non-diagonal preconditioning
 61 is also analyzed [5, 6], where appropriate preconditioners are given for specific
 62 applications, and competitive numerical performances are observed. This work inverts
 63 one of the preconditioners (not needed in our method). Recently, preconditioning for
 64 strongly convex problems has also been discussed [18].

65 **1.2. Contributions.** The contributions of this paper are three-fold.

66 First, we provide a criterion for choosing preconditioners based on an ergodic
 67 convergence result. We also show that ADMM corresponds to a special choice of
 68 preconditioners in Preconditioned PDHG (PPDHG).

69 Second, we show that PPDHG converges when its subproblem is solved inexactly
 70 to a specified relative-error condition. Remarkably, this condition does not need to
 71 be checked since it is naturally satisfied when one applies a *fixed number* of inner
 72 loops using any of several common subproblem solvers including proximal gradient
 73 descent, FISTA with restart, proximal block coordinate descent, as well as even faster
 74 block-coordinate-gradient-descent (BCGD) methods (e.g., [24, 1, 19]).

75 Third, the diverse choice of subproblem solution methods, especially the BCGD
 76 methods, lets us deal with the difficult subproblems that arise when we apply precon-
 77 ditioners. With appropriate preconditioners and subproblem solvers, both PDHG and
 78 ADMM can be accelerated and their total running time significantly reduced. The
 79 efficiency of our algorithm is demonstrated by numerical experiments.

80 It is worth mentioning that our fixed number of inner loops is different from the
 81 “finite inner loops” claimed in [5, 6]. In their settings, one subproblem is essentially
 82 applying a preconditioned proximal operator to a convex quadratic function at the
 83 current iterate, which has a closed form solution. The same operator is applied n times
 84 starting at the current iterate for any $n \geq 1$, and convergence can still be established,
 85 they call these n operations finite inner loops.

86 **1.3. Organization.** The rest of this paper is organized as follows: Section 2
 87 establishes notation and reviews some basic results. In Section 3, we first provide a
 88 criterion for choosing preconditioners of PDHG, then introduce the bounded relative

89 error condition and show that it is satisfied by a fixed number of inner loops. We es-
 90 tablish convergence of the inexact preconditioned PDHG. Section 4 provides numerical
 91 examples on several popular applications. Finally, Section 5 concludes the paper.

92 **2. Preliminaries.** In this section, we review some basic concepts, introduce our
 93 notation, and state some known results. For the sake of brevity, we omit proofs and
 94 direct references. We refer the reader to the textbook [3].

95 We use $\|\cdot\|$ for the ℓ_2 -norm, and $\langle \cdot, \cdot \rangle$ for the usual dot product. $M \succ 0$ denotes
 96 a symmetric positive definite matrix M , and $M \succeq 0$ denotes a symmetric positive
 97 semidefinite matrix M .

98 $\lambda_{\min}(M)$ and $\lambda_{\max}(M)$ stand for the smallest eigenvalue and the largest eigenvalues
 99 of M , respectively. $\kappa(M) = \frac{\lambda_{\max}(M)}{\lambda_{\min}(M)}$ is the condition number of M . For $M \succeq 0$, let
 100 $\|\cdot\|_M$ and $\langle \cdot, \cdot \rangle_M$ denote the (semi-)norm and inner product induced by M , respectively.

For a proper closed convex function $\phi : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$, its subdifferential at
 $x \in \mathbf{dom} f$ is defined by

$$\partial\phi(x) = \{v \in \mathbb{R}^n \mid \phi(z) \geq \phi(x) + \langle v, z - x \rangle \forall z \in \mathbb{R}^n\},$$

and its convex conjugate is

$$\phi^*(y) = \sup_{x \in \mathbb{R}^n} \{\langle y, x \rangle - \phi(x)\},$$

101 we have $y \in \partial\phi(x)$ if and only if $x \in \partial\phi^*(y)$.

102 For any symmetric $M \succ 0$, we define the extended proximal operator of ϕ to be

$$103 \quad (2.1) \quad \text{Prox}_{\phi}^M(x) := \arg \min_{y \in \mathbb{R}^n} \{\phi(y) + \frac{1}{2} \|y - x\|_M^2\},$$

105 When $M = \gamma^{-1}I$ where $\gamma > 0$, it reduces to the classic proximal operator.

106 For the extended proximal operator (2.1) we also have the following generalization
 107 of Moreau's Identity:

108 **LEMMA 2.1** ([10], Theorem 3.1(ii)). *For any proper closed convex function ϕ and*
 109 *$M \succ 0$, we have*

$$110 \quad (2.2) \quad x = \text{Prox}_{\phi}^M(x) + M^{-1} \text{Prox}_{\phi^*}^{M^{-1}}(Mx).$$

111 A proper closed function is said to be a Kurdyka-Łojasiewicz (KL) function, if
 112 for each $x_0 \in \mathbf{dom} f$, there exist $\eta \in (0, \infty]$, a neighborhood U of x_0 and a continuous
 113 concave function $\varphi : [0, \eta) \rightarrow \mathbb{R}_+$ such that:

- 114 1. $\varphi(0) = 0$,
- 115 2. φ is C^1 on $(0, \eta)$,
- 116 3. for all $s \in (0, \eta)$, $\varphi'(s) > 0$,
- 117 4. for all $x \in U \cap \{x \mid f(x_0) < f(x) < f(x_0) + \eta\}$, the KL inequality holds:

$$118 \quad \varphi'(f(x) - f(x_0)) \text{dist}(0, \partial f(x)) \geq 1.$$

119 **3. Acceleration of PDHG.** Throughout this section, the following regularity
 120 assumption is assumed:

121 **ASSUMPTION 1.**

- 122 1. $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$, $g : \mathbb{R}^m \rightarrow \mathbb{R} \cup \{+\infty\}$ are proper closed convex.

123 2. A primal-dual solution pair (x^*, z^*) of (1.1) and (1.2) exists, i.e.,

$$124 \quad \mathbf{0} \in \partial f(x^*) + A^T z^*, \quad \mathbf{0} \in \partial g(Ax^*) - z^*.$$

125 The problem (1.1) also has the following convex-concave saddle-point formulation:

$$126 \quad (3.1) \quad \min_{x \in \mathbb{R}^n} \max_{z \in \mathbb{R}^m} \varphi(x, z) := f(x) + \langle Ax, z \rangle - g^*(z)$$

127
128 A primal-dual solution pair (x^*, z^*) is a solution of (3.1) and vice versa.

129 **3.1. Preconditioned PDHG.** The method of Primal-Dual Hybrid Gradient
130 (PDHG) [39, 7] for solving (1.1) uses the iteration

$$131 \quad (3.2) \quad \begin{aligned} x^{k+1} &= \text{Prox}_{\tau f}(x^k - \tau A^T z^k), \\ z^{k+1} &= \text{Prox}_{\sigma g^*}(z^k + \sigma A(2x^{k+1} - x^k)). \end{aligned}$$

132
133 Convergence of (3.2) to a primal-dual solution pair of (1.1) is established when
134 $\frac{1}{\tau\sigma} \geq \|A\|^2$ [7]. In order to achieve faster convergence by exploiting the structure of
135 subproblems, we can apply preconditioners $M_1, M_2 \succ 0$ (their choices are discussed
136 below) to obtain Preconditioned PDHG (PPDHG):

$$137 \quad (3.3) \quad \begin{aligned} x^{k+1} &= \text{Prox}_f^{M_1}(x^k - M_1^{-1} A^T z^k), \\ z^{k+1} &= \text{Prox}_{g^*}^{M_2}(z^k + M_2^{-1} A(2x^{k+1} - x^k)), \end{aligned}$$

138
139 where the extended proximal operators $\text{Prox}_f^{M_1}$ and $\text{Prox}_{g^*}^{M_2}$ are defined in (2.1).

140 Note that there is no need to form M_1^{-1} and M_2^{-1} since (3.3) is equivalent to

$$141 \quad (3.4) \quad \begin{aligned} x^{k+1} &= \arg \min_{x \in \mathbb{R}^n} \{f(x) + \langle x - x^k, A^T z^k \rangle + \frac{1}{2} \|x - x^k\|_{M_1}^2\}, \\ z^{k+1} &= \arg \min_{z \in \mathbb{R}^m} \{g^*(z) - \langle z - z^k, A(2x^{k+1} - x^k) \rangle + \frac{1}{2} \|z - z^k\|_{M_2}^2\}. \end{aligned}$$

142
143 **3.2. Choice of preconditioners by an ergodic convergence result.** The
144 convergence of PPDHG is not new. In fact, PPDHG is a special case of a general primal-
145 dual algorithm considered in [8]. In this section, we discuss how to select appropriate
146 preconditioners M_1 and M_2 based on an ergodic convergence result from [8]. In
147 particular, we show that ADMM corresponds to the choice $M_1 = \frac{1}{\tau} I_{n \times n}$, $M_2 = \tau A A^T$,
148 which has faster convergence than PDHG in terms of outer iterations.

149 Let us start with the following lemma which characterizes primal-dual solution
150 pairs of (1.1) and (1.2).

151 **LEMMA 3.1.** *Under Assumption 1, (X, Z) is a primal-dual solution pair of (1.1)*
152 *if and only if $\varphi(X, z) - \varphi(x, Z) \leq 0$ for any $(x, z) \in \mathbb{R}^{n+m}$.*

153 *Proof.* If (X, Z) is a primal-dual solution pair of (1.1), then

$$154 \quad -A^T Z \in \partial f(X), \quad AX \in \partial g^*(Z).$$

155 As a result, for any $(x, z) \in \mathbb{R}^{n+m}$ we have

$$156 \quad f(x) \geq f(X) + \langle -A^T Z, x - X \rangle, \quad g^*(z) \geq g^*(Z) + \langle AX, z - Z \rangle,$$

157 adding them together gives $\varphi(X, z) - \varphi(x, Z) \leq 0$.

158 On the other hand, if $\varphi(X, z) - \varphi(x, Z) \leq 0$ for any $(x, z) \in \mathbb{R}^{n+m}$, then

159 $\langle AX, z \rangle + f(X) - g^*(z) - \langle Ax, Z \rangle - f(x) + g^*(Z) \leq 0$ for any $(x, z) \in \mathbb{R}^{n+m}$.

160 Taking $x = X$ yields $\langle AX, z - Z \rangle - g^*(z) + g^*(Z) \leq 0$, so $AX \in \partial g^*(Z)$; Similarly,
 161 taking $z = Z$ gives $\langle AX - Ax, Z \rangle + f(X) - f(x) \leq 0$, so $-A^T Z \in \partial f(X)$. As a result,
 162 (X, Z) is a primal-dual solution pair of (1.1). \square

163 On the other hand, we have the following ergodic convergence result, which is
 164 adapted from Theorem 1 of [8].

165 **THEOREM 3.2.** *Let $(x^k, z^k), n = 0, 1, \dots, N$ be a sequence generated by PPDHG*
 166 *(3.3). Under Assumption 1, if in addition*

167 (3.5)
$$\tilde{M} := \begin{pmatrix} M_1 & -A^T \\ -A & M_2 \end{pmatrix} \succeq 0,$$

168 then, for any $x \in \mathbb{R}^n$ and $z \in \mathbb{R}^m$, it holds that

169 (3.6)
$$\varphi(X^N, z) - \varphi(x, Z^N) \leq \frac{1}{2N} (x - x^0, z - z^0) \begin{pmatrix} M_1 & -A^T \\ -A & M_2 \end{pmatrix} \begin{pmatrix} x - x^0 \\ z - z^0 \end{pmatrix},$$

170 where $X^N = \frac{1}{N} \sum_{i=1}^N x_i$ and $Z^N = \frac{1}{N} \sum_{i=1}^N z_i$.

171 *Proof.* This follows from Theorem 1 and Remark 3 of [8], by setting $L_f = 0$,
 172 $\frac{1}{\tau} D_x(x, x_0) = \frac{1}{2} \|x - x^0\|_{M_1}^2$, $\frac{1}{\sigma} D_z(z, z_0) = \frac{1}{2} \|z - z^0\|_{M_2}^2$, and $K = A$. \square

173 In view of Lemma 3.1 and Theorem 3.2, in order to accelerate the ergodic conver-
 174 gence, the preconditioners M_1 and M_2 should be chosen such that (3.5) is satisfied,
 175 and the right hand side of (3.6) should be small. In view of this, we can obtain some
 176 useful criteria for choosing M_1 and M_2 .

177 First, by Schur complement lemma, the condition (3.5) is equivalent to $M_2 \succeq$
 178 $AM_1^{-1}A^T$. Hence, for a given M_1 , the optimal M_2 is $AM_1^{-1}A^T$.

179 Second, PDHG (3.2) corresponds to $M_1 = \frac{1}{\tau} I_{n \times n}$, $M_2 = \frac{1}{\sigma} I_{m \times m}$ with $\frac{1}{\sigma} \geq \|A\|^2$.
 180 On the other hand, one can show that ADMM applied to (1.1) corresponds to $M_1 =$
 181 $\frac{1}{\tau} I_{n \times n}$, $M_2 = \tau AA^T$ (see Appendix A, this is also implicitly shown in [7, Sec. 4.3]),
 182 where M_2 is optimal for M_1 since $AM_1^{-1}A^T = \tau AA^T = M_2$. In this regard, ADMM
 183 corresponds to a better choice of preconditioners than PDHG, which explains why
 184 ADMM uses fewer outer iterations than PDHG in practice. This is also verified in our
 185 numerical experiments in Section 4.

186 Finally, there might be better choices of preconditioners than that of ADMM.
 187 This can bring us faster algorithms and is left as future work.

188 **3.3. PPDHG with fixed finite inner loops.** Solving the subproblems in (3.3)
 189 exactly or nearly so is wasteful. Choosing the number of inner loops based on a
 190 condition requires checking the condition. It is convenient if we can simply fix the
 191 number of inner loops.

192 In this subsection, we describe the ‘‘bounded relative error’’ of the z -subproblem
 193 in (3.3) and then show that this can be satisfied by running a fixed number of inner
 194 loops, uniformly for every outer loop.

195 **DEFINITION 3.3.** *Given x^k, x^{k+1} and z^k , the z -subproblem in PPDHG (3.3) is*
 196 *said to be solved with bounded relative error if there is a constant $c > 0$ such that*

197 (3.7)
$$\mathbf{0} \in \partial g^*(z^{k+1}) + M_2(z^{k+1} - z^k - M_2^{-1}A(2x^{k+1} - x^k)) + \varepsilon^{k+1},$$

198 (3.8)
$$\|\varepsilon^{k+1}\| \leq c \|z^{k+1} - z^k\|.$$

Remarkably, this condition does not need to be checked at any iteration. For a given $c > 0$, it can be satisfied by a fixed number of inner loops using proximal gradient descent (see Theorem 3.4). One can also use faster solvers for the z -subproblem, e.g., FISTA with restart [27], and solvers that suit the subproblem structure, e.g., cyclic proximal BCD (see Theorem 3.6). Although the error in solving z -subproblems appears to be neither summable nor square summable at first glance, convergence can still be established. We summarize our algorithm in Algorithm 3.1.

Algorithm 3.1 Inexact preconditioned PDHG

Input: $f : \mathbb{R}^n \rightarrow \mathbb{R}, g : \mathbb{R}^m \rightarrow \mathbb{R}, A \in \mathbb{R}^{m \times n}$, preconditioners M_1 and M_2 , initial (x_0, z_0) , subproblem solver S for the z -subproblem in (3.3), fixed inner iterations p , max outer iterations K .

Output: (x^K, z^K)

```

1: for  $k \leftarrow 0, 1, \dots, K - 1$  do
2:    $x^{k+1} = \text{Prox}_f^{M_1}(x^k - M_1^{-1}A^T z^k)$ ;
3:    $z_0^{k+1} = z^k$ ;
4:   for  $i \leftarrow 0, 1, \dots, p - 1$  do
5:      $z_{i+1}^{k+1} = S(z_i^{k+1}, x^{k+1}, x^k)$ ;
6:   end for
7:    $z^{k+1} = z_p^{k+1}$ ;            $\triangleright$  approximate  $\text{Prox}_{g^*}^{M_2}(z^k + M_2^{-1}A(2x^{k+1} - x^k))$ 
8: end for

```

THEOREM 3.4. Under Assumption 1, if $p \geq 1$ iterations of proximal gradient descent with stepsize $\gamma \in (0, \frac{2\lambda_{\min}(M_2)}{\lambda_{\max}^2(M_2)})$ are applied to solve the z -subproblem in (3.3), and is initialized with the last iterate z^k , then the subproblem is solved with bounded relative error with the following constant for (3.8)

$$(3.9) \quad c = c(p) = \frac{\frac{1}{\gamma} + \lambda_{\max}(M_2)}{1 - \tau^p} (\tau^p + \tau^{p-1}),$$

where $\tau = \sqrt{1 - \gamma(2\lambda_{\min}(M_2) - \gamma\lambda_{\max}^2(M_2))} < 1$.

Proof. The z -subproblem in (3.4) is of the form

$$(3.10) \quad \min_{z \in \mathbb{R}^m} h_1(z) + h_2(z),$$

where

$$h_1(z) = g^*(z),$$

$$h_2(z) = \frac{1}{2} \|z - z^k - M_2^{-1}A(2x^{k+1} - x^k)\|_{M_2}^2.$$

In Algorithm 3.1, an inexact z^{k+1} is given by

$$z_0^{k+1} = z^k,$$

$$z_{i+1}^{k+1} = \text{Prox}_{\gamma h_1}(z_i^{k+1} - \gamma \nabla h_2(z_i^{k+1})), \quad i = 0, 1, \dots, p - 1,$$

$$z^{k+1} = z_p^{k+1}.$$

The optimality condition of the last iteration above reads

$$\mathbf{0} \in \partial h_1(z_p^{k+1}) + \nabla h_2(z_{p-1}^{k+1}) + \frac{1}{\gamma}(z_p^{k+1} - z_{p-1}^{k+1}),$$

228 compare this with (3.7) and use $z_p^{k+1} = z^{k+1}$, we have

$$229 \quad \varepsilon^{k+1} = \frac{1}{\gamma}(z_p^{k+1} - z_{p-1}^{k+1}) + \nabla h_2(z_{p-1}^{k+1}) - \nabla h_2(z_p^{k+1}),$$

230 we need to show that ε^{k+1} satisfies (3.8).

231 Let z_\star^{k+1} be the solution of (3.10), $\alpha = \lambda_{\min}(M_2)$, and $\beta = \lambda_{\max}(M_2)$, then $h_1(z)$
232 is convex and $h_2(z)$ is α -strongly convex and β -Lipschitz differentiable. As a result, [3,
233 Prop. 26.16(ii)] gives

$$234 \quad \|z_i^{k+1} - z_\star^{k+1}\| \leq \tau^i \|z_0^{k+1} - z_\star^{k+1}\|, \quad \forall i = 0, 1, \dots, p,$$

235 where $\tau = \sqrt{1 - \gamma(2\alpha - \gamma\beta^2)}$.

236 Let $a_i = \|z_i^{k+1} - z_\star^{k+1}\|$, then $a_i \leq \tau^i a_0$. Therefore,

$$237 \quad (3.11) \quad \|\varepsilon^{k+1}\| \leq \left(\frac{1}{\gamma} + \beta\right) \|z_p^{k+1} - z_{p-1}^{k+1}\|$$

$$238 \quad (3.12) \quad \leq \left(\frac{1}{\gamma} + \beta\right) (a_p + a_{p-1})$$

$$239 \quad (3.13) \quad \leq \left(\frac{1}{\gamma} + \beta\right) (\tau^p + \tau^{p-1}) a_0.$$

240
241 On the other hand, we have

$$242 \quad \|z^{k+1} - z^k\| \geq a_0 - a_p$$

$$243 \quad (3.14) \quad \geq (1 - \tau^p) a_0.$$

245 Combining (3.11) and (3.14) gives

$$246 \quad \|\varepsilon^{k+1}\| \leq c \|z^{k+1} - z^k\|,$$

247 where c is given in (3.9). □

248 *Remark 3.5.* Similarly to proximal gradient descent, one can show that finite
249 iterations of FISTA with restart also satisfies the bounded relative error condition in
250 Def. 3.3. The proof is omitted due to space limitation.

251 **THEOREM 3.6.** *Let Assumption 1 holds and g be block separable, i.e.,*
252 $z = (z_1, z_2, \dots, z_l)$ *and $g(z) = \sum_{i=1}^l g_i(z_i)$. Let the stepsize of cyclic proximal BCD be*
253 γ , *which is small enough such that*

$$254 \quad 0 < \gamma \leq \min \left\{ \frac{2\lambda_{\min}(M_2)}{\lambda_{\max}^2(M_2)}, \frac{1 - \sqrt{1 - \gamma(2\lambda_{\min}(M_2) - \gamma\lambda_{\max}^2(M_2))}}{4\sqrt{2}\gamma\lambda_{\max}(M_2)}, \right.$$

$$255 \quad \left. \frac{1}{4l\lambda_{\max}(M_2)}, \frac{2l\lambda_{\max}(M_2)}{17l\lambda_{\max}(M_2) + 2\left(\frac{1 - \sqrt{1 - \gamma(2\lambda_{\min}(M_2) - \gamma\lambda_{\max}^2(M_2))}}{\gamma}\right)^2} \right\}.$$

257 *Then, if $p \geq 1$ epochs of cyclic proximal BCD are applied to solve the z -subproblem*
258 *in (3.3), and is initialized with the last iterate z^k , then the subproblem is solved with*
259 *bounded relative error with*

$$260 \quad (3.15) \quad c = c(p) = \frac{(l\lambda_{\max}(M_2) + \frac{1}{\gamma})(\rho^p + \rho^{p-1})}{1 - \rho^p},$$

261
262 where $\rho = 1 - \frac{(1 - \sqrt{1 - \gamma(2\lambda_{\min}(M_2) - \gamma\lambda_{\max}^2(M_2))})^2}{2\gamma} < 1$.

263 *Proof.* See Appendix B. □

264 **3.4. Convergence of inexact PPDHG.** In this subsection, we proceed to
 265 establish the convergence of Algorithm 3.1. First, we transform Algorithm 3.1 into
 266 an ADMM-like algorithm in Proposition 3.8 below. Then, in Theorems 3.11 and 3.12
 267 below, we prove the convergence of Algorithm 3.1 by using a generalized augmented
 268 Lagrangian of this ADMM-like algorithm.

269 First, let us show that PPDHG (3.3) is equivalent to an ADMM-like algorithm
 270 applied on the dual problem (1.2), this is similar to the equivalence of PDHG (3.2)
 271 and Linearized ADMM applied to the dual problem (1.2) shown in [14]. Therefore, we
 272 call this ADMM-like algorithm Preconditioned Linearized ADMM (PLADMM):

$$\begin{aligned} z^{k+1} &= \text{Prox}_{g^*}^{M_2}(z^k + M_2^{-1}AM_1^{-1}(-A^T z^k - y^k + u^k)), \\ (3.16) \quad y^{k+1} &= \text{Prox}_{f^*}^{M_1^{-1}}(u^k - A^T z^{k+1}), \\ u^{k+1} &= u^k - A^T z^{k+1} - y^{k+1}. \end{aligned}$$

276 *Remark 3.7.* When $M_1 = \frac{1}{\tau}I$, $M_2 = \lambda I$, PLADMM (3.16) is Linearized ADMM,
 277 or Split Inexact Uzawa [38].

278 Furthermore, the Inexact PPDHG in Algorithm 3.1 is equivalent to Inexact
 279 PLADMM, which is summarized in Algorithm 3.2.

280 Let us also define the following generalized augmented Lagrangian for PLADMM:

$$(3.17) \quad L(z, y, u) = g^*(z) + f^*(y) + \langle -A^T z - y, M_1^{-1}u \rangle + \frac{1}{2}\|A^T z + y\|_{M_1^{-1}}^2.$$

283 Inspired by the framework of [36], this generalized augmented Lagrangian will serve
 284 as a Lyapunov function to establish convergence of Algorithm 3.2 and 3.1.

Algorithm 3.2 Inexact preconditioned linearized ADMM

Input: $f^* : \mathbb{R}^n \rightarrow \mathbb{R}$, $g^* : \mathbb{R}^m \rightarrow \mathbb{R}$, $A \in \mathbb{R}^{m \times n}$, preconditioners M_1 and M_2 ,
 initial vector (z_0, y_0, u_0) , subproblem solver S for the z -subproblem in (3.16), number
 of inner loops p , number of outer iterations K .

Output: (z^K, y^K, u^K)

```

1: for  $k \leftarrow 0, 1, \dots, K - 1$  do
2:    $z_0^{k+1} = z^k$ ;
3:   for  $i \leftarrow 0, 1, \dots, p - 1$  do
4:      $z_{i+1}^{k+1} = S(z_i^{k+1}, y^k, u^k)$ ;
5:   end for
6:    $z^{k+1} = z_p^{k+1}$ ;  $\triangleright$  approximate  $\text{Prox}_{g^*}^{M_2}(z^k + M_2^{-1}AM_1^{-1}(-A^T z^k - y^k + u^k))$ .
7:    $y^{k+1} = \text{Prox}_{f^*}^{M_1^{-1}}(u^k - A^T z^{k+1})$ ;
8:    $u^{k+1} = u^k - A^T z^{k+1} - y^{k+1}$ ;
9: end for
```

285 **PROPOSITION 3.8.** *Under Assumption 1 and the transforms $u^k = M_1 x^k$, $y^{k+1} =$
 286 $u^k - A^T z^k - u^{k+1}$, PPDHG (3.3) is equivalent to PLADMM (3.16), and the Inexact
 287 PPDHG in Algorithm 3.1 is equivalent to the Inexact PLADMM in Algorithm 3.2.*

288 *Proof.* First, let us transform PPDHG in (3.3) to PLADMM (3.16).
 289 Set $u^k = M_1 x^k$, $y^{k+1} = u^k - A^T z^k - u^{k+1}$, then (2.2) and (3.3) gives

$$290 \quad y^{k+1} = M_1 x^k - A^T z^k - M_1 x^{k+1} = \text{Prox}_{f^*}^{M_1^{-1}}(u^k - A^T z^k),$$

292 and we also have

$$293 \quad u^{k+1} = u^k - A^T z^k - y^{k+1},$$

$$294 \quad z^{k+1} = \text{Prox}_{g^*}^{M_2}(z^k + M_2^{-1} A M_1^{-1} (-A^T z^k - y^{k+1} + u^{k+1})).$$

296 If z -update is performed first, then we arrive at PLADMM (3.16).

297 Notice that for the inexact PPDHG in Algorithm 3.1, we are solving the
298 z -subproblem of PPDHG (3.3) with bounded relative error as in Definition 3.3,
299 therefore we are essentially doing the same to the z -subproblem of PLADMM (3.16),
300 which gives Algorithm 3.2. \square

301 In order to establish convergence of Algorithm 3.1, we also need the following
302 Assumption 2, in addition to Assumption 1.

303 ASSUMPTION 2.

- 304 1. $f(x)$ is μ_f -strongly convex.
- 305 2. $g^*(z) + f^*(-A^T z)$ is coercive, i.e.,

$$306 \quad \lim_{\|z\| \rightarrow \infty} g^*(z) + f^*(-A^T z) = \infty.$$

- 307 3. $g^*(z)$ is a KL function.

308 THEOREM 3.9. Under Assumptions 1 and 2. Choose the preconditioners M_1, M_2
309 and the number of inner loops p in Algorithm 3.1 such that

$$310 \quad C_1 = \frac{1}{2} M_1^{-1} - \frac{\|M_1\|}{\mu_f^2} I \succ 0,$$

$$311 \quad C_2 = M_2 - \frac{1}{2} A M_1^{-1} A^T - c(p) I \succ 0,$$

313 where $c(p)$ is related to the z -subproblem solver S and M_2 (see, e.g., (3.9) and (3.15)).
314 Define $L^k := L(z^k, y^k, u^k)$, then the inexact PLADMM in Algorithm 3.2 satisfies the
315 following sufficient descent and lower boundedness properties:

$$316 \quad (3.18) \quad L^k - L^{k+1} \geq \|y^k - y^{k+1}\|_{C_1}^2 + \|z^k - z^{k+1}\|_{C_2}^2,$$

$$317 \quad (3.19) \quad L^k \geq g^*(z^*) + f^*(-A^T z^*) > -\infty.$$

319 *Proof.* Since the z -subproblem of Algorithm 3.2 is solved with bounded relative
320 error in Def. 3.3, we have

$$321 \quad (3.20) \quad \mathbf{0} \in \partial g^*(z^{k+1}) + M_2(z^{k+1} - z^k - M_2^{-1} A M_1^{-1} (-A^T z^k - y^k + u^k)) + \varepsilon^{k+1},$$

323 where ε^{k+1} satisfies (3.8):

$$324 \quad (3.21) \quad \|\varepsilon^{k+1}\| \leq c(p) \|z^{k+1} - z^k\|.$$

326 The y and u updates gives

$$327 \quad (3.22) \quad \mathbf{0} = \nabla f^*(y^{k+1}) + M_1^{-1}(y^{k+1} - u^k + A^T z^{k+1}) = \nabla f^*(y^{k+1}) - M_1^{-1} u^{k+1},$$

$$328 \quad (3.23) \quad u^{k+1} = u^k - A^T z^{k+1} - y^{k+1}.$$

330 In order to show (3.18), let us write

$$\begin{aligned}
331 \quad & g^*(z^k) \geq g^*(z^{k+1}) \\
332 \quad & \quad + \langle M_2(z^k - z^{k+1}) + AM_1^{-1}(-A^T z^k - y^k + u^k) - \varepsilon^{k+1}, z^k - z^{k+1} \rangle, \\
333 \quad & f^*(y^k) \geq f^*(y^{k+1}) + \langle M_1^{-1}u^{k+1}, y^k - y^{k+1} \rangle, \\
334
\end{aligned}$$

335 Assembling these inequalities with (3.21) gives us

$$\begin{aligned}
336 \quad & L^k - L^{k+1} \geq \|z^k - z^{k+1}\|_{M_2 - c(p)I}^2 \\
337 \quad & \quad + \langle AM_1^{-1}(-A^T z^k - y^k + u^k), z^k - z^{k+1} \rangle + \langle M_1^{-1}u^{k+1}, y^k - y^{k+1} \rangle \\
338 \quad & \quad + \langle -A^T z^k - y^k, M_1^{-1}u^k \rangle - \langle A^T z^{k+1} - y^{k+1}, M_1^{-1}(u^k - A^T z^{k+1} - y^{k+1}) \rangle \\
339 \quad & \quad + \frac{1}{2} \|A^T z^k + y^k\|_{M_1^{-1}}^2 - \frac{1}{2} \|A^T z^{k+1} + y^{k+1}\|_{M_1^{-1}}^2 \\
340 \quad & = \|z^k - z^{k+1}\|_{M_2 - c(p)I}^2 \\
341 \quad (A) \quad & \quad + \langle AM_1^{-1}(-A^T z^k - y^k), z^k - z^{k+1} \rangle + \langle M_1^{-1}u^{k+1}, y^k - y^{k+1} \rangle \\
342 \quad (B) \quad & \quad + \langle -y^k, M_1^{-1}u^k \rangle - \langle -y^{k+1}, M_1^{-1}u^k \rangle \\
343 \quad & \quad + \frac{1}{2} \|A^T z^k + y^k\|_{M_1^{-1}}^2 - \frac{3}{2} \|A^T z^{k+1} + y^{k+1}\|_{M_1^{-1}}^2, \\
344
\end{aligned}$$

345 where the terms in (A) and (B) simplify to

$$346 \quad (3.24) \quad \langle AM_1^{-1}(-A^T z^k - y^k), z^k - z^{k+1} \rangle + \langle M_1^{-1}(-A^T z^{k+1} - y^{k+1}), y^k - y^{k+1} \rangle.$$

348 Now we will use the following cosine rule on the two inner products above:

$$349 \quad \langle a - b, a - c \rangle_{M_1^{-1}} = \frac{1}{2} \|a - b\|_{M_1^{-1}}^2 + \frac{1}{2} \|a - c\|_{M_1^{-1}}^2 - \frac{1}{2} \|b - c\|_{M_1^{-1}}^2.$$

350 Set $a = A^T z^k$, $c = A^T z^{k+1}$, and $b = -y^k$ to obtain

$$\begin{aligned}
351 \quad & \langle AM_1^{-1}(-A^T z^k - y^k), z^k - z^{k+1} \rangle = -\frac{1}{2} \|A^T z^k + y^k\|_{M_1^{-1}}^2 - \frac{1}{2} \|A^T z^k - A^T z^{k+1}\|_{M_1^{-1}}^2 \\
352 \quad (3.25) \quad & \quad + \frac{1}{2} \|y^k + A^T z^{k+1}\|_{M_1^{-1}}^2. \\
353
\end{aligned}$$

354 Set $a = y^{k+1}$, $c = y^k$, and $b = -A^T z^{k+1}$ to obtain

$$\begin{aligned}
355 \quad & \langle M_1^{-1}(-A^T z^{k+1} - y^{k+1}), y^k - y^{k+1} \rangle = \frac{1}{2} \|A^T z^{k+1} + y^{k+1}\|_{M_1^{-1}}^2 + \frac{1}{2} \|y^k - y^{k+1}\|_{M_1^{-1}}^2 \\
356 \quad (3.26) \quad & \quad - \frac{1}{2} \|A^T z^{k+1} + y^k\|_{M_1^{-1}}^2. \\
357
\end{aligned}$$

358 Combining (3.24), (3.25), and (3.26) yields

$$\begin{aligned}
359 \quad & L^k - L^{k+1} \geq \|z^k - z^{k+1}\|_{M_2 - \frac{1}{2}AM_1^{-1}A^T - c(p)I}^2 + \|y^k - y^{k+1}\|_{\frac{1}{2}M_1^{-1}}^2 \\
360 \quad (3.27) \quad & \quad - \|A^T z^{k+1} + y^{k+1}\|_{M_1^{-1}}^2. \\
361
\end{aligned}$$

362 Since f is μ_f -strongly convex, we know that ∇f^* is $\frac{1}{\mu_f}$ -Lipschitz continuous. Conse-
363 quently,

$$\begin{aligned}
364 \quad & \|A^T z^{k+1} + y^{k+1}\|_{M_1^{-1}}^2 = \|u^k - u^{k+1}\|_{M_1^{-1}}^2 \leq \frac{1}{\lambda_{\min}(M_1^{-1})} \|M_1^{-1}(u^k - u^{k+1})\|^2 \\
365 \quad (3.28) \quad & \leq \stackrel{(3.22)}{\frac{\|M_1\|}{\mu_f^2}} \|y^k - y^{k+1}\|^2. \\
366
\end{aligned}$$

367 Combining (3.27) and (3.28) gives (3.18).

368 Now, to show (3.19), we use (3.22) to get

$$369 \quad f^*(y^k) \geq f^*(-A^T z^k) + \langle M_1^{-1} u^k, y^k + A^T z^k \rangle,$$

370 And, thus,

$$371 \quad L^k = g^*(z^k) + f^*(y^k) + \langle -A^T z^k - y^k, M_1^{-1} u^k \rangle + \frac{1}{2} \|A^T z^k + y^k\|_{M_1^{-1}}^2$$

$$372 \quad (3.29) \quad \geq g^*(z^k) + f^*(-A^T z^k) + \frac{1}{2} \|A^T z^k + y^k\|_{M_1^{-1}}^2,$$

$$373$$

374 and finally (3.19). \square

375 *Remark 3.10.* In order for $C_2 > 0$, we can set $M_2 = AM_1^{-1}A^T$ as suggested by
 376 subsection 3.2, since $c(p) \propto \alpha^p$ for some $0 < \alpha < 1$ in (3.9) and (3.15), we know that
 377 there exists $p_0 \geq 1$ such that $C_2 > 0$ for any $p \geq p_0$. In our numerical experiments,
 378 Algorithm 3.1 always converges for $p \geq 1$.

379 We conclude this section by showing the convergence of (x^k, z^k) in Algorithm 3.1
 380 to a primal-dual solution pair of (1.1) and (1.2).

381 **THEOREM 3.11.** *Let the assumptions in Theorem 3.9 hold. Then, (x^k, z^k) in*
 382 *Algorithm 3.1 is bounded, and any cluster point of $\{x^k, z^k\}$ is a primal-dual solution*
 383 *pair of (1.1) and (1.2).*

384 *Proof.* According to Theorem 3.8, We just need to show that $\{M_1^{-1}u^k, z^k\}$ is
 385 bounded and its cluster points are primal-dual solution pairs of (1.1).

386 Since L^k is nonincreasing, (3.29) tells us that

$$387 \quad g^*(z^k) + f^*(-A^T z^k) + \frac{1}{2} \|A^T z^k + y^k\|_{M_1^{-1}}^2 \leq L^0 < +\infty.$$

388 Since $g^*(z) + f^*(-A^T z)$ is coercive, we get that $\{z^k\}$ is bounded, and from the
 389 boundedness of $\{A^T z^k + y^k\}$, the boundedness of $\{y^k\}$. Furthermore, (3.22) gives us

$$390 \quad \|M_1^{-1}(u^k - u^0)\| \leq \frac{1}{\mu_f} \|y^k - y^0\|.$$

391 Therefore, $\{M_1^{-1}u^k\}$ is also bounded.

392 Suppose (z^c, y^c, u^c) is a cluster point of $\{z^k, y^k, u^k\}$. Let us show that (z^c, y^c, u^c)
 393 is saddle point of $L(z, y, u)$, i.e.,

$$394 \quad (3.30) \quad \mathbf{0} \in \partial L(z^c, y^c, u^c),$$

396 or equivalently,

$$397 \quad \mathbf{0} \in \partial g^*(z^c) - AM_1^{-1}u^c,$$

$$398 \quad \mathbf{0} = \nabla f^*(y^c) - M_1^{-1}u^c,$$

$$399 \quad \mathbf{0} = A^T z^c + y^c,$$

401 which ensures $(M_1^{-1}u^c, z^c)$ as a primal-dual solution pair of (1.1).

402 In order to show (3.30), we first notice that (3.17) gives

$$403 \quad \partial_x L(z^{k+1}, y^{k+1}, u^{k+1}) = \partial g^*(z^{k+1}) - AM_1^{-1}u^{k+1} + AM_1^{-1}(A^T z^{k+1} + y^{k+1}),$$

$$404 \quad \nabla_y L(z^{k+1}, y^{k+1}, u^{k+1}) = \nabla f^*(y^{k+1}) - M_1^{-1}u^{k+1} + M_1^{-1}(A^T z^{k+1} + y^{k+1}),$$

$$405 \quad \nabla_u L(z^{k+1}, y^{k+1}, u^{k+1}) = M_1^{-1}(-A^T z^{k+1} - y^{k+1}).$$

407 Compare these with the optimality conditions (3.20), (3.22), and (3.23), we have

$$408 \quad d^{k+1} = (d_z^{k+1}, d_y^{k+1}, d_u^{k+1}) \in \partial L(z^{k+1}, y^{k+1}, u^{k+1}),$$

409 where

$$410 \quad d_z^{k+1} = M_2(z^k - z^{k+1}) + 2AM_1^{-1}(u^k - u^{k+1}) - AM_1^{-1}(u^{k-1} - u^k) - \varepsilon^{k+1},$$

$$411 \quad d_y^{k+1} = M_1^{-1}(u^k - u^{k+1}),$$

$$412 \quad d_u^{k+1} = M_1^{-1}(u^{k+1} - u^k).$$

414 Since (3.18), and (3.19) implies $z^k - z^{k+1}, y^k - y^{k+1} \rightarrow \mathbf{0}$, (3.22) gives $u^k - u^{k+1} \rightarrow \mathbf{0}$.
415 Combine these with (3.8), we have $d^k \rightarrow \mathbf{0}$.

416 Finally, let us take a subsequence $\{z^{k_s}, y^{k_s}, u^{k_s}\} \rightarrow (z^c, y^c, u^c)$, since $d^{k_s} \rightarrow \mathbf{0}$
417 as $s \rightarrow +\infty$, [33, Def. 8.3] and [33, Prop. 8.12] yield (3.30), which tells us that
418 $(M_1^{-1}u^c, z^c)$ is a primal-dual solution pair of (1.1). \square

419 Also, we can show that (x^k, z^k) in Algorithm 3.1 actually converges. Since the
420 proof consists of a standard technique of using the KL property in Assumption 2,
421 which is not very relevant to the main idea of this subsection, we leave it to Appendix
422 C.

423 **THEOREM 3.12.** *Let the assumptions in Theorem 3.9 hold, then the $\{x^k, z^k\}$ in*
424 *Algorithm 3.1 converges to a primal-dual solution pair of (1.1).*

425 *Proof.* See Appendix C. \square

426 **4. Numerical experiments.** In this section, we compare our inexact preconditioned PDHG in algorithm 3.1 with PDHG (3.2) and PDHG with diagonal preconditioning [29]. We consider three popular applications of PDHG: TV-L¹ denoising, graph cuts, and estimation of earth mover's distance. Although they do not satisfy all the assumptions in our theory, we still observe significant speedup compared to other algorithms.

432 When we write these examples in the form of (1.1), A is one of the following:

433 **Case 1:** The 2D discrete gradient operator $D : \mathbb{R}^{M \times N} \rightarrow \mathbb{R}^{2M \times N}$:

434 Let the images be of size $M \times N$, and h be the length of discretization interval,
435 then

$$436 \quad (Du)_{i,j} = \begin{pmatrix} (Du)_{i,j}^1 \\ (Du)_{i,j}^2 \end{pmatrix},$$

438 where

$$439 \quad (Du)_{i,j}^1 = \begin{cases} \frac{1}{h}(u_{i+1,j} - u_{i,j}) & \text{if } i < M, \\ 0 & \text{if } i = M, \end{cases}$$

$$440 \quad (Du)_{i,j}^2 = \begin{cases} \frac{1}{h}(u_{i,j+1} - u_{i,j}) & \text{if } j < N, \\ 0 & \text{if } j = N. \end{cases}$$

442 **Case 2:** The weighted gradient operator $D_w : \mathbb{R}^{M \times N} \rightarrow \mathbb{R}^{2M \times N}$:

$$443 \quad D_w = \text{diag}(w)D,$$

444 where $w \in (\mathbb{R}^+)^{2MN}$ is a weight vector.

445 **Case 3:** The 2D discrete divergence operator $\text{div}: \mathbb{R}^{2M \times N} \rightarrow \mathbb{R}^{M \times N}$:

446 (4.1)
$$\text{div}(p)_{i,j} = h(p_{i,j}^1 - p_{i-1,j}^1 + p_{i,j}^2 - p_{i,j-1}^2),$$

448 where $p = (p^1, p^2)^T \in \mathbb{R}^{2M \times N}$, $p_{0,j}^1 = p_{M,j}^1 = 0$ and $p_{i,0}^2 = p_{i,N}^2 = 0$ for
449 $i = 1, \dots, M, j = 1, \dots, N$.

450 In view of the special structures of these operators, we choose cyclic proximal
451 Block Coordinate Descent (BCD) as the z -subproblem solver in Algorithm 3.1. In
452 particular, we split $\{1, 2, \dots, m\}$ into 2 blocks (in case 3) or 4 blocks (in cases 1 and 2)
453 according to Claims 4.1 and 4.2, which are inspired by the popular red-black ordering
454 [34] for solving sparse linear system.

455 According to Theorem 3.6, finite inner loops of cyclic proximal BCD satisfy the
456 bounded relative error condition in Def.3.3, and we can expect that this solver brings
457 faster overall convergence. The intuition is that when g^* is linear (or equivalently,
458 g is a δ function), the z -subproblem in Alg.3.1 reduces to a linear system with a
459 structured sparse matrix AA^T . As a result, Gradient Descent amounts to Richardson
460 method [32, 34] and cyclic BCD becomes Gauss-Seidel method [16, 34]. In view of the
461 special structures of A , Gauss-Seidel is faster (see Chapter 4 of [34]). Therefore, we
462 can anticipate a faster convergence when cyclic proximal BCD is used.

463 Furthermore, the following two claims tell us that under special block designs
464 inspired by red-black ordering, the two subproblems have closed-form solutions, which
465 are easy to implement and compute. Furthermore, updates within each block can be
466 implemented in parallel.

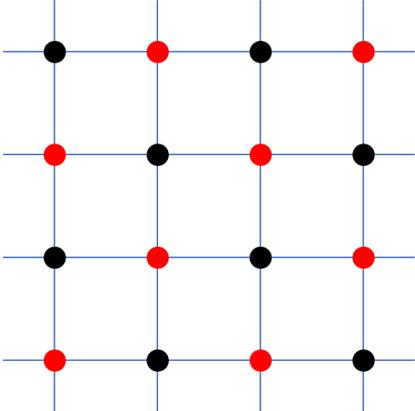


FIG. 1. two-block ordering in Claim 4.1

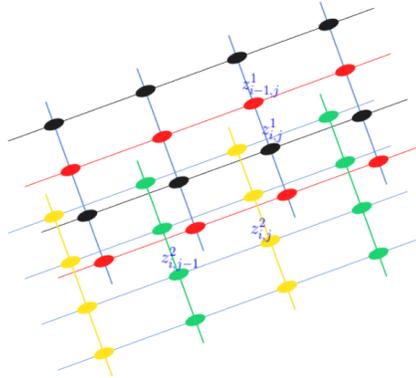


FIG. 2. four-block ordering in Claim 4.2

467 CLAIM 4.1. When $A = \text{div}$ (i.e. $A^T = -D$) and $M_2 = \tau AA^T$, for $z \in \mathbb{R}^{M \times N}$, we
468 separate z into two block z_b, z_r where

469
$$z_b := \{z_{i,j} \mid i+j \text{ is even}\}, z_r := \{z_{i,j} \mid i+j \text{ is odd}\},$$

470 for $1 \leq i \leq M, 1 \leq j \leq N$. If $g(z) = \sum_{i,j} g_{i,j}(z_{i,j})$ and $\text{prox}_{\lambda g_{i,j}^*}$ have closed-form
471 solutions for all $1 \leq i \leq M, 1 \leq j \leq N$ and $\lambda > 0$, then BCD subproblems on the
472 z -subproblem of Algorithm 3.1 have closed-form solutions, and updates within each
473 block can be implemented in parallel.

474 *Proof.* As illustrated in Fig. 1, on the z -subproblem, the update of every black
 475 node depends only on its neighbor red nodes, thus all the black nodes can be updated
 476 in parallel and with closed-form solutions. The same argument applies to the red
 477 nodes. See Appendix D for a complete explanation. \square

478 CLAIM 4.2. When $A = D$ or $A = D_w$ (i.e. $A^T = -\text{div}$ or $A^T = -\text{div diag}(w)$)
 479 and $M_2 = \tau AA^T$, for $z = (z^1, z^2)^T \in \mathbb{R}^{2M \times N}$, we separate z into four blocks $z_b, z_r,$
 480 z_y and z_g , where

$$481 \quad z_b = \{z_{i,j}^1 \mid i \text{ is odd}\}, \quad z_r = \{z_{i,j}^1 \mid i \text{ is even}\},$$

$$482 \quad z_y = \{z_{i,j}^2 \mid j \text{ is odd}\}, \quad z_g = \{z_{i,j}^2 \mid j \text{ is even}\},$$

484 for $1 \leq i \leq M, 1 \leq j \leq N$. If $g(z) = \sum_{i,j} g_{i,j}(z_{i,j})$ and $\text{prox}_{\lambda g_{i,j}^*}$ have closed-form
 485 solutions for all $1 \leq i \leq M, 1 \leq j \leq N$ and $\lambda > 0$, then BCD subproblems on the
 486 z -subproblem of Algorithm 3.1 have closed-form solutions, and updates within each
 487 block can be implemented in parallel.

488 *Proof.* In Figure 2, the 4 blocks are in 4 different colors. Nodes with the same
 489 color can be updated in parallel with closed-form solutions, as within one color nodes
 490 are independent with each other during the updates. See Appendix D for details. \square

491 In Table 1, Table 2 and Figure 7, PDHG denotes the PDHG in (3.2); DP-PDHG
 492 denotes the diagonal preconditioned PDHG in [29], PPDHG denotes PPDHG in
 493 (3.3) where the $(k+1)$ th z -subproblem is solved until $\frac{\|z^k - z^{k+1}\|_2}{\max\{1, \|z^{k+1}\|_2\}} < 10^{-5}$ using
 494 the TFOCS [4] implementation of FISTA with restart; Alg. 3.1, BCD denotes our
 495 inexact preconditioned PDHG in Algorithm 3.1, where the inner loop solver S is cyclic
 496 proximal BCD. Except for DP-PDHG, only the best runtime over certain choices of
 497 parameters is presented.

498 Comparison of PDHG and DP-PDHG have already been presented for TV-L¹
 499 denoising and graph cuts in [29], and PDHG is proposed to estimate the earth mover's
 500 distance in [22]. In order to provide a direct comparison, we use their problem
 501 formulations.

502 **4.1. Total variation based image denoising.** The (discrete) TV-L¹ model
 503 for image denoising can be expressed as

$$504 \quad \text{minimize} \quad \Phi(u) = \|Du\|_1 + \lambda \|u - f\|_1,$$

505 where D is the 2D discrete gradient operator with $h = 1$, $u \in \mathbb{R}^{M \times N}$ is the sought
 506 solution, $f \in \mathbb{R}^{M \times N}$ is a noisy input image, and λ is a regularization parameter. In
 507 our experiment we input a 1024×1024 image with noise level 0.15 and set $\lambda = 1$, see
 508 Fig. 3. We run the algorithms until $\delta^k := \frac{|\Phi^k - \Phi^*|}{|\Phi^*|} < 10^{-6}$, where Φ^k is the objective
 509 value at k th iteration and Φ^* is the optimal objective value obtained by calling CVX.

510 Our Numerical results on TV-L¹ model are summarized in Table 1, where the best
 511 results for $\tau \in \{10, 1, 0.1, 0.01, 0.001\}$ and $p \in \{1, 2, 3\}$ are presented. Our Algorithm
 512 3.1 is significantly faster than the other three algorithms.

513 Remarkably, our algorithm's number of outer iterations is less than that of PPDHG
 514 with the stopping criterion $\frac{\|z^k - z^{k+1}\|_2}{\max\{1, \|z^{k+1}\|_2\}} < 10^{-5}$, as this kind of stopping criteria
 515 may become looser as z^k is closer to z^* . In this example, $\frac{\|z^k - z^{k+1}\|_2}{\max\{1, \|z^{k+1}\|_2\}} < 10^{-5}$ only
 516 requires 1 inner iteration of FISTA when Outer Iter ≥ 368 , while as high as 228 inner
 517 loops on average during the first 100 outer iterations. In comparison, our algorithm
 518 achieves both less outer iterations and cheaper cost per outer iteration.

519 In addition, the diagonal preconditioner designed in [29] has little effects when
 520 $A = D$. In fact, $M_1 = \text{diag}(\Sigma_i |A_{i,j}|)$ will be $4I_n$ and $M_2 = \text{diag}(\Sigma_j |A_{i,j}|)$ will be $2I_m$
 521 if we ignore the Neumann boundary condition. With these fixed and almost the same
 522 parameters, DP-PDHG performs even worse than PDHG.

Method	Parameters	Outer Iter	Runtime(s)
PDHG	$\tau = 0.01, M_1 = \tau^{-1}I_n, M_2 = \tau\ D\ ^2I_m$	2990	114.2576
DP-PDHG	$M_1 = \text{diag}(\Sigma_i D_{i,j}), M_2 = \text{diag}(\Sigma_j D_{i,j})$	8856	329.7890
PPDHG (3.3)	$\tau = 0.1, M_1 = \tau^{-1}I_n, M_2 = \tau DD^T$	963	5.9777×10^3
Alg. 3.1, BCD	$\tau = 0.01, M_1 = \tau^{-1}I_n, M_2 = \tau DD^T, p = 1$	541	26.2704

TABLE 1
 TV- L^1 denoising.



FIG. 3. Noisy image



FIG. 4. Denoised image (Alg. 3.1, BCD)

523 **4.2. Graph cuts.** The total-variation-based graph cut model is to minimize the
 524 follow weighted TV energy:

$$\begin{aligned}
 & \text{minimize} && \|D_w u\|_1 + \langle u, \omega^u \rangle \\
 & \text{subject to} && 0 \leq u \leq 1,
 \end{aligned}$$

526 where $w^u \in \mathbb{R}^{M \times N}$ is a vector of unary weights, $w^b \in \mathbb{R}^{2MN}$ is a vector of binary
 527 weights, and $D_w = \text{diag}(w^b)D$, where D is the 2D discrete gradient operator with
 528 $h = 1$. Specifically, $w_{i,j}^u = \alpha(\|I_{i,j} - \mu_f\|^2 - \|I_{i,j} - \mu_b\|^2)$, $w_{i,j}^{b,1} = \exp(-\beta|I_{i+1,j} - I_{i,j}|)$
 529 and $w_{i,j}^{b,2} = \exp(-\beta|I_{i,j+1} - I_{i,j}|)$. In our experiment the image is of the size 660×720 ,
 530 and we set $\alpha = 1/2$, $\beta = 10$, $\mu_f = [0; 0; 1]$ (for the blue foreground) and $\mu_b = [0; 1; 0]$
 531 (for the green background). We run all algorithms until $\delta^k := \frac{|\Phi^k - \Phi^*|}{|\Phi^*|} < 10^{-8}$, where
 532 Φ^k is the objective value at k th iteration and Φ^* is the optimal objective value obtained
 533 by calling CVX.

534 The best results of $\tau \in \{10, 1, 0.1, 0.01, 0.001\}$ and $p \in \{1, 2, 3\}$ are summarized in
 535 Table 2, where we can see that our algorithm yields the best performance on runtime.
 536 Also, our algorithm's number of outer iterations is close to that of PPDHG.

Method	Parameters	Outer Iter	Runtime(s)
PDHG	$\tau = 1, M_1 = \tau^{-1}I_n, M_2 = \tau\ D_w\ ^2I_m$	5529	140.5777
DP-PDHG	$M_1 = \text{diag}(\Sigma_i D_{w_{i,j}}), M_2 = \text{diag}(\Sigma_j D_{w_{i,j}})$	3572	108.3573
PPDHG (3.3)	$\tau = 10, M_1 = \tau^{-1}I_n, M_2 = \tau D_w D_w^T$	282	938.3787
Alg. 3.1, BCD	$\tau = 10, M_1 = \tau^{-1}I_n, M_2 = \tau D_w D_w^T, p = 2$	411	14.9663

TABLE 2
Graph cuts



FIG. 5. Input image



FIG. 6. Graph cut (Alg. 3.1, BCD)

537 **4.3. Estimation of earth mover's distance.** We consider the estimation of
 538 earth mover's distance, which is a popular model in image processing, computer vision
 539 and statistics [20, 25, 28]. From [22] we know that the problem can be formulated as

$$540 \quad (4.2) \quad \begin{aligned} & \text{minimize} && \|m\|_{1,2} \\ & \text{subject to} && \text{div}(m) + \rho^1 - \rho^0 = 0, \end{aligned}$$

541 where $m \in \mathbb{R}^{2M \times N}$ is the sought flux vector on the $M \times N$ grid, and ρ^0, ρ^1 represents
 542 two mass distributions on the $M \times N$ grid. The setting in our experiment here is the
 543 same with that in [22], i.e. $M = N = 256, h = \frac{N-1}{4}$, and for ρ^0 and ρ^1 see Fig. 8.

544 Since the iterates m^k may not satisfy the linear constraint, the objective $\Phi(m) =$
 545 $I_{\{m|\text{div}(m)=\rho^0-\rho^1\}} + \|m\|_{1,2}$ is not comparable. Instead, we compare $\|m^k\|_{1,2}$ and
 546 the constraint violation until $k = 100000$ outer iterations in Fig. 7, where we set
 547 $\tau = 3 \times 10^{-6}$ as in [22], and $\sigma = \frac{1}{\tau\|\text{div}\|^2}$. In Fig. 7, we can see that our algorithm
 548 provides much lower constraint violation as well as much better estimation for the
 549 earth mover's distance $\|m\|_{1,2}$. Fig. 8 shows the solution obtained by Alg. 3.1, where
 550 m is the flux that moves the standing cat ρ^1 into the crouching cat ρ^0 . DP-PDHG
 551 and PPDHG are extremely slow in this example. Similar to 4.1, when $A = \text{div}$,
 552 the diagonal preconditioners proposed in [29] are approximately equivalent to fixed
 553 constant parameters $\tau = \frac{1}{2h}, \sigma = \frac{1}{4h}$ and they lead to extremely slow convergence. As
 554 for PPDHG, it suffers from the expensive cost per outer iteration as in the previous
 555 two experiments.

556 It is worth mentioning that unlike [22], the algorithms in our experiments are not
 557 implemented in a parallel fashion. On the other hand, in our Algorithm 3.1 with cyclic
 558 proximal BCD as the inner loop solver, coordinates in each block in the block designs
 559 of Fig. 1 and 2 can be updated in parallel. Therefore, one can expect a further speed
 560 up by a parallel implementation.

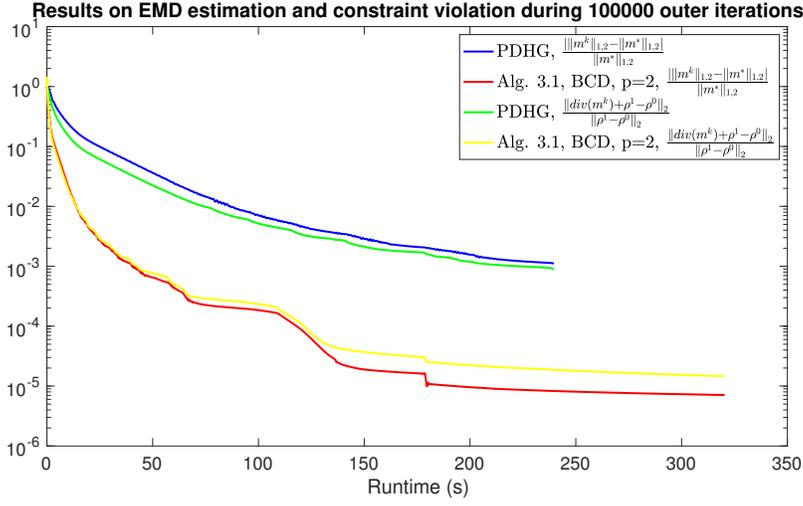


FIG. 7. For PDHG, $\tau = 3 \times 10^{-6}$, $\sigma = \frac{1}{\tau \|\text{div}\|^2}$; For Alg. 3.1, BCD, $\tau = 3 \times 10^{-6}$, $M_1 = \tau^{-1} I_n$, $M_2 = \tau \text{divdiv}^T$, $p = 2$. $\|m^*\|_{1,2} = 0.6718$ is given by gurobi of CVX.

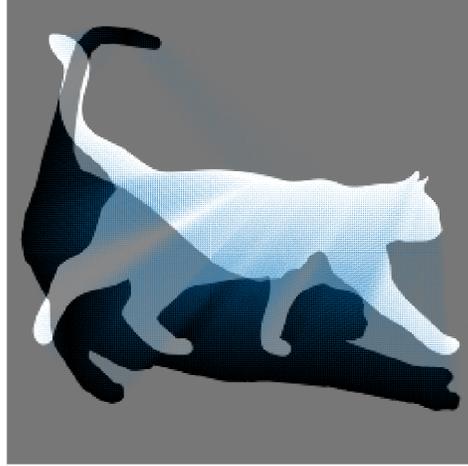


FIG. 8. ρ^0 , ρ^1 are the white standing cat, and the black crouching cat, respectively. The images are of the size 256×256 , and the earth mover's distance between ρ^0 and ρ^1 is 0.6718.

561 **5. Concluding Remarks.** In this paper, We provide an algorithmic framework
 562 for apply preconditioning and fast subproblem solvers on PDHG and ADMM with
 563 convergence guarantees. Remarkably, we allow a fixed number of inner iterations for
 564 one of the subproblems. Although the examples in our numerical experiments do
 565 not satisfy all the assumptions, significant accelerations in both outer iteration and

runtime are observed when proper preconditioners and subproblem solvers are applied.

There are still some interesting questions, which need to be addressed in the future: (a) According to Theorem 3.2, there may be better preconditioners than $M_1 = \frac{1}{\tau}I_{n \times n}, M_2 = \tau AA^T$, which lead to ADMM iterations. (b) It is possible that convergence of Algorithm 3.1 can also be established for even faster accelerated subproblem solvers like APCG [23], NU_ACDM [1], and A2BCD [19]. (c) It is possible that a broad class of algorithms can be accelerated by integrating preconditioning, fixed number of inner loops, and suitable subproblem solvers. We hope our framework can be applied on more algorithms with faster convergence guarantees.

Appendix A. ADMM as a special case of PPDHG.

In this section we show that if we choose $M_1 = \frac{1}{\tau}$ and $M_2 = \tau AA^T$ in PPDHG (3.3), then it is equivalent to ADMM on the primal problem (1.1).

By Theorem 1 of [37], we know that ADMM is primal-dual equivalent, in the sense that one can recover primal iterates from dual iterates and vice versa. Therefore, it suffices to show that $M_1 = \frac{1}{\tau}$ and $M_2 = \tau AA^T$ in PPDHG (3.3) on the primal problem is equivalent to ADMM on the dual problem (1.2).

In Theorem 3.8 we have shown that, under an appropriate change of variables, PPDHG on the primal is equivalent to PLADMM in (3.16) on the dual. As a result, we just need to demonstrate that PLADMM on the dual is exactly ADMM on the dual when $M_1 = \frac{1}{\tau}I_{n \times n}$ and $M_2 = \tau AA^T$.

For the z -update in (3.16), we have

$$\begin{aligned}
z^{k+1} &= \arg \min_{z \in \mathbb{R}^m} \{g^*(z) - \tau \langle z - z^k, A(-A^T z^k - y^k + u^k) \rangle + \frac{\tau}{2} \|z - z^k\|_{AA^T}^2\} \\
&= \arg \min_{z \in \mathbb{R}^m} \{g^*(z) - \tau \langle z - z^k, A(-y^k + u^k) \rangle + \frac{\tau}{2} \|z\|_{AA^T}^2\} \\
&= \arg \min_{z \in \mathbb{R}^m} \{g^*(z) + \tau \langle z, A(y^k - u^k) \rangle + \frac{\tau}{2} \|A^T z\|^2\} \\
&= \arg \min_{z \in \mathbb{R}^m} \{g^*(z) + \tau \langle A^T z, -u^k \rangle + \frac{\tau}{2} \|A^T z + y^k\|^2\} \\
\text{(A.1)} \quad &= \arg \min_{z \in \mathbb{R}^m} \{g^*(z) + \tau \langle -A^T z - y^k, u^k \rangle + \frac{\tau}{2} \|A^T z + y^k\|^2\}.
\end{aligned}$$

and for the y -update we have

$$\begin{aligned}
y^{k+1} &= \text{Prox}_{f^*}^{M_1^{-1}}(u^k - A^T z^{k+1}) \\
&= \arg \min_{y \in \mathbb{R}^n} \{f^*(y) + \frac{\tau}{2} \|y - u^k + A^T z^{k+1}\|^2\} \\
\text{(A.2)} \quad &= \arg \min_{y \in \mathbb{R}^n} \{f^*(y) + \tau \langle -A^T z^{k+1} - y, u^k \rangle + \frac{\tau}{2} \|A^T z^{k+1} + y\|^2\}.
\end{aligned}$$

Define $v^k = \tau u^k$, (A.1), (A.2), and the u -update in (3.16) become

$$\begin{aligned}
z^{k+1} &= \arg \min_{z \in \mathbb{R}^m} \{g^*(z) + \langle -A^T z - y^k, v^k \rangle + \frac{\tau}{2} \|A^T z + y^k\|^2\}, \\
y^{k+1} &= \arg \min_{y \in \mathbb{R}^n} \{f^*(y) + \langle -A^T z^{k+1} - y, v^k \rangle + \frac{\tau}{2} \|A^T z^{k+1} + y\|^2\}, \\
v^{k+1} &= v^k - \tau(A^T z^{k+1} + y^{k+1}),
\end{aligned}$$

604 which are ADMM iterations on the dual problem (1.2).

605 **Appendix B. Proof of Theorem 3.6: Cyclic proximal BCD satisfies**
 606 **bounded relative error condition.**

607 The z -subproblem in (3.3) is of the form

$$608 \quad \min_{z \in \mathbb{R}^m} h_1(z) + h_2(z),$$

609 where

$$610 \quad h_1(z) = g^*(z) = \sum_{i=1}^l g_i^*(z_i),$$

$$611 \quad h_2(z) = \frac{1}{2} \|z - z^k - M_2^{-1}A(2x^{k+1} - x^k)\|_{M_2}^2.$$

613 And z^{k+1} is given by

$$614 \quad z_0^{k+1} = z^k,$$

$$615 \quad z_{i+1}^{k+1} = S(z_i^{k+1}, x^{k+1}, x^k), \quad i = 0, 1, \dots, p-1,$$

$$616 \quad z_p^{k+1} = z_p^{k+1}.$$

618 Here the inner loop solver S is cyclic proximal BCD.

619 Let us define

$$620 \quad T(z) = \text{Prox}_{\gamma g^*(z)}(z - \gamma \nabla h_2(z)),$$

$$621 \quad B(z) = \frac{1}{\gamma}(z - T(z)),$$

623 and the i th coordinate operator of B :

$$624 \quad B_i(z) = (0, \dots, (B(z))_i, \dots, 0).$$

625 Then

$$626 \quad z_{i+1}^{k+1} = S(z_i^{k+1}, x^{k+1}, x^k) = (I - \gamma B_i)(I - \gamma B_2) \dots (I - \gamma B_1) z_i^{k+1}.$$

627 By [3, Prop. 26.16(ii)], we know that $T(z)$ is a contraction with coefficient $\theta =$
 628 $\sqrt{1 - \gamma(2\lambda_{\min}(M_2) - \gamma\lambda_{\max}^2(M_2))}$. Together with [3,], we know that for $\forall z_1, z_2 \in \mathbb{R}^m$
 629 we have,

$$630 \quad \langle B(z_1) - B(z_2), z_1 - z_2 \rangle = \frac{1}{\gamma} \|z_1 - z_2\|^2 - \frac{1}{\gamma} \langle T(z_1) - T(z_2), z_1 - z_2 \rangle$$

$$631 \quad \geq \mu \|z_1 - z_2\|^2,$$

633 where $\mu = \frac{1-\theta}{\gamma}$.

634 Let $z_\star^{k+1} = \arg \min_{z \in \mathbb{R}^m} \{h_1(z) + h_2(z)\}$. By [9, Thm 3.5], we know that

$$635 \quad (\text{B.1}) \quad \|z_i^{k+1} - z_\star^{k+1}\| \leq \rho^i \|z_0^{k+1} - z_\star^{k+1}\|, \quad \forall i = 1, 2, \dots, p.$$

637 where $\rho = 1 - \frac{\gamma\mu^2}{2}$.

638 Let $y_j = (I - \gamma B_j) \dots (I - \gamma B_1) z_{p-1}^{k+1}$ for $j = 1, \dots, l$ and $y_0 = z_{p-1}^{k+1}$. Note that
 639 $(z_p^{k+1})_j = (y_j)_j$ for $j = 1, 2, \dots, l$, and the blocks of y_j satisfies

$$640 \quad (y_j)_t = \begin{cases} \left(\text{Prox}_{\gamma g^*} (y_{j-1} - \gamma \nabla h_2(y_{j-1})) \right)_t, & \text{if } t = j \\ (y_{j-1})_t, & \text{otherwise.} \end{cases}$$

642 On the other hand, we have

$$643 \quad \text{Prox}_{\gamma g^*} (y_{j-1} - \gamma \nabla h_2(y_{j-1})) = \arg \min_{y \in \mathbb{R}^m} \{g^*(y) + \frac{1}{2\gamma} \|y - y_{j-1} + \gamma \nabla h_2(y_{j-1})\|^2\}.$$

644 Since g^* and $\|\cdot\|^2$ are separable, we obtain

$$645 \quad \mathbf{0} \in \partial g_j^*((y_j)_j) + \frac{1}{\gamma} \left((y_j)_j - (y_{j-1})_j + \gamma (\nabla h_2(y_{j-1}))_j \right), \quad \forall j = 1, 2, \dots, l,$$

646 or equivalently,

$$647 \quad \mathbf{0} \in \partial g_j^*((z_p^{k+1})_j) + \frac{1}{\gamma} \left((z_p^{k+1})_j - (z_{p-1}^{k+1})_j + \gamma (\nabla h_2(y_{j-1}))_j \right), \quad \forall j = 1, 2, \dots, l.$$

648 As a result,

$$649 \quad \mathbf{0} \in \partial g^*(z_p^{k+1}) + \frac{1}{\gamma} (z_p^{k+1} - z_{p-1}^{k+1} + \gamma \xi_p), \quad \forall j = 1, 2, \dots, l,$$

650 where $(\xi_p)_j = (\nabla h_2(y_{j-1}))_j$ for $j = 1, 2, \dots, l$. Compare this with (3.7), we know that

$$651 \quad \varepsilon^{k+1} = \xi_p - \nabla h_2(z_p^{k+1}) + \frac{1}{\gamma} (z_p^{k+1} - z_{p-1}^{k+1}).$$

652 Notice that the first $j-1$ blocks of y_{j-1} are the same with those of $y_l = z_p^{k+1}$, and
 653 the rest of the blocks are the same with those of $y_0 = z_{p-1}^{k+1}$, so we have

$$654 \quad \begin{aligned} \|\varepsilon^{k+1}\| &\leq \sum_{j=1}^l \lambda_{\max}(M_2) \|y_{j-1} - z_p^{k+1}\| + \frac{1}{\gamma} \|z_p^{k+1} - z_{p-1}^{k+1}\| \\ 655 \quad &\leq l \lambda_{\max}(M_2) \|z_{p-1}^{k+1} - z_p^{k+1}\| + \frac{1}{\gamma} \|z_p^{k+1} - z_{p-1}^{k+1}\| \\ 656 \quad &\leq (l \lambda_{\max}(M_2) + \frac{1}{\gamma}) (\|z_p^{k+1} - z_{\star}^{k+1}\| + \|z_{p-1}^{k+1} - z_{\star}^{k+1}\|) \end{aligned}$$

658 Combine this with (B.1), we arrive at

$$659 \quad (\text{B.2}) \quad \|\varepsilon^{k+1}\| \leq (l \lambda_{\max}(M_2) + \frac{1}{\gamma}) (\rho^p + \rho^{p-1}) \|z_0^{k+1} - z_{\star}^{k+1}\|.$$

661 We also have

$$662 \quad \begin{aligned} \|z^{k+1} - z^k\| &= \|z_p^{k+1} - z_0^{k+1}\| \\ 663 \quad &\geq \|z_0^{k+1} - z_{\star}^{k+1}\| - \|z_p^{k+1} - z_{\star}^{k+1}\| \\ 664 \quad &\geq (1 - \rho^p) \|z_0^{k+1} - z_{\star}^{k+1}\| \end{aligned}$$

666 Combine this with (B.2), we obtain

$$667 \quad \|\varepsilon^{k+1}\| \leq \frac{(\lambda_{\max}(M_2) + \frac{1}{\gamma})(\rho^p + \rho^{p-1})}{1 - \rho^p} \|z^{k+1} - z^k\|.$$

668 **Appendix C. Proof of Theorem 3.12: KL property gives sequence con-**
669 **vergence.**

670 According to Theorem 3.8, We just need to show that $\{M_1^{-1}u^k, z^k\}$ converges to
671 a primal-dual solution pair of (1.1).

672 By Theorem 3.11, we can take $\{z^{k_s}, y^{k_s}, u^{k_s}\} \rightarrow (z^c, y^c, u^c)$. Note that
673 $L(z^{k_s}, y^{k_s}, u^{k_s})$ is monotonic nonincreasing and lower bounded due to Theorem 3.9,
674 which implies the convergence of $L(z^{k_s}, y^{k_s}, u^{k_s})$. Since L is lower semicontinuous, we
675 have

$$676 \quad (C.1) \quad L(z^c, y^c, u^c) \leq \lim_{s \rightarrow \infty} L(z^{k_s}, y^{k_s}, u^{k_s}).$$

678 Since the only potentially discontinuous terms in L is g^* , we have

$$679 \quad (C.2) \quad \lim_{s \rightarrow \infty} L(z^{k_s}, y^{k_s}, u^{k_s}) - L(z^c, y^c, u^c) \leq \limsup_{s \rightarrow \infty} g^*(z^{k_s}) - g^*(z^c).$$

681 By (3.20), we know that

$$682 \quad g^*(z^c) \geq g^*(z^{k_s}) \\ 683 \quad + \langle M_2(z^{k_s-1} - z^{k_s}) + AM_1^{-1}(-A^T z^{k_s-1} - y^{k_s-1} + u^{k_s-1}) - \varepsilon^{k_s}, z^c - z^{k_s} \rangle,$$

685 By Theorem 3.9, we know that $z^{k_s-1} - z^{k_s} \rightarrow \mathbf{0}$. Since $z^{k_s} \rightarrow z^c$ and $\{z^k, y^k, u^k\}$ is
686 bounded, we obtain

$$687 \quad \limsup_{s \rightarrow \infty} g^*(z^{k_s}) - g^*(z^c) \leq 0.$$

688 Combine this with (C.1) and (C.2), we conclude that $\lim_{s \rightarrow \infty} L(z^{k_s}, y^{k_s}, u^{k_s}) =$
689 $L(z^c, y^c, u^c)$.

690 Since g^* is a KL function, L is also KL. As a result, similar to Theorem 2.9 of [2],
691 we can claim the convergence of $\{z^k, y^k, u^k\}$ to $\{z^c, y^c, u^c\}$.

692 **Appendix D. Two-block ordering in Claim 4.1 and Four-block ordering**
693 **in Claim 4.2.**

694 According to (3.4), when $M_2 = \tau AA^T$, the z -subproblem of Algorithm 3.1 is

$$695 \quad (D.1) \quad z^{k+1} = \arg \min_{z \in \mathbb{R}^m} \{g^*(z) - \langle z - z^k, A(2x^{k+1} - x^k) \rangle + \frac{\tau}{2} \|A^T(z - z^k)\|_2^2\}.$$

697 Let us prove Claim 4.1 first.

698 In claim 4.1, $A = \text{div} \in \mathbb{R}^{MN \times 2MN}$ and $z \in \mathbb{R}^{MN}$. Following the definition of the
699 sets z_b and z_r in Claim 4.1, we separate the MN columns of $A^T = -D$ into two blocks
700 L_b, L_r associated with z_b and z_r , respectively. Therefore, we have $A^T z = L_b z_b + L_r z_r$
701 for any $z \in \mathbb{R}^{MN}$.

702 By the red-black ordering in Fig. 1, different columns of L_b are orthogonal to each
703 other, therefore, $L_b^T L_b$ is diagonal. Similarly, $L_r^T L_r$ is also diagonal.

704 Let b be the set of black nodes and r the set of red nodes, then we can rewrite
705 (D.1) as

$$706 \quad (D.2) \quad z^{k+1} = \arg \min_{z_b, z_r \in \mathbb{R}^{MN/2}} \{g_b^*(z_b) + g_r^*(z_r) + \langle z_b + z_r, c^k \rangle \\ 707 \quad + \frac{\tau}{2} \|L_b(z_b - z_b^k) + L_r(z_r - z_r^k)\|_2^2\},$$

708

709 where $g_b^*(z_b) = \sum_{(i,j) \in b} g_{i,j}^*(z_{i,j})$, $g_r^*(z_r) = \sum_{(i,j) \in r} g_{i,j}^*(z_{i,j})$, and $c^k = -A(2x^{k+1} -$
 710 $x^k)$.

711 The cyclic proximal BCD applied on black and red blocks is then

$$712 \quad (D.3) \quad z_b^{k+\frac{t+1}{p}} = \operatorname{prox}^{\tau L_b^T L_b}_{g_b^*(\cdot) + \langle \cdot, \tau L_b^T L_r(z_r^{k+\frac{t}{p}} - z_r^k) + c_b^k \rangle} \left(z_b^{k+\frac{t}{p}} \right),$$

$$713 \quad (D.4) \quad z_r^{k+\frac{t+1}{p}} = \operatorname{prox}^{\tau L_r^T L_r}_{g_r^*(\cdot) + \langle \cdot, \tau L_r^T L_b(z_b^{k+\frac{t+1}{p}} - z_b^k) + c_r^k \rangle} \left(z_r^{k+\frac{t}{p}} \right).$$

715 for $t = 0, 1, \dots, p-1$, where p is the number of inner loops as in Algorithm 3.1.

716 These updates have closed-form solutions since $L_b^T L_b$ and $L_r^T L_r$ are diagonal, and
 717 all $\operatorname{prox}_{g_{i,j}^*}$ have closed form solutions. Furthermore, updates within each block can
 718 be implemented in parallel.

719 The proof of Claim 4.2 follows in a similar way. When $A = D$ or $A = D_w$, we
 720 separate the columns of A^T into four blocks L_b, L_r, L_y, L_g associated with $z_b, z_r,$
 721 z_y, z_g , respectively. Therefore, we have $A^T z = L_b z_b + L_r z_r + L_y z_y + L_g z_g$ for all
 722 $z \in \mathbb{R}^{2MN}$. Similarly, by the block design in Fig. 2, we know that cyclic proximal
 723 BCD iterations on the z -subproblem have closed-form solutions, and updates within
 724 each block can be implemented in a parallel fashion.

725

REFERENCES

- 726 [1] Z. ALLEN-ZHU, Z. QU, P. RICHTÁRIK, AND Y. YUAN, *Even faster accelerated coordinate*
 727 *descent using non-uniform sampling*, in International Conference on Machine Learning,
 728 2016, pp. 1110–1119.
- 729 [2] H. ATTOUCH, J. BOLTE, AND B. F. SVAITER, *Convergence of descent methods for semi-algebraic*
 730 *and tame problems: proximal algorithms, forward-backward splitting, and regularized gauss-*
 731 *seidel methods*, *Mathematical Programming*, 137 (2013), pp. 91–129.
- 732 [3] H. H. BAUSCHKE, P. L. COMBETTES, ET AL., *Convex Analysis and Monotone Operator Theory*
 733 *in Hilbert Spaces*, vol. 2011, Springer, 2017.
- 734 [4] S. R. BECKER, E. J. CANDÈS, AND M. C. GRANT, *Templates for convex cone problems with*
 735 *applications to sparse signal recovery*, *Mathematical programming computation*, 3 (2011),
 736 p. 165.
- 737 [5] K. BREDIES AND H. SUN, *Preconditioned douglas–rachford splitting methods for convex-concave*
 738 *saddle-point problems*, *SIAM Journal on Numerical Analysis*, 53 (2015), pp. 421–444.
- 739 [6] K. BREDIES AND H. SUN, *A proximal point analysis of the preconditioned alternating direction*
 740 *method of multipliers*, *Journal of Optimization Theory and Applications*, 173 (2017),
 741 pp. 878–907.
- 742 [7] A. CHAMBOLLE AND T. POCK, *A first-order primal-dual algorithm for convex problems with*
 743 *applications to imaging*, *Journal of mathematical imaging and vision*, 40 (2011), pp. 120–145.
- 744 [8] A. CHAMBOLLE AND T. POCK, *On the ergodic convergence rates of a first-order primal-dual*
 745 *algorithm*, *Mathematical Programming*, 159 (2016), pp. 253–287.
- 746 [9] Y. T. CHOW, T. WU, AND W. YIN, *Cyclic coordinate-update algorithms for fixed-point*
 747 *problems: Analysis and applications*, *SIAM Journal on Scientific Computing*, 39 (2017),
 748 pp. A1280–A1300.
- 749 [10] P. L. COMBETTES AND N. N. REYES, *Moreau’s decomposition in banach spaces*, *Mathematical*
 750 *Programming*, 139 (2013), pp. 103–114.
- 751 [11] J. ECKSTEIN AND D. P. BERTSEKAS, *On the douglas–rachford splitting method and the*
 752 *proximal point algorithm for maximal monotone operators*, *Mathematical Programming*, 55
 753 (1992), pp. 293–318.
- 754 [12] J. ECKSTEIN AND W. YAO, *Relative-error approximate versions of douglas–rachford splitting*
 755 *and special cases of the admm*, *Mathematical Programming*, pp. 1–28.
- 756 [13] J. ECKSTEIN AND W. YAO, *Approximate admm algorithms derived from lagrangian splitting*,
 757 *Computational Optimization and Applications*, 68 (2017), pp. 363–405.
- 758 [14] E. ESSER, X. ZHANG, AND T. F. CHAN, *A general framework for a class of first order primal-*
 759 *dual algorithms for convex optimization in imaging science*, *SIAM Journal on Imaging*
 760 *Sciences*, 3 (2010), pp. 1015–1046.

- 761 [15] D. FEIJER AND F. PAGANINI, *Stability of primal–dual gradient dynamics and applications to*
762 *network optimization*, Automatica, 46 (2010), pp. 1974–1981.
- 763 [16] C. F. GAUSS, *Werke (in German)*, 9, Göttingen: Königlichen Gesellschaft der Wissenschaften,
764 1903.
- 765 [17] P. GISELSSON AND S. BOYD, *Diagonal scaling in douglas-rachford splitting and admm*, in
766 Decision and Control (CDC), 2014 IEEE 53rd Annual Conference on, IEEE, 2014, pp. 5033–
767 5039.
- 768 [18] P. GISELSSON AND S. BOYD, *Linear convergence and metric selection for douglas-rachford*
769 *splitting and admm*, IEEE Transactions on Automatic Control, 62 (2017), pp. 532–544.
- 770 [19] R. HANNAH, F. FENG, AND W. YIN, *A2bcd: An asynchronous accelerated block coordinate*
771 *descent algorithm with optimal complexity*, arXiv preprint arXiv:1803.05578, (2018).
- 772 [20] E. LEVINA AND P. BICKEL, *The earth mover’s distance is the mallows distance: Some insights*
773 *from statistics*, in null, IEEE, 2001, p. 251.
- 774 [21] M. LI, L.-Z. LIAO, AND X. YUAN, *Inexact alternating direction methods of multipliers*
775 *with logarithmic–quadratic proximal regularization*, Journal of Optimization Theory and
776 Applications, 159 (2013), pp. 412–436.
- 777 [22] W. LI, E. K. RYU, S. OSHER, W. YIN, AND W. GANGBO, *A parallel method for earth mover’s*
778 *distance*, UCLA Comput. Appl. Math. Pub.(CAM) Rep, (2017), pp. 17–12.
- 779 [23] Q. LIN, Z. LU, AND L. XIAO, *An accelerated proximal coordinate gradient method*, in Advances
780 in Neural Information Processing Systems, 2014, pp. 3059–3067.
- 781 [24] Q. LIN, Z. LU, AND L. XIAO, *An accelerated randomized proximal coordinate gradient method*
782 *and its application to regularized empirical risk minimization*, SIAM Journal on Optimiza-
783 tion, 25 (2015), pp. 2244–2273.
- 784 [25] L. MÉTIVIER, R. BROSSIER, Q. MÉRIGOT, E. OUDET, AND J. VIRIEUX, *Measuring the misfit*
785 *between seismograms using an optimal transport distance: application to full waveform*
786 *inversion*, Geophysical Supplements to the Monthly Notices of the Royal Astronomical
787 Society, 205 (2016), pp. 345–377.
- 788 [26] M. K. NG, F. WANG, AND X. YUAN, *Inexact alternating direction methods for image recovery*,
789 SIAM Journal on Scientific Computing, 33 (2011), pp. 1643–1668.
- 790 [27] B. O’DONOGHUE AND E. CANDÈS, *Adaptive restart for accelerated gradient schemes*, Foundations
791 of computational mathematics, 15 (2015), pp. 715–732.
- 792 [28] O. PELE AND M. WERMAN, *Fast and robust earth mover’s distances.*, in ICCV, vol. 9, 2009,
793 pp. 460–467.
- 794 [29] T. POCK AND A. CHAMBOLLE, *Diagonal preconditioning for first order primal-dual algorithms*
795 *in convex optimization*, in Computer Vision (ICCV), 2011 IEEE International Conference
796 on, IEEE, 2011, pp. 1762–1769.
- 797 [30] T. POCK, D. CREMERS, H. BISCHOF, AND A. CHAMBOLLE, *An algorithm for minimizing the*
798 *mumford-shah functional*, in Computer Vision, 2009 IEEE 12th International Conference
799 on, IEEE, 2009, pp. 1133–1140.
- 800 [31] J. RASCH AND A. CHAMBOLLE, *Inexact first-order primal-dual algorithms*, arXiv preprint
801 arXiv:1803.10576, (2018).
- 802 [32] L. F. RICHARDSON, *ix. the approximate arithmetical solution by finite differences of physical*
803 *problems involving differential equations, with an application to the stresses in a masonry*
804 *dam*, Phil. Trans. R. Soc. Lond. A, 210 (1911), pp. 307–357.
- 805 [33] R. T. ROCKAFELLAR AND R. J.-B. WETS, *Variational Analysis*, vol. 317, Springer Science &
806 Business Media, 2009.
- 807 [34] Y. SAAD, *Iterative methods for sparse linear systems*, vol. 82, siam, 2003.
- 808 [35] T. VALKONEN, *A primal–dual hybrid gradient method for nonlinear operators with applications*
809 *to mri*, Inverse Problems, 30 (2014), p. 055012.
- 810 [36] Y. WANG, W. YIN, AND J. ZENG, *Global convergence of ADMM in nonconvex nonsmooth*
811 *optimization*, arXiv preprint arXiv:1511.06324, (2015).
- 812 [37] M. YAN AND W. YIN, *Self equivalence of the alternating direction method of multipliers*, in
813 Splitting Methods in Communication, Imaging, Science, and Engineering, Springer, 2016,
814 pp. 165–194.
- 815 [38] X. ZHANG, M. BURGER, AND S. OSHER, *A unified primal-dual algorithm framework based on*
816 *bregman iteration*, Journal of Scientific Computing, 46 (2011), pp. 20–46.
- 817 [39] M. ZHU AND T. CHAN, *An efficient primal-dual hybrid gradient algorithm for total variation*
818 *image restoration*, UCLA CAM Report, 34 (2008).