

INERTIAL, CORRECTED, PRIMAL-DUAL PROXIMAL SPLITTING*

TUOMO VALKONEN†

Abstract. We study inertial versions of primal-dual proximal splitting, also known as the Chambolle–Pock method. Our starting point is the preconditioned proximal point formulation of this method. By adding correctors corresponding to the antisymmetric part of the relevant monotone operator, using a FISTA-style gap unrolling argument, we are able to derive gap estimates instead of merely ergodic gap estimates. Moreover, based on adding a diagonal component to this corrector, we are able to combine strong convexity based acceleration with inertial acceleration. We test our proposed method on image processing and inverse problems, obtaining convergence improvements for sparse Fourier inversion and positron emission tomography.

Key words. inertia, primal-dual, proximal point, Chambolle–Pock, splitting, acceleration

AMS subject classifications. 49M29, 65K10, 65K15, 90C30, 90C47

DOI. 10.1137/18M1182851

1. Introduction. For convex, proper, and lower semicontinuous $G : X \rightarrow \overline{\mathbb{R}}$ and $F^* : Y \rightarrow \overline{\mathbb{R}}$, and a bounded linear operator $K \in \mathcal{L}(X; Y)$ on Hilbert spaces X and Y , we will derive inertial primal-dual optimization methods for the problem

$$(1.1) \quad \min_{x \in X} G(x) + F(Kx).$$

If K is the identity, and F is smooth, a classical algorithm for the iterative solution of (1.1) is the forward-backward splitting method $x^{i+1} := \text{prox}_{\tau G}(x^i - \tau \nabla F(x^i))$, where $\tau L < 1$ for L the Lipschitz factor of ∇F . That is, we take proximal steps with respect to G and gradient steps with respect to F . The proximal step needs to be efficiently realizable, i.e., G needs to be “prox-simple.” If no strong convexity is present, the iterates of the forward-backward splitting generally converge weakly, and the function values at the rate $O(1/N)$. By applying *inertia*, the latter can be improved to $O(1/N^2)$. This in essence consists of rebasing the algorithm at an inertial variable \bar{x}^i :

$$(1.2) \quad x^{i+1} := \text{prox}_{\tau G}(\bar{x}^i - \tau \nabla F(\bar{x}^i)), \quad \text{where} \quad \bar{x}^i := (1 + \alpha_i)x^i - \alpha_i x^{i-1},$$

for suitable inertial parameters $\{\alpha_i\}_{i \in \mathbb{N}}$. In FISTA [3], which itself is an extension of Nesterov’s accelerated gradient descent [20], one would take

$$(1.3) \quad \alpha_{i+1} := \lambda_{i+1}(\lambda_i^{-1} - 1) \quad \text{for} \quad \lambda_{i+1}^{-1} := \sqrt{\lambda_i^{-2} + 1/4} + 1/2.$$

For this scheme, no convergence rates of the iterates themselves are known, although weak convergence can be obtained with small modifications [7]. Several studies have

*Received by the editors April 24, 2018; accepted for publication (in revised form) February 13, 2020; published electronically May 14, 2020.
<https://doi.org/10.1137/18M1182851>

Funding: This research has been supported by the EPSRC First Grant EP/P021298/1, “PARTIAL Analysis of Relations in Tasks of Inversion for Algorithmic Leverage,” as well as by Academy of Finland grants 314701 and 320022.

†Department of Mathematics and Statistics, University of Helsinki, Finland, and ModeMat, Escuela Politécnica Nacional, Quito, Ecuador. Previously: Department of Mathematical Sciences, University of Liverpool, Liverpool L69 3BX, United Kingdom (tuomo.valkonen@iki.fi).

sought to further optimize the inertial parameters; we refer merely to a few of the most recent works [16, 2] and references therein.

If F is nonsmooth, but G is smooth, we can apply forward-backward splitting or FISTA with the roles of the two functions exchanged. However, if K is not the identity, $F \circ K$ is rarely prox-simple, so these methods are not practically applicable. Nevertheless, denoting by F^* the Fenchel conjugate of F , we can reformulate (1.1) as

$$(1.4) \quad \min_{x \in X} \max_{y \in Y} G(x) + \langle Kx, y \rangle - F^*(y).$$

A popular iterative method for this class of problems is the *primal-dual proximal splitting* (PDPS), commonly known as the Chambolle–Pock method [8]. It takes alternate proximal steps with respect to the primal and dual variables x and y :

$$(1.5) \quad \begin{cases} x^{i+1} := \text{prox}_{\tau_i G}(x^i - \tau_i K^* y^i), \\ x_\omega^{i+1} := \omega_i (x^{i+1} - x^i) + x^i, \\ y^{i+1} := \text{prox}_{\sigma_{i+1} F^*}(y^i + \sigma_{i+1} K x_\omega^{i+1}). \end{cases}$$

In the basic version the overrelaxation parameter $\omega_i \equiv 1$, and the primal and dual step lengths $\tau_i \equiv \tau_0$, $\sigma_i \equiv \sigma_0$ with $\tau_0 \sigma_0 \|K\|^2 < 1$. This yields an $O(1/N)$ convergence rate for an ergodic gap functional and weak convergence of the iterates. If G is strongly convex with factor $\gamma > 0$, an accelerated version updates $\tau_{i+1} := \omega_i \tau_i$ and $\sigma_{i+1} := \sigma_i / \omega_i$ for $\omega_i := 1/\sqrt{1+\gamma\tau_i}$. This yields $O(1/N^2)$ rates for the ergodic gap as well as $\|x^N - \hat{x}\|^2$.

Several recent works [26, 11, 9, 10, 18, 1] have applied inertia and closely related overrelaxation [15, 13] to the basic method (1.5). For inertia, writing $u^i = (x^i, y^i)$ and $\bar{u}^i = (\bar{x}^i, \bar{y}^i)$, similarly to (1.2), one rebases the algorithm at $\bar{u}^i := (1 + \alpha_i)u^i - \alpha_i u^{i-1}$ in place of u^i . In [9], $O(1/N)$ convergence of an *ergodic* gap functional is shown for this method. No $O(1/N^2)$ results are known to us, or results for a nonergodic gap or iterates. *In this work, we want to improve upon these convergence rate results, possibly by modifying the algorithm.*

A crucial ingredient for inertia to work in (1.2) is a gap unrolling argument. To demonstrate this argument, we take for simplicity $F = 0$. Then (1.2) implies $q^{i+1} := -\tau^{-1}(x^{i+1} - \bar{x}^i) \in \partial G(x^{i+1})$. Defining the auxiliary sequence $\zeta^{i+1} := \lambda_i^{-1} x^{i+1} - (\lambda_i^{-1} - 1)x^i$, for all $\hat{x} \in X$ one then has¹

$$(1.6) \quad C_0 := \frac{1}{2\tau} \|\zeta^0 - \hat{x}\|^2 \geq -\tau^{-1} \sum_{i=0}^{N-1} \langle \zeta^{i+1} - \zeta^i, \zeta^{i+1} - \hat{x} \rangle = \sum_{i=0}^{N-1} \lambda_i^{-1} \langle q^{i+1}, \zeta^{i+1} - \hat{x} \rangle.$$

If we do not apply inertia, that is, $\lambda_i \equiv 1$, we have $\zeta^{i+1} = x^{i+1}$, so by convexity and Jensen's inequality

$$(1.7) \quad C_0 \geq \sum_{i=0}^{N-1} (G(x^{i+1}) - G(\hat{x})) \geq N (G(\tilde{x}^N) - G(\hat{x})), \quad \text{where} \quad \tilde{x}^N := \frac{1}{N} \sum_{i=0}^{N-1} x^{i+1}.$$

Due to the variable \tilde{x}^N , this $O(1/N)$ estimate is ergodic. If, on the other hand, we update λ_i as in (1.3), we can unroll the ergodicity: Since

¹For the inequality apply Pythagoras' identity to convert the right-hand inner product into norms squared. Then cancel repeated terms and estimate the remaining negative terms. The details in abstract form can also be found in Theorem 2.3.

$$\lambda_i (\zeta^{i+1} - \hat{x}) = \lambda_i (x^{i+1} - \hat{x}) + (1 - \lambda_i) (x^{i+1} - x^i),$$

we estimate from (1.6) by rearrangements and the definition of the subdifferential that

$$\begin{aligned} C_0 &\geq \sum_{i=0}^{N-1} \lambda_i^{-2} \left[\lambda_i \langle q^{i+1}, x^{i+1} - \hat{x} \rangle + (1 - \lambda_i) \langle q^{i+1}, x^{i+1} - x^i \rangle \right] \\ (1.8a) \quad &\geq \sum_{i=0}^{N-1} \lambda_i^{-2} \left[\lambda_i (G(x^{i+1}) - G(\hat{x})) + (1 - \lambda_i) (G(x^{i+1}) - G(x^i)) \right] \\ &= \sum_{i=0}^{N-1} \left[\lambda_i^{-2} (G(x^{i+1}) - G(\hat{x})) - (\lambda_i^{-2} - \lambda_i^{-1}) (G(x^i) - G(\hat{x})) \right]. \end{aligned}$$

Telescoping and the recurrence $\lambda_i^{-2} = \lambda_{i+1}^{-2} - \lambda_{i+1}^{-1}$ established from (1.3) now yield

$$(1.8b) \quad C_0 + \lambda_0^{-2} (1 - \lambda_0) (G(x^0) - G(\hat{x})) \geq \lambda_{N-1}^{-2} (G(x^N) - G(\hat{x})).$$

Since the recurrence also implies that λ_N is of the order $O(1/N^2)$ [3], this gives the improved convergence rate. Similar arguments can be applied to the forward step component F , as we will demonstrate in section 3 in a more general setting.

How could such argumentation be applied to the PDPS (1.5)? It was discovered in [15] that the method can be written as the “preconditioned proximal point method”

$$(1.9) \quad 0 \in H(u^{i+1}) + W_{i+1}^{-1} M_{i+1} (u^{i+1} - u^i)$$

in the space $U := X \times Y$ with the general notation $u = (x, y)$ for the monotone operator $H : U \rightrightarrows U$, the linear preconditioner $M_{i+1} \in \mathcal{L}(U; U)$, and the step length operator $W_{i+1} \in \mathcal{L}(U; U)$ defined as

$$(1.10) \quad H(u) := \begin{pmatrix} \partial G(x) + K^* y \\ \partial F^*(y) - Kx \end{pmatrix}, \quad M_{i+1} := \begin{pmatrix} I & -\tau_i K^* \\ -\sigma_{i+1} \omega_i K & I \end{pmatrix}, \quad W_{i+1} := \begin{pmatrix} \tau_i I & 0 \\ 0 & \sigma_{i+1} I \end{pmatrix}.$$

The overrelaxation parameter ω_i and the step length parameters τ_i and σ_{i+1} are as after (1.5). Clearly $H(u) = \partial \hat{G}(u) + \Gamma u$ for the convex function $\hat{G}(u) := G(x) + F^*(y)$ and an antisymmetric operator Γ . We can thus apply (1.8) to \hat{G} . However, the antisymmetric operator Γ does not arise as a subdifferential, so similar arguments do not apply to it. Indeed, given some inertial parameters $\lambda_i > 0$ and the primal-dual auxiliary sequence $z^{i+1} := \lambda_i^{-1} u^{i+1} - (\lambda_i^{-1} - 1) u^i$, corresponding to ζ^{i+1} above, it does not seem possible to develop a useful estimate out of $\sum_{i=0}^{N-1} \lambda_i^{-1} \langle H(u^{i+1}), z^{i+1} - \hat{u} \rangle$ alone, unless $\lambda_i \equiv 1$. In section 2, we are therefore going to correct the inertial scheme against the antisymmetry of Γ . We do this in the context of general proximal point methods for the solution of the variational inclusion $0 \in H(\hat{u})$. We demonstrate how to test for convergence rates of such methods based on the ideas introduced in [22] for noninertial methods.

We adapt and improve the inertial unrolling argument (1.8) to the setting of this testing theory and corrected inertial methods in section 3. With an eye toward convergence rate proofs, we also develop parameter growth estimates, and briefly demonstrate the theory by application to FISTA. Based on the general results

of sections 2 and 3, we then develop our proposed *inertial, corrected, primal-dual proximal splitting* (IC-PDPS) in section 4. Using the corrector, we will also be able to incorporate strong convexity based acceleration into the inertial method. We finish with conclusions and numerical experience in section 5. Readers wishing to simply implement our proposed method can find it in an explicit and mostly self-contained form near the end in Algorithm 4.1. Only the step length rules need to be taken from a choice of theorems given in the algorithm description.

Notation. We write $\overline{\mathbb{R}} := [-\infty, \infty]$ for the extended reals and $\mathcal{L}(X; Y)$ for the space of bounded linear operators between Hilbert spaces X and Y . The identity operator in any space is I . For $T, S \in \mathcal{L}(X; X)$, we write $T \geq S$ when $T - S$ is positive semidefinite. Also for possibly non-self-adjoint T , we introduce the inner product $\langle x, z \rangle_T := \langle Tx, z \rangle$, and, for positive semidefinite T , the seminorm $\|x\|_T := \sqrt{\langle x, x \rangle_T}$. For a set $A \subset \mathbb{R}$ and a scalar $c \in \mathbb{R}$, we write $A \geq c$ if every element $t \in A$ satisfies $t \geq c$. We write $H : X \rightrightarrows Y$ for H being a set-valued map from X to Y .

2. General inertial methods, correctors. We will now study the application of inertia to general preconditioned proximal point schemes, of which the PDPS is an instance. As we discussed in the introduction, [15] showed that the PDPS (1.5) can be written as solving $0 \in H(u^{i+1}) + W_{i+1}^{-1} M_{i+1} (u^{i+1} - u^i)$ for u^{i+1} with the choices (1.10). To simplify the analysis of accelerated methods, [22, 25] moved the step length operator at the front, rewriting the method with $\tilde{H}_{i+1} := W_{i+1} H$ as an instance of the more general scheme

$$(PP) \quad 0 \in \tilde{H}_{i+1} (u^{i+1}) + M_{i+1} (u^{i+1} - u^i)$$

that can also model forward steps by suitable definitions of \tilde{H}_{i+1} .

In this section, we will study the application of inertia to (PP). We start by formulating a simple extension of (1.2) to (PP). As in (1.3), writing the inertial parameter as $\alpha_{i+1} = \lambda_{i+1}(\lambda_i^{-1} - 1)$, we now take an invertible linear operator Λ_{i+1} as our fundamental inertial parameter. Given an initial iterate $u^0 = \bar{u}^0 \in U$, we then rebase u^i in (PP) to \bar{u}^i to obtain the method

$$(2.1) \quad \begin{cases} 0 \in \tilde{H}_{i+1} (u^{i+1}) + M_{i+1} (u^{i+1} - \bar{u}^i), \\ \bar{u}^{i+1} := u^{i+1} + \Lambda_{i+2} (\Lambda_{i+1}^{-1} - I) (u^{i+1} - u^i). \end{cases}$$

We assume that $\tilde{H}_{i+1} : U \rightrightarrows U$ and $M_{i+1}, \Lambda_{i+1} \in \mathcal{L}(U; U)$ on a Hilbert space U .

Remark 2.1. The operator Λ_{i+1} has the index $i + 1$ off-by-one compared to λ_i in (1.2) and (1.3). This is for consistency with the historical development of the PDPS (1.5) into the form (1.9) or (PP): compare (1.10), where primal step lengths within the step length operator W_{i+1} have index i , and dual step lengths index $i + 1$. This will generally be the case: operator indices agree with dual parameter indices, while primal parameter indices will be one less. We have not reindexed the parameters to maintain the property $\sigma_i \tau_i = \sigma_0 \tau_0$ of the PDPS.

We want to correct for any antisymmetric or otherwise challenging components $\Gamma_{i+1} \in \mathcal{L}(U; U)$ of \tilde{H}_{i+1} . We therefore introduce the *corrector*

$$(2.2) \quad \check{M}_{i+1} := \Gamma_{i+1} (\Lambda_{i+1}^{-1} - I)$$

and modify (2.1) into the *general corrected inertial method*

$$(PP-I) \quad \begin{cases} 0 \in \tilde{H}_{i+1}(u^{i+1}) + M_{i+1}(u^{i+1} - \bar{u}^i) + \check{M}_{i+1}(u^{i+1} - u^i), \\ \bar{u}^{i+1} := u^{i+1} + \Lambda_{i+2}^{-1}(\Lambda_{i+1}^{-1} - I)(u^{i+1} - u^i). \end{cases}$$

In this section, our task is to develop general convergence estimates for (PP-I), which we will then use to prove convergence rates of more specific instances of the general method, in particular the IC-PDPS in section 4. To interpret the main assumption of our abstract convergence estimate, and for later use, we recall the following three-point inequality.

LEMMA 2.2. *Let $F : X \rightarrow \overline{\mathbb{R}}$ be proper, convex, lower semicontinuous with ∇F L -Lipschitz. Then*

$$\langle \nabla F(z), x - \hat{x} \rangle \geq F(x) - F(\hat{x}) - \frac{L}{2} \|x - z\|^2 \quad (\hat{x}, z, x \in X).$$

Proof. Since F has L -Lipschitz gradient, it is smooth in the sense of convex analysis (also known as satisfying the *descent inequality*), $F(z) - F(x) \geq \langle \nabla F(z), z - x \rangle - \frac{L}{2} \|x - z\|^2$. By convexity $F(\hat{x}) - F(z) \geq \langle \nabla F(z), \hat{x} - z \rangle$. Summing these two estimates, we obtain the claim. \square

We will develop our convergence estimates following the testing framework of [22]. The idea introduced there was to pick a suitably designed testing operator $Z_{i+1} \in \mathcal{L}(U; U)$, and then apply the testing functional $u \mapsto \langle uu^{i+1} - \hat{u} \rangle_{Z_{i+1}}$ to both sides of (PP). An almost trivial argument based on a simple assumption on \tilde{H}_{i+1} and Pythagoras' (three-point) identity would then show that $Z_{i+1}M_{i+1}$ forms a local metric that measures convergence rates. However, presently, we cannot in general obtain estimates on the principal sequence $\{u^i\}_{i \in \mathbb{N}}$. Rather, we will obtain estimates on the auxiliary sequence $\{z^i\}_{i \in \mathbb{N}}$, defined through

$$(2.3) \quad z^0 := u^0 \quad \text{and} \quad z^{i+1} := \Lambda_{i+1}^{-1}u^{i+1} - (\Lambda_{i+1}^{-1} - I)u^i \quad (i \in \mathbb{N}).$$

This adds some additional complexity to the main condition (2.4) of the next theorem. We will motivate the condition after the proof.

THEOREM 2.3. *On a Hilbert space U , for $i = 0, \dots, N-1$, let $\tilde{H}_{i+1} : U \rightrightarrows U$ as well as $M_{i+1}, Z_{i+1}, \Gamma_{i+1}, \Lambda_i \in \mathcal{L}(U; U)$ with Λ_i invertible. Given an initial iterate $u^0 = \bar{u}^0 \in U$, let $\{u^{i+1}\}_{i \in \mathbb{N}}$ be defined through the solution of (PP-I), and the auxiliary sequence $\{z^i\}_{i \in \mathbb{N}}$ by (2.3). Suppose, for $i = 1, \dots, N-1$, that $Z_{i+1}M_{i+1} \geq 0$ is self-adjoint, and for some $\hat{u} \in U$ and a placeholder real value $\mathcal{V}_{i+1}(\hat{u}) \in \mathbb{R}$ that*

$$(2.4) \quad \begin{aligned} & \left\langle \tilde{H}_{i+1}(u^{i+1}) - \Gamma_{i+1}(u^{i+1} - \hat{u}), z^{i+1} - \hat{u} \right\rangle_{\Lambda_{i+1}^* Z_{i+1}} \\ & \geq \mathcal{V}_{i+1}(\hat{u}) - \frac{1}{2} \|z^{i+1} - z^i\|_{\Lambda_{i+1}^* Z_{i+1} M_{i+1} \Lambda_{i+1}}^2 \end{aligned}$$

and

$$(2.5) \quad \Lambda_{i+1}^* Z_{i+1} (M_{i+1} \Lambda_{i+1} + 2\Gamma_{i+1}) \geq \Lambda_{i+2}^* Z_{i+2} M_{i+2} \Lambda_{i+2}.$$

Then

$$(2.6) \quad \frac{1}{2} \|z^N - \hat{u}\|_{\Lambda_{N+1}^* Z_{N+1} M_{N+1} \Lambda_{N+1}}^2 + \sum_{i=0}^{N-1} \mathcal{V}_{i+1}(\hat{u}) \leq \frac{1}{2} \|z^0 - \hat{u}\|_{\Lambda_1^* Z_1 M_1 \Lambda_1}^2 \quad (N \geq 1).$$

Proof. Application of $\langle \cdot, Z_{i+1}^* \Lambda_{i+1} (z^{i+1} - \hat{u}) \rangle$ to the main inclusion of (PP-I) yields for some $q^{i+1} \in \tilde{H}_{i+1}(u^{i+1})$ that

$$(2.7) \quad 0 = \langle q^{i+1} + M_{i+1}(u^{i+1} - \bar{u}^i) + \check{M}_{i+1}(u^{i+1} - u^i), \Lambda_{i+1}(z^{i+1} - \hat{u}) \rangle_{Z_{i+1}}.$$

By (2.3) and (2.2), we deduce that

$$\check{M}_{i+1}(u^{i+1} - u^i) = \Gamma_{i+1}(\Lambda_{i+1}^{-1} - I)(u^{i+1} - u^i) = \Gamma_{i+1}(z^{i+1} - u^{i+1}).$$

By the definition \bar{u}^i in (PP-I), and of the auxiliary sequence $\{z^{i+1}\}_{i \in \mathbb{N}}$ in (2.3), taking $u^{-1} := u^0$, moreover,

$$(2.8) \quad \begin{aligned} \Lambda_{i+1}(z^{i+1} - z^i) &= u^{i+1} - (I - \Lambda_{i+1})u^i - \Lambda_{i+1}[\Lambda_i^{-1}u^i - (\Lambda_i^{-1} - I)u^{i-1}] \\ &= u^{i+1} - [I - \Lambda_{i+1} + \Lambda_{i+1}\Lambda_i^{-1}]u^i - \Lambda_{i+1}(I - \Lambda_i^{-1})u^{i-1} \\ &= u^{i+1} - \bar{u}^i. \end{aligned}$$

Therefore, we transform (2.7) into

$$(2.9) \quad 0 = \langle q^{i+1} + M_{i+1}\Lambda_{i+1}(z^{i+1} - z^i) + \Gamma_{i+1}(z^{i+1} - u^{i+1}), \Lambda_{i+1}(z^{i+1} - \hat{u}) \rangle_{Z_{i+1}}.$$

Writing for brevity $A := \Lambda_{i+1}^* Z_{i+1} M_{i+1} \Lambda_{i+1}$, which by assumption is self-adjoint and positive semidefinite, the standard three-point formula or Pythagoras' identity states

$$\langle z^{i+1} - z^i, z^{i+1} - \hat{u} \rangle_A = \frac{1}{2} \|z^{i+1} - z^i\|_A^2 - \frac{1}{2} \|z^i - \hat{u}\|_A^2 + \frac{1}{2} \|z^{i+1} - \hat{u}\|_A^2.$$

We therefore transform (2.9) into

$$\begin{aligned} 0 &= \langle q^{i+1} + \Gamma_{i+1}(z^{i+1} - u^{i+1}), z^{i+1} - \hat{u} \rangle_{\Lambda_{i+1}^* Z_{i+1}} + \frac{1}{2} \|z^{i+1} - z^i\|_{\Lambda_{i+1}^* Z_{i+1} M_{i+1} \Lambda_{i+1}}^2 \\ &\quad - \frac{1}{2} \|z^i - \hat{u}\|_{\Lambda_{i+1}^* Z_{i+1} M_{i+1} \Lambda_{i+1}}^2 + \frac{1}{2} \|z^{i+1} - \hat{u}\|_{\Lambda_{i+1}^* Z_{i+1} M_{i+1} \Lambda_{i+1}}^2. \end{aligned}$$

Using (2.4), we obtain

$$\begin{aligned} 0 &\geq \mathcal{V}_{i+1}(\hat{u}) - \frac{1}{2} \|z^i - \hat{u}\|_{\Lambda_{i+1}^* Z_{i+1} M_{i+1} \Lambda_{i+1}}^2 \\ &\quad + \frac{1}{2} \|z^{i+1} - \hat{u}\|_{\Lambda_{i+1}^* Z_{i+1} M_{i+1} \Lambda_{i+1}}^2 + \langle \Gamma_{i+1}(z^{i+1} - \hat{u}), z^{i+1} - \hat{u} \rangle_{\Lambda_{i+1}^* Z_{i+1}}. \end{aligned}$$

Using (2.5) and summing over $i = 0, \dots, N-1$ establishes (2.6). \square

Remark 2.4. Consider $\Lambda_{i+1} = Z_{i+1} = M_{i+1} = I$. Then (2.4) reads

$$(2.10) \quad \langle \tilde{H}_{i+1}(u^{i+1}), u^{i+1} - \hat{u} \rangle \geq \mathcal{V}_{i+1}(\hat{u}) + \langle u^{i+1} - \hat{u}, u^{i+1} - \hat{u} \rangle_{\Gamma_{i+1}} - \frac{1}{2} \|u^{i+1} - u^i\|^2.$$

With $\tau L \leq 1$ and $\Gamma_{i+1} = 0$, take first $\tilde{H}_{i+1}(u) = \tau \nabla F(u^i)$, fixing ∇F to be evaluated at the previous iteration to obtain a gradient descent method. Then it is easy to see how this estimate with $\mathcal{V}_{i+1}(\hat{u}) = \tau[F(u^{i+1}) - F(\hat{u})]$ follows from Lemma 2.2. Thus \mathcal{V}_{i+1} measures function value differences. Similarly, if $\tilde{H}_{i+1}(u) = \tau \partial G(u)$ for nonsmooth but (γ -strongly) convex G , we can take $\Gamma_{i+1} = \tau \gamma I$ in (2.10). We can also combine $\tilde{H}_{i+1}(u) = \tau[\partial G(u) + \nabla F(u^i)]$ to obtain forward-backward splitting. In other words, (2.4) is an operator-relative inertia-aware convexity and smoothness condition, where the variable $\mathcal{V}_{i+1}(\hat{u})$ can be used to model function value and other

gap estimates. We will need this operator-relativity to apply distinct inertial and testing parameters on the primal and dual variables of the PDPS; compare the block structure of (1.10).

Minding this interpretation of (2.4), the claim (2.6) of the theorem with additional positivity and growth assumptions can thus be used to show the convergence of the sum of the value estimates $\mathcal{V}_{i+1}(\hat{u})$ to zero, as well as $z^N \rightarrow \hat{u}$. Our task in the following is to obtain that growth and to unroll the sum into a simple estimate.

3. Unrolling and parameter growth estimates. To develop the IC-PDPS method, we will seek to satisfy the conditions of Theorem 2.3 for an algorithm inspired by the proximal point interpretation of the PDPS. Before we do this, in this section, we will prove general inertial unrolling arguments (subsection 3.2), refining (1.8) to the testing framework. We also prove parameter growth estimates with an eye toward converge rate proofs (subsection 3.3), and demonstrate how our corrector term allows combining inertia with strong convexity based acceleration (3.4). We finish by applying these estimates to the FISTA to demonstrate how it fits into our overall approach (subsection 3.4). This also demonstrates how our approach works without the additional challenges of the primal-dual setup.

3.1. Scalar parameter choices. To place the general estimates that make up the major part of this section into context, we start by specializing Theorem 2.3 to $\Lambda_{i+1} = \lambda_i I$, $Z_{i+1} = \phi_i I$, $W_{i+1} = \tau_i I$, $M_{i+1} = I$, and $\Gamma_{i+1} := \gamma \tau_i I$ for some scalars $\lambda_i, \phi_i, \tau_i > 0$, and $\gamma \geq 0$. Also taking $\tilde{H}_{i+1}(x) := \tau_i(\partial G(x) + \nabla F(x^i))$, we immediately rewrite (PP-I), with change of symbol² u into x , as

$$\begin{aligned} \text{(PP-i)} \quad & \begin{cases} 0 \in \tau_i [\partial G(x^{i+1}) + \nabla F(x^i)] + (x^{i+1} - \bar{x}^i) + \gamma \tau_i (\lambda_i^{-1} - 1) (x^{i+1} - x^i), \\ \bar{x}^{i+1} := x^{i+1} + \lambda_{i+1} (\lambda_i^{-1} - 1) (x^{i+1} - x^i). \end{cases} \end{aligned}$$

Moreover, with the change of symbol² of z into ζ , the auxiliary sequence defined in (2.3) becomes

$$(3.1) \quad \zeta^0 := x^0 \quad \text{and} \quad \zeta^{i+1} := \lambda_i^{-1} x^{i+1} - (\lambda_i^{-1} - 1) x^i \quad (i \in \mathbb{N}).$$

Immediately, Theorem 2.3 specializes into the following:

COROLLARY 3.1. *On a Hilbert space X , let $G : X \rightarrow \overline{\mathbb{R}}$ and $F : X \rightarrow \mathbb{R}$ be convex, proper, and lower semicontinuous, with F differentiable. Let $\lambda_i, \phi_i, \tau_i > 0$, ($i = 0, \dots, N-1$), and $\gamma \geq 0$. For an initial iterate $x^0 = \bar{x}^0 \in X$, let $\{x^{i+1}\}_{i=0}^{N-1}$ be generated by (PP-i) and the auxiliary sequence $\{\zeta^i\}_{i=0}^{N-1}$ by (3.1). For each $i = 0, \dots, N-1$, suppose for some $\hat{x} \in X$ and a placeholder value $\mathcal{V}_{i+1}(\hat{x}) \in \mathbb{R}$, we have the estimate*

$$(3.2) \quad \phi_i \lambda_i \tau_i \langle \partial G(x^{i+1}) + \nabla F(x^i) - \gamma(x^{i+1} - \hat{x}), \zeta^{i+1} - \hat{x} \rangle \geq \mathcal{V}_{i+1}(\hat{x}) - \frac{\lambda_i^2 \phi_i}{2} \|\zeta^{i+1} - \zeta^i\|^2$$

and the inequality

$$(3.3) \quad \lambda_{i+1}^2 \phi_{i+1} \leq \lambda_i^2 \phi_i (1 + 2\gamma \lambda_i^{-1} \tau_i).$$

²We reserve the symbols u and z for the abstract (section 2) and primal-dual (section 4) problems. In the latter we take primal-dual pairs $u = (x, y)$ and $z = (\zeta, \eta)$, so the primal variables match the symbols of this section.

Then

$$(3.4) \quad \frac{\phi_N \lambda_N^2}{2} \|\zeta^N - \hat{x}\|^2 + \sum_{i=0}^{N-1} \mathcal{V}_{i+1}(\hat{x}) \leq \frac{\phi_0 \lambda_0^2}{2} \|\zeta^0 - \hat{x}\|^2 \quad (N \geq 1).$$

3.2. Inertial unrolling. We start by refining the proximal step inertial unrolling argument (1.8). We will in subsection 3.3 see that the recurrence inequality (3.5) assumed by the next lemma generalizes the recurrence $\lambda_i^{-2} = \lambda_{i+1}^{-2} - \lambda_{i+1}^{-1}$ from the introduction, satisfied by the FISTA.

LEMMA 3.2. *Let $G : X \rightarrow \overline{\mathbb{R}}$ be convex, proper, and lower semicontinuous. Suppose $\lambda_i \in (0, 1]$ and $\phi_i, \tau_i > 0$ satisfy the recurrence inequality*

$$(3.5) \quad \phi_{i+1} \tau_{i+1} (1 - \lambda_{i+1}) \leq \phi_i \tau_i \quad (i = 0, \dots, N-1).$$

For any given $\{x^i\}_{i=0}^N$, let the auxiliary variables $\{\zeta^i\}_{i=0}^N$ be generated by (3.1). Assume $\partial G(x^{i+1})$ to be nonempty for $i = 0, \dots, N-1$, and $\hat{x} \in [\partial G]^{-1}(0)$. Then

$$(3.6) \quad \sum_{i=0}^{N-1} \inf_{q^{i+1} \in \partial G(x^{i+1})} \phi_i \tau_i \lambda_i \langle q^{i+1}, \zeta^{i+1} - \hat{x} \rangle \geq \phi_{N-1} \tau_{N-1} (G(x^N) - G(\hat{x})) \\ - \phi_0 \tau_0 (1 - \lambda_0) (G(x^0) - G(\hat{x})).$$

Proof. For all $i = 0, \dots, N-1$, pick $q^{i+1} \in \partial G(x^{i+1})$, and define

$$s_N^G = \sum_{i=0}^{N-1} \phi_i \tau_i \lambda_i \langle q^{i+1}, \zeta^{i+1} - \hat{x} \rangle.$$

Then we need to show that

$$(3.7) \quad s_N^G \geq \phi_{N-1} \tau_{N-1} (G(x^N) - G(\hat{x})) - \phi_0 \tau_0 (1 - \lambda_0) (G(x^0) - G(\hat{x})).$$

Observe that the auxiliary variables $\{\zeta^{i+1}\}_{i=0}^{N-1}$ satisfy

$$(3.8) \quad \lambda_i (\zeta^{i+1} - \hat{x}) = \lambda_i (x^{i+1} - \hat{x}) + (1 - \lambda_i) (x^{i+1} - x^i).$$

With this and the convexity of G , we estimate

$$(3.9) \quad s_N^G = \sum_{i=0}^{N-1} \phi_i \tau_i \left[\lambda_i \langle q^{i+1}, x^{i+1} - \hat{x} \rangle + (1 - \lambda_i) \langle q^{i+1}, x^{i+1} - x^i \rangle \right] \\ \geq \sum_{i=0}^{N-1} \phi_i \tau_i \left[\lambda_i (G(x^{i+1}) - G(\hat{x})) + (1 - \lambda_i) (G(x^{i+1}) - G(x^i)) \right] \\ = \sum_{i=0}^{N-1} \left[\phi_i \tau_i (G(x^{i+1}) - G(\hat{x})) - \phi_i \tau_i (1 - \lambda_i) (G(x^i) - G(\hat{x})) \right].$$

Since $G(x^i) \geq G(\hat{x})$, the recurrence inequality (3.5) together with a telescoping argument now gives (3.7). \square

We can also include a forward step in the unrolling argument.

LEMMA 3.3. Let $G, F : X \rightarrow \overline{\mathbb{R}}$ be convex, proper, and lower semicontinuous. Suppose F has L -Lipschitz gradient and that $\lambda_i \in (0, 1]$ and $\phi_i, \tau_i > 0$ satisfy the recurrence inequality (3.5) for $i = 0, \dots, N-1$. For any given $\{x^i\}_{i=0}^N$, let the auxiliary variables $\{\zeta^i\}_{i=0}^N$ be generated by (3.1). Assume $\partial G(x^{i+1})$ to be nonempty for all $i = 0, \dots, N-1$ and that $\hat{x} \in [\partial G + \nabla F]^{-1}(0)$. Then

$$(3.10) \quad \sum_{i=0}^{N-1} \inf_{q^{i+1} \in \partial G(x^{i+1})} \left[\phi_i \tau_i \lambda_i \langle q^{i+1} + \nabla F(\bar{x}^i), \zeta^{i+1} - \hat{x} \rangle + \frac{\phi_i \tau_i \lambda_i^2 L}{2} \|\zeta^{i+1} - \zeta^i\|^2 \right] \\ \geq \phi_{N-1} \tau_{N-1} \left[(G+F)(x^N) - (G+F)(\hat{x}) \right] - \phi_0 \tau_0 (1 - \lambda_0) \left[(G+F)(x^0) - (G+F)(\hat{x}) \right].$$

Proof. Similarly to (2.8), $\lambda_i(\zeta^{i+1} - \zeta^i) = (x^{i+1} - \bar{x}^i)$. By (3.8) and Lemma 2.2, therefore

$$s_N^F := \sum_{i=0}^{N-1} \left[\phi_i \tau_i \lambda_i \langle \nabla F(\bar{x}^i), \zeta^{i+1} - \hat{x} \rangle + \frac{\phi_i \tau_i \lambda_i^2 L}{2} \|\zeta^{i+1} - \zeta^i\|^2 \right] \\ = \sum_{i=0}^{N-1} \phi_i \tau_i \left[\lambda_i \langle \nabla F(\bar{x}^i), x^{i+1} - \hat{x} \rangle + (1 - \lambda_i) \langle \nabla F(\bar{x}^i), x^{i+1} - x^i \rangle + \frac{L}{2} \|x^{i+1} - \bar{x}^i\|^2 \right] \\ \geq \sum_{i=0}^{N-1} \phi_i \tau_i \left[\lambda_i (F(x^{i+1}) - F(\hat{x})) + (1 - \lambda_i) (F(x^{i+1}) - F(x^i)) \right] \\ = \sum_{i=0}^{N-1} \left[\phi_i \tau_i (F(x^{i+1}) - F(\hat{x})) - \phi_i \tau_i (1 - \lambda_i) (F(x^i) - F(\hat{x})) \right].$$

Picking $q^{i+1} \in \partial G(x^{i+1})$, ($i = 0, \dots, N-1$), and summing with the estimate (3.9) for G , we deduce

$$s_N^G + s_N^F \geq \sum_{i=0}^{N-1} \phi_i \tau_i \left[\left[(G+F)(x^{i+1}) - (G+F)(\hat{x}) \right] - (1 - \lambda_i) \left[(G+F)(x^i) - (G+F)(\hat{x}) \right] \right].$$

Since $(G+F)(x^i) \geq (G+F)(\hat{x})$, the recurrence inequality (3.5) together with a telescoping argument now gives the claim. \square

3.3. Parameter growth estimates. As suggested by the unrolled estimates (3.6) and (3.10), we want to make $\phi_{N-1} \tau_{N-1}$ grow as fast as possible while satisfying (3.5) and $\lambda_i \in (0, 1]$. We now develop such estimates through a series of lemmas. The first of these lemmas with $\epsilon = 0$ is the FISTA rate argument [3, Lemma 4.3].

LEMMA 3.4. With $\lambda_0 = 1$, suppose $\lambda_i^{-2} - \epsilon \lambda_i^{-1} = \lambda_{i+1}^{-2} - \lambda_{i+1}^{-1}$ for some $\epsilon \in [-1, 1]$ and all $i = 0, \dots, N-1$. Then $\{\lambda_i^{-1}\}_{i \in \mathbb{N}}$ is nondecreasing, $\lambda_N^{-1} \geq 1 + (1 - \epsilon)N/2$, and we equivalently define λ_{i+1} through

$$(3.11) \quad \lambda_{i+1} = \frac{2}{1 + \sqrt{1 + 4(\lambda_i^{-2} - \epsilon \lambda_i^{-1})}}.$$

Proof. The update (3.11) is a simple solution of the quadratic equation $\lambda_i^{-2} - \epsilon \lambda_i^{-1} = \lambda_{i+1}^{-2} - \lambda_{i+1}^{-1}$. The latter also rearranges as

$$(3.12) \quad \lambda_{i+1}^{-2} - \lambda_{i+1}^{-1} = \lambda_i^{-2} - \lambda_i^{-1} + (1 - \epsilon) \lambda_i^{-1}.$$

This shows that $\{t_i\}_{i \in \mathbb{N}}$ for $t_i := \lambda_i^{-2} - \lambda_i^{-1}$ is nondecreasing. Since $\lambda_i^{-1} = \frac{1}{2}(1 + \sqrt{1 + 4t_i})$, we obtain the claim that $\{\lambda_i^{-1}\}_{i \in \mathbb{N}}$ is nondecreasing.

We still need to prove the rate-of-growth claim. Since $\lambda_0 = 1$, (3.12) also yields

$$\lambda_N^{-2} - \lambda_N^{-1} = \sum_{i=0}^{N-1} (1 - \epsilon) \lambda_i^{-1}.$$

Let us make the inductive assumption that $\lambda_i^{-1} \geq 1 + (1 - \epsilon)i/2$ for $i = 0, \dots, N - 1$. Clearly this holds for $N = 0$ by the choice $\lambda_0 = 1$, taking care of the inductive base. For the inductive step, we get from above and $\epsilon \in [-1, 1]$ that

$$\lambda_N^{-2} - \lambda_N^{-1} \geq (1 - \epsilon)N + \frac{(1 - \epsilon)^2}{4} N(N - 1) \geq \frac{1 - \epsilon}{2} N + \left(\frac{1 - \epsilon}{2} \right)^2 N^2.$$

This quadratic inequality together with $\lambda_N > 0$ implies

$$\lambda_N^{-1} \geq \frac{1 + \sqrt{1 + 2(1 - \epsilon)N + (1 - \epsilon)^2 N^2}}{2} = 1 + \frac{1 - \epsilon}{2} N,$$

which verifies the inductive step and establishes the claim. \square

LEMMA 3.5. *The sequence $\{\phi_i \tau_i\}_{i \in \mathbb{N}}$ is nondecreasing and the conditions (3.3) and (3.5) hold, more precisely*

$$(3.13) \quad \lambda_{i+1}^2 \phi_{i+1} = \lambda_i^2 \phi_i (1 + 2\gamma \lambda_i^{-1} \tau_i) \quad \text{and} \quad \phi_{i+1} \tau_{i+1} (1 - \lambda_{i+1}) = (1 - \epsilon \lambda_i) \phi_i \tau_i$$

for all $i \in \mathbb{N}$ for some $\epsilon \in [0, 1]$ in the following cases:

- (i) *If $\gamma = 0$ and we take $\tau_i \equiv \tau$ for any $\tau > 0$; $\phi_i := \lambda_i^{-2}$; $\phi_0 = \lambda_0 = 1$, and update λ_{i+1} for any $\epsilon \in [0, 1]$ according to (3.11). Then also*

$$\phi_N \tau_N \geq (1 - \epsilon)^2 N^2 \tau / 4 \quad \text{and} \quad \lambda_N^2 \phi_N = 1 \quad (N \in \mathbb{N}).$$

- (ii) *If $\gamma > 0$ and we take $\lambda_i \equiv \lambda \in (0, 1)$ and $\tau_i \equiv \tau := \lambda^2 / [2\gamma(1 - \lambda)]$ constants, and $\phi_{i+1} = c\phi_i$ with $c := (1 - \epsilon\lambda) / (1 - \lambda) > 1$ for any $\epsilon \in [0, 1]$ and $\phi_0 > 0$. Then also*

$$\phi_N \tau_N \geq \phi_0 \tau c^N \quad \text{and} \quad \lambda_N^2 \phi_N \geq \lambda^2 \phi_0 c^N \quad (N \in \mathbb{N}).$$

- (iii) *If we are constrained to have $\phi_i = c_0 \tau_i^{-2}$ for some constant $c_0 > 0$, and with $\lambda_0 = 1$, $\tau_0 > 0$ and $\epsilon \in [0, 1]$ update*

$$(3.14) \quad \tau_{i+1} = \frac{1 - \lambda_{i+1}}{1 - \epsilon \lambda_i} \tau_i \quad \text{and} \quad \lambda_{i+1} = \frac{\sqrt{\lambda_i^2 + 2\gamma \lambda_i \tau_i}}{1 - \epsilon \lambda_i + \sqrt{\lambda_i^2 + 2\gamma \lambda_i \tau_i}}.$$

Then, for some constants $c, c' > 0$, for all $N \in \mathbb{N}$, also

$$\begin{aligned} \phi_N \tau_N &\geq c' N^2 & \text{and} & & \lambda_N^2 \phi_N &\geq c N^2 & (\gamma > 0), \\ \phi_N \tau_N &\geq (1 - \epsilon) \tau_0^{-1} N & \text{and} & & \lambda_N^2 \phi_N &= c_0 \tau_0^{-2} & (\gamma = 0). \end{aligned}$$

The choice $\epsilon = 0$ in (iii) would be the simplest and also optimal in the sense that both (3.3) and (3.5) would hold as equalities. However, we will see that a nonzero choice performs significantly better in practice—with the same asymptotic guarantees.

Proof. It is clear that the inequalities (3.3) and (3.5) follow from (3.13) and $\epsilon \geq 0$.

(i) Since $\gamma = 0$, the first part of (3.13) holds when $\phi_i \lambda_i^2 = \phi_0 \lambda_0^2$. This follows from our choices $\phi_0 = \lambda_0 = 1$ and $\phi_i = \lambda_i^{-2}$. Inserting $\phi_i = \lambda_i^{-2}$ and $\tau_i \equiv \tau$, the second part of (3.13) reduces to $\lambda_i^{-2} - \epsilon \lambda_i^{-1} = \lambda_{i+1}^{-2} - \lambda_{i+1}^{-1}$. Lemma 3.4 shows that $\phi_N \tau_N = \lambda_N^{-2} \tau \geq \tau(1 - \epsilon)^2 N^2 / 4$. This is the claimed estimate. Moreover, since $\{\lambda_i^{-1}\}_{i \in \mathbb{N}}$ is nondecreasing by Lemma 3.4, we see that $\{\phi_i \tau_i\}_{i \in \mathbb{N}}$ is nondecreasing.

(ii) The second part of (3.13) agrees with the chosen update rule for ϕ_{i+1} . Inserting this rule into the first part of (3.13) and using the fact that also $\tau_i \equiv \tau$, we see the latter to be satisfied if $(1 - \lambda)(1 + 2\gamma\lambda^{-1}\tau) = 1$. This is satisfied by our chosen $\tau = \lambda^2/[2\gamma(1 - \lambda)]$. Since τ and λ are constants, the claimed growth estimates follow from $\phi_N = \phi_0 c^N$. Clearly $\{\phi_i \tau_i\}_{i \in \mathbb{N}} = \{\phi_0 \tau c^i\}_{i \in \mathbb{N}}$ is nondecreasing.

(iii) Finally, with $\phi_i = c_0 \tau_i^{-2}$, (3.13) holds if

$$(3.15) \quad \lambda_{i+1}^2 \tau_{i+1}^{-2} = \lambda_i^2 \tau_i^{-2} (1 + 2\gamma \lambda_i^{-1} \tau_i) \quad \text{and} \quad \tau_{i+1}^{-1} (1 - \lambda_{i+1}) = (1 - \epsilon \lambda_i) \tau_i^{-1}.$$

The latter agrees with our update rule for τ_{i+1} . Inserting this, the former holds if $\lambda_{i+1} (1 - \epsilon \lambda_i) = (1 - \lambda_{i+1}) \sqrt{\lambda_i^2 + 2\gamma \lambda_i \tau_i}$. This is satisfied by our update rule for λ_{i+1} .

To derive the growth estimates, suppose first that $\gamma > 0$. With $\theta_i := \lambda_i \tau_i^{-1} = c_0^{-1/2} \lambda_i \phi_i^{1/2}$, the first part of (3.15) reads $\theta_{i+1}^2 = \theta_i^2 (1 + 2\gamma \theta_i)$. Thus $\{\theta_i\}_{i \in \mathbb{N}}$ is nondecreasing, moreover, the recurrence is of the same form as the standard acceleration rule for the PDPS, where we would have $\phi_i = c_0 \tau_i^{-2}$ in place of θ_i ; compare section 1 and [8, 25]. Hence $\theta_i \geq ci$ for some constant $c > 0$. Since $\lambda_i^2 \phi_i = c_0 \theta_i^2$, this gives one of the claimed rates. From the second part of (3.15),

$$(3.16) \quad \tau_{i+1}^{-1} = (1 - \epsilon \lambda_i) \tau_i^{-1} + \lambda_{i+1} \tau_{i+1}^{-1} = \tau_i^{-1} - \epsilon \theta_i + \theta_{i+1}.$$

Repeating this recursively, since $\theta_0 = \tau_0^{-1}$, for some constant $c' > 0$,

$$(3.17) \quad \tau_N^{-1} = \tau_0^{-1} - \epsilon \theta_0 + \theta_N + \sum_{i=1}^{N-1} (1 - \epsilon) \theta_i \geq (1 - \epsilon) \tau_0^{-1} + cN + \sum_{i=0}^{N-1} (1 - \epsilon) ci \geq c' N^2.$$

Therefore also $\phi_N \tau_N = c_0 \tau_N^{-1}$ has the claimed growth estimate.

Suppose then that $\gamma = 0$. We obtain (3.16) as above, however now with constant $\theta_i \equiv \lambda_i \tau_i^{-1} = \lambda_0 \tau_0^{-1} = \tau_0^{-1}$. Therefore, arguing as in (3.17) only yields $\tau_N^{-1} \geq (1 - \epsilon) \tau_0^{-1} N$. This is again the claimed growth estimate.

Finally, since $\{\theta_i\}_{i \in \mathbb{N}}$ is in both cases ($\gamma = 0$ or $\gamma > 0$) nondecreasing, we see from (3.16) that $\{\tau_i^{-1}\}_{i \in \mathbb{N}}$ and consequently $\{\phi_i \tau_i\}_{i \in \mathbb{N}} = \{c_0 \tau_i^{-1}\}_{i \in \mathbb{N}}$ are nondecreasing. \square

3.4. Combining inertia with strong convexity. Let $G : X \rightarrow \overline{\mathbb{R}}$ be proper, lower semicontinuous, and (strongly) convex with parameter $\gamma \geq 0$. We now demonstrate with a simple inertial proximal point method, (PP-i) with $F = 0$, how we are able to incorporate strong convexity based acceleration with inertia. We do this by considering the convex functions

$$(3.18) \quad G_\gamma(x; \hat{x}) := G(x) - \frac{\gamma}{2} \|x - \hat{x}\|^2.$$

Indeed, $0 \in \partial G(\hat{x})$ if and only if $0 \in \partial G_\gamma(\hat{x}; \hat{x})$ with $\partial G_\gamma(x; \hat{x}) = \partial G(x) - \gamma(x - \hat{x})$. Lemma 3.2 applied to $G_\gamma(\cdot; \hat{x})$ thus shows

$$(3.19) \quad \sum_{i=0}^{N-1} \mathcal{V}_{i+1}(\hat{x}) \geq \phi_{N-1} \tau_{N-1} [G_\gamma(x^N; \hat{x}) - G_\gamma(\hat{x}; \hat{x})] - \phi_0 \tau_0 (1 - \lambda_0) [G_\gamma(x^0; \hat{x}) - G_\gamma(\hat{x}; \hat{x})]$$

for

$$\mathcal{V}_{i+1}(\hat{x}) := \inf_{q^{i+1} \in \partial G(x^{i+1})} \phi_i \tau_i \lambda_i \langle q^{i+1} - \gamma(x^{i+1} - \hat{x}), \zeta^{i+1} - \hat{x} \rangle.$$

This choice of $\mathcal{V}_{i+1}(\hat{x})$ by definition verifies (3.2). It therefore follows from Corollary 3.1 with $F \equiv 0$ that

$$\begin{aligned} & \frac{\phi_N \lambda_N^2}{2} \|\zeta^N - \hat{x}\|^2 + \phi_{N-1} \tau_{N-1} [G_\gamma(x^N; \hat{x}) - G_\gamma(\hat{x}; \hat{x})] \\ & \leq \frac{\phi_0 \lambda_0^2}{2} \|\zeta^0 - \hat{x}\|^2 + \phi_0 \tau_0 (1 - \lambda_0) [G_\gamma(x^0; \hat{x}) - G_\gamma(\hat{x}; \hat{x})] \quad (N \geq 1). \end{aligned}$$

If we choose our parameters according to Lemma 3.5(ii), then $\phi_N \tau_N$ and $\lambda_N^2 \phi_N$ grow exponentially. Crucially $G_\gamma(x^N; \hat{x}) \geq G_\gamma(\hat{x}; \hat{x})$, so this implies linear convergence of $G_\gamma(x^N; \hat{x}) \rightarrow G_\gamma(\hat{x}; \hat{x})$ and of the auxiliary sequence $\zeta^N \rightarrow \hat{x}$.

3.5. Connections. We now discuss how our results relate to known algorithms.

Example 3.6 (FISTA). Let $G : X \rightarrow \overline{\mathbb{R}}$ and $F : X \rightarrow \mathbb{R}$ be convex, proper, and lower semicontinuous with ∇F existing and L -Lipschitz. If $\gamma = 0$, take $\tau_i \equiv \tau \in (0, 1/L]$, and λ_{i+1} by (3.11) for $\lambda_0 = 1$ and $\epsilon = 0$. Then given initial iterates $\bar{x}^0 = x^0$, (PP-i) becomes the inertial forward-backward splitting or FISTA

$$\begin{cases} x^{i+1} := \text{prox}_{\tau G}(\bar{x}^i - \tau \nabla F(\bar{x}^i)), \\ \bar{x}^{i+1} := x^{i+1} + \lambda_{i+1} (\lambda_i^{-1} - 1) (x^{i+1} - x^i). \end{cases}$$

We have $G(x^N) + F(x^N) \rightarrow G(\hat{x}) + F(\hat{x})$ at the rate $O(1/N^2)$ for any minimizer \hat{x} of $G + F$.

Demonstration. The algorithm is clear from (PP-i). Lemma 3.5(i) shows that (3.3) and (3.5) are satisfied, $\phi_N \tau_N \geq \tau N^2/4$, and $\lambda_N^2 \phi_N = 1$. Clearly (3.2) holds with

$$\begin{aligned} \mathcal{V}_{i+1}(\hat{x}) &:= \inf_{q^{i+1} \in \partial G(x^{i+1})} \phi_i \tau_i \lambda_i \langle q^{i+1} + \nabla F(\bar{x}^i) - \gamma(x^{i+1} - \hat{x}), \zeta^{i+1} - \hat{x} \rangle \\ &\quad + \frac{\lambda_i^2 \phi_i}{2} \|\zeta^{i+1} - \zeta^i\|^2. \end{aligned}$$

Using in (3.10) the fact that $\tau_i L \leq 1$, similarly to subsection 3.4, we may therefore refer to Lemma 3.3 and Corollary 3.1 to verify the estimate

$$\begin{aligned} (3.20) \quad & \frac{\phi_N \lambda_N^2}{2} \|\zeta^N - \hat{x}\|^2 + \phi_{N-1} \tau_{N-1} [(\hat{G}_\gamma + F)(x^N) - (\hat{G}_\gamma + F)(\hat{x})] \\ & \leq \frac{\phi_0 \lambda_0^2}{2} \|\zeta^0 - \hat{x}\|^2 + \phi_0 \tau_0 (1 - \lambda_0) [(\hat{G}_\gamma + F)(x^0) - (\hat{G}_\gamma + F)(\hat{x})], \end{aligned}$$

where we use the shorthand notation $\hat{G}_\gamma := G_\gamma(\cdot; \hat{x})$. Inserting our choice $\lambda_0 = 1$ with $\gamma = 0$, $\zeta^0 = x^0$ and the growth estimates from above, we obtain

$$(G + F)(x^N) \leq (G + F)(\hat{x}) + \frac{2}{(N-1)^2 \tau} \|x^0 - \hat{x}\|^2 \quad (N \geq 2).$$

This verifies the claimed convergence rates. \square

Example 3.7 (FISTA combined with strong convexity). In Example 3.6, suppose in addition that G is strongly convex with parameter $\gamma > 0$. Take $0 < \lambda \leq$

$\sqrt{L^{-2}\gamma^2 + 2L^{-1}\gamma} - L^{-1}\gamma$ and $\tau := \lambda^2/[2\gamma(1-\lambda)]$. Also let $\tilde{\tau} := \tau/[1 + (\lambda^{-1} - 1)\gamma\tau]$. Then (PP-i) with $\lambda_i \equiv \lambda$ and $\tau_i \equiv \tau$ becomes

$$\begin{cases} x^{i+1} := \text{prox}_{\tilde{\tau}G}(\tilde{x}^i - \tilde{\tau}\nabla F(\tilde{x}^i)), \\ \tilde{x}^{i+1} := x^{i+1} + \lambda(\lambda^{-1} - 1)(x^{i+1} - x^i), \\ \hat{x}^{i+1} := (\tilde{x}^{i+1} + \gamma\tau(\lambda^{-1} - 1)x^{i+1}) / (1 + \gamma\tau(\lambda^{-1} - 1)). \end{cases}$$

Both $G_\gamma(x^N; \hat{x}) + F(x^N) \rightarrow G_\gamma(\hat{x}; \hat{x}) + F(\hat{x})$ and $\zeta^N \rightarrow \hat{x}$ at the linear rate $O((1-\lambda)^N)$.

Demonstration. To derive the claimed algorithm, we divide the first step of (PP-i) by $1 + \gamma\tau(\lambda^{-1} - 1)$. This yields $\tilde{x}^i - \tilde{\tau}\nabla F(\tilde{x}^i) \in \tilde{\tau}\partial G(x^{i+1}) + x^{i+1}$; the rest follows from the definition of \tilde{x}^i and the proximal map. Observe that $\lambda \in (0, 1)$. Lemma 3.5(ii) with $\epsilon = 0$ and $\phi_0 = 1$ thus proves (3.3) and (3.5) and shows that $\phi_N\tau_N \geq \tau/(1-\lambda)^N$ and $\lambda_N^2\phi_N \geq \lambda^2/(1-\lambda)^N$. As in Example 3.6, we now obtain (3.20) provided $\tau L \leq 1$. With our chosen $\tau = \lambda^2/[2\gamma(1-\lambda)]$, this constraint resolves as our assumed upper bound $\lambda \leq \sqrt{L^{-2}\gamma^2 + 2L^{-1}\gamma} - L^{-1}\gamma$. With the growth rates above and $\tau = \lambda^2/[2\gamma(1-\lambda)]$, (3.20) after some rearrangements yields

$$\begin{aligned} (1-\lambda)\|\zeta^N - \hat{x}\|^2 + \gamma^{-1}[(\hat{G}_\gamma + F)(x^N) - (\hat{G}_\gamma + F)(\hat{x})] \\ \leq (1-\lambda)^{N+1} \left(\|x^0 - \hat{x}\|^2 + \gamma^{-1}[(\hat{G}_\gamma + F)(x^0) - (\hat{G}_\gamma + F)(\hat{x})] \right). \end{aligned}$$

Thus the claimed convergence rates hold. \square

Remark 3.8. Douglas–Rachford splitting for the problem $\min_{x \in X} F(x) + G(x)$ reads

$$\begin{cases} x^{i+1} = \text{prox}_{\gamma F}(v^i), \\ y^{i+1} = \text{prox}_{\gamma G}(2x^{i+1} - v^i), \\ v^{i+1} = v^i + y^{i+1} - x^{i+1}. \end{cases}$$

It can be presented in the form (PP) with $\tilde{H}_{i+1} = H$ for $u = (x, y, v)$ and

$$H(u) := \begin{pmatrix} \tau\partial F(x) + y - v \\ \tau\partial G(y) + v - x \\ x - y \end{pmatrix} \quad \text{and} \quad M_{i+1} := \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & I \end{pmatrix}.$$

With

$$\Gamma_{i+1} := \begin{pmatrix} 0 & I & -I \\ -I & 0 & I \\ I & -I & 0 \end{pmatrix}, \quad \text{taking instead} \quad M_{i+1} := \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & \lambda_{i+1}^{-1}I \end{pmatrix},$$

our approach can be used to construct a corrected inertial Douglas–Rachford splitting. We will, however, not pursue this. Instead, in the next section we take the primal-dual proximal splitting as an example of an algorithm with a nontrivial corrector and Γ_{i+1} .

Inertial Douglas–Rachford splitting has previously been studied in [21]. The “corrected” algorithm derived from our approach will be different. Another accelerated approach is considered in [5]. They apply Douglas–Rachford splitting to H defined in (1.10) by writing it in the form $H(u) = \partial\hat{G}(u) + \Gamma u$ for the convex function $\hat{G}(u) := G(x) + F^*(y)$ and an antisymmetric operator Γ , as we did in section 1. What this ingenious approach yields is essentially a doubly overrelaxed PDPS.

4. Inertial primal-dual proximal splitting. We now return to the saddle point problem (1.4). We suppose $G : X \rightarrow \overline{\mathbb{R}}$ and $F^* : Y \rightarrow \overline{\mathbb{R}}$ are (strongly) convex with factors $\gamma, \rho \geq 0$, and $K \in \mathcal{L}(X; Y)$. Recalling (1.10), for some step length, testing, and inertial parameters $\tau_i, \sigma_{i+1}, \phi_i, \psi_{i+1}, \lambda_i, \mu_{i+1} > 0$, we then take $\tilde{H}_{i+1} = W_{i+1}^{-1} H$ as well as

(4.1a)

$$H(u) := \begin{pmatrix} \partial G(x) + K^* y \\ \partial F^*(y) - Kx \end{pmatrix}, \quad W_{i+1} := \begin{pmatrix} \tau_i I & 0 \\ 0 & \sigma_{i+1} I \end{pmatrix}, \quad Z_{i+1} := \begin{pmatrix} \phi_i I & 0 \\ 0 & \psi_{i+1} I \end{pmatrix},$$

(4.1b)

$$\Gamma_{i+1} := \begin{pmatrix} \gamma \tau_i I & \tau_i K^* \\ -\sigma_{i+1} K & \rho \sigma_{i+1} I \end{pmatrix}, \quad \Lambda_{i+1} := \begin{pmatrix} \lambda_i I & 0 \\ 0 & \mu_{i+1} I \end{pmatrix}, \quad \text{and}$$

(4.1c)

$$M_{i+1} = \begin{pmatrix} I & -\mu_{i+1}^{-1} \tau_i K^* \\ -\lambda_i^{-1} \sigma_{i+1} \omega_i K & I \end{pmatrix} \quad \text{for} \quad \omega_i := \frac{\lambda_i \phi_i \tau_i}{\lambda_{i+1} \phi_{i+1} \tau_{i+1}}.$$

We then observe from (2.2) that the corrector

$$(4.2) \quad \tilde{M}_{i+1} = \Gamma_{i+1} (\Lambda_{i+1}^{-1} - I) = \begin{pmatrix} \gamma \tau_i (\lambda_i^{-1} - 1) I & \tau_i (\mu_{i+1}^{-1} - 1) K^* \\ -\sigma_{i+1} (\lambda_i^{-1} - 1) K & \rho \sigma_{i+1} (\mu_{i+1}^{-1} - 1) I \end{pmatrix}.$$

We need to satisfy the conditions of Theorem 2.3 for this setup, in particular (2.4) for some $\mathcal{V}_{k+1}(\hat{u})$ and $\hat{u} \in H^{-1}(0)$, and show that the estimate (2.6) is useful, in particular that $Z_{i+1} M_{i+1}$ and $\sum_{i=0}^{N-1} \mathcal{V}_{i+1}(\hat{u})$ are positive, that the former grows at a good rate, and that the latter becomes a useful gap functional. In the first instance, we intend to develop it into (a multiple of) the Lagrangian (duality) gap

$$(4.3) \quad \mathcal{G}(x, y; \hat{x}, \hat{y}) := (G(x) + \langle \hat{y}, Kx \rangle - F^*(\hat{y})) - (G(\hat{x}) + \langle y, K\hat{x} \rangle - F^*(y)).$$

Recalling G_γ and $(F^*)_\rho$ defined in (3.18), we also introduce the strong convexity adjusted gap

(4.4)

$$\mathcal{G}_{\gamma, \rho}(x, y; \hat{x}, \hat{y}) := (G_\gamma(x; \hat{x}) + \langle \hat{y}, Kx \rangle - (F^*)_\rho(\hat{y}; \hat{y})) - (G_\gamma(\hat{x}; \hat{x}) + \langle y, K\hat{x} \rangle - (F^*)_\rho(y; \hat{y})).$$

Since the problem $\min_x \max_y G_\gamma(x; \hat{x}) + \langle Kx, y \rangle - (F^*)_\rho(y; \hat{y})$ has the solution (\hat{x}, \hat{y}) , it is clear that $\mathcal{G}_{\gamma, \rho}$ is nonnegative, and zero at (\hat{x}, \hat{y}) . Before proving convergence we derive an explicit algorithm from (4.1).

4.1. Algorithm derivation. With the structure (4.1) fixed, we are ready to develop the skeleton of an explicit algorithm out of (PP-I). Since (PP-I) updates

$$(\bar{x}^i, \bar{y}^i) = \bar{u}^i := u^i + \Lambda_{i+1}(\Lambda_i^{-1} - I)(u^i - u^{i-1}),$$

we have

(4.5)

$$\bar{x}^i = x^i + \lambda_i(\lambda_{i-1}^{-1} - 1)(x^i - x^{i-1}), \quad \text{and} \quad \bar{y}^i = y^i + \mu_{i+1}(\mu_i^{-1} - 1)(y^i - y^{i-1}).$$

Using (4.2) and (4.1) we expand (PP-I) as

$$\begin{cases} 0 \in \tau_i \partial G(x^{i+1}) + \tau_i K^* y^{i+1} + (x^{i+1} - \bar{x}^i) - \mu_{i+1}^{-1} \tau_i K^* (y^{i+1} - \bar{y}^i) \\ \quad + \gamma \tau_i (\lambda_i^{-1} - 1) (x^{i+1} - x^i) + \tau_i (\mu_{i+1}^{-1} - 1) K^* (y^{i+1} - y^i), \\ 0 \in \sigma_{i+1} \partial F^*(y^{i+1}) - \sigma_{i+1} K x^{i+1} - \lambda_i^{-1} \sigma_{i+1} \omega_i K (x^{i+1} - \bar{x}^i) + (y^{i+1} - \bar{y}^i) \\ \quad + \rho \sigma_{i+1} (\mu_{i+1}^{-1} - 1) (y^{i+1} - y^i) - \sigma_{i+1} (\lambda_i^{-1} - 1) K (x^{i+1} - x^i). \end{cases}$$

The second line in both inclusions comes from the corrector term. Collecting all instances of the same iterate together, this can be simplified as

$$(4.6) \quad \begin{cases} 0 \in \tau_i \partial G(x^{i+1}) + [1 + \gamma \tau_i (\lambda_i^{-1} - 1)] x^{i+1} - [\bar{x}^i + \gamma \tau_i (1 - \lambda_i) x^i] \\ \quad + \tau_i K^* [\mu_{i+1}^{-1} \bar{y}^i - (\mu_{i+1}^{-1} - 1) y^i], \\ 0 \in \sigma_{i+1} \partial F^*(y^{i+1}) + [1 + \rho \sigma_{i+1} (\mu_{i+1}^{-1} - 1)] y^{i+1} - [\bar{y}^i + \rho \sigma_{i+1} (\mu_{i+1}^{-1} - 1) y^i] \\ \quad - \lambda_i^{-1} \sigma_{i+1} (1 + \omega_i) K x^{i+1} + \lambda_i^{-1} \sigma_{i+1} \omega_i K \bar{x}^i + \sigma_{i+1} (\lambda_i^{-1} - 1) K x^i. \end{cases}$$

Using (4.5) we can write

$$\begin{aligned} \bar{y}^i &:= \mu_{i+1}^{-1} \bar{y}^i - (\mu_{i+1}^{-1} - 1) y^i = \mu_{i+1}^{-1} y^i + (\mu_i^{-1} - 1) (y^i - y^{i-1}) - (\mu_{i+1}^{-1} - 1) y^i \\ &= y^i + (\mu_i^{-1} - 1) (y^i - y^{i-1}). \end{aligned}$$

Similarly, defining

$$\tilde{x}^{i+1} := x^{i+1} + (\lambda_i^{-1} - 1) (x^{i+1} - x^i),$$

we can write

$$\begin{aligned} x_\omega^{i+1} &:= \lambda_i^{-1} (1 + \omega_i) x^{i+1} - \lambda_i^{-1} \omega_i \bar{x}^i - (\lambda_i^{-1} - 1) x^i \\ &= [\lambda_i^{-1} x^{i+1} - (\lambda_i^{-1} - 1) x^i] + \omega_i [\lambda_i^{-1} x^{i+1} - \lambda_i^{-1} \bar{x}^i] \\ &= \tilde{x}^{i+1} + \omega_i [\lambda_i^{-1} x^{i+1} - (\lambda_i^{-1} - 1) x^i + (\lambda_i^{-1} - 1) x^i - \lambda_i^{-1} \bar{x}^i] \\ &= \tilde{x}^{i+1} + \omega_i (\tilde{x}^{i+1} - \tilde{x}^i). \end{aligned}$$

Also introducing

$$\tilde{\tau}_i := \tau_i / [1 + \gamma \tau_i (\lambda_i^{-1} - 1)] \quad \text{and} \quad \tilde{\sigma}_{i+1} := \sigma_{i+1} / [1 + \rho \sigma_{i+1} (\mu_{i+1}^{-1} - 1)],$$

we now rewrite (4.6) as

$$\begin{cases} 0 \in \tau_i \partial G(x^{i+1}) + (\tau_i / \tilde{\tau}_i) x^{i+1} - [\bar{x}^i + \gamma \tau_i (\lambda_i^{-1} - 1) x^i] + \tau_i K^* \tilde{y}^i, \\ 0 \in \sigma_{i+1} \partial F^*(y^{i+1}) + (\sigma_{i+1} / \tilde{\sigma}_{i+1}) y^{i+1} - [\bar{y}^i + \rho \sigma_{i+1} (\mu_{i+1}^{-1} - 1) y^i] - \sigma_{i+1} K x_\omega^{i+1}. \end{cases}$$

Multiplying, respectively, by $\tilde{\tau}_i / \tau_i$ and $\tilde{\sigma}_{i+1} / \sigma_{i+1}$, and recalling (4.5), we obtain the proximal updates of Algorithm 4.1, which we have written somewhat more compactly by additionally introducing the iterates x_γ^i and y_ρ^i . The updates of \bar{x}^{i+1} , \tilde{x}^{i+1} , \bar{y}^{i+1} , and \tilde{y}^{i+1} in the main step of the method are simply the definitions from above. The step length parameters will still need to be determined from one of the theorems referenced in Algorithm 4.1. Observe how the “corrected” inertial variables \tilde{x}^{i+1} and \tilde{y}^{i+1} differ from the standard inertial variables \bar{x}^{i+1} and \bar{y}^{i+1} .

Before developing specific rules for the step lengths and inertial parameters, we still need to provide the estimate (2.4). This process will produce additional conditions on the parameters.

4.2. Basic conditions. We now verify the basic conditions of Theorem 2.3 and the positivity of $Z_{i+1} M_{i+1}$.

LEMMA 4.1. *With the setup (4.1), the condition (2.5) holds if*

$$(4.7a) \quad \lambda_i^2 \phi_i (1 + 2\gamma \tau_i \lambda_i^{-1}) \geq \lambda_{i+1}^2 \phi_{i+1},$$

$$(4.7b) \quad \mu_{i+1}^2 \psi_{i+1} (1 + 2\rho \sigma_{i+1} \mu_{i+1}^{-1}) \geq \mu_{i+2}^2 \psi_{i+2},$$

$$(4.7c) \quad \lambda_i \phi_i \tau_i = \mu_i \psi_i \sigma_i.$$

Algorithm 4.1. IC-PDPS.

Require: On Hilbert spaces X and Y , a linear operator $K \in \mathcal{L}(X; Y)$, and convex, proper, and lower semicontinuous $G : X \rightarrow \mathbb{R}$ and $F^* : Y \rightarrow \mathbb{R}$ with factors $\gamma, \rho \geq 0$ of (strong) convexity.

- 1: Determine step length and inertial parameters $\{(\tau_i, \sigma_{i+1}, \lambda_i, \mu_{i+1}, \omega_i)\}_{i \in \mathbb{N}}$ from a suitable one among Theorems 4.5 to 4.8.
- 2: Pick initial iterates $\tilde{x}^0 := \bar{x}^0 := x^0 \in X$, and $\tilde{y}^0 := \bar{y}^0 := y^0 \in \text{dom } \partial F^*$.

3: Let $i := 0$.

4: **repeat** $\begin{cases} \tilde{\tau}_i := \tau_i / [1 + \gamma \tau_i (\lambda_i^{-1} - 1)], \\ \tilde{\sigma}_{i+1} := \sigma_{i+1} / [1 + \rho \sigma_{i+1} (\mu_{i+1}^{-1} - 1)]. \end{cases}$

5: Let

6: Compute

$$\begin{cases} x_\gamma^i := [\bar{x}^i + \gamma \tau_i (\lambda_i^{-1} - 1) x^i] / [1 + \gamma \tau_i (\lambda_i^{-1} - 1)], \\ x^{i+1} := \text{prox}_{\tilde{\tau}_i G}(x_\gamma^i - \tilde{\tau}_i K^* \tilde{y}^i), \\ \bar{x}^{i+1} := x^{i+1} + \lambda_{i+1} (\lambda_i^{-1} - 1) (x^{i+1} - x^i), \\ \tilde{x}^{i+1} := x^{i+1} + (\lambda_i^{-1} - 1) (x^{i+1} - x^i), \\ y_\rho^i := [\bar{y}^i + \rho \sigma_{i+1} (\mu_{i+1}^{-1} - 1) y^i] / [1 + \rho \sigma_{i+1} (\mu_{i+1}^{-1} - 1)], \\ y^{i+1} := \text{prox}_{\tilde{\sigma}_{i+1} F^*}(y_\rho^i + \sigma_{i+1} K[\tilde{x}^{i+1} + \omega_i (\bar{x}^{i+1} - \bar{x}^i)]), \\ \bar{y}^{i+1} := y^{i+1} + \mu_{i+2} (\mu_{i+1}^{-1} - 1) (y^{i+1} - y^i), \\ \tilde{y}^{i+1} := y^{i+1} + (\mu_{i+1}^{-1} - 1) (y^{i+1} - y^i). \end{cases}$$

7: Update $i := i + 1$.

8: **until** a stopping criterion is satisfied.

Observe: If $\gamma = 0$, then $\tilde{\tau}_i = \tau_i$ and $x_\gamma^i = \bar{x}^i$. If $\rho = 0$, then $\tilde{\sigma}_{i+1} = \sigma_{i+1}$ and $y_\rho^i = \bar{y}^i$.

Proof. Inserting the operators from (4.1), the condition (2.5) reads

$$\begin{pmatrix} \phi_i \lambda_i (\lambda_i + 2\gamma \tau_i) I & \phi_i \lambda_i \tau_i K^* \\ -\psi_{i+1} \mu_{i+1} \sigma_{i+1} (2 + \omega_i) K & \psi_{i+1} \mu_{i+1} (\mu_{i+1} + 2\rho \sigma_{i+1}) I \end{pmatrix} \\ \geq \begin{pmatrix} \phi_{i+1} \lambda_{i+1}^2 I & -\phi_{i+1} \lambda_{i+1} \tau_{i+1} K^* \\ -\psi_{i+2} \mu_{i+2} \sigma_{i+2} \omega_{i+1} K & \psi_{i+2} \mu_{i+2}^2 I \end{pmatrix}.$$

Further inserting ω_i and ω_{i+1} from (4.1c), and using (4.7c), the off-diagonal components cancel out. The diagonal components that are left are simply (4.7a) and (4.7b). \square

LEMMA 4.2. Let $i \in \mathbb{N}$. If (4.1) and (4.7) hold, then $Z_{i+1} M_{i+1}$ is self-adjoint. If, moreover,

$$(4.8) \quad (1 - \kappa) \mu_{i+1}^2 \psi_{i+1} \geq \phi_i \tau_i^2 \|K\|^2 \quad \text{for some } \kappa \in [0, 1),$$

then $Z_{i+1} M_{i+1}$ is positive definite, more precisely

$$(4.9) \quad Z_{i+1} M_{i+1} \geq \delta Z_{i+1} \quad \text{for } \delta := 1 - \sqrt{1 - \kappa}.$$

Proof. For now, take arbitrary $\delta \in [0, \kappa]$. From (4.1c), using Cauchy's inequality

$$Z_{i+1} M_{i+1} = \begin{pmatrix} \phi_i I & -\mu_{i+1}^{-1} \phi_i \tau_i K^* \\ -\mu_{i+1}^{-1} \phi_i \tau_i K & \psi_{i+1} I \end{pmatrix} \geq \begin{pmatrix} \delta \phi_i I & 0 \\ 0 & \psi_{i+1} I - (1 - \delta)^{-1} \mu_{i+1}^{-2} \phi_i \tau_i^2 K K^* \end{pmatrix}.$$

Clearly then $Z_{i+1}M_{i+1}$ is self-adjoint. Using (4.8), we have

$$\psi_{i+1}I - (1-\delta)^{-1}\mu_{i+1}^{-2}\phi_i\tau_i^2KK^* \geq \psi_{i+1}I - (1-\delta)^{-1}(1-\kappa)\psi_{i+1}I = (\kappa-\delta)(1-\delta)^{-1}\psi_{i+1}I.$$

To make a specific choice of δ , we equate $\delta = (\kappa - \delta)(1 - \delta)^{-1}$. This gives the quadratic equation $2\delta - \delta^2 - \kappa = 0$ with the solution $\delta = 1 - \sqrt{1 - \kappa}$. The rest is trivial. \square

4.3. Gap unrolling and alignment. We now derive a basic convergence estimate using Theorem 2.3. This involves verifying (2.4) for some $\mathcal{V}_{i+1}(\hat{u})$ and estimating the sum of the latter from below to yield a useful gap estimate. For the statement of the next lemma, we recall the definition of the strong convexity adjusted functions G_γ and $(F^*)_\rho$ from subsection 3.4, and the corresponding gap functional defined in (4.4).

LEMMA 4.3. *Suppose (4.1) and (4.7) hold. Take*

$$(4.10a) \quad (1 - \lambda_{i+1})\phi_{i+1}\tau_{i+1} \leq \phi_i\tau_i, \quad \lambda_0 = 1,$$

$$(4.10b) \quad (1 - \mu_{i+1})\psi_{i+1}\sigma_{i+1} \leq \psi_i\sigma_i \leq \psi_{i+1}\sigma_{i+1}, \quad \mu_0 = 1, \quad \text{and}$$

$$(4.10c) \quad \phi_i\tau_i = \psi_i\sigma_i \quad (i \in \mathbb{N}).$$

Then for any $\hat{u} = (\hat{x}, \hat{y}) \in H^{-1}(0)$, the iterates generated by Algorithm 4.1 satisfy

$$(4.11) \quad \frac{1}{2}\|z^N - \hat{u}\|_{\Lambda_{N+1}^*Z_{N+1}M_{N+1}\Lambda_{N+1}}^2 + \phi_{N-1}\tau_{N-1}\mathcal{G}_{\gamma,\rho}(u^N; \hat{u}) \leq C_0(\hat{u}) \quad (N \geq 1),$$

where for any $w^0 \in \partial(F^*)_\rho(y^0)$ we set

$$C_0(\hat{u}) := \frac{1}{2}\|z^0 - \hat{u}\|_{\Lambda_1^*Z_1M_1\Lambda_1}^2 + \psi_0\sigma_0 \langle w^0 - K\hat{x}, y^0 - \hat{y} \rangle.$$

Proof. Observe that Algorithm 4.1 explicitly requires $y^0 \in \text{dom } \partial F^*$, so some $w^0 \in \partial(F^*)_\rho(y^0; \hat{y})$ exists. The proximal steps moreover ensure $y^{i+1} \in \text{dom } \partial F^*$ and $x^{i+1} \in \text{dom } \partial G$ for all $i \in \mathbb{N}$.

By the defining (2.3), the auxiliary sequence $\{z^i\}_{i \in \mathbb{N}} \subset X \times Y$ for $z^i = (\zeta^i, \eta^i)$ satisfies

$$(4.12) \quad \lambda_i\zeta^{i+1} = x^{i+1} - (1 - \lambda_i)x^i \quad \text{and} \quad \mu_{i+1}\eta^{i+1} = y^{i+1} - (1 - \mu_{i+1})y^i \quad (i \in \mathbb{N}).$$

Since $\mu_0 = 1$ and $\eta^0 = y^0$, the latter also works for $i = -1$ and any, superfluous, $y^{-1} \in Y$. We observe that

$$\tilde{H}_{i+1}(u^{i+1}) - \Gamma_{i+1}(u^{i+1} - \hat{u}) = \begin{pmatrix} \tau_i[\partial G(x^{i+1}) - \gamma(x^{i+1} - \hat{x}) + K^*\hat{y}] \\ \sigma_{i+1}[\partial F^*(y^{i+1}) - \rho(y^{i+1} - \hat{y}) - K\hat{x}] \end{pmatrix}.$$

Let us define (recall subsection 3.4)

$$(4.13a) \quad \bar{G}_\gamma(x; \hat{u}) := G(x) - \frac{\gamma}{2}\|x - \hat{x}\|^2 + \langle K^*\hat{y}, x - \hat{x} \rangle = G_\gamma(x; \hat{x}) + \langle K^*\hat{y}, x - \hat{x} \rangle, \quad \text{and}$$

$$(4.13b) \quad (\bar{F}^*)_\rho(y; \hat{u}) := F^*(y) - \frac{\rho}{2}\|y - \hat{y}\|^2 - \langle K\hat{x}, y - \hat{y} \rangle = (F^*)_\rho(y; \hat{y}) - \langle K\hat{x}, y - \hat{y} \rangle.$$

Then \bar{G}_γ and $(\bar{F}^*)_\rho$ are convex with $\bar{G}_\gamma(x; \hat{u}) \geq \bar{G}_\gamma(\hat{x}; \hat{u})$, and $(\bar{F}^*)_\rho(y; \hat{u}) \geq (\bar{F}^*)_\rho(\hat{y}; \hat{u})$ for all $x \in X$ and $y \in Y$. Moreover, (2.4) holds with

$$\begin{aligned} \mathcal{V}_{i+1}(\hat{u}) &:= \inf \left\langle \tilde{H}_{i+1}(u^{i+1}) - \Gamma_{i+1}(u^{i+1} - \hat{u}), z^{i+1} - \hat{u} \right\rangle_{\Lambda_{i+1}^*Z_{i+1}} \\ &= \inf \left[\phi_i\tau_i\lambda_i \left\langle \partial \bar{G}_\gamma(x^{i+1}; \hat{u}), \zeta^{i+1} - \hat{x} \right\rangle + \psi_{i+1}\sigma_{i+1}\mu_{i+1} \left\langle \partial (\bar{F}^*)_\rho(y^{i+1}; \hat{u}), \eta^{i+1} - \hat{y} \right\rangle \right]. \end{aligned}$$

For each $i \in \mathbb{N}$, let $\bar{q}^{i+1} \in \partial \bar{G}_\gamma(x^{i+1}; \hat{u})$, and let $w^i \in Y$ be such that $\bar{w}^i = w^i - K\hat{x} \in \partial(\bar{F}^*)_\rho(y^i; \hat{u})$. Define

$$(4.14) \quad s_N := s_N^G + s_N^{F^*} := \sum_{i=0}^{N-1} \phi_i \tau_i \lambda_i \langle \bar{q}^{i+1}, \zeta^{i+1} - \hat{x} \rangle + \sum_{i=0}^{N-1} \psi_{i+1} \sigma_{i+1} \mu_{i+1} \langle \bar{w}^{i+1}, \eta^{i+1} - \hat{y} \rangle.$$

Since we have assumed (4.1) and (4.7), we may use Lemmas 4.1 and 4.2 to verify (2.5) and the self-adjointness of $Z_{i+1}M_{i+1}$. We may therefore use Theorem 2.3 to establish (4.11) if we further show that

$$(4.15) \quad s_N \geq \phi_{N-1} \tau_{N-1} \mathcal{G}_{\gamma, \rho}(u^N; \hat{u}) - c_0$$

for

$$c_0 := \psi_0 \sigma_0 \mu_0 \langle \bar{w}^0, \eta^0 - \hat{y} \rangle = \psi_0 \sigma_0 \langle w^0 - K\hat{x}, y^0 - \hat{y} \rangle.$$

Indeed, this establishes the right-hand side as a lower bound on $\sum_{i=0}^{N-1} \mathcal{V}_{i+1}(\hat{u})$.

The difficulty in working with s_N is that unless $\gamma = \rho = 0$, our algorithm will give $\phi_i \tau_i = \psi_i \sigma_i$, not $\phi_i \tau_i = \psi_{i+1} \sigma_{i+1}$. We therefore have to realign variables. Using the assumption $\lambda_0 = 1$, (4.10a), and (4.12), by Lemma 3.2 and (4.13a) we have

$$(4.16) \quad \begin{aligned} s_N^G &\geq \phi_{N-1} \tau_{N-1} [\bar{G}_\gamma(x^N; \hat{u}) - \bar{G}_\gamma(\hat{x}; \hat{u})] \\ &= \phi_{N-1} \tau_{N-1} [G_\gamma(x^N; \hat{x}) - G_\gamma(\hat{x}; \hat{x})] + \phi_{N-1} \tau_{N-1} \langle K^* \hat{y}, x^N - \hat{x} \rangle. \end{aligned}$$

Regarding $s_N^{F^*}$, by the second inequality of (4.10b) we have $\psi_N \sigma_N \geq \psi_{N-1} \sigma_{N-1}$. Moreover, \hat{y} minimizes $(\bar{F}^*)_\rho(\cdot; \hat{u})$. Using these two facts after an application analogous to (4.16) of Lemma 3.2, we get

$$(4.17) \quad \begin{aligned} s_N^{F^*} &= \sum_{i=0}^N \psi_i \sigma_i \mu_i \langle \bar{w}^i, \eta^i - \hat{y} \rangle - c_0 \\ &\geq \psi_N \sigma_N [(\bar{F}^*)_\rho(y^N; \hat{u}) - (\bar{F}^*)_\rho(\hat{y}; \hat{u})] - c_0 \\ &\geq \psi_{N-1} \sigma_{N-1} [(\bar{F}^*)_\rho(y^N; \hat{u}) - (\bar{F}^*)_\rho(\hat{y}; \hat{u})] - c_0 \\ &= \psi_{N-1} \sigma_{N-1} [(F^*)_\rho(y^N; \hat{y}) - (F^*)_\rho(\hat{y}; \hat{y})] - \psi_{N-1} \sigma_{N-1} \langle K\hat{x}, y^N - \hat{y} \rangle - c_0. \end{aligned}$$

Combining (4.14), (4.16), and (4.17), thus

$$\begin{aligned} s_N &\geq \phi_{N-1} \tau_{N-1} [G_\gamma(x^N; \hat{x}) - G_\gamma(\hat{x}; \hat{x})] + \psi_{N-1} \sigma_{N-1} ((F^*)_\rho(y^N; \hat{y}) - (F^*)_\rho(\hat{y}; \hat{y})) \\ &\quad + \phi_{N-1} \tau_{N-1} \langle K^* \hat{y}, x^N - \hat{x} \rangle - \psi_{N-1} \sigma_{N-1} \langle K\hat{x}, y^N - \hat{y} \rangle - c_0. \end{aligned}$$

Now (4.10c) establishes (4.15). \square

4.4. Step length and inertial parameter rules. We now consider several cases of the factors of (strong) convexity ρ and γ being zero or positive. Throughout, as in the proof of (4.3), we write $z^i = (\zeta^i, \eta^i)$, for the auxiliary sequence $\{z^i\}_{i \in \mathbb{N}}$ defined in (2.3). We first summarize the various lemmas and their conditions from above.

LEMMA 4.4. With $\lambda_0 = 1$ and $\tau_0, \sigma_0, \phi_0, \psi_0 > 0$, suppose that $\mu_i = \lambda_i$ as well as

$$(4.18a) \quad \psi_i \sigma_i = \phi_i \tau_i, \quad \omega_i \lambda_{i+1} \phi_{i+1} \tau_{i+1} = \lambda_i \phi_i \tau_i,$$

$$(4.18b)$$

$$\lambda_i^2 \phi_i (1 + 2\gamma \tau_i \lambda_i^{-1}) \geq \lambda_{i+1}^2 \phi_{i+1}, \quad (1 - \lambda_{i+1}) \phi_{i+1} \tau_{i+1} \leq \phi_i \tau_i \leq \phi_{i+1} \tau_{i+1},$$

$$(4.18c)$$

$$\lambda_i^2 \psi_i (1 + 2\rho \sigma_i \lambda_i^{-1}) \geq \lambda_{i+1}^2 \psi_{i+1}, \quad (1 - \kappa) \lambda_{i+1}^2 \psi_{i+1} \geq \phi_i \tau_i^2 \|K\|^2 \quad (i \in \mathbb{N})$$

for some $\kappa \in [0, 1)$. Then the iterates generated by Algorithm 4.1 and the auxiliary sequence generated by (2.3) satisfy with $\delta := 1 - \sqrt{1 - \kappa}$ for any $N \geq 1$ and any $\hat{u} \in H^{-1}(0)$ the estimate

$$(4.19)$$

$$\frac{\delta \phi_N \lambda_N^2}{2} \|\zeta^N - \hat{x}\|^2 + \frac{\delta \psi_{N+1} \lambda_{N+1}^2}{2} \|\eta^N - \hat{y}\|^2 + \phi_{N-1} \tau_{N-1} \mathcal{G}_{\gamma, \rho}(u^N; \hat{u}) \leq C_0(\hat{u}).$$

Proof. We first show that the setup (4.1) and the conditions (4.7), (4.8), and (4.10) hold. Indeed, the second part of (4.18a) is simply the choice of ω_i in (4.1c), while the rest of (4.1c) follows from the derivation of Algorithm 4.1 from this structural setup in subsection 4.1. Moreover, since $\mu_i = \lambda_i$, the first part of (4.18a) implies (4.7c) and (4.10c). Likewise, (4.18b) implies (4.7a) and (4.10a). The conditions (4.18c) in turn imply (4.7b) and (4.8). Together (4.18a) and (4.18b) imply (4.10b). Therefore (4.7), (4.8), and (4.10) hold in their entirety. We can thus apply Lemmas 4.1 and 4.3 to obtain the estimate (4.11). By application of Lemma 4.2 we then derive (4.19) from (4.11). \square

THEOREM 4.5. Suppose $\gamma = 0$ and $\rho = 0$. Take $\tau_0, \sigma_0 > 0$ with $\tau_0 \sigma_0 \|K\|^2 < 1$, $\lambda_0 = \mu_0 = 1$, $\epsilon \in [0, 1)$, and update $(i \in \mathbb{N})$

$$(4.20a) \quad \tau_{i+1} := \tau_i \lambda_i^{-1} \lambda_{i+1}, \quad \omega_i := 1,$$

$$(4.20b) \quad \sigma_{i+1} := \sigma_i \lambda_i^{-1} \lambda_{i+1}, \quad \mu_{i+1} := \lambda_{i+1} := \lambda_i / (1 + (1 - \epsilon) \lambda_i).$$

Then the iterates generated by Algorithm 4.1 satisfy $\mathcal{G}(u^N; \hat{u}) \rightarrow 0$ at the rate $O(1/N)$ for any $\hat{u} \in H^{-1}(0)$.

Proof. We will use Lemma 4.4, for which we need to verify (4.18). We use Lemma 3.5(iii) to verify (4.18b) for $\phi_i = \tau_i^{-2}$, $\tau_{i+1} = \tau_i(1 - \lambda_{i+1})/(1 - \epsilon \lambda_i)$, and λ_{i+1} as in (3.14). With $\gamma = 0$, the latter agrees with the expression for λ_{i+1} in (4.20). With ϕ_i and $\rho = 0$ inserted, the rest of (4.18) reads

$$(4.21a) \quad \psi_i \sigma_i = \tau_i^{-1}, \quad \omega_i \lambda_{i+1} \tau_{i+1}^{-1} = \lambda_i \tau_i^{-1},$$

$$(4.21b) \quad \lambda_i^2 \psi_i \geq \lambda_{i+1}^2 \psi_{i+1}, \quad (1 - \kappa) \lambda_{i+1}^2 \psi_{i+1} \geq \|K\|^2 \quad (i \in \mathbb{N}).$$

Clearly the first part of (4.21b) holds by taking $\psi_i = \lambda_i^{-2} \tau_0^{-1} \sigma_0^{-1}$ for all $i \in \mathbb{N}$. Let us assume

$$(4.22)$$

$$\omega_i = (\lambda_{i+1}^{-1} - 1)/(\lambda_i^{-1} - \epsilon), \quad \tau_{i+1} = \tau_i \lambda_i^{-1} \lambda_{i+1} \omega_i, \quad \text{and} \quad \sigma_{i+1} := \sigma_i \lambda_i^{-1} \lambda_{i+1} / \omega_i.$$

Inserting $\tau_{i+1} = \tau_i(1 - \lambda_{i+1})/(1 - \epsilon \lambda_i)$ from above and ω_i from (4.22), the second part of (4.21a) holds. It therefore only remains to secure the first part of (4.21a) and the second part of (4.21b). With ψ_i inserted, this is to say

$$\sigma_i \tau_i = \sigma_0 \tau_0 \lambda_i^2 \quad \text{and} \quad (1 - \kappa) \geq \tau_0 \sigma_0 \|K\|^2.$$

The second condition is simply our initial condition on the step lengths. The first condition holds if $\sigma_i = \tau_i^{-1} \lambda_i^2 \sigma_0 \tau_0$. Using that $\tau_{i+1} = \tau_i \lambda_i^{-1} \lambda_{i+1} \omega_i$, this holds when σ_{i+1} as in (4.22). We have therefore proved (4.18) to hold when (4.22) does, $\phi_i = \tau_i^{-2}$, $\psi_i = \lambda_i^{-2} \tau_0^{-1} \sigma_0^{-1}$, and $\tau_0 \sigma_0 \|K\|^2 < 1$.

Take now as stated $\omega_i = 1$, and observe that the update rule for λ_{i+1} in (4.20) gives the rule for ω_i in (4.22). Moreover, the rules for τ_{i+1} and σ_{i+1} in (4.20) are consistent with (4.22). Therefore (4.22), consequently (4.18), holds under the conditions of the theorem and the choices of the testing parameters ϕ_i and ψ_i in the previous paragraph. Lemma 4.4 thus yields (4.19). Since $\mathcal{G}_{\gamma, \rho}(u^N; \hat{u}) \geq 0$ when $\hat{u} \in H^{-1}(0)$, the growth estimate of Lemma 3.5(iii) in the case $\gamma = 0$ applied in (4.19) establishes the claimed convergence rate. \square

Thus, without any strong convexity, inertia and correction improve the ergodic $O(1/N)$ convergence of the gap for the PDPS to nonergodic convergence.

THEOREM 4.6. *Suppose $\gamma > 0$ and $\rho = 0$. Take $\epsilon \in [0, 1)$, $\tau_0, \sigma_0 > 0$ with $\tau_0 \sigma_0 \|K\|^2 < 1$, initialize $\lambda_0 := \mu_0 := 1$, and update ($i \in \mathbb{N}$)*

$$(4.23a) \quad \tau_{i+1} := \tau_i \lambda_i^{-1} \lambda_{i+1} \omega_i, \quad \omega_i := (\lambda_{i+1}^{-1} - 1) / (\lambda_i^{-1} - \epsilon),$$

$$(4.23b) \quad \sigma_{i+1} := \sigma_i \lambda_i^{-1} \lambda_{i+1} / \omega_i, \quad \mu_{i+1} := \lambda_{i+1} := \frac{\sqrt{\lambda_i^2 + 2\gamma \lambda_i \tau_i}}{1 - \epsilon \lambda_i + \sqrt{\lambda_i^2 + 2\gamma \lambda_i \tau_i}}.$$

Then the iterates generated by Algorithm 4.1 satisfy both $\mathcal{G}_{\gamma, 0}(u^N; \hat{u}) \rightarrow 0$ and $\|\zeta^N - \hat{x}\|^2 \rightarrow 0$ at the rate $O(1/N^2)$ for any $\hat{u} \in H^{-1}(0)$.

Proof. Note that λ_{i+1} given in (3.14) agrees with that in (4.23). Also note that in the proof of Theorem 4.5, we did not involve the choices (4.20) until the final paragraph. Therefore, we may follow the proof of Theorem 4.5 to see (4.18) to hold when (4.22) does, $\phi_i = \tau_i^{-2}$, $\psi_i = \lambda_i^{-2} \tau_0^{-1} \sigma_0^{-1}$, and $\tau_0 \sigma_0 \|K\|^2 < 1$.

Observe now that (4.22) is consistent with the updates of ω_i , τ_{i+1} , and σ_{i+1} in (4.23). By taking the testing parameters ϕ_i and ψ_i as above, we have therefore verified (4.18), so Lemma 4.4 yields (4.19). The growth estimate of Lemma 3.5(iii) in the case $\gamma > 0$ applied in (4.19) establishes the claimed convergence rates. \square

THEOREM 4.7. *Suppose $\gamma = 0$ and $\rho > 0$. Take $\tau_0 > 0$ and $\epsilon \in [0, 1/2]$ with $\tau_0 \|K\|^2 < 2\rho$, initialize $\lambda_0 := \mu_0 := 1$, and update ($i \in \mathbb{N}$)*

$$\begin{aligned} \tau_i &:= \tau_0, & \omega_i &:= (\lambda_{i+1}^{-1} - 1) / (\lambda_i^{-1} - \epsilon) = \lambda_{i+1} \lambda_i^{-1}, \\ \sigma_{i+1} &:= \frac{\lambda_i^2}{2\rho}, & \mu_{i+1} := \lambda_{i+1} &:= \frac{2}{1 + \sqrt{1 + 4(\lambda_i^{-2} - \epsilon \lambda_i^{-1})}}, \end{aligned}$$

Then the iterates generated by Algorithm 4.1 satisfy both $\mathcal{G}_{0, \rho}(u^N; \hat{u}) \rightarrow 0$ and $\|\eta^N - \hat{y}\|^2 \rightarrow 0$ at the rate $O(1/N^2)$ for any $\hat{u} \in H^{-1}(0)$.

Proof. We use Lemma 3.5(i) to verify (4.18b) for $\phi_i = \lambda_i^{-2}$ as well as $\tau_i = \tau_0$ and λ_{i+1} as stated. Inserting the ϕ_i and τ_i , the rest of (4.18) now reduces to

$$\begin{aligned} \psi_i \sigma_i &= \lambda_i^{-2} \tau_0, & \omega_i \lambda_{i+1}^{-1} &= \lambda_i^{-1}, \\ \lambda_i^2 \psi_i (1 + 2\rho \sigma_i \lambda_i^{-1}) &\geq \lambda_{i+1}^2 \psi_{i+1}, & \text{and} & \quad (1 - \kappa) \lambda_{i+1}^2 \psi_{i+1} \geq \lambda_i^{-2} \tau_0^2 \|K\|^2. \end{aligned}$$

The second condition is one version of our update rule for ω_i . We still need to show that the two versions of the rule are equal. If we take $\psi_i := 2\rho \tau_0 \lambda_{i-1}^{-2} \lambda_i^{-2}$ and σ_i

as stated, introducing the new variable λ_{-1} , not used in the algorithm, the rest becomes

$$\lambda_{i-1}^{-2} + \lambda_i^{-1} \geq \lambda_i^{-2} \quad \text{and} \quad (1 - \kappa)2\rho \geq \tau_0 \|K\|^2 \quad (i \in \mathbb{N}).$$

The latter condition is satisfied by our initial step length assumption for some $\kappa \in (0, 1)$. Since $\lambda_0 = 1$, the first condition holds for $i = 0$ for any $\lambda_{-1} > 0$. By Lemma 3.4, $\lambda_{i-1}^{-2} - \epsilon \lambda_{i-1}^{-1} = \lambda_i^{-2} - \lambda_i^{-1}$ for $i \in \mathbb{N}$. Therefore the first condition holds, and the two expressions for ω_i are equivalent.

We have thus verified (4.18), so Lemma 4.4 gives the estimate (4.19). The growth estimate of Lemma 3.5 (i) applied there establishes the claimed gap convergence rate. The convergence rate of the dual auxiliary variable is determined by the rate of growth of $\lambda_{N+1}^2 \psi_{N+1} = 2\rho\tau_0 \lambda_N^{-2}$. By the same Lemma 3.5(i), $\phi_N \tau_N = \lambda_N^{-2} \phi_0$ grows at the rate $\Theta(N^2)$, so we get the claimed $O(1/N^2)$ convergence. \square

THEOREM 4.8. *Suppose $\gamma > 0$ and $\rho > 0$. Take $\lambda \in (0, 1)$ and $\epsilon \in [0, 1)$ with $\|K\|^2 < 4\gamma\rho(\lambda^{-1} - \epsilon)(\lambda^{-1} - 1)$. Update $(i \in \mathbb{N})$*

$$\begin{aligned} \tau_i &:= \lambda^2 / [2\gamma(1 - \lambda)], & \omega_i &:= (\lambda^{-1} - 1) / (\lambda^{-1} - \epsilon), \\ \sigma_i &:= \lambda^2 / [2\rho(1 - \lambda)], & \mu_{i+1} &:= \lambda_{i+1} := \lambda. \end{aligned}$$

Then the iterates generated by Algorithm 4.1 satisfy both $\mathcal{G}_{\gamma,\rho}(u^N; \hat{u}) \rightarrow 0$ and $\|z^N - \hat{u}\|^2 \rightarrow 0$ at a linear rate for (the unique) $\hat{u} \in H^{-1}(0)$.

Proof. We use Lemma 3.5(ii) to verify (4.18b) for $\phi_i = c^i$ for $c := (1 - \epsilon\lambda)/(1 - \lambda) > 1$ as well as $\tau_i \equiv \tau_0$ and $\lambda_{i+1} \equiv \lambda$ as stated. The rest of (4.18) now reduces to

$$\begin{aligned} \psi_i \sigma_i &= \phi_i \tau_0, & \omega_i c &= 1, \\ \psi_i(1 + 2\rho\sigma_i\lambda^{-1}) &\geq \psi_{i+1}, & \text{and} & \quad (1 - \kappa)\lambda^2\psi_{i+1} \geq \phi_i \tau_0^2 \|K\|^2. \end{aligned}$$

Clearly our choice of $\omega_i = 1/c$ satisfies the second condition. Taking $\psi_i = \phi_i \rho / \gamma$ and, as stated, $\sigma_i = \gamma\tau_0 / \rho = \lambda^2 / [2\rho(1 - \lambda)]$, the first condition is also satisfied, while the third condition becomes $\phi_{i+1}(1 + 2\gamma\tau_0\lambda^{-1}) \geq \phi_i$. As $\lambda_i \equiv \lambda$, this is the first part of (4.18b), which we have already verified. The last condition becomes $(1 - \kappa)\lambda^2\rho\gamma^{-1}c \geq \tau_0^2 \|K\|^2$, which with c and τ_0 expanded is $4(1 - \kappa)\gamma\rho(1 - \epsilon\lambda)(1 - \lambda) \geq \lambda^2 \|K\|^2$. This is secured by our assumed bound on $\|K\|$.

We have thus verified (4.18), so Lemma 4.4 yields the estimate (4.19). The growth estimate of Lemma 3.5(ii) applied there establishes the claimed gap and primal variable convergence rates. Since $\psi_i = \phi_i \rho / \gamma$ and $\mu_i = \lambda_i$, the dual variable converges at the same rate as the primal variable. \square

Remark 4.9 (partial gaps). Convergence of the Lagrangian gap is weak compared to the true duality gap. Suppose the product set $B_x \times B_y \subset X \times Y$ is bounded and contains some $\hat{u} \in H^{-1}(0)$. In the literature, *partial gaps* are considered,

$$0 \leq \mathcal{G}(u; B_x, B_y) := \sup_{\bar{x} \in B_x} \inf_{\bar{y} \in B_y} \mathcal{G}(u; (\bar{x}, \bar{y})).$$

Ergodic partial gap estimates can be derived in the noninertial setting (see [8]) because of the fact that $\hat{u} \in H^{-1}(0)$ is never actually needed in the proofs; $\hat{u} \in X \times Y$ can be any element. Indeed, even in our work, the main reason we need to assume \hat{u} to be a solution are the final phases of the unrolling Lemmas 3.2 and 3.3. However, because of this, we cannot derive partial gap estimates.

5. Numerical experience. We study the performance of the proposed algorithm on three image processing and inverse problems: denoising, sparse Fourier inversion, and positron emission tomography (PET), all with total variation (TV) regularization. Denoising is the most basic image processing task, while sparse Fourier inversion is used for magnetic resonance image reconstruction; see, e.g., [4, 17]. These two problems are of the form

$$(5.1) \quad \min_{x \in \mathbb{R}^{n_1 n_2}} \frac{1}{2} \|z - Tx\|_2^2 + \beta \|Dx\|_{2,1},$$

where $n_1 \times n_2$ is the size of the unknown image x in pixels, $z \in \mathbb{R}^m$ is the corrupted data, and $\beta > 0$ a regularization parameter. The matrix $D \in \mathbb{R}^{2n_1 n_2 \times n_1 n_2}$ is a discretization of the gradient operator, and $\|g\|_{2,1} := \sum_{i=1}^{n_1 n_2} \sqrt{g_{i,1}^2 + g_{i,2}^2}$ for $g = (g_{\cdot,1}, g_{\cdot,2}) \in \mathbb{R}^{2n_1 n_2}$. We take D as forward-differences with Neumann boundary conditions.

The operator $T \in \mathbb{R}^{k \times n_1 n_2}$ depends on the problem in question: for denoising, $T = I$ is the identity, and for sparse Fourier inversion it is the composition $T = S\mathcal{F}$ with a subsampling operator $S \in \mathbb{R}^{k \times n_1 n_2}$ and the discrete Fourier transform \mathcal{F} . For denoising $k = n_1 n_2$, while for sparse Fourier reconstruction, $k \ll n_1 n_2$.

To implement variants of the PDPS, we note that (5.1) can in all three cases be written in the saddle point form

$$\min_{x \in \mathbb{R}^{n_1 n_2}} \max_{y \in \mathbb{R}^{2n_1 n_2}} \frac{1}{2} \|z - Tx\|_2^2 + \langle Dx, y \rangle - \delta_{\beta B}(y),$$

where $B = B_{\mathbb{R}^2}^{n_1 n_2}$ for $B_{\mathbb{R}^2}$ the Euclidean unit ball in \mathbb{R}^2 . Since T is in both cases related to a unitary operator, we can easily compute the proximal map of $G(x) := \frac{1}{2} \|z - Tx\|_2^2$.

The PET problem is slightly different. We take as T a discrete Radon transform, each $[Tx]_j$ being the integral of the image x over a line with angle parameter θ_j and displacement r_j . As the efficient and precise realization of such an operator in general cases is outside the scope of the present work, in our simplified setting, we consider only the four angles $\theta_j \in \{0^\circ, 45^\circ, 90^\circ, 135^\circ\}$ and displacements r_j such that Tx consists of all row sums, all column sums, and all diagonal sums of x rewritten as a $n_1 \times n_2$ matrix. We also change the first fidelity term in (5.1) to model, instead of Gaussian noise, Poisson noise. Finally, we need to force $x \geq 0$. That is, our problem is

$$\min_{x \in [0, \infty)^{n_1 n_2}} \langle Tx, \mathbb{1} \rangle - \langle b, \log(Tx + c) \rangle + \beta \|Dx\|_{2,1},$$

where $\mathbb{1} := (1, \dots, 1) \in \mathbb{R}^k$, $b \in (0, \infty)^k$ is the measured data, and $c \in (0, \infty)^k$ is a background intensity, assumed known. The logarithm is applied componentwise.

Computing the proximal step with respect to the fidelity term is challenging due to the structure of T . We therefore also write this term as a conjugate, observing that $g_j(z) := z - b_j \log(z + c_j)$ has the conjugate $g_j^*(\phi_j) = -b_j + c_j(1 - \phi_j) + b_j \log(b_j / (1 - \phi_j))$. Introducing the additional upper bound $x \leq 1$, this leads to

$$\min_{x \in \mathbb{R}^{n_1 n_2}} \max_{(\phi, y) \in \mathbb{R}^k \times \mathbb{R}^{2n_1 n_2}} \delta_{[0,1]^{n_1 n_2}}(x) + \langle (Tx, Dx), (\phi, y) \rangle - \left(\delta_{\beta B}(y) + \sum_{j=1}^k g_j^*(\phi_j) \right),$$

Without the additional upper bound, this problem arranged as the prototype problem (1.4) would have $G = \delta_{[0, \infty)^{n_1 n_2}}$, which has the conjugate $G^* = \delta_{(-\infty, 0]^{n_1 n_2}}$. Although

our algorithms guarantee $x^{i+1} \in [0, \infty)^{n_1 n_2}$, the conjugate will cause the true (non-Lagrangian) duality gap

$$(5.2) \quad \tilde{\mathcal{G}}(x, y) := G(x) + F(Kx) + G^*(-K^*y) + F^*(y) \geq \mathcal{G}(x, y; \hat{x}, \hat{y})$$

to be infinite in practice. However, we wish to report the true duality gap instead of the Lagrangian duality gap, as it does not depend on knowing a solution (\hat{x}, \hat{y}) . This is why we have added the upper bound $x \leq 1$. Any greater upper bound would also work, giving a slightly different duality gap.

5.1. Data. We performed the numerical experiments on the first two of our models on the parrot image (#23) from the free Kodak image suite, depicted in Figure 5.1(a) together with the corrupted data and restored images for the test problems. We also performed some experiments (see Figure 5.8) on all 24 images of this image suite. However, the effect of the exact image on the ranking of the tested algorithms is generally small. The size of all the images is $n_1 \times n_2 = 768 \times 512$. To study scalability, we also scaled it down to $n_1 \times n_2 = 192 \times 128$ pixels. Together with the dual variable, the problem dimensions are therefore $768 \cdot 512 \cdot 3 = 1179648 \simeq 10^6$ and $128 \cdot 128 \cdot 3 = 49152 \approx 4 \cdot 10^4$.

For the denoising problem we added Gaussian noise with standard deviation 51 (−13.9dB) to the original test image. To remove the noise, we first choose $\beta = 0.2$ (low regularization parameter), and then $\beta = 1$ (high regularization parameter). Following [12], we scale this parameter by the factor 0.25 for operations on the downscaled image. We also added noise in the other test problems to avoid *inverse crimes* [19]. For sparse Fourier inversion, we used the same level of noise as for denoising. The sparse Fourier inversion experiments are only performed on the original nondownscaled image with the regularization parameter $\beta = 0.1$ (sparse Fourier inversion).

For the PET problem, instead of photographs, we use the Shepp–Logan phantom in Figure 5.2. This is because the limited number of angles encoded in T (reduction of

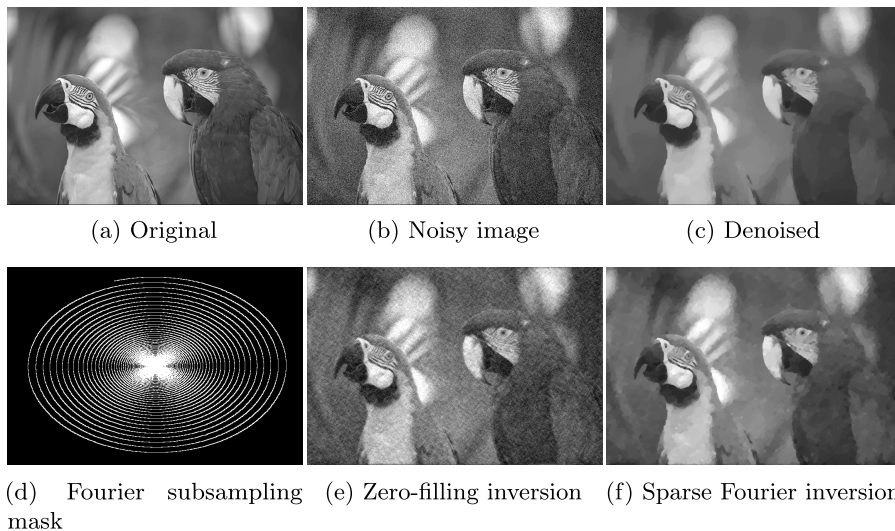


FIG. 5.1. Input data and reconstructions. The original image is #23 from the free Kodak image suite, available online at the time of writing at <http://r0k.us/graphics/kodak/>. Since raw data z for the sparse Fourier inversion is not visually informative, (e) displays the naïve zero-filling inversion $\mathcal{F}^* S^* z$ for the subsampling operator S corresponding to the spiral mask in (d).

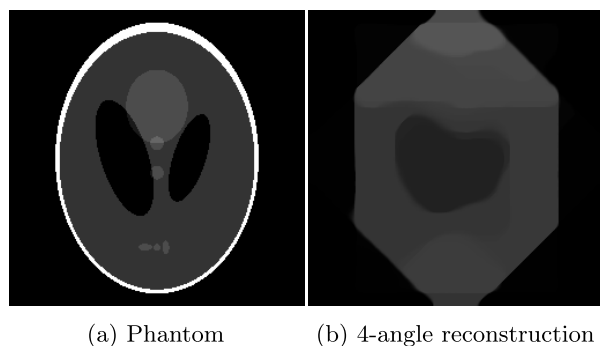


FIG. 5.2. Shepp-Logan brain phantom and its reconstruction from simulated 4-angle (0° , 45° , 90° , 135°) PET. The 4-angle tomography of a 256×256 image consists of 1534 data points, meaning the reconstruction is achieved with just 2.3% of data.

data to a mere 2.3% for the phantom) would not give a recognizable reconstruction of a more complex image. Moreover, the phantom is more relevant to the problem in question. As the resolution, we take $n_1 \times n_2 = 256 \times 256$. To obtain the simulated measurement data b , we apply Poisson noise to the row, column, and diagonal sums in Tx , and then add the background $c := \mathbb{1}$.

5.2. Algorithmic setup. We compare our algorithm (IC-PDPS) to the basic PDPS of [8] and the basic inertial (I-PDPS) and overrelaxed (R-PDPS) variants from [9]. The latter is essentially the Vũ-Comdat algorithm. We do not include FISTA and other non-primal-dual algorithms in our comparisons: of all of our example problems, they apply easily only to TV denoising in its dual form. Similarly, the basic ADMM [14] requires difficult inversions for our problems. Its more efficient preconditioned variant [27], on the other hand, is equivalent to the PDPS [4].

We use the same initial choices of $\tau_0 = 9.9/L$ and $\sigma_0 = 0.1/L$ with $L := \sqrt{8} \geq \|D\|$ [6] for all algorithms and the denoising and sparse Fourier inversion model problems. For the PET problem we take $\sigma_0 = 30/L'$ and $\tau_0 = 0.033/L'$ for an estimate $L' \geq \sqrt{\|T\|^2 + L^2}$. The ratio between τ_0 and σ_0 has been hand-optimized for the baseline PDPS. For the R-PDPS we take the additional overrelaxation parameter $\rho = 1.5$. For the I-PDPS we use fixed inertial parameter $\alpha = 0.9/3$: according to [9], the sequence of parameters $\{\alpha_i\}_{i \in \mathbb{N}}$ has to be nondecreasing with $\alpha_i < 1/3$. We also tested the FISTA rule, which did in practice yield better results for TV denoising but completely failed for the other problems. Hence we use the provably convergent fixed parameter.

The denoising problem is strongly convex with factor $\gamma = 1$, so we include results for both the unaccelerated and accelerated versions of the PDPS and IC-PDPS (Theorems 4.5 and 4.6). We also apply the rules of Theorem 4.7 to the problem with the primal and dual variables exchanged. This is denoted “dual IC-PDPS.” The R-PDPS and the I-PDPS cannot with provable convergence be combined with strong convexity based acceleration: trying to do so was the starting point of our research. For acceleration we use $\gamma = 0.5 < 1$, which is the maximal value for which the ergodic gap is known to converge at the rate $O(1/N^2)$ for the PDPS ($\gamma = 1$ only yields convergence of the iterates; see [8, 22, 23, 25]). For IC-PDPS $\gamma = 1$ is allowed, and provably yields convergence of the gap, but in practice yields worse results than $\gamma = 0.5$.

The IC-PDPS has one further parameter: $\epsilon \in [0, 1)$. For denoising and sparse Fourier inversion we generally take $\epsilon = 0.7$, and for PET, $\epsilon = 0.9$. We also report the denoising convergence behavior for $\epsilon = 0.5$ and $\epsilon = 0$ later in Figure 5.5.

For our reporting, we computed a target optimal solution \hat{x} by taking one million iterations of the basic PDPS. However, the convergence of the basic PDPS for sparse Fourier inversion appears to be very slow: judging by the gap in Figure 5.6(a), the IC-PDPS converges much faster, while both the PDPS and I-PDPS flatten out. We therefore computed the target solution for sparse Fourier inversion by taking one million iterations of the IC-PDPS. Note that the target solution is not used to compute the gap; instead of the Lagrangian duality gap (4.3), we report the true duality gap given in (5.2), as this does not depend on knowing a solution (\hat{x}, \hat{y}) .

We report the distance to \hat{x} in decibels $10 \log_{10}(\|x^i - \hat{x}\|^2 / \|\hat{x}\|^2)$, as well as the duality gap, again in decibels relative to the initial gap as $10 \log_{10}(\tilde{\mathcal{G}}(x^i, y^i)^2 / \tilde{\mathcal{G}}(x^0, y^0)^2)$. For the initial iterates we always took $x^0 = 0$ and $y^0 = 0$. The hardware we used was a MacBook Pro with 16 GB RAM and a 2.8 GHz Intel Core i5 CPU. The codes were written in MATLAB+C-MEX.

5.3. Results. The results for TV denoising of the downsampled image are in Figure 5.3 and for the original image in Figure 5.4 and Table 5.1. The latter includes both the high and low values of the regularization parameter β . For the downsampled experiments we only report the lower value of β . The comparison for different values of ϵ for IC-PDPS is moreover in Figure 5.5, for the higher value of β . The results for sparse Fourier inversion are in Figure 5.6 and Table 5.2(a) and for PET in Figure 5.7 and Table 5.2(b). Finally, Figure 5.8 displays for denoising and sparse Fourier inversion the minimum and maximum intervals for the duality gap over all 24 images in the image suite. We have excluded R-PDPS from these results to avoid overcrowding; its performance is comparable to I-PDPS, as can be gleaned from the other figures.

For TV denoising, the unaccelerated IC-PDPS is clearly the worst algorithm, while I-PDPS and R-PDPS slightly improve upon the basic PDPS. As expected from the $O(1/N)$ versus $O(1/N^2)$ convergence rates, all of these methods are significantly worse than the accelerated PDPS, the accelerated IC-PDPS, and the accelerated dual IC-PDPS. For the downsampled image and for low β for the original resolution image, they are all comparable for the gap, but accelerated IC-PDPS somewhat surprisingly has asymptotically better iterate convergence. Of course, judging by the timings in Table 5.1 in particular, the iterations of the IC-PDPS are somewhat more costly, so the basic accelerated PDPS appears the best choice in this case.

For high β , the results are initially similar, but both variants of the accelerated IC-PDPS are asymptotically better than the accelerated PDPS. This suggests that the IC-PDPS might perform better when there is “more work to be done.” This is somewhat confirmed by the results for sparse Fourier inversion, which is a significantly

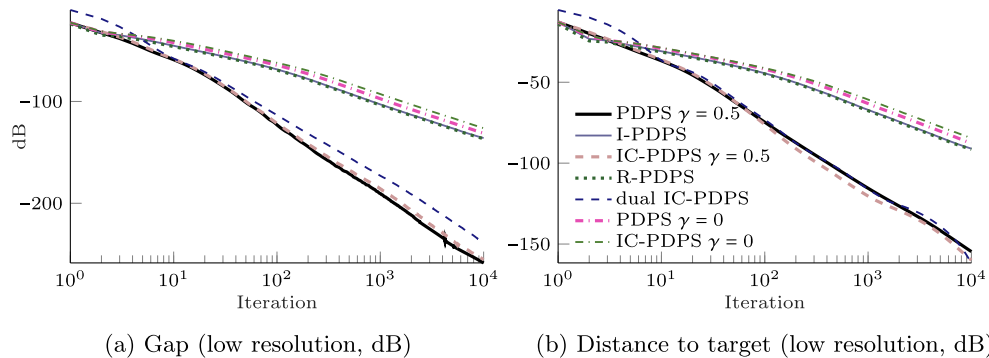


FIG. 5.3. Denoising convergence behavior for low-resolution image.

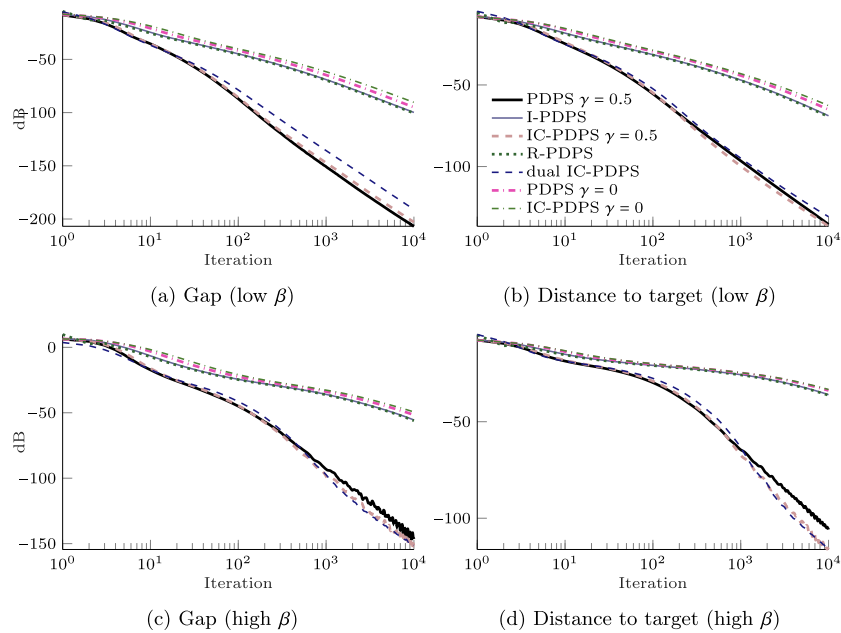


FIG. 5.4. Denoising convergence behavior.

TABLE 5.1

Denoising performance: CPU time and number of iterations (at a resolution of 10 after 100 iterations) to reach given duality gap and distance to target. The dashes indicate that the algorithm never reached (within the maximum number of iterations) the corresponding quality.

(a) Low regularization								
Method	gap ≤ -40 dB		gap ≤ -90 dB		tgt ≤ -40 dB		tgt ≤ -90 dB	
	iter	time	iter	time	iter	time	iter	time
PDPS $\gamma = 0.5$	14	0.47s	120	4.27s	38	1.33s	690	24.74s
I-PDPS	58	2.39s	4870	203.75s	400	16.70s	—	—
IC-PDPS $\gamma = 0.5$	14	0.75s	120	6.87s	40	2.25s	590	33.99s
R-PDPS	55	2.36s	4630	202.46s	380	16.58s	—	—
dual IC-PDPS	13	0.64s	160	8.53s	43	2.25s	750	40.18s
PDPS $\gamma = 0$	82	2.87s	6950	245.99s	560	19.79s	—	—
IC-PDPS $\gamma = 0$	99	5.00s	9710	494.99s	650	33.09s	—	—

(b) High regularization								
Method	gap ≤ -40 dB		gap ≤ -90 dB		tgt ≤ -40 dB		tgt ≤ -90 dB	
	iter	time	iter	time	iter	time	iter	time
PDPS $\gamma = 0.5$	70	2.35s	890	30.34s	250	8.50s	4330	147.73s
I-PDPS	1810	70.74s	—	—	—	—	—	—
IC-PDPS $\gamma = 0.5$	73	3.85s	740	39.51s	270	14.38s	2740	146.43s
R-PDPS	1720	68.07s	—	—	—	—	—	—
dual IC-PDPS	91	4.78s	770	40.85s	330	17.48s	2560	135.94s
PDPS $\gamma = 0$	2580	84.57s	—	—	—	—	—	—
IC-PDPS $\gamma = 0$	3240	157.19s	—	—	—	—	—	—

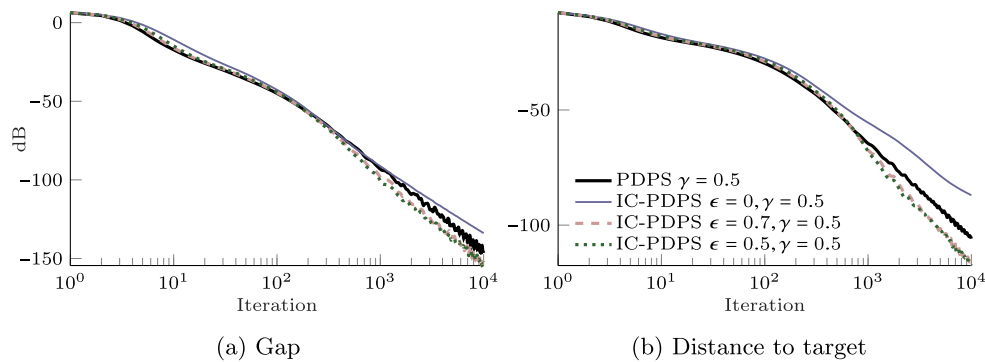


FIG. 5.5. Effect of ϵ on denoising convergence behavior.

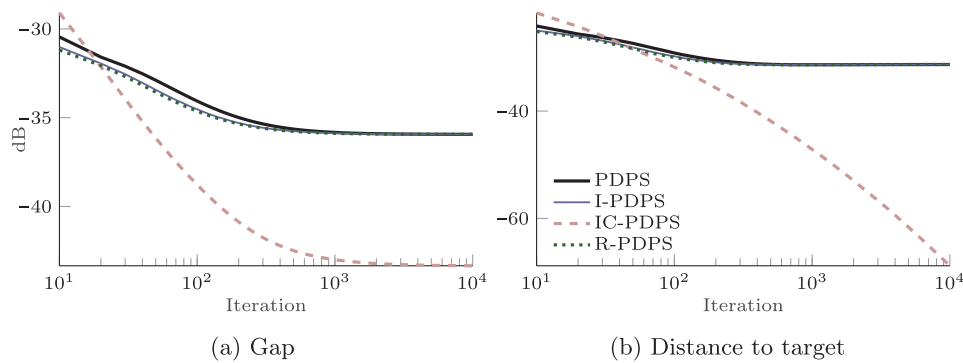


FIG. 5.6. Sparse Fourier inversion convergence behaviour. In (b), the target is computed by taking one million iterations of IC-PDPS instead of PDPS; see subsection 5.2.

TABLE 5.2

Sparse Fourier inversion and PET performance: CPU time and number of iterations (at a resolution of 10) to reach given duality gap and distance to target. The dashes indicate that the algorithm never reached (within the maximum number of iterations) the corresponding quality.

(a) Sparse Fourier inversion					(b) PET				
Method	gap ≤ -35 dB		tgt ≤ -35 dB		Method	gap ≤ -40 dB		tgt ≤ -30 dB	
	iter	time	iter	time		iter	time	iter	time
PDPS	210	15.96s	—	—	PDPS	3740	43.31s	6190	71.68s
I-PDPS	150	11.87s	—	—	I-PDPS	3220	48.72s	4380	66.28s
IC-PDPS	40	3.90s	180	17.88s	IC-PDPS	2180	30.94s	6010	85.32s
R-PDPS	140	12.12s	—	—	R-PDPS	—	—	4150	57.32s

more difficult problem than TV denoising. There the gap convergence performance of IC-PDPS is significantly better than PDPS or I-PDPS: according to Table 5.2(a), compared to the PDPS only 75% of the computational time is required to obtain -35 dB gap reduction.

For the PET problem, Figure 5.7 and Table 5.2(b) indicate that IC-PDPS has good gap convergence behavior, taking 30% less time than the PDPS to reach -40 dB, but has primal variable convergence behavior comparable to the PDPS. This indicates that the IC-PDPS has good convergence of the dual variable.

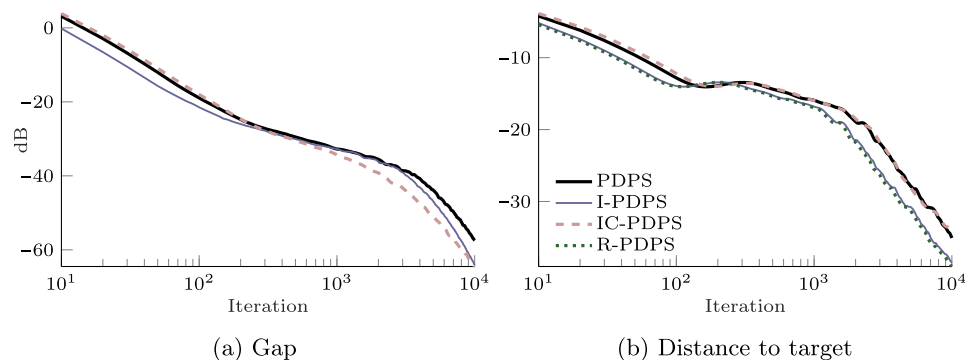


FIG. 5.7. Convergence behavior for the PET example problem.

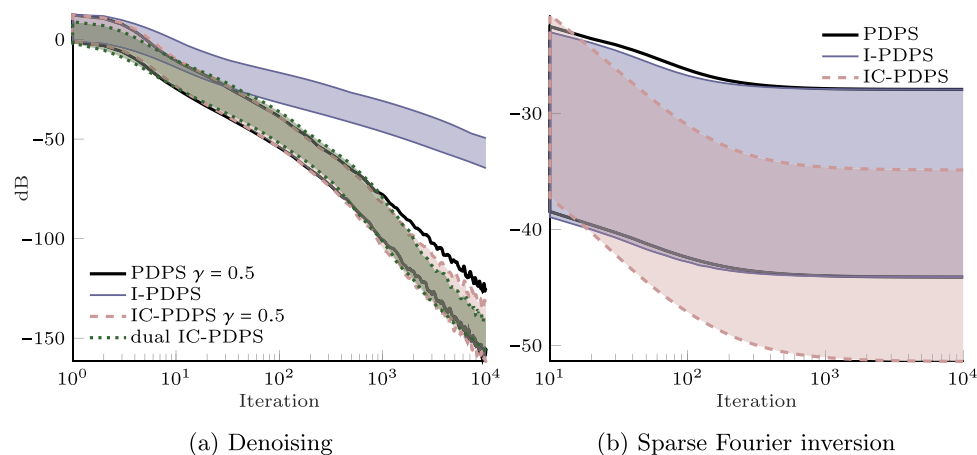


FIG. 5.8. Gap convergence behavior over multiple images (24). The filled areas indicate on each iteration the minimum and maximum gaps (dB) over all the images.

From Figure 5.8 we can see that the exact image does not significantly alter the rankings of the algorithms, with IC-PDPS performing significantly better than the other methods for sparse Fourier inversion.

5.4. Conclusion. While our proposed IC-PDPS does not always improve upon the basic, inertial, and overrelaxed PDPS, it never does significantly worse by iteration count. For some problems, such as sparse Fourier inversion and PET, it offers improved performance. Moreover, we have theoretically guaranteed the $O(1/N)$ convergence of the Lagrangian gap functional or the $O(1/N^2)$ convergence of the strong convexity adjusted gap $\mathcal{G}_{\gamma,\rho}$. This is better than the merely ergodic convergence known of the PDPS and the basic inertial and overrelaxed variants.

A data statement for the EPSRC. Our algorithm implementations and the publicly available test images have been archived on Zenodo at [24].

REFERENCES

- [1] F. ALVAREZ AND H. ATTOUCH, *An inertial proximal method for maximal monotone operators via discretization of a nonlinear oscillator with damping*, Set-Valued Var. Anal., 9 (2001), pp. 3–11, <https://doi.org/10.1023/A:1011253113155>.

- [2] H. ATTOUCH AND A. CABOT, *Convergence rates of inertial forward-backward algorithms*, SIAM J. Optim., 28 (2018), pp. 849–874, <https://doi.org/10.1137/17M1114739>.
- [3] A. BECK AND M. TEOULLE, *A fast iterative shrinkage-thresholding algorithm for linear inverse problems*, SIAM J. Imaging Sci., 2 (2009), pp. 183–202, <https://doi.org/10.1137/080716542>.
- [4] M. BENNING, F. KNOLL, C.-B. SCHÖNLIEB, AND T. VALKONEN, *Preconditioned ADMM with nonlinear operator constraint*, in System Modeling and Optimization: 27th IFIP TC 7 Conference, L. Bociu, J.-A. Désidéri, and A. Habbal, eds., Springer, New York, 2016, pp. 117–126, https://doi.org/10.1007/978-3-319-55795-3_10.
- [5] K. BREDIES AND H. P. SUN, *Preconditioned Douglas–Rachford algorithms for TV- and TGV-regularized variational imaging problems*, J. Math. Imaging Vision, 52 (2015), pp. 317–344, <https://doi.org/10.1007/s10851-015-0564-1>.
- [6] A. CHAMBOLLE, *An algorithm for total variation minimization and applications*, J. Math. Imaging Vision, 20 (2004), pp. 89–97, <https://doi.org/10.1023/B:JMIV.0000011325.36760.1e>.
- [7] A. CHAMBOLLE AND C. DOSSAL, *On the convergence of the iterates of “FISTA,”* J. Optim. Theory Appl., 16 (2015), <https://hal.inria.fr/hal-01060130v3>.
- [8] A. CHAMBOLLE AND T. POCK, *A first-order primal-dual algorithm for convex problems with applications to imaging*, J. Math. Imaging Vision, 40 (2011), pp. 120–145, <https://doi.org/10.1007/s10851-010-0251-1>.
- [9] A. CHAMBOLLE AND T. POCK, *On the ergodic convergence rates of a first-order primal-dual algorithm*, Math. Program., 159 (2016), pp. 253–287, <https://doi.org/10.1007/s10107-015-0957-3>.
- [10] R. H. CHAN, S. MA, AND J. YANG, *Inertial Primal-Dual Algorithms for Structured Convex Optimization*, <https://arxiv.org/abs/1409.2992>, 2014.
- [11] L. CONDAT, *A primal-dual splitting method for convex optimization involving lipschitzian, proximable and linear composite terms*, J. Optim. Theory Appl., 158 (2013), pp. 460–479, <https://doi.org/10.1007/s10957-012-0245-9>.
- [12] J. C. DE LOS REYES, C.-B. SCHÖNLIEB, AND T. VALKONEN, *Bilevel parameter learning for higher-order total variation regularisation models*, J. Math. Imaging Vision, 57 (2017), pp. 1–25, <https://doi.org/10.1007/s10851-016-0662-8>.
- [13] J. ECKSTEIN AND D. BERTSEKAS, *On the Douglas–Rachford splitting method and the proximal point algorithm for maximal monotone operators*, Math. Program., 55 (1992), pp. 293–318, <https://doi.org/10.1007/BF01581204>.
- [14] D. GABAY, *Applications of the method of multipliers to variational inequalities*, in Augmented Lagrangian Methods: Applications to the Numerical Solution of Boundary-Value Problems, M. Fortin and R. Glowinski, eds., Stud. Math. Appl. 15, North-Holland, Amsterdam, 1983, pp. 299–331.
- [15] B. HE AND X. YUAN, *Convergence analysis of primal-dual algorithms for a saddle-point problem: From contraction perspective*, SIAM J. Imaging Sci., 5 (2012), pp. 119–149, <https://doi.org/10.1137/100814494>.
- [16] D. KIM AND J. A. FESSLER, *Another look at the fast iterative shrinkage/thresholding algorithm (FISTA)*, SIAM J. Optim., 28 (2018), pp. 223–250, <https://doi.org/10.1137/16M108940X>.
- [17] F. KNOLL, K. BREDIES, T. POCK, AND R. STOLLBERGER, *Second order total generalized variation (TGV) for MRI*, Magnetic Resonance Medicine, 65 (2011), pp. 480–491, <https://doi.org/10.1002/mrm.22595>.
- [18] D. LORENZ AND T. POCK, *An inertial forward-backward algorithm for monotone inclusions*, J. Math. Imaging Vision, 51 (2015), pp. 311–325, <https://doi.org/10.1007/s10851-014-0523-2>.
- [19] J. L. MUELLER AND S. SILTANEN, *Linear and Nonlinear Inverse Problems with Practical Applications*, Comput. Sci. Eng. 10, SIAM, Philadelphia, PA, 2012, <https://doi.org/10.1137/1.9781611972344>.
- [20] Y. NESTEROV, *A method of solving a convex programming problem with convergence rate $O(1/k^2)$* , Soviet Math. Dokl., 27 (1983), pp. 372–376.
- [21] P. PATRINOS, L. STELLA, AND A. BEMPORAD, *Douglas–Rachford splitting: Complexity estimates and accelerated variants*, in Proceedings of the 53rd IEEE Conference on Decision and Control, 2014, pp. 4234–4239, <https://doi.org/10.1109/CDC.2014.7040049>.
- [22] T. VALKONEN, *Testing and non-linear preconditioning of the proximal point method*, Appl. Math. Optim. (2018), <https://doi.org/10.1007/s00245-018-9541-6>.
- [23] T. VALKONEN, *Block-proximal methods with spatially adapted acceleration*, Electron. Trans. Numer. Anal., 51 (2019), pp. 15–49, <https://doi.org/10.1553/etna.vol51s15>.
- [24] T. VALKONEN, *Source Codes for “Inertial, Corrected, Primal-Dual Proximal Splitting,”* <https://doi.org/10.5281/zenodo.3531934> (2019).
- [25] T. VALKONEN AND T. POCK, *Acceleration of the PDHGM on partially strongly convex functions*, J. Math. Imaging Vision, 59 (2017), pp. 394–414, <https://doi.org/10.1007/s10851-016-0692-2>.

- [26] B. C. VŮ, *A splitting algorithm for dual monotone inclusions involving cocoercive operators*, Adv. Comput. Math., 38 (2013), pp. 667–681, <https://doi.org/10.1007/s10444-011-9254-8>.
- [27] X. ZHANG, M. BURGER, AND S. OSHER, *A unified primal-dual algorithm framework based on Bregman iteration*, J. Sci. Comput., 46 (2011), pp. 20–46, <https://doi.org/10.1007/s10915-010-9408-8>.