

**PRECONDITIONED ACCELERATED GRADIENT DESCENT  
METHODS FOR LOCALLY LIPSCHITZ SMOOTH OBJECTIVES  
WITH APPLICATIONS TO THE SOLUTION OF NONLINEAR  
PDES**

JEA-HYUN PARK, ABNER J. SALGADO, AND STEVEN M. WISE

ABSTRACT. We develop a theoretical foundation for the application of Nesterov’s accelerated gradient descent method (AGD) to the approximation of solutions of a wide class of partial differential equations (PDEs). This is achieved by proving the existence of an *invariant set* and exponential convergence rates when its preconditioned version (PAGD) is applied to minimize *locally Lipschitz* smooth, strongly convex objective functionals. We introduce a second-order ordinary differential equation (ODE) with a preconditioner built-in and show that PAGD is an explicit time-discretization of this ODE, which requires a natural time step restriction for energy stability. At the continuous time level, we show an exponential convergence of the ODE solution to its steady state using a simple energy argument. At the discrete level, assuming the aforementioned step size restriction, the existence of an invariant set is proved and a matching exponential rate of convergence of the PAGD scheme is derived by mimicking the energy argument and the convergence at the continuous level. Applications of the PAGD method to numerical PDEs are demonstrated with certain nonlinear elliptic PDEs using pseudo-spectral methods for spatial discretization, and several numerical experiments are conducted. The results confirm the global geometric and mesh size-independent convergence of the PAGD method, with an accelerated rate that is improved over the preconditioned gradient descent (PGD) method.

## 1. INTRODUCTION

The purpose of this work is to broaden the context in which a well-known and efficient algorithm for unconstrained convex minimization, the so-called *Nesterov’s accelerated gradient descent* (AGD) scheme, can be utilized and, further, to shed some light on its convergence properties. This method iteratively finds approximations to the solution of the following optimization problem: given  $G : \mathbb{H} \rightarrow \mathbb{R}$ , find

$$x^* = \operatorname{argmin} \{G(x) \mid x \in \mathbb{H}\}.$$

Here, and in what follows,  $\mathbb{H}$  is a real, separable Hilbert space with inner product  $(\cdot, \cdot)_{\mathbb{H}}$  and the so-called *objective* functional,  $G$ , is assumed to be strongly convex and *locally Lipschitz* smooth; see Section 2 for definitions and notation. We immediately comment that the assumptions on the objective guarantee the existence and uniqueness of a minimizer (e.g., [12, Theorem 7.4-4, Theorem 8.2-2]).

---

2010 *Mathematics Subject Classification.* 65B99, 65J08, 65N35, 65K10.

*Key words and phrases.* Preconditioning, Nesterov Acceleration, Momentum Method, Convex Optimization, Nonlinear Elliptic Partial Differential Equations, Pseudo-Spectral Methods, Lyapunov.

Convex minimization is ubiquitous, and our main interest in this problem comes from the fact that many important nonlinear partial differential equations (PDEs) can be viewed as the Euler equations of certain convex objective functions. For example, the classical minimal surface problem (see [14]) and the  $p$ -Laplacian equation (see [4]) have this structure, just to name a few. However, the method discussed in this work is particularly powerful for semilinear PDEs. For example, time discretizations of many important models in material science often end up involving such problems, e.g., via the convex splitting technique (see [16]). In a related context, the current explosion of interest in statistical learning has drawn the attention of practitioners to so-called first order schemes, i.e., those that only require knowledge of first order derivatives, which is suitable in dealing with large data sets. These considerations are important for solving nonlinear PDE as well. One of the main thrusts of this research is to show that Nesterov's accelerated schemes, which are popular in statistical learning, can be utilized as fast solvers for nonlinear PDE once we resolve some nontrivial, technical difficulties specific to such problems.

The first and most naïve approach to find  $x^*$  would be to appeal directly to the first order necessary (and, in this context, sufficient) optimality condition, namely, the Euler equation:

$$(1.1) \quad G'(x^*) = 0,$$

where  $G'$  denotes the Fréchet derivative of  $G$ . In the examples that we have in mind, however, this requires the simultaneous solution to a very large number of nonlinear equations, and the direct solution of the system is not feasible in practice. Other approaches better suited for minimization must be constructed. According to [7], iterative methods for minimizing functionals date back in 1847 when Cauchy proposed the so-called gradient descent method (GD). The solution to (1.1) can be seen as the steady state of the *gradient flow*

$$(1.2) \quad X(0) = x_0, \quad \dot{X}(t) = -G'(X(t)), \quad t > 0.$$

Here  $x_0 \in \mathbb{H}$  is arbitrary and, in the second equation, we are implicitly identifying the dual space of  $\mathbb{H}$ , denoted by  $\mathbb{H}'$ , with  $\mathbb{H}$  itself. Under the assumptions we have imposed on the objective  $G$ , it is possible to show that this flow satisfies  $X(t) \rightarrow x^*$  as  $t \rightarrow \infty$ , see [30, Theorem 2.4]. The idea of GD is to approximate the solution to this flow via a *forward Euler* time discretization with a fixed step size  $s$ : given  $x_0 \in \mathbb{H}$ , for  $k \geq 0$ , find  $x_{k+1}$  satisfying

$$(1.3) \quad x_{k+1} = x_k - sG'(x_k).$$

While this idea seems straightforward, more in-depth discussions on this method started only in the 1960s, where some practical step size rules and convergence analyses were established. It was shown that if the objective functional is convex and Lipschitz smooth, then GD converges to the minimizer,  $x^*$ , and it exhibits a first order rate of convergence in the objective. Here, and in what follows, by an  $n^{\text{th}}$  *order (algebraic) convergence* in the objective, we mean that  $G(x_k) - G^* \leq O(1/k^n)$ , as  $k \rightarrow \infty$ , where  $G^* = G(x^*)$  is the minimum of  $G$ . By an *exponential* or a *geometric convergence* in the objective we mean that  $G(x_k) - G^* \leq O(r^k)$ , as  $k \rightarrow \infty$ , for some  $r \in (0, 1)$ . In the latter case, we call  $r$  the rate of (exponential) convergence. It can further be shown that, if the objective is, in addition, strongly convex, then the rate of convergence is exponential, and that it matches the rate of convergence of the solution of (1.2) to  $x^*$ . (See [22, Theorem 2.1.15] or Remark 5.11). Some

physical intuition for the evolution of the solution to (1.2) is provided in Section 4. See, in particular, Remark 4.1.

Evidently, all considerations regarding convergence are subject to the norm  $\|\cdot\|_{\mathbb{H}}$ . It is possible to improve the convergence rate by using an equivalent norm, through which the level sets of the objective  $G$  look “more circular.” In the numerical linear algebra and numerical PDE communities, this is commonly known as *preconditioning*. In the context of (1.2) and GD, this is achieved by introducing an operator  $\mathcal{L} : \mathbb{H} \rightarrow \mathbb{H}'$  and considering the evolution of  $\dot{X}(t) = -\mathcal{L}^{-1}G'(X(t))$ . Notice that we no longer implicitly identify  $\mathbb{H}'$  with  $\mathbb{H}$ . The time-discrete counterpart of (1.3) is known as the *preconditioned gradient descent method* (PGD) and is as follows: given  $x_0 \in \mathbb{H}$ , for  $k \geq 0$ , find  $x_{k+1}$  such that

$$x_{k+1} = x_k - s\mathcal{L}^{-1}G'(x_k).$$

If the preconditioner is suitably chosen, then the convergence rate of GD can be substantially improved (see [15]). Note that we will tacitly assume in the sequel that  $\mathcal{L}$  is independent of the iteration index  $k$ . We remark that Newton’s method may be viewed as a kind of *generalized* preconditioned gradient descent method if we assume that  $G$  is twice Fréchet differentiable and allow for the possibility that the preconditioner can change at each iteration. In particular, Newton’s method is expressed as

$$G''(x_k)(x_{k+1} - x_k) = -G'(x_k) =: r_k,$$

where  $G''(x_k)$  is the second Fréchet derivative of  $G$ , and  $r_k$  is the so-called residual. Then, Newton’s method is a generalized preconditioned gradient descent method for which the preconditioner satisfies  $s\mathcal{L}_k = G''(x_k)$ . One of the difficulties with Newton’s method is that the preconditioner constantly changes, in general, and must be recomputed and re-inverted at each iteration step, which can prove quite costly. Furthermore,  $G''$  may not exist in all applications of interest. Indeed, in the sequel, we will not assume that  $G''$  exists.

To improve the convergence rate of GD, Nesterov [21] suggested a scheme that accelerates the GD method. For convex and Lipschitz smooth objectives the Nesterov’s accelerated gradient descent (AGD) scheme achieves a second order convergence rate. Later, he showed that if the objective is, in addition, strongly convex, then AGD achieves a faster exponential convergence rate than GD (see [22, Theorem 2.1.15]).

However, while the GD scheme has a strong physical intuition behind it, it is not completely clear what mechanism is at play to provide an acceleration in the AGD scheme. Some attempts have been made to understand this in the literature. Attouch et al. [3] studied asymptotic behaviors of the solutions to a heavy ball system (similar to (4.1)) and showed their convergence to minimizers (if they exist) of locally Lipschitz objectives that are bounded below at the continuous time level. Goudou and Munier [17] looked into a similar system with quasiconvex, locally Lipschitz objectives at the continuous time level and the convergence of its implicit discretization, what they call *proximal inertial algorithm*. Apparently, [29] is the first work that explains the acceleration happening in AGD both quantitatively and intuitively and inspired many researchers including us. For convex, Lipschitz smooth objectives, they were able to show that the solutions to a second order ODE  $\ddot{X} + \frac{3}{t}\dot{X} + \nabla G(X) = 0$  converges to the set of minimizers of  $G$  quadratically fast, as  $t \rightarrow \infty$ , and a matching, discrete convergence rate was established for a version

of AGD. Wibisono et al. [30] took a similar approach in more generality in the language of *Bregman Lagrangian flow*. However, these two works did not explain the exponential acceleration for the strongly convex objectives. This limitation was one of the motivations of our work and removing it is one of the goals of this paper. Recently, there have appeared more works that address the same issue and provide more general or unifying frameworks. For (globally)  $L$ -smooth,  $\mu$ -strongly convex objectives (see Section 2 for definitions and notation), the best known convergence rates of AGD and the associated ODE are  $G(x_k) - G^* \leq O((1 - \sqrt{\mu/L})^k)$  as  $k \rightarrow \infty$  and  $G(X(t)) - G^* \leq O(e^{-\sqrt{\mu}t})$  as  $t \rightarrow \infty$  respectively. Similarly, for (globally)  $L$ -smooth, convex objectives,  $G(x_k) - G^* \leq O(1/k^2)$  as  $k \rightarrow \infty$  and  $G(X(t)) - G^* \leq O(1/t^2)$  as  $t \rightarrow \infty$  are the best known convergence rates respectively. Within the same framework as in [30], but using a different Lyapunov function, Wilson et al. [31] showed the best convergence rates for both convex and strongly convex cases at both continuous and discrete level. Shi et al. [27] looked into what they call *high-resolution ODE* and provided a finer understanding about the momentum-type schemes and obtained similar results. The convergence rates for the strongly convex case that they derived were not the best. However, they were able to explain the difference in the performance of Polyak's momentum method and AGD at the continuous level. Siegel [28] also analyzed a system of ODEs to study a version of AGD and obtained the best known rate of convergence for the strongly convex case at both continuous and discrete levels. He also studied non-smooth but still strongly convex objectives and stochastic versions. Luo and Chen [20] obtained the same best convergence rates for all the four cases mentioned above using a single ODE system but using a time rescaling argument when dealing with the convex case. Laborde and Oberman [19] studied perturbed ODE systems and the corresponding version of AGD in the stochastic framework. As a byproduct, they obtained the same best convergence rates for the perturbed version of AGD with strongly convex objectives in the deterministic setting. There are also related works from a PDE point of view rather than convex optimization. Schaeffer and Hou [25] studied accelerated methods for nonlinear elliptic operators, which may not have a variational structure, i.e., the PDE may not have an appropriate objective. They also proposed similar methods for viscosity solutions. Benyamini et al. [6] and Calder and Yezzi [9] studied *PDE accelerations* that are similar to [30] in spirit and applied them to image processing and minimal surface obstacle problems respectively.

The work contained herein includes the following important contributions, in particular, from a numerical PDE point of view, which can be seen more clearly from the literature comparison of Appendix C. See Table 2.

1. We prove all of our results under the more general assumption that the objective functional is *locally Lipschitz* smooth. Almost all earlier works assume that the objective is *globally Lipschitz* smooth (see e.g., [22, 30, 29, 2, 31, 19, 25, 20, 28, 27]). This is too restrictive to approximate solutions of nonlinear PDEs. If the objective functional associated with the PDE of interest grows just a bit faster than quadratic functionals (i.e., those of very mild nonlinearity), it violates the global Lipschitz condition and is beyond the theoretical guarantee. On the other hand, the local condition does not require anything outside of a certain bounded set so that much more nonlinear PDEs can be dealt with. Only a few works from a dynamical system point of view (e.g., [3, 17]) assume local Lipschitz condition. However, those works address only convergence itself at the continuous time

level and did not discuss discrete level analysis. To the best of our knowledge, this is the first work that provides convergence rates under the local Lipschitz smoothness assumption at the continuous or discrete level.

2. We prove the existence of an invariant set  $\mathfrak{B}$  of the PAGD method. That is, every sequence generated by the scheme stays in a certain bounded set. The local Lipschitz assumption is meaningful when it is furnished with an invariant set so that we have no restriction in exploiting the Lipschitz condition. We emphasize that this is not a trivial technicality. Unlike the gradient descent method, the accelerated methods are not descent methods. In fact, they oscillate. Thus, a simple sublevel set argument does not work. Even worse, they involve extrapolations of the main iterates. Consequently, a naïve attempt to obtain an invariant set leads to an impasse: to control the extrapolations, one wants to use the Lipschitz condition, but under the local Lipschitz condition, one cannot use it before proving that they are in a fixed bounded set. Again, to the best of our knowledge, our work is the first that addresses and resolves this issue.
3. We provide a detailed discrete analysis for a nonlinear PDE. All existing works mentioned before either did not discuss numerical examples or did not explain how concrete numerical examples fit the abstract framework, and they omitted whether their numerical examples satisfy the assumptions that they imposed. In contrast, we show that our examples satisfy all the necessary assumptions.
4. We provide an intuitive explanation for the acceleration mechanism behind AGD for strongly convex and locally Lipschitz smooth objectives. Inspired by [29], we view AGD as a discretization of a certain second order ordinary differential equation (ODE) — we present how to discretize this ODE to obtain AGD — and show that the solution to this ODE converges exponentially fast to its stationary point, which is the minimizer of  $G$ . We also provide an energy based proof of the exponential convergence rate of PAGD. This proof mimics the analysis of the continuous counterpart that is previously developed and shows what dissipation mechanisms are at play to achieve the aforementioned acceleration. We also show that the rates of convergence of the ODE model and AGD match and those rates at the continuous and discrete level are both the best known rates.
5. We build a preconditioner into the problem itself (even at the continuous time level) to analyze the scheme with a preconditioner in an explicit way. This seems deceptively simple. After all, preconditioning is nothing but using a different norm, hence a numerical analysis in one norm implicitly suggests the possibility of a similar analysis in another norm. On the other hand, determining an effective preconditioning strategy is a nontrivial matter. Furthermore, it is our observation that preconditioning in the classical optimization setting, especially for problems related to data analysis and machine learning, is underutilized and is a potential growth area in the future. Likewise, preconditioning strategies related to spectral collocation methods applied to nonlinear PDE are uncommon, but, as we shall see, are effective and efficient solver tools.

This paper is organized as follows. In Section 2, we summarize the notation, assumptions, and main tools that we will use. In Section 3, we introduce several numerical schemes that are closely related to our discussion and summarize their convergence rates. In Section 4, we explore the connection between PAGD and a second order ODE and how this connection can help understand the acceleration behind PAGD intuitively. In Section 5, we prove the existence of an invariant set

for the PAGD scheme and its exponential convergence. We take an ODE inspired approach, whose intuition lies in the developments of Section 4. In Section 6, we illustrate the application of the PAGD method to the solution of some numerical PDEs. These numerical experiments show the improvement in convergence by both acceleration and preconditioning. Finally, in the Appendices, we provide the derivation of the initial value problem (IVP) which corresponds to the limiting case of PAGD and a specific discretization of the IVP that leads to PAGD.

## 2. PRELIMINARIES

Let us begin by introducing the setting, assumptions, and some basic properties of the objects that we are interested in. By  $\mathbb{H}$ , we denote a real and separable Hilbert space with inner product  $(\cdot, \cdot)_{\mathbb{H}}$  and associated norm  $\|\cdot\|_{\mathbb{H}}$ . Since we will use other inner products and norms on  $\mathbb{H}$ , for clarity, we will refer to  $(\cdot, \cdot)_{\mathbb{H}}$  and  $\|\cdot\|_{\mathbb{H}}$  as the *canonical inner product* and *canonical norm*, respectively. The dual of  $\mathbb{H}$  is denoted by  $\mathbb{H}'$ . Its canonical operator norm is denoted by  $\|\cdot\|_{\mathbb{H}'}$ . For  $v \in \mathbb{H}$  and  $f \in \mathbb{H}'$ , the symbol  $\langle f, v \rangle$  represents their duality pairing, that is,  $\langle f, v \rangle = f(v) \in \mathbb{R}$ .

A *preconditioner* is defined by a linear operator  $\mathcal{L} : \mathbb{H} \rightarrow \mathbb{H}'$ . Such an operator induces a bilinear form: for  $x, y \in \mathbb{H}$ ,

$$(2.1) \quad (x, y)_{\mathcal{L}} = \langle \mathcal{L}x, y \rangle = \mathcal{L}[x](y).$$

We further assume that the bilinear form defined in (2.1) satisfies the following properties: there exist  $C_1, C_2 > 0$  such that, for any  $x, y \in \mathbb{H}$ ,

$$(2.2) \quad (x, y)_{\mathcal{L}} = (y, x)_{\mathcal{L}}, \quad (x, y)_{\mathcal{L}} \leq C_2 \|x\|_{\mathbb{H}} \|y\|_{\mathbb{H}}, \quad C_1 \|x\|_{\mathbb{H}}^2 \leq (x, x)_{\mathcal{L}}.$$

Let us state some immediate, but important consequences without proof for the sake of brevity.

**Proposition 2.1** (properties of  $\mathcal{L}$ ). *Let  $\mathbb{H}$  be a real, separable Hilbert space with inner product  $(\cdot, \cdot)_{\mathbb{H}}$ , and suppose that  $\mathcal{L} : \mathbb{H} \rightarrow \mathbb{H}'$  is a linear mapping that satisfies (2.2). Then,  $(\cdot, \cdot)_{\mathcal{L}}$  is an inner product on  $\mathbb{H}$  and the object*

$$\|x\|_{\mathcal{L}} = \sqrt{(x, x)_{\mathcal{L}}}, \quad \forall x \in \mathbb{H},$$

*is a norm, which is, in fact, equivalent to the canonical norm,  $\|\cdot\|_{\mathbb{H}}$ . By the Riesz Representation Theorem,  $\mathcal{L}$  is invertible. The inverse is continuous and, in fact, it is just the Riesz Map with respect to the  $\mathcal{L}$ -inner product, denoted  $\mathfrak{R}_{\mathcal{L}}$ . We write  $\mathcal{L}^{-1} = \mathfrak{R}_{\mathcal{L}} : \mathbb{H}' \rightarrow \mathbb{H}$ . The object*

$$(2.3) \quad (f, g)_{\mathcal{L}^{-1}} = \langle f, \mathcal{L}^{-1}g \rangle, \quad \forall f, g \in \mathbb{H}',$$

*is an inner product on the Hilbert space  $\mathbb{H}'$  and the object*

$$(2.4) \quad \|f\|_{\mathcal{L}^{-1}} = \sqrt{(f, f)_{\mathcal{L}^{-1}}} = \sqrt{\langle f, \mathcal{L}^{-1}f \rangle}, \quad \forall f \in \mathbb{H}',$$

*is a norm. The new norm on  $\mathbb{H}'$  is an operator norm in the sense that*

$$(2.5) \quad \|f\|_{\mathcal{L}^{-1}} = \sup_{0 \neq x \in \mathbb{H}} \frac{\langle f, x \rangle}{\|x\|_{\mathcal{L}}} = \sup_{\substack{x \in \mathbb{H} \\ \|x\|_{\mathcal{L}}=1}} \langle f, x \rangle, \quad \forall f \in \mathbb{H}'.$$

*Finally, we have*

$$(2.6) \quad \|\mathcal{L}^{-1}f\|_{\mathcal{L}} = \|f\|_{\mathcal{L}^{-1}}, \quad \forall f \in \mathbb{H}' \quad \text{and} \quad \|\mathcal{L}x\|_{\mathcal{L}^{-1}} = \|x\|_{\mathcal{L}}, \quad \forall x \in \mathbb{H}.$$

**Remark 2.2** (no preconditioning). By setting  $\mathcal{L} = \mathfrak{R}_{\mathbb{H}}^{-1}$ , we can remove the preconditioning, where  $\mathfrak{R}_{\mathbb{H}} : \mathbb{H}' \rightarrow \mathbb{H}$  is the canonical Riesz map. Hence, PAGD is a generalization of AGD.  $\blacksquare$

Our objective  $G : \mathbb{H} \rightarrow \mathbb{R}$  will be assumed to be Fréchet differentiable at every point in  $\mathbb{H}$ . We denote by  $G'(x) \in \mathbb{H}'$  the Fréchet derivative of  $G$  at the point  $x \in \mathbb{H}$ . Since the definition of Fréchet differentiability involves a norm, the actual derivative is possibly norm dependent. The following result shows that, actually, the definition is invariant as long as the norms are equivalent.

**Proposition 2.3** (equivalent norms). *Let  $\mathbb{H}$  be a real and separable Hilbert space with norm  $\|\cdot\|_{\mathbb{H}}$ , and  $G : \mathbb{H} \rightarrow \mathbb{R}$  be Fréchet differentiable at  $x \in \mathbb{H}$ . Assume that  $\|\|\cdot\|\|_{\mathbb{H}}$  is another norm on  $\mathbb{H}$ . If  $\|\|\cdot\|\|_{\mathbb{H}}$  is equivalent to  $\|\cdot\|_{\mathbb{H}}$ , then  $G$  is also Fréchet differentiable at  $x$  with respect to  $\|\|\cdot\|\|_{\mathbb{H}}$ . Furthermore, the derivatives coincide.*

Notice that nothing is said about continuity in the previous statement. For convex functions, the continuity of the derivatives is automatic once the Fréchet differentiability is guaranteed (see [23, p. 20 Corollary]).

**Proposition 2.4** (continuity). *Let  $\mathbb{H}$  be a real and separable Hilbert space and  $D \subset \mathbb{H}$  be open and convex. If  $G : D \rightarrow \mathbb{R}$  is convex and Fréchet differentiable, then  $x \mapsto G'(x)$  is norm continuous on  $D$ .*

The following two definitions provide a framework to describe the geometry of the graph of our objective functional.

**Definition 2.5** (Lipschitz smoothness). Let  $\mathbb{H}$  be a real and separable Hilbert space, and  $G : \mathbb{H} \rightarrow \mathbb{R}$  be Fréchet differentiable at every point. We say that  $G$  is *locally Lipschitz smooth* (with respect to  $\mathcal{L}$ -norm) iff, for every bounded, convex set  $B \subset \mathbb{H}$ , there exists a constant  $L_B > 0$  such that

$$(2.7) \quad \langle G'(x) - G'(y), x - y \rangle \leq L_B \|y - x\|_{\mathcal{L}}^2 \quad \forall x, y \in B.$$

For brevity, we say that  $G$  is  $L_B$ -smooth on  $B$ . If the constant  $L_B = L > 0$  can be chosen to be independent of  $B$ , then we say that  $G$  is *globally Lipschitz smooth with a constant  $L$* , or simply  $L$ -smooth.

**Remark 2.6** (terminology). The above definition is a weaker notion than the local Lipschitz continuity of the Fréchet derivative of  $G$ , which is given by

$$(2.8) \quad \|G'(x) - G'(y)\|_{\mathcal{L}^{-1}} \leq L_B \|x - y\|_{\mathcal{L}} \quad \forall x, y \in B,$$

for some  $L_B > 0$ . Of course, this implies the local Lipschitz smoothness of  $G$  (2.7). In this paper, to avoid confusion, whenever (2.8) holds, we will say that  $G$  is *locally Lipschitz smooth in the strong sense* or that  $G'$  is *locally Lipschitz in the strong sense*. We need this stronger condition when we conduct the continuous level analysis (Section 4). Note, however, for convex functions, the global versions of the two definitions are equivalent. That is, if  $B = \mathbb{H}$ , (2.7) implies (2.8) (see [22, Theorem 2.1.5 (2.1.8)]).  $\blacksquare$

**Definition 2.7** (strong convexity). Let  $G : \mathbb{H} \rightarrow \mathbb{R}$  be Fréchet differentiable. We say that  $G$  is  $\mu$ -strongly convex (with respect to  $\mathcal{L}$ -norm) iff there exists a constant  $\mu > 0$  such that

$$(2.9) \quad \langle G'(x) - G'(y), x - y \rangle \geq \mu \|y - x\|_{\mathcal{L}}^2 \quad \forall x, y \in \mathbb{H}.$$

We now state an equivalent characterization of these notions.

**Theorem 2.8** (equivalence). *Let  $\mathbb{H}$  be a real and separable Hilbert space, and  $G : \mathbb{H} \rightarrow \mathbb{R}$  be Fréchet differentiable.  $G$  is  $L_B$ -smooth on the bounded convex set  $B \subset \mathbb{H}$  if and only if*

$$(2.10) \quad G(y) - G(x) - \langle G'(x), y - x \rangle \leq \frac{L_B}{2} \|y - x\|_{\mathcal{L}}^2 \quad \forall x, y \in B.$$

Similarly,  $G$  is  $\mu$ -strongly convex if and only if

$$(2.11) \quad G(y) - G(x) - \langle G'(x), y - x \rangle \geq \frac{\mu}{2} \|y - x\|_{\mathcal{L}}^2 \quad \forall x, y \in \mathbb{H}.$$

*Proof.* These results follow from Taylor's Theorem with integral remainder. See also [22, Theorem 2.1.5, Theorem 2.1.9].  $\square$

Among the two characterizations of Lipschitz smoothness and strong convexity stated above, we will call (2.10) and (2.11) the upper and the lower *quadratic trap* of  $G$ , respectively. The constant  $\frac{L_B}{\mu}$  is called the (*local*) *condition number* of the objective functional  $G$  with respect to the  $\mathcal{L}$ -norm. In what follows, we will use its reciprocal, denoted by  $\rho = \frac{\mu}{L_B} \in (0, 1]$ , to quantify rates of convergence. Note that the condition number crucially depends on the norm that is used to describe the geometry of the graph of  $G$ . Choosing a good preconditioner,  $\mathcal{L}$ , is at the heart of much of scientific computing.

We conclude this section by stating a pair of well-known identities which we will use frequently. For any  $A, B \in \mathbb{H}$ ,

$$(2.12) \quad (A, B)_{\mathcal{L}} = \frac{1}{2} \|A\|_{\mathcal{L}}^2 + \frac{1}{2} \|B\|_{\mathcal{L}}^2 - \frac{1}{2} \|A - B\|_{\mathcal{L}}^2$$

$$(2.13) \quad = \frac{1}{2} \|A + B\|_{\mathcal{L}}^2 - \frac{1}{2} \|A\|_{\mathcal{L}}^2 - \frac{1}{2} \|B\|_{\mathcal{L}}^2.$$

### 3. OPTIMIZATION SCHEMES

Here we briefly review several algorithms that are closely related to our main algorithm of interest. To focus on the main differences between the schemes of interest, we will not pay attention to choices of step size and stopping criteria of the algorithms. For those readers who are interested in these details, we refer, for instance, to [22, 8, 5, 11].

#### **Algorithm 1:** Preconditioned gradient descent method (PGD)

**Data:**  $G$ : The objective

**Data:**  $s > 0$ : The step size

**Data:**  $x_0 \in \mathbb{H}$ : The initial guess

**Result:** The sequence  $\{x_k\}_{k \geq 1}$  that approximates  $x^*$ , the minimizer of  $G$

**for**  $k \geq 0$  **do**

  |  $x_{k+1} = x_k - s\mathcal{L}^{-1}G'(x_k)$ ;

**end**

We begin by presenting the PGD scheme in Algorithm 1 and describing its convergence properties. To do so, we introduce

$$B = \{x \in \mathbb{H} \mid G(x) \leq G(x_0)\},$$



which is a bounded, convex set containing the minimizer. Then, assuming that  $G$  is  $L_B$ -smooth on  $B$  and  $\mu$ -strongly convex, and that the step size satisfies  $s \in (0, 2/(L_B + \mu)]$ , it is possible to show that  $x_k \in B$  for all  $k \geq 0$ . Moreover, in this setting, the scheme converges exponentially fast to the minimizer (see [15, 11, 22]). In particular, if  $s = 2/(L_B + \mu)$ , then

$$(3.1) \quad \|x_k - x^*\|_{\mathcal{L}} \leq \left(\frac{1-\rho}{1+\rho}\right)^k \|x_0 - x^*\|_{\mathcal{L}}.$$

**Algorithm 2:** Preconditioned accelerated gradient descent method (PAGD)

**Data:**  $G$ : The objective

**Data:**  $\eta > 0$ : The friction coefficient

**Data:**  $s > 0$ : The step size

**Data:**  $x_0 \in \mathbb{H}$ : The initial guess

**Result:** The sequence  $\{x_k\}_{k \geq 1}$  that approximates  $x^*$ , the minimizer of  $G$

**Define:**  $\theta = \eta\sqrt{s}$  and  $\lambda = \frac{1-\theta}{1+\theta}$ ;

**Set:**  $x_{-1} = v_0 = x_0 \in \mathbb{H}$ ;

**for**  $k \geq 0$  **do**

$$(3.2) \quad y_k = x_k + \lambda(x_k - x_{k-1}),$$

$$(3.3) \quad x_{k+1} = y_k - s\mathcal{L}^{-1}G'(y_k),$$

$$(3.4) \quad v_{k+1} = x_k + \frac{1}{\theta}(x_{k+1} - x_k).$$

**end**

To improve on the convergence of GD (Algorithm 1 with  $\mathcal{L} = \mathfrak{R}_{\mathbb{H}}^{-1}$ ), Nesterov [21] devised an algorithm, which “accelerates” the rate of convergence of GD. The improved algorithm is commonly known as *Nesterov’s accelerated gradient descent method* (AGD). The preconditioned version of this scheme, PAGD, is presented in Algorithm 2. Roughly speaking, it computes an extrapolation, (3.2), takes a gradient step there, (3.3), and repeats the same process. Notice that an actual implementation does not need to compute the sequence  $\{v_k\}_{k \geq 0}$ . We need it for the theoretical analysis. As we will see in Section 5, for convergence, the algorithm must satisfy the condition  $s \leq 1/L_B$  and  $\eta \leq \sqrt{\mu}$  where  $L_B > 0$  is the (local) Lipschitz smoothness constant of  $G$  with respect to a bounded convex neighborhood of the minimizer,  $B$ , and  $\mu$  is the strong convexity constant.

It must be noted that PAGD, as presented in Algorithm 2, is practical only if the objective functional is  $\mu$ -strongly convex ( $\mu > 0$ ). Otherwise, a convergence result may not be available. There exists a more general scheme, which one may call *accelerated gradient descent method with variable weights* (see [22, p. 78]), that is applicable to merely convex objectives. We do not discuss this case here.

Let us now compare the performances of GD and AGD (Algorithm 1 and Algorithm 2 with  $\mathcal{L} = \mathfrak{R}_{\mathbb{H}}^{-1}$  respectively) by comparing  $G(x_k) - G^*$ , where  $k$  is the number of iterations and  $G^* = G(x^*)$  is the minimum of  $G$ . To the best of our knowledge, the existing results on AGD are established under the assumption that the objective is *globally* Lipschitz smooth. Thus, for the rest of the summary of this section, the objective  $G$  is assumed to be (globally)  $L$ -smooth. If GD is applied

to a (merely) convex,  $L$ -smooth objective functional with a step size condition  $0 < s \leq 1/L$ , then we have a first order convergence in the objective functional, i.e.,  $G(x_k) - G^* \leq O(1/k)$  as  $k \rightarrow \infty$  (see [22, Corollary 2.1.2]). On the other hand, AGD with variable weights (the more general version mentioned above) provides a second order convergence, that is,  $G(x_k) - G^* \leq O(1/k^2)$  as  $k \rightarrow \infty$ . If the objective is, in addition,  $\mu$ -strongly convex, the convergence rates of the two schemes become exponential. Specifically, estimate (3.1) and the quadratic traps show that the convergence rate of GD is  $G(x_k) - G^* \leq O\left(\left(\frac{1-\rho}{1+\rho}\right)^{2k}\right)$  as  $k \rightarrow \infty$ , where we recall that  $\rho = \mu/L$ . This is in contrast to AGD, which converges with a rate of  $G(x_k) - G^* \leq O\left((1 - \sqrt{\rho})^k\right)$  as  $k \rightarrow \infty$ ; see [22, Theorem 2.2.3]. If  $\rho \ll 1$ , this acceleration can be significant. As we will see later, PAGD achieves the same rate of exponential convergence even if the objective is locally Lipschitz smooth instead of the Lipschitz smoothness being imposed globally.

#### 4. AN ODE MODEL FOR PAGD

We study a continuous time analogue of PAGD, a second order ODE, inspired by [29]. As we will see, the discussion in this section turns out to be informative. It not only provides an intuitive understanding of Nesterov's acceleration, but also guides us to important results at the discrete level. As mentioned in the introduction, there are recent works that arrive at the same conclusion for some parts of our results using similar ideas. However, our unique contributions rely on our specific layout of various quantities and calculations. Thus, we include such details in a condensed manner while referring to existing work otherwise.

To streamline the discussion, we start by directly introducing the initial value problem (IVP) whose certain discretization leads to PAGD:

$$(4.1) \quad \ddot{X}(t) + 2\eta\dot{X}(t) + \mathcal{L}^{-1}G'(X(t)) = 0, \quad t > 0, \quad X(0) = x_0, \quad \dot{X}(0) = 0.$$

Interestingly, this is the same system as what Poljak [24] had in mind when he proposed the *heavy ball* method. See Appendix A for its derivation. Conversely, PAGD can be viewed as a discretization of this IVP although not every choice in the process can be seen as natural or intuitive. It is given in Appendix B. Note that it involves some ingredients appearing in Section 4.2.

**Remark 4.1** (physical interpretation). This IVP (4.1) describes the motion of a ball of unit mass in the potential  $G$  with friction coefficient  $2\eta$  which starts from the initial position  $x_0$  at rest; see Figure 1. Our physical intuition suggests that the ball will converge to its minimal point as it exhausts the initial energy under the action of friction. If the friction, quantified by  $\eta$ , is too small it will oscillate much as it reaches the minimal point and will converge only after a long travel. On the other hand, if the friction is too large, it will not move sufficiently rapidly, and this, in turn, will also lead to a slow convergence.

Let us compare this with another physical system. We can interpret the physics of the gradient flow as a limiting case of the same dynamics. The gradient flow  $\dot{X}(t) = -\mathcal{L}^{-1}G'(X(t))$  can be viewed, up to a constant factor  $2\eta$ , as a massless limit of the IVP (4.1). That is, a physical thought experiment suggests that the surroundings hold the particle back as soon as it gets accelerated since it is so light. A real life example of this kind is a very viscous fluid, such as honey, flowing down a bowl. Our physical experience suggests that it will not oscillate and will flow along

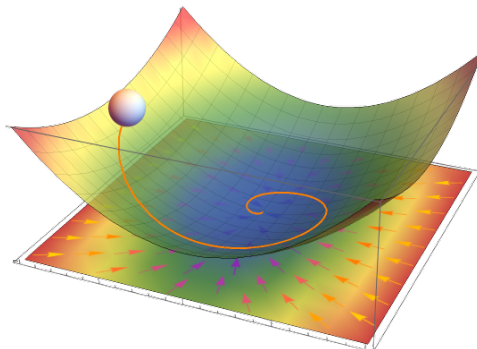


FIGURE 1. A rolling ball system. The IVP (4.1) describes a ball of unit mass rolling down a bowl-shaped potential landscape with a constant friction coefficient.

the steepest descent direction every moment. However, it will reach the bottom slower than the rolling ball will if the friction is appropriately strong. ■

**4.1. Analysis of the IVP.** As one can expect from the fact that the IVP (4.1) describes a concrete physical situation, its solution possesses good properties. In this and the following section, however, we need a slightly stronger Lipschitz condition on  $G'$  than in the discrete level discussion.

**Lemma 4.2** (existence and uniqueness). *Suppose that  $G : \mathbb{H} \rightarrow \mathbb{R}$  is  $\mu$ -strongly convex and locally Lipschitz smooth in the strong sense, i.e., (2.8) holds. Then, for any  $T > 0$ , there exists a unique solution  $X \in C^2(0, T; \mathbb{H})$  to the initial value problem (4.1) and the solution obeys the following energy identity*

$$(4.2) \quad \frac{1}{2} \|\dot{X}(t)\|_{\mathcal{L}}^2 + G(X(t)) - G^* = G(x_0) - G^* - 2\eta \int_0^t \|\dot{X}(\tau)\|_{\mathcal{L}}^2 d\tau \quad \forall t \geq 0.$$

*Consequently, the solution exists for all  $t \in [0, \infty)$  and it is twice continuously differentiable.*

*Proof.* See [3, Theorem 3.1 and Proposition 4.2] for the existence, uniqueness, and smoothness. For the energy law, take the  $\mathcal{L}$ -inner product of the first equation of (4.1) with  $\dot{X}$  and integrate over time  $\tau \in [0, t]$ . □

**Remark 4.3** (smoothness of the solution). The fact that  $X \in C^2((0, \infty); \mathbb{H}) \cap C([0, \infty); \mathbb{H})$  justifies the manipulations we will carry out when we derive the IVP (4.1) in Appendix A. ■

**4.2. Convergence to equilibrium.** We now wish to prove that the solution to the IVP (4.1) with  $G$  being locally Lipschitz smooth and  $\mu$ -strongly convex converges to its attractive steady state solution as  $t \rightarrow \infty$ , which is the minimal point in this case, at a matching rate with that of PAGD. This is one of the highlights of this work. To this end, we introduce an auxiliary variable

$$V(t) = X(t) - x^* + \frac{1}{\eta} \dot{X}(t)$$

so that the first equation of the IVP (4.1) can be rewritten

$$(4.3) \quad \eta \dot{V}(t) + \eta \dot{X}(t) + \mathcal{L}^{-1}G'(X(t)) = 0.$$

We also introduce an energy

$$(4.4) \quad E(X, V) = \frac{\eta}{2} \|V\|_{\mathcal{L}}^2 + \frac{1}{\eta} (G(X) - G^*),$$

where we recall  $G^* = G(x^*) = \min_{x \in \mathbb{H}} G(x)$ . We will show that  $E$  is a *Lyapunov* function for the IVP (4.1). For notational convenience, set  $E_0 = E(x_0, x_0) = \frac{1}{\eta} (G(x_0) - G^*) + \frac{\eta}{2} \|x_0 - x^*\|_{\mathcal{L}}^2$ .

**Theorem 4.4** (exponential decay). *Let  $G : \mathbb{H} \rightarrow \mathbb{R}$  be locally Lipschitz smooth in the strong sense and  $\mu$ -strongly convex. Denote by  $X$  the unique solution to the IVP (4.1). If  $\eta^2 \leq \mu$ , the exponentially inflated energy  $\mathcal{E}(t) = e^{\eta t} E(X(t), V(t))$  is nonincreasing. Consequently, the Lyapunov function (4.4) decays to zero at an exponential rate:*

$$(4.5) \quad E(X(t), V(t)) = \frac{\eta}{2} \|V(t)\|_{\mathcal{L}}^2 + \frac{1}{\eta} (G(X(t)) - G^*) \leq e^{-\eta t} E_0.$$

*Proof.* Existence and uniqueness of  $X$  is guaranteed by Lemma 4.2. Let us now prove the estimate (4.5). Taking the inner product of (4.3) with  $V(t)$ , and using the identity (2.12), we obtain, suppressing the time variable,

$$(4.6) \quad \begin{aligned} 0 &= \eta(V, \dot{V})_{\mathcal{L}} + \eta(\dot{X}, X - x^* + \frac{1}{\eta} \dot{X})_{\mathcal{L}} + (\mathcal{L}^{-1}G'(X), X - x^*)_{\mathcal{L}} \\ &\quad + \frac{1}{\eta} (\mathcal{L}^{-1}G'(X), \dot{X})_{\mathcal{L}} \\ &= \eta(V, \dot{V})_{\mathcal{L}} + \frac{1}{2} \|\dot{X}\|_{\mathcal{L}}^2 + \frac{\eta^2}{2} \|V\|_{\mathcal{L}}^2 - \frac{\eta^2}{2} \|X - x^*\|_{\mathcal{L}}^2 + \langle G'(X), X - x^* \rangle \\ &\quad + \frac{1}{\eta} \langle G'(X), \dot{X} \rangle. \end{aligned}$$

The lower quadratic trap, (2.11), implies

$$(4.7) \quad G(X) - G^* - \langle G'(X), X - x^* \rangle \leq -\frac{\mu}{2} \|X - x^*\|_{\mathcal{L}}^2.$$

Substituting (4.6) into the time derivative of the inflated energy and then using the above estimate (4.7), we have

$$(4.8) \quad \begin{aligned} \dot{\mathcal{E}}(t) &= e^{\eta t} \left[ \frac{\eta^2}{2} \|V\|_{\mathcal{L}}^2 + \eta(V, \dot{V})_{\mathcal{L}} + (G(X) - G^*) + \frac{1}{\eta} \langle G'(X), \dot{X} \rangle \right] \\ &= e^{\eta t} \left[ -\frac{1}{2} \|\dot{X}\|_{\mathcal{L}}^2 + \frac{\eta^2}{2} \|X - x^*\|_{\mathcal{L}}^2 - \langle G'(X), X - x^* \rangle + G(X) - G^* \right] \\ &\leq -\frac{1}{2} e^{\eta t} \left[ \|\dot{X}\|_{\mathcal{L}}^2 + \frac{\mu - \eta^2}{2} \|X - x^*\|_{\mathcal{L}}^2 \right]. \end{aligned}$$

The last term is always nonpositive provided  $\eta^2 \leq \mu$ , and this implies  $\mathcal{E}(t) \leq \mathcal{E}(0) = E_0$ . This completes the proof.  $\square$

**Remark 4.5** (physical interpretation). We can rigorously explain the physical intuition given in Remark 4.1 through Theorem 4.4 and its proof. If the friction, quantified by  $\eta$ , is too small the decay to the attraction point is slow as  $\eta$  governs

the decay rate  $e^{-\eta t}$ . On the other hand, If the friction is too large, say  $\eta > \sqrt{\mu}$ , then we cannot guarantee the boundedness of  $\mathcal{E}(t)$ . ■

## 5. AN ENERGY APPROACH TO CONVERGENCE OF PAGD

In this section, we prove the existence of an invariant set of PAGD and its exponential convergence in the objective as well as in the residual when it is applied to a strongly convex, locally Lipschitz smooth objective. We follow the ODE arguments developed in Section 4. Throughout this section, we assume  $\eta = \sqrt{\mu}$ , the optimal choice for the friction coefficient in view of Theorem 4.4. Note that this does not undermine generality. If  $\tilde{\mu}$  is the largest strong convexity constant of  $G$ , that is, the supremum of  $\mu$ 's that satisfies the strong convexity, (2.9), then any  $\mu \in (0, \tilde{\mu}]$  can be taken as a (non-optimal) strong convexity constant. Thus, the general case  $\eta^2 \leq \tilde{\mu}$  corresponds to  $\eta^2 = \mu \leq \tilde{\mu}$ , the optimal friction coefficient associated with a non-optimal strong convexity constant.

As a first step, we show that the assumption of the local Lipschitz smoothness is sufficient for our analysis, as the iterates lie within a bounded set. We first show that, for every  $k \geq 0$ , the  $y_k$  iterate of PAGD lies in the segment between  $x_k$  and  $v_k$ . This is used frequently in the convergence proof.

**Lemma 5.1** (convex hull). *For every  $k \geq 0$ , the iterates constructed in PAGD, described in Algorithm 2, satisfy  $y_k \in \overline{x_k v_k}$ . Specifically, they satisfy the following four equivalent equations:*

$$(5.1) \quad \begin{cases} y_k = \frac{1}{1+\theta}x_k + \frac{\theta}{1+\theta}v_k, & x_k = (1+\theta)y_k - \theta v_k, \\ v_k = \left(1 + \frac{1}{\theta}\right)y_k - \frac{1}{\theta}x_k, & x_k - y_k = \theta(y_k - v_k). \end{cases}$$

*Proof.* If  $k = 0$  this is trivial since  $x_0 = y_0 = v_0$ . For  $k \geq 1$ , we eliminate  $x_{k-1}$  from (3.2) and (3.4) with the index being  $k-1$  to get

$$\left(1 - \frac{1}{\theta}\right)y_k + \lambda v_k = \left((1+\lambda)\left(1 - \frac{1}{\theta}\right) + \frac{\lambda}{\theta}\right)x_k = -\frac{\lambda}{\theta}x_k.$$

Rearranging terms and using  $\lambda = \frac{1-\theta}{1+\theta}$ , we obtain the equalities that are listed above. □

We now show that there is an invariant set for the iterates of PAGD.

**Lemma 5.2** (invariant set). *Assume that the objective  $G : \mathbb{H} \rightarrow \mathbb{R}$  is  $\mu$ -strongly convex and locally Lipschitz smooth. Define*

$$(5.2) \quad \mathfrak{B} = \{x \in \mathbb{H} \mid \|x - x^*\|_{\mathcal{L}} \leq R\},$$

where  $R = R_1 + \frac{1}{\eta}R_2$ ,  $R_1 = \sqrt{\frac{2}{\mu}(G(x_0) - G^*)}$ ,  $R_2 = \sqrt{2r(G(x_0) - G^*)}$ , and  $r > 1$ . Let PAGD, as described in Algorithm 2, be implemented with a step size rule

$$(5.3) \quad s \in \left(0, \min \left\{ L_B^{-1}, \left(\frac{r-1}{r+1}\right)^2 \mu^{-1} \right\} \right],$$

where  $L_B$  the local Lipschitz smoothness constant of  $G$  associated to the set  $\mathfrak{B}$ . Then, for all  $k \geq 0$ , we have that  $\|x_k - x^*\|_{\mathcal{L}} \leq R_1$ , hence  $x_k \in \mathfrak{B}$ , and  $y_k, v_k \in \mathfrak{B}$ .

*Proof.* The outline of this proof is simple although it is long. We mimic the energy law developed in Section 4.1 to obtain a bound on the distance between the main iterates and the minimizer and that on the speed. Once we get the bounds, it is easy to prescribe an appropriate ball, which will be our invariant set.

We will prove the statement by induction. For  $k = 0$ , the statement is trivial since  $x_0 = y_0 = v_0$  and the strong convexity implies  $\|x_0 - x^*\|_{\mathcal{L}} \leq R_1$ . Suppose that  $\|x_k - x^*\|_{\mathcal{L}} \leq R_1$  (hence  $x_k \in \mathfrak{B}$ ) and  $y_k, v_k \in \mathfrak{B}$  are true for  $k = 0, 1, 2, \dots, N$ . We need to show that  $\|x_{N+1} - x^*\|_{\mathcal{L}} \leq R_1$  (hence  $x_{N+1} \in \mathfrak{B}$ ) and  $v_{N+1} \in \mathfrak{B}$ , then Lemma 5.1 implies  $y_{N+1} \in \overline{x_{N+1}v_{N+1}} \subset \mathfrak{B}$  since  $\mathfrak{B}$  is a convex set as a sublevel set of a convex function.

Note that the condition  $s \leq (\frac{r-1}{r+1})^2 \mu^{-1}$ , which is implied by (5.3), ensures  $\lambda^{-1}$  to be bounded above since

$$(5.4) \quad \frac{1}{\lambda} = \frac{1 + \sqrt{s\mu}}{1 - \sqrt{s\mu}} \leq r.$$

First, a similar argument to [11, Proposition 4.6] shows that the  $x_{N+1}$  update from  $y_N$  is a descent step in terms of  $G$ . That is, the section of  $G$  across the line  $\overleftrightarrow{y_N x_{N+1}}$  also inherits the strong convexity and the local Lipschitz smoothness with the same constants on the one-dimensional affine subset

$$B_{N+1} = \{x = y_N - \tau \mathcal{L}^{-1} G'(y_N) \in \mathbb{H} \mid \tau \in \mathbb{R}\}.$$

Let  $S(\tau) = G(y_N - \tau \mathcal{L}^{-1} G'(y_N))$  denote the section.

Since we know that  $y_N \in \mathfrak{B}$ , we can bound  $S$  in a neighborhood of  $\tau = 0$  using the upper quadratic trap

$$\begin{aligned} U(\tau) &:= G(y_N) + \langle G'(y_N), -\tau \mathcal{L}^{-1} G'(y_N) \rangle + \frac{L_B}{2} \|\tau \mathcal{L}^{-1} G'(y_N)\|_{\mathcal{L}}^2 \\ &= G(y_N) - \tau \|G'(y_N)\|_{\mathcal{L}^{-1}}^2 + \frac{L_B \tau^2}{2} \|G'(y_N)\|_{\mathcal{L}^{-1}}^2. \end{aligned}$$

Observe that  $S(0) = U(0) = G(y_N)$ , that  $U(\tau)$  is decreasing around  $\tau = 0$  since  $dU/d\tau(0) = -\|G'(y_N)\|_{\mathcal{L}^{-1}}^2 \leq 0$ , and that the optimal step size to minimize  $U$  is  $1/L_B$  since  $dU/d\tau(1/L_B) = 0$ . This implies that  $S(s) \leq U(s) \leq U(0)$  for any  $s \in [0, 2/L_B]$ . Moreover, for  $s \in [0, 1/L_B]$ , we have

$$(5.5) \quad \begin{aligned} G(x_{N+1}) &= S(s) \leq U(s) = G(y_N) - s \|G'(y_N)\|_{\mathcal{L}^{-1}}^2 + \frac{L_B s^2}{2} \|G'(y_N)\|_{\mathcal{L}^{-1}}^2 \\ &\leq G(y_N) - \frac{s}{2} \|G'(y_N)\|_{\mathcal{L}^{-1}}^2, \end{aligned}$$

which is the desired descent property in  $G$  from  $y_N$  to  $x_{N+1}$ .

Now, we want to mimic the energy argument that we carried out in Section 4.1. Substitute (3.2) into (3.3), and add and subtract  $x_k - x_{k-1}$ , to obtain the discrete counterpart of (4.1)

$$(5.6) \quad x_{k+1} - 2x_k + x_{k-1} + (1 - \lambda)(x_k - x_{k-1}) + s \mathcal{L}^{-1} G'(y_k) = 0.$$

Note that defining  $x_{-1} := x_0$  allows us to extend this equality to the case  $k = 0$ . Take the  $\mathcal{L}$ -inner product of this identity with  $x_k - x_{k-1}$  and add for  $0 \leq k \leq N$ .

Then, using (2.13), the first term telescopes to simplify

$$\begin{aligned}
& \sum_{k=0}^N (x_{k+1} - 2x_k + x_{k-1}, x_k - x_{k-1})_{\mathcal{L}} \\
&= \frac{1}{2} \sum_{k=0}^N \left( \|x_{k+1} - x_k\|_{\mathcal{L}}^2 - \|x_{k+1} - 2x_k + x_{k-1}\|_{\mathcal{L}}^2 - \|x_k - x_{k-1}\|_{\mathcal{L}}^2 \right) \\
&= \frac{1}{2} \|x_{N+1} - x_N\|_{\mathcal{L}}^2 - \frac{1}{2} \sum_{k=0}^N \|x_{k+1} - 2x_k + x_{k-1}\|_{\mathcal{L}}^2.
\end{aligned}$$

We leave the second term as it is. For the third term, using (3.2),  $G(y_k) - G(x_k) \leq \langle G'(y_k), y_k - x_k \rangle$  from convexity, and (5.5), it follows

$$\begin{aligned}
(5.7) \quad & s \sum_{k=0}^N \langle G'(y_k), x_k - x_{k-1} \rangle = \frac{s}{\lambda} \sum_{k=0}^N \langle G'(y_k), y_k - x_k \rangle \\
& \geq \frac{s}{\lambda} \sum_{k=0}^N (G(y_k) - G(x_k)) \geq \frac{s}{\lambda} \sum_{k=0}^N \left( G(x_{k+1}) - G(x_k) + \frac{s}{2} \|G'(y_k)\|_{\mathcal{L}^{-1}}^2 \right) \\
& = \frac{s}{\lambda} G(x_{N+1}) - \frac{s}{\lambda} G(x_0) + \frac{s^2}{2\lambda} \sum_{k=0}^N \|G'(y_k)\|_{\mathcal{L}^{-1}}^2.
\end{aligned}$$

Gathering all the three terms together and rearranging, we get

$$\begin{aligned}
(5.8) \quad & \frac{1}{2} \|x_{N+1} - x_N\|_{\mathcal{L}}^2 + \frac{s}{\lambda} G(x_{N+1}) \leq \frac{s}{\lambda} G(x_0) + \frac{1}{2} \sum_{k=0}^N \|x_{k+1} - 2x_k + x_{k-1}\|_{\mathcal{L}}^2 \\
& \quad - \frac{s^2}{2\lambda} \sum_{k=0}^N \|G'(y_k)\|_{\mathcal{L}^{-1}}^2 - (1-\lambda) \sum_{k=0}^N \|x_k - x_{k-1}\|_{\mathcal{L}}^2.
\end{aligned}$$

Similarly, take the  $\mathcal{L}$ -inner product of (5.6) with  $x_{k+1} - x_k$  and sum over  $0 \leq k \leq N$ . This time, use (2.12) for the first term to get

$$\begin{aligned}
& \sum_{k=0}^N (x_{k+1} - 2x_k + x_{k-1}, x_{k+1} - x_k)_{\mathcal{L}} \\
&= \frac{1}{2} \|x_{N+1} - x_N\|_{\mathcal{L}}^2 + \frac{1}{2} \sum_{k=0}^N \|x_{k+1} - 2x_k + x_{k-1}\|_{\mathcal{L}}^2.
\end{aligned}$$

For the second term, using Cauchy-Schwarz and Young's inequality, we have

$$\begin{aligned}
& (1-\lambda) \sum_{k=0}^N (x_k - x_{k-1}, x_{k+1} - x_k)_{\mathcal{L}} \\
& \geq -\frac{1-\lambda}{2} \sum_{k=0}^N \left( \|x_k - x_{k-1}\|_{\mathcal{L}}^2 + \|x_{k+1} - x_k\|_{\mathcal{L}}^2 \right).
\end{aligned}$$

For the third term, use (3.3) and argue as in (5.7) to get

$$\begin{aligned}
& s \sum_{k=0}^N \langle G'(y_k), x_{k+1} - x_k \rangle = s \sum_{k=0}^N \langle G'(y_k), y_k - s\mathcal{L}^{-1}G'(y_k) - x_k \rangle \\
& \geq -s^2 \sum_{k=0}^N \|G'(y_k)\|_{\mathcal{L}^{-1}}^2 + s \sum_{k=0}^N \langle G'(y_k), y_k - x_k \rangle \\
& \geq sG(x_{N+1}) - sG(x_0) - \frac{s^2}{2} \sum_{k=0}^N \|G'(y_k)\|_{\mathcal{L}^{-1}}^2.
\end{aligned}$$

Gathering all these estimates we get

$$\begin{aligned}
(5.9) \quad & \frac{1}{2} \|x_{N+1} - x_N\|_{\mathcal{L}}^2 + sG(x_{N+1}) \leq sG(x_0) - \frac{1}{2} \sum_{k=0}^N \|x_{k+1} - 2x_k + x_{k-1}\|_{\mathcal{L}} \\
& + \frac{s^2}{2} \sum_{k=0}^N \|G'(y_k)\|_{\mathcal{L}^{-1}}^2 + \frac{1-\lambda}{2} \sum_{k=0}^N \left( \|x_k - x_{k-1}\|_{\mathcal{L}}^2 + \|x_{k+1} - x_k\|_{\mathcal{L}}^2 \right).
\end{aligned}$$

Add (5.8) and (5.9) and rearrange, then after some cancellations, it follows

$$\begin{aligned}
(5.10) \quad & \frac{1+\lambda}{2} \|x_{N+1} - x_N\|_{\mathcal{L}}^2 + s \left( 1 + \frac{1}{\lambda} \right) (G(x_{N+1}) - G^*) \\
& \leq s \left( 1 + \frac{1}{\lambda} \right) (G(x_0) - G^*) - \frac{s^2}{2} \left( \frac{1}{\lambda} - 1 \right) \sum_{k=0}^N \|G'(y_k)\|_{\mathcal{L}^{-1}}^2
\end{aligned}$$

$$(5.11) \quad \leq s \left( 1 + \frac{1}{\lambda} \right) (G(x_0) - G^*),$$

since  $0 < \lambda < 1$ . By removing the kinetic term from this estimate, strong convexity leads to

$$G(x_0) - G^* \geq G(x_{N+1}) - G^* \geq \frac{\mu}{2} \|x_{N+1} - x^*\|_{\mathcal{L}}^2,$$

which implies

$$(5.12) \quad \|x_{N+1} - x^*\|_{\mathcal{L}} \leq R_1,$$

which, in turn, proves  $x_{N+1} \in \mathfrak{B}$ . Similarly, discarding the potential term, dividing through  $\frac{s(1+\lambda)}{2}$ , and using (5.4), we obtain

$$\left\| \frac{x_{N+1} - x_N}{\sqrt{s}} \right\|_{\mathcal{L}} \leq \sqrt{\frac{2}{\lambda} (G(x_0) - G^*)} \leq R_2.$$

Then, from the definition of  $v_{N+1}$  (3.4),

$$\|v_{N+1} - x^*\|_{\mathcal{L}} \leq \|x_N - x^*\|_{\mathcal{L}} + \frac{1}{\eta} \left\| \frac{x_{N+1} - x_N}{\sqrt{s}} \right\|_{\mathcal{L}} \leq R_1 + \frac{1}{\eta} R_2 = R,$$

which implies  $v_{N+1} \in \mathfrak{B}$ . This completes the proof.  $\square$

**Remark 5.3** (step size restriction). The additional condition  $s \leq \left(\frac{r-1}{r+1}\right)^2 \mu^{-1}$  on the step size is not restrictive at all in practice. For example, if  $r = 3$ , we require that  $s \leq 1/4\mu$ . The purpose of this condition is to bound  $\lambda^{-1}$  as explained in the proof. However,  $\lambda^{-1}$  becomes unbounded when  $s\mu$  is close to 1. If we set  $s = 1/L_B$ ,  $s\mu$  is the (inverse) condition number and the (inverse) condition number being close to 1 makes the problem more amenable because it means that  $G$  is almost quadratic.



Moreover, even from a theoretical point of view, as  $r$  increases, the invariant set  $\mathfrak{B}$  gets larger, which means  $L_B^{-1}$  gets smaller, while  $(\frac{r-1}{r+1})^2 \mu^{-1}$  approaches  $\mu^{-1}$ . Since  $L_B^{-1} < \mu^{-1}$  (unless  $G$  is perfectly quadratic), the second argument of the minimum in (5.3) eventually becomes of no effect.  $\blacksquare$

**Corollary 5.4** (convergence of residuals). *Assume that  $G : \mathbb{H} \rightarrow \mathbb{R}$  is  $\mu$ -strongly convex and locally Lipschitz smooth. Suppose that PAGD, as described in Algorithm 2, is implemented with a step size that obeys condition (5.3), where  $r > 1$ ,  $\mathfrak{B}$  is the invariant set given by (5.2), and  $L_B$  is the Lipschitz smoothness constant associated with  $\mathfrak{B}$ . In this setting, the residuals  $\{G'(y_k)\}_{k \geq 0}$  converge to zero in the  $\mathcal{L}^{-1}$ -norm at least  $\ell^2$ -fast. In other words,*

$$\sum_{k=0}^{\infty} \|G'(y_k)\|_{\mathcal{L}^{-1}}^2 < \infty.$$

*Proof.* Moving the summation term of (5.10) to the left hand side and dropping the other nonnegative terms, we have

$$\sum_{k=0}^N \|G'(y_k)\|_{\mathcal{L}^{-1}}^2 \leq \frac{2(1+\lambda)}{s(1-\lambda)} (G(x_0) - G^*).$$

Letting  $N \rightarrow \infty$  completes the proof.  $\square$

Of course, this result is far from optimal. An exponential convergence of the residuals in the  $\mathcal{L}^{-1}$ -norm will be proved in Corollary 5.9.

We can now begin the proof of convergence *per se*. We begin with an estimate for the discrete time derivative of the potential energy.

**Lemma 5.5** (discrete derivative of potential energy). *Let the objective  $G : \mathbb{H} \rightarrow \mathbb{R}$  be  $\mu$ -strongly convex and locally Lipschitz smooth. Suppose that the step size in Algorithm 2 satisfies (5.3), where  $r > 1$ ,  $\mathfrak{B}$  is the invariant set given by (5.2), and  $L_B$  is the Lipschitz smoothness constant associated with  $\mathfrak{B}$ . Then, we have that*

$$\frac{1}{\eta} \frac{G(x_{k+1}) - G(x_k)}{\sqrt{s}} \leq -\frac{\sqrt{s}}{2\eta} \|G'(y_k)\|_{\mathcal{L}^{-1}}^2 + \frac{1}{\theta} \langle G'(y_k), y_k - x_k \rangle - \frac{\eta}{2\sqrt{s}} \|x_k - y_k\|_{\mathcal{L}}^2.$$

*Proof.* Since we have an invariant set,  $\mathfrak{B}$ , we can utilize the local Lipschitz smoothness with respect to it. Combined with  $sL_B \leq 1$ , it leads to

$$(5.13) \quad G(x_{k+1}) = G(y_k - s\mathcal{L}^{-1}G'(y_k)) \leq G(y_k) - \frac{s}{2} \|G'(y_k)\|_{\mathcal{L}^{-1}}^2.$$

Using this estimate, the strong convexity of  $G$  yields

$$\begin{aligned} G(x_k) &\geq G(y_k) + \langle G'(y_k), x_k - y_k \rangle + \frac{\mu}{2} \|x_k - y_k\|_{\mathcal{L}}^2 \\ &\geq G(x_{k+1}) + \frac{s}{2} \|G'(y_k)\|_{\mathcal{L}^{-1}}^2 + \langle G'(y_k), x_k - y_k \rangle + \frac{\mu}{2} \|x_k - y_k\|_{\mathcal{L}}^2. \end{aligned}$$

Rearranging the last estimate, multiplying through by  $1/\theta$ , and recalling  $\theta = \sqrt{s\mu}$ ,  $\eta = \sqrt{\mu}$ , we obtain the desired result.  $\square$

We also need an analogue of (4.6), a certain relation derived from the scheme.

**Lemma 5.6** (discrete analogue of (4.6)). *The iterates constructed by PAGD, as described in Algorithm 2, satisfy*

$$(5.14) \quad \frac{\eta}{2\sqrt{s}} \left( \|v_{k+1} - x^*\|_{\mathcal{L}}^2 - \|v_k - x^*\|_{\mathcal{L}}^2 \right) + \frac{1}{\theta} \langle G'(y_k), y_k - x_k \rangle + \frac{\eta^2}{2} \|v_k - x^*\|_{\mathcal{L}}^2 \\ - \frac{\eta}{2\sqrt{s}} \|v_{k+1} - v_k\|_{\mathcal{L}}^2 + \frac{1}{2s} \|y_k - x_k\|_{\mathcal{L}}^2 - \frac{\eta^2}{2} \|y_k - x^*\|_{\mathcal{L}}^2 + \langle G'(y_k), y_k - x^* \rangle = 0.$$

*Proof.* Substituting (3.3) in (3.4), and using the relations (5.1), we have

$$v_{k+1} = x_k + \frac{1}{\theta}(y_k - x_k) - \frac{\sqrt{s}}{\eta} \mathcal{L}^{-1} G'(y_k) = x_k + v_k - y_k - \frac{\sqrt{s}}{\eta} \mathcal{L}^{-1} G'(y_k).$$

Rearranging, and multiplying through by  $\frac{\eta}{\sqrt{s}}$ , we obtain the discrete analogue of the ODE (4.3)

$$(5.15) \quad \eta \frac{v_{k+1} - v_k}{\sqrt{s}} + \eta \frac{y_k - x_k}{\sqrt{s}} + \mathcal{L}^{-1} G'(y_k) = 0.$$

The discrete analogue of  $V(t)$  is  $v_k - x^*$ , so following the proof of Theorem 4.4, we now take the  $\mathcal{L}$ -inner product of (5.15) with  $v_k - x^*$  to obtain

$$(5.16) \quad \frac{\eta}{\sqrt{s}} (v_{k+1} - v_k, v_k - x^*)_{\mathcal{L}} + \frac{\eta}{\sqrt{s}} (y_k - x_k, v_k - x^*)_{\mathcal{L}} + \langle G'(y_k), v_k - x^* \rangle = 0.$$

Using (2.13), the first term can be rewritten as

$$\frac{\eta}{\sqrt{s}} (v_{k+1} - v_k, v_k - x^*)_{\mathcal{L}} = \\ \frac{\eta}{2\sqrt{s}} \left( \|v_{k+1} - x^*\|_{\mathcal{L}}^2 - \|v_k - x^*\|_{\mathcal{L}}^2 \right) - \frac{\eta}{2\sqrt{s}} \|v_{k+1} - v_k\|_{\mathcal{L}}^2.$$

For the second term in (5.16), we use relations (5.1), and then the identity (2.12) to get

$$\frac{\eta}{\sqrt{s}} (y_k - x_k, v_k - x^*)_{\mathcal{L}} = \frac{\eta}{\sqrt{s}} \frac{1}{\theta} (y_k - x_k, \theta v_k - \theta x^*)_{\mathcal{L}} \\ = \frac{1}{s} (y_k - x_k, y_k - x_k + \theta(y_k - x^*))_{\mathcal{L}} \\ = \frac{1}{2s} \|y_k - x_k\|_{\mathcal{L}}^2 + \frac{\eta^2}{2} \|v_k - x^*\|_{\mathcal{L}}^2 - \frac{\eta^2}{2} \|y_k - x^*\|_{\mathcal{L}}^2.$$

Finally, for the third term of (5.16), we use (5.1) similarly to the above, then it follows

$$\langle G'(y_k), v_k - x^* \rangle = \frac{1}{\theta} \langle G'(y_k), \theta v_k - \theta x^* \rangle \\ = \frac{1}{\theta} \langle G'(y_k), y_k - x_k \rangle + \langle G'(y_k), y_k - x^* \rangle.$$

Then, the desired result follows upon combining the last three identities.  $\square$

We need one more relation between the iterates.

**Lemma 5.7** (relation between iterates). *The iterates constructed by PAGD, as described in Algorithm 2, satisfy*

$$\frac{\eta}{2\sqrt{s}} \|v_{k+1} - v_k\|_{\mathcal{L}}^2 = \frac{\eta}{2\sqrt{s}} \|x_k - y_k\|_{\mathcal{L}}^2 + \frac{\sqrt{s}}{2\eta} \|G'(y_k)\|_{\mathcal{L}^{-1}}^2 + \langle G'(y_k), y_k - x_k \rangle.$$

*Proof.* Combine (3.4) and the relations (5.1), and then use (3.3) to obtain

$$\begin{aligned} v_{k+1} - v_k &= x_k + \frac{1}{\theta}(x_{k+1} - x_k) + \frac{1}{\theta}x_k - \left(1 + \frac{1}{\theta}\right)y_k \\ &= x_k - y_k + \frac{1}{\theta}(x_{k+1} - y_k) = x_k - y_k - \frac{\sqrt{s}}{\eta}\mathcal{L}^{-1}G'(y_k). \end{aligned}$$

Take  $\mathcal{L}$ -norm square on both sides and then multiply by  $\frac{\eta}{2\sqrt{s}}$ .  $\square$

We are now in a position to prove the main result of this section, the exponential convergence of PAGD using energy arguments. The following result and its consequences are another of the main contributions of this work. To state it, we recall that the Lyapunov function of (4.1)  $E : \mathbb{H}^2 \rightarrow \mathbb{R}$  is defined in (4.4) and that  $E_0$  is its value at the initial state as mentioned in Theorem 4.4.

**Theorem 5.8** (exponential decay). *Let the objective  $G : \mathbb{H} \rightarrow \mathbb{R}$  be locally Lipschitz smooth and  $\mu$ -strongly convex. If PAGD, as described in Algorithm 2, is applied to approximate  $x^* = \operatorname{argmin}_{x \in \mathbb{H}} G(x)$  with a step size satisfying (5.3), where  $r > 1$ ,  $\mathfrak{B}$  is the invariant set given by (5.2), and  $L_B$  is the Lipschitz smoothness constant associated with  $\mathfrak{B}$ , then the Lyapunov function (4.4) decays exponentially along the iterates  $\{x_k\}_{k \geq 0}$ . More specifically, for  $k \geq 0$ , we have*

$$(5.17) \quad E(x_{k+1}, v_{k+1} - x^*) \leq (1 - \theta)E(x_k, v_k - x^*), \quad E(x_k, v_k - x^*) \leq (1 - \theta)^k E_0.$$

*Proof.* Define, for  $k \geq 0$ ,  $\mathcal{E}_k = (1 - \theta)^{-k} E(x_k, v_k - x^*)$ , which is the discrete analogue of the exponentially inflated energy in the the proof of Theorem 4.4. To simplify notation, we set  $C_{\theta, k} = (1 - \theta)^{-(k+1)} > 0$ . Then, similarly to the ODE case, one can show the discrete time derivative of  $\mathcal{E}_k$  is nonpositive as follows. First, we simply use the forward difference time derivative, rearrange, and use Lemma 5.5 to get

$$\begin{aligned} \frac{\mathcal{E}_{k+1} - \mathcal{E}_k}{\sqrt{s}} &= \frac{1}{\sqrt{s}} \left[ (1 - \theta)^{-(k+1)} \left( \frac{1}{\eta} (G(x_{k+1}) - G^*) + \frac{\eta}{2} \|v_{k+1} - x^*\|_{\mathcal{L}}^2 \right) \right. \\ &\quad \left. - (1 - \theta)^{-k} \left( \frac{1}{\eta} (G(x_k) - G^*) + \frac{\eta}{2} \|v_k - x^*\|_{\mathcal{L}}^2 \right) \right] \\ &= C_{\theta, k} \left[ \frac{1}{\eta} \frac{G(x_{k+1}) - G(x_k)}{\sqrt{s}} + (G(x_k) - G^*) \right. \\ &\quad \left. + \frac{\eta}{2\sqrt{s}} (\|v_{k+1} - x^*\|_{\mathcal{L}}^2 - \|v_k - x^*\|_{\mathcal{L}}^2) + \frac{\eta^2}{2} \|v_k - x^*\|_{\mathcal{L}}^2 \right] \\ &\leq C_{\theta, k} \left[ -\frac{\sqrt{s}}{2\eta} \|G'(y_k)\|_{\mathcal{L}^{-1}}^2 - \frac{\eta}{2\sqrt{s}} \|x_k - y_k\|_{\mathcal{L}}^2 + (G(x_k) - G^*) \right. \\ &\quad \left. + \frac{1}{\theta} \langle G'(y_k), y_k - x_k \rangle + \frac{\eta}{2\sqrt{s}} (\|v_{k+1} - x^*\|_{\mathcal{L}}^2 - \|v_k - x^*\|_{\mathcal{L}}^2) \right. \\ &\quad \left. + \frac{\eta^2}{2} \|v_k - x^*\|_{\mathcal{L}}^2 \right]. \end{aligned}$$

We continue by using Lemma 5.6 and then Lemma 5.7, then it follows

$$\begin{aligned} \frac{\mathcal{E}_{k+1} - \mathcal{E}_k}{\sqrt{s}} &\leq C_{\theta,k} \left[ -\frac{\sqrt{s}}{2\eta} \|G'(y_k)\|_{\mathcal{L}^{-1}}^2 - \left(\frac{\eta}{2\sqrt{s}} + \frac{1}{2s}\right) \|x_k - y_k\|_{\mathcal{L}}^2 + (G(x_k) - G^*) \right. \\ &\quad \left. + \frac{\eta}{2\sqrt{s}} \|v_{k+1} - v_k\|_{\mathcal{L}}^2 + \frac{\eta^2}{2} \|y_k - x^*\|_{\mathcal{L}}^2 + \langle G'(y_k), x^* - y_k \rangle \right] \\ &= C_{\theta,k} \left[ -\frac{1}{2s} \|y_k - x_k\|_{\mathcal{L}}^2 + (G(x_k) - G^*) + \langle G'(y_k), y_k - x_k \rangle \right. \\ &\quad \left. + \langle G'(y_k), x^* - y_k \rangle + \frac{\eta^2}{2} \|y_k - x^*\|_{\mathcal{L}}^2 \right]. \end{aligned}$$

Finally, add and subtract  $G(y_k)$  from the last expression, and use the following estimates, which are simple rearrangements of the lower and upper quadratic traps,

$$\begin{aligned} G(y_k) - G^* + \langle G'(y_k), x^* - y_k \rangle &\leq -\frac{\mu}{2} \|y_k - x^*\|_{\mathcal{L}}^2 \\ G(x_k) - G(y_k) + \langle G'(y_k), y_k - x_k \rangle &\leq \frac{L_B}{2} \|x_k - y_k\|_{\mathcal{L}}^2, \end{aligned}$$

then we arrive at

$$\begin{aligned} \frac{\mathcal{E}_{k+1} - \mathcal{E}_k}{\sqrt{s}} &\leq C_{\theta,k} \left[ G(x_k) - G(y_k) + \langle G'(y_k), y_k - x_k \rangle + G(y_k) - G^* \right. \\ &\quad \left. + \langle G'(y_k), x^* - y_k \rangle - \frac{1}{2s} \|y_k - x_k\|_{\mathcal{L}}^2 + \frac{\eta^2}{2} \|y_k - x^*\|_{\mathcal{L}}^2 \right] \\ &\leq C_{\theta,k} \left[ \frac{1}{2} \left( L_B - \frac{1}{s} \right) \|y_k - x_k\|_{\mathcal{L}}^2 \right]. \end{aligned}$$

The step size condition forces the last term to be nonpositive. Therefore, we conclude that  $\{\mathcal{E}_k\}_{k \geq 0}$  is nonincreasing, from which we obtain (5.17).  $\square$

The following estimates are evident.

**Corollary 5.9** (rate of convergence). *In the setting of Theorem 5.8, we have that the iterates of PAGD, as described in Algorithm 2, converge to  $x^*$ , the minimizer of  $G$ , at an exponential rate. More specifically, for a suitable  $r > 1$  the step size can be set  $s = 1/L_B$  and, in this case, for  $k \geq 0$ ,*

$$(5.18) \quad \frac{1}{\eta} (G(x_k) - G^*) + \frac{\eta}{2} \|v_k - x^*\|_{\mathcal{L}}^2 \leq (1 - \sqrt{\rho})^k E_0,$$

which implies

$$(5.19) \quad G(x_k) - G^* \leq (1 - \sqrt{\rho})^k \eta E_0, \quad \|x_k - x^*\|_{\mathcal{L}} \leq (1 - \sqrt{\rho})^{\frac{k}{2}} \sqrt{\frac{2E_0}{\eta}}.$$

Furthermore, we have exponential convergence in the  $\mathcal{L}^{-1}$ -norm of the residuals: for  $k \geq 0$ ,

$$(5.20) \quad \|G'(y_k)\|_{\mathcal{L}^{-1}} \leq 3L_B \sqrt{\frac{2E_0}{\eta}} (1 - \sqrt{\rho})^{\frac{k-1}{2}}.$$

*Proof.* We can choose an appropriate  $r > 1$  so that the step size condition (5.3) reduces to  $s \in (0, L_B^{-1}]$ ; see Remark 5.3. Estimate (5.18) and the first estimate of

(5.19) follow from (5.17) upon setting  $s = 1/L_B$ . The second estimate of (5.19) follows by applying strong convexity of  $G$  to the first estimate of (5.19).

Next, from the estimate (5.13), one obtains

$$(5.21) \quad G^* \leq G(x_{k+1}) \leq G(y_k) - \frac{1}{2L_B} \|G'(y_k)\|_{\mathcal{L}^{-1}}^2,$$

from which, one obtains the following by rearranging and then using the upper quadratic trap

$$\|G'(y_k)\|_{\mathcal{L}^{-1}} \leq \sqrt{2L_B(G(y_k) - G^*)} \leq L_B \|y_k - x^*\|_{\mathcal{L}}.$$

In addition, we also have, from the definition of  $y_k$  and  $0 < \lambda < 1$ ,

$$\|y_k - x^*\|_{\mathcal{L}} = \|x_k - x^* + \lambda(x_k - x_{k-1} \pm x^*)\|_{\mathcal{L}} \leq 2\|x_k - x^*\|_{\mathcal{L}} + \|x_{k-1} - x^*\|_{\mathcal{L}}.$$

Combining the last two estimates and using (5.19), we obtain (5.20).  $\square$

**Remark 5.10** (total energy). The exponential decrease of the “total energy” at every step, given in (5.18), does not imply that the “potential energy”  $G(x_{k+1}) - G^*$  or the “kinetic energy”  $\frac{\mu}{2} \|v_{k+1} - x^*\|_{\mathcal{L}}^2$  decay monotonically by themselves. Corollary 5.9 only asserts exponential bounds. The same is true for the decay of the  $\mathcal{L}^{-1}$ -norm of the residuals  $\|G'(y_k)\|_{\mathcal{L}^{-1}}$ . In fact, the numerical illustrations of Section 6 show that these quantities may oscillate.  $\blacksquare$

**Remark 5.11** (matching convergence rates). As discussed in Section 3, in the case of  $G$  being locally Lipschitz smooth,  $\mu$ -strongly convex, the (best) contraction factor for PGD is  $(\frac{1-\rho}{1+\rho})^2$  while we have  $1 - \sqrt{\rho}$  for PAGD (see Theorem 5.8), where we recall  $\rho = \mu/L_B$  and  $L_B > 0$  is the Lipschitz smoothness constant on some appropriate invariant set  $\mathfrak{B}$ . It must be pointed out that this rate for PGD is achieved by choosing a “particularly good” step size that is only available to PGD:  $s = \frac{2}{L_B + \mu}$  (see [22, Theorem 2.1.15]). More specifically, we have a contraction factor  $1 - s \frac{2\mu L_B}{L_B + \mu}$  for PGD provided  $0 < s \leq \frac{2}{L_B + \mu}$ . If one uses the step size  $s = 1/L_B$ , then the contraction factor for PGD turns out to be  $\frac{1-\rho}{1+\rho}$ . This choice makes it easier to see the rate match the continuous time model. Setting  $p = 2$  in [30, SI (Supplement Information) Theorem H.2] we see that the gradient flow  $\dot{X} = -\mathcal{L}^{-1}G'(X)$  has convergence rate  $G(X(t)) - G^* \leq (G(X(0)) - G^*)e^{-\mu t}$ . However, using an estimate available to  $\mu$ -strongly convex functions (see [22, Theorem 2.1.10 (2.1.19)]), we can do better to get  $G(X(t)) - G^* \leq (G(X(0)) - G^*)e^{-2\mu t}$ . Then, we see that setting  $t = sk$  and  $s = 1/L_B$  for the gradient flow, and assuming  $\rho \ll 1$ , the contraction factor can be approximated by

$$e^{-2\mu s} \approx 1 - 2\mu s = 1 - 2\rho,$$

which is close to  $\frac{1-\rho}{1+\rho}$ . Similarly, setting  $t = \sqrt{s}k$  and  $s = 1/L_B$  in (4.5) and referring to Corollary 5.9, we have the contraction factors

$$e^{-\sqrt{\mu s}} \approx 1 - \sqrt{\mu s} = 1 - \sqrt{\rho}$$

for the IVP (4.1), which matches that of PAGD.  $\blacksquare$

## 6. NUMERICAL EXPERIMENTS

In this section, we carry out a series of numerical experiments aimed at illustrating the theory that we have developed. In all our examples, we approximate the solution to the nonlinear PDE (6.1) by iteratively minimizing an energy related to this PDE. The approximate solution is computed using a pseudo-spectral method (see [10, 26]), which was implemented in an in-house Matlab R2016a© code. This pseudo-spectral code heavily uses the built-in `fft` and `ifft` Matlab internal routines to invert preconditioners and apply residuals.

The energy minimization is carried out with GD, AGD, PGD, or PAGD, where the algorithm terminates if one of the following cases is true:

- (a) the  $\infty$ -norm (when the true solution is unknown) or the  $\mathcal{L}_N$ -norm (when the true solution is known) of the search direction is smaller than a certain tolerance, which we will call *convergence*;
- (b) the norm being measured is larger than a certain upper tolerance, which we will call *blow up*;
- (c) the number of iterations reaches a certain number, which we will call *no convergence*.

In the conditions above, we mean by “search directions” the residual if the scheme does not involve a preconditioner. If the scheme involves a (discrete) preconditioner  $\mathcal{L}_N$  (see (6.14) for definition), the search direction is the solution to  $\mathcal{L}_N s = r$ , where  $r$  is the residual. In all implementations, the initial guess is always zero.

**6.1. The continuous problem.** We approximate the solution to the following “nonlocal” PDE:

$$(6.1) \quad (-\Delta)^\alpha u + |u|^{p-2}u + tu = f \quad \text{in } \Omega = (0, 1)^2 \subset \mathbb{R}^2,$$

supplemented with periodic boundary conditions, where  $\alpha > 0$ ,  $p \geq 2$ , and  $t > 0$ . Here and in what follows, all functions are real-valued except for the exponential functions appearing in Fourier series and Fourier coefficients. The nonlocal operator  $(-\Delta)^\alpha$  is the *spectral* fractional Laplacian, which is defined via Fourier series as

For every  $v \in L^2_{\text{per}}(\Omega)$ , we have that  $v(\mathbf{x}) = \sum_{\mathbf{m} \in \mathbb{Z}^2} \hat{v}_{\mathbf{m}} e^{2\pi i \mathbf{m} \cdot \mathbf{x}}$ , where the equality is in the  $L^2(\Omega)$ -sense,  $\mathbf{x} = (x, y) \in \bar{\Omega}$ ,  $\mathbf{i} = \sqrt{-1}$ , and  $\hat{v}_{\mathbf{m}} = \int_{\Omega} v(\mathbf{x}) e^{-2\pi i \mathbf{m} \cdot \mathbf{x}} d\mathbf{x}$ ,  $\mathbf{m} \in \mathbb{Z}^2$ .

Thus, we define

$$(-\Delta)^\alpha v(\mathbf{x}) = \sum_{\mathbf{m} \in \mathbb{Z}^2} (4\pi^2 |\mathbf{m}|^2)^\alpha \hat{v}_{\mathbf{m}} e^{2\pi i \mathbf{m} \cdot \mathbf{x}},$$

provided that the sum is finite.

With this definition at hand, it is not difficult to see that (6.1) in its weak form, can be seen as the Euler-Lagrange equation for the functional

$$(6.2) \quad G(u) = \int_{\Omega} \left( \frac{1}{2} |(-\Delta)^{\frac{\alpha}{2}} u|^2 + \frac{1}{p} |u|^p + \frac{t}{2} |u|^2 - fu \right) d\mathbf{x},$$

over the space  $\mathbb{H} = H^{\alpha}_{\text{per}}(\Omega) \cap L^p(\Omega)$ . It is well known that  $\{e^{2\pi i \mathbf{m} \cdot \mathbf{x}}\}_{\mathbf{m} \in \mathbb{Z}^2}$  is an orthonormal basis of  $L^2_{\text{per}}(\Omega)$ . Then,  $\mathbb{H}$  can be equivalently defined via

$$\mathbb{H} = \left\{ v \in L^p_{\text{per}}(\Omega) \left| \sum_{\mathbf{m} \in \mathbb{Z}^2} |\mathbf{m}|^\alpha |\hat{v}_{\mathbf{m}}|^2 < \infty \right. \right\}.$$

The existence and uniqueness of a weak solution to (6.1) is guaranteed for any  $f \in L^p_{\text{per}}(\Omega)$ , where  $1/p + 1/p' = 1$ , since, in this case, the energy is well-defined, strictly convex, and coercive.

For the space  $\mathbb{H}$  to possess a Hilbert structure, a restriction on  $p$  must be imposed depending on  $\alpha$ . For ease of notation, let  $(\cdot, \cdot)$  and  $\|\cdot\|$  denote the  $L^2(\Omega)$ -inner product and  $L^2(\Omega)$ -norm respectively. A natural inner product on  $H^\alpha_{\text{per}}(\Omega)$  is given by  $(v, w)_{H^\alpha_{\text{per}}(\Omega)} = ((-\Delta)^{\frac{\alpha}{2}} v, (-\Delta)^{\frac{\alpha}{2}} w) + (v, w)$ ,

and its associated norm by  $\|v\|_{H^\alpha_{\text{per}}(\Omega)} = \sqrt{(v, v)_{H^\alpha_{\text{per}}(\Omega)}}$ . The following is a standard Sobolev embedding result. For a proof, see, e.g., [1, Theorem 7.34].

**Proposition 6.1** (Sobolev embedding). *Let  $\alpha \in (0, 1]$ . For all  $p \in [2, p^*]$  with  $p^* = \frac{2}{1-\alpha}$  if  $\alpha < 1$  or  $p \in [2, \infty)$  if  $\alpha = 1$ , there exists  $C_{\text{emb}} = C_{\text{emb}}(p, \alpha) > 0$  such that, for all  $v \in H^\alpha_{\text{per}}(\Omega)$ ,*

$$(6.3) \quad \|v\|_{L^p(\Omega)} \leq C_{\text{emb}} \|v\|_{H^\alpha_{\text{per}}(\Omega)}.$$

We introduce the preconditioner  $\mathcal{L}u = (-\Delta)^\alpha u + \nu u$ ,

where  $\nu \geq 0$  is a free parameter, which induces a inner product

$$(6.4) \quad \langle \mathcal{L}u, v \rangle = \int_{\Omega} ((-\Delta)^{\frac{\alpha}{2}} u (-\Delta)^{\frac{\alpha}{2}} v + \nu uv) \, d\mathbf{x}.$$

**Remark 6.2** (notation). As it is clear from its definition, the Lipschitz constant of  $G'$  depends on the norm being used. Thus, we will make a difference between the case with preconditioner and without it.  $\hat{L}$  denotes the Lipschitz constant with respect to the preconditioner-induced norm  $\|\cdot\|_{\mathcal{L}}$ , while  $L$  is the constant with respect to the original norm  $\|\cdot\|_{\mathbb{H}}$ . ■

We investigate the properties of  $G$  in the following result.

**Proposition 6.3** (properties of  $G$ ). *Let  $G$  be given by (6.2) and the preconditioner  $\mathcal{L}$  by (6.4). Then,  $G$  is strongly convex with respect to  $\mathcal{L}$ -norm. If, in addition,  $p$  satisfies the conditions of Proposition 6.1, then  $G$  is locally Lipschitz smooth with respect to  $\mathcal{L}$ -norm.*

*Proof.* First, the action of  $G'$  is characterized by the following: for  $v, w \in \mathbb{H}$ ,

$$\langle G'(v), w \rangle = ((-\Delta)^{\frac{\alpha}{2}} v, (-\Delta)^{\frac{\alpha}{2}} w) + t(v, w) + (|v|^{p-2} v, w) - (f, w).$$

Note also that the following estimates hold, which are a special case of [4, Lemma 2.1]: for  $p > 1$ , there exist  $C_{p1}, C_{p2} > 0$ , which depend only on  $p$ , such that for all  $\xi, \eta \in \mathbb{R}$ ,

$$(6.5) \quad \left| |\xi|^{p-2} \xi - |\eta|^{p-2} \eta \right| \leq C_{p1} |\xi - \eta| (|\xi| + |\eta|)^{p-2},$$

$$(6.6) \quad (|\xi|^{p-2} \xi - |\eta|^{p-2} \eta)(\xi - \eta) \geq C_{p2} |\xi - \eta|^2 (|\xi| + |\eta|)^{p-2}.$$

Thus, using (6.6)

$$\langle G'(v) - G'(w), v - w \rangle \geq \|(-\Delta)^{\frac{\alpha}{2}}(v - w)\|^2 + t \|v - w\|^2 \geq \hat{\mu} \|v - w\|_{\mathcal{L}}^2,$$

where

$$(6.7) \quad \hat{\mu} = \min\{1, t/\nu\}.$$

Observe that this holds without referring to Sobolev embedding. Note also that this implies the coercivity of  $G$  with respect to  $\mathcal{L}$ -norm, that is,  $\lim_{\|v\|_{\mathcal{L}} \rightarrow \infty} G(v) = \infty$ .

Next, thanks to coercivity, for any bounded, convex set  $B \subset \mathbb{H}$  there exists  $M_B \in \mathbb{R}$  such that  $B \subset \{x \in \mathbb{H} \mid G(x) \leq M_B\}$ . Hence, for each  $v \in B$ , using Cauchy-Schwarz inequality and Young's inequality, it follows that there exists  $\varepsilon > 0$  such that

$$(6.8) \quad \begin{aligned} M_B \geq G(v) &= \frac{1}{2} \|(-\Delta)^{\frac{\alpha}{2}} v\|^2 + \frac{1}{p} \|v\|_{L^p(\Omega)}^p + \frac{t}{2} \|v\|^2 - (f, v) \\ &\geq \frac{1}{2} \|(-\Delta)^{\frac{\alpha}{2}} v\|^2 + \frac{1}{p} \|v\|_{L^p(\Omega)}^p + \frac{t}{4} \|v\|^2 - \frac{1}{2\varepsilon} \|f\|^2. \end{aligned}$$

Rearranging this, we see that there exists  $C_{f,t,p,B} > 0$  such that

$$(6.9) \quad \|v\|_{L^p(\Omega)} \leq C_{f,t,p,B} \quad \forall v \in B.$$

On the other hand, using (6.5), Hölder's inequality, and (6.9), we have, for all  $v, w \in B$ ,

$$\begin{aligned} &\langle G'(v) - G'(w), v - w \rangle \\ &= \|(-\Delta)^{\frac{\alpha}{2}}(v - w)\|^2 + t\|v - w\|^2 + (|v|^{p-2}v - |w|^{p-2}w, v - w) \\ &\leq \|(-\Delta)^{\frac{\alpha}{2}}(v - w)\|^2 + t\|v - w\|^2 + C_{p1} \int_{\Omega} |v - w|^2 (|v| + |w|)^{p-2} d\mathbf{x} \\ &\leq \|(-\Delta)^{\frac{\alpha}{2}}(v - w)\|^2 + t\|v - w\|^2 + C_{p3} \|v - w\|_{L^p(\Omega)}^2 \left( \|v\|_{L^p(\Omega)}^{p-2} + \|w\|_{L^p(\Omega)}^{p-2} \right) \\ &\leq \|(-\Delta)^{\frac{\alpha}{2}}(v - w)\|^2 + t\|v - w\|^2 + 2C_{f,t,p,B}^{p-2} C_{p3} \|v - w\|_{L^p(\Omega)}^2, \end{aligned}$$

where  $C_{p3} > 0$  is a constant reflecting the equivalence between  $(|v| + |w|)^{p-2}$  and  $|v|^{p-2} + |w|^{p-2}$ .

Finally, owing to the restriction on  $p$ , Proposition 6.1 guarantees that

$$\|v - w\|_{L^p(\Omega)}^2 \leq C_{emb}^2 \|v - w\|_{H_{\text{per}}^{\alpha}(\Omega)}^2,$$

so that

$$\langle G'(v) - G'(w), v - w \rangle \leq \hat{L}_B \|(v - w)\|_{\mathcal{L}}^2,$$

with  $\hat{L}_B = \max\{1, t/\nu, 2C_{f,t,p,B}^{p-2} C_{p3} C_{emb}^2\}$ .  $\square$

**Remark 6.4** (strong Lipschitz smoothness). The proof of Proposition 6.3 can be easily modified to show that  $G$  is locally Lipschitz smooth in the strong sense, i.e., (2.8) holds.  $\blacksquare$

**6.2. Discretization.** We discretize the model problem (6.1) by introducing a uniform grid of points. To simplify the presentation, we choose  $N \in \mathbb{N}$  with  $N = 2K + 1$  for some integer  $K \geq 1$ . (The details for the case that  $N$  is even are only slightly more complicated.) Define  $h = 1/N$ , and introduce the grid domain  $\Omega_N = \{(x_{\ell}, y_m) \in [0, 1]^2 \mid x_{\ell} = \ell h, y_m = m h, 0 \leq \ell, m \leq N\}$ . For ease of notation, let us introduce  $\mathbb{N}_N^2 = \{\mathbf{m} = (m_1, m_2) \in \mathbb{Z}^2 \mid 1 \leq m_1, m_2 \leq N\}$  and  $\mathbb{Z}_K^2 = \{\mathbf{r} = (r_1, r_2) \in \mathbb{Z}^2 \mid -K \leq r_1, r_2 \leq K\}$ . Then, for  $\mathbf{m} \in \mathbb{N}_N^2$ , we can denote  $\mathbf{x}_{\mathbf{m}} = (x_{m_1}, y_{m_2}) \in \Omega_N$ . This notation must not be confused with that of the iterates of PAGD. Define the space of periodic grid functions

$$(6.10) \quad \begin{aligned} \mathbb{H}_N &= \{v_N : \Omega_N \rightarrow \mathbb{R} \mid v_N(0, hm) = v_N(hN, hm), v_N(h\ell, 0) = v_N(h\ell, hN), \\ &\quad 0 \leq m, \ell \leq N\}, \end{aligned}$$



endowed with the  $L_N^2$ -inner product  $(v_N, w_N)_N = h^2 \sum_{\mathbf{m} \in \mathbb{N}_N^2} v_N(\mathbf{x}_\mathbf{m}) w_N(\mathbf{x}_\mathbf{m})$ .

More generally, for  $p \geq 1$ , define  $\|w_N\|_{N,p} = \left( h^2 \sum_{\mathbf{m} \in \mathbb{N}_N^2} |w_N(\mathbf{x}_\mathbf{m})|^p \right)^{\frac{1}{p}}$ . Given  $w_N \in \mathbb{H}_N$ , its *discrete Fourier transform* (DFT) is

$$\hat{w}_K(\mathbf{r}) = h^2 \sum_{\mathbf{s} \in \mathbb{N}_N^2} w_N(\mathbf{x}_\mathbf{s}) e^{-2\pi i \mathbf{r} \cdot \mathbf{x}_\mathbf{s}}, \quad \mathbf{r} \in \mathbb{Z}_K^2.$$

The *discrete fractional Laplacian*  $(-\Delta_N)^\alpha : \mathbb{H}_N \rightarrow \mathbb{H}_N$  is defined by

$$(6.11) \quad [(-\Delta_N)^\alpha w_N](\mathbf{x}_\mathbf{m}) = \sum_{\mathbf{r} \in \mathbb{Z}_K^2} (4\pi^2 |\mathbf{r}|^2)^\alpha \hat{w}_K(\mathbf{r}) e^{2\pi i \mathbf{r} \cdot \mathbf{x}_\mathbf{m}}.$$

Finally, for  $v_N, w_N \in \mathbb{H}_N$ , the  $H_N^\alpha$ -inner product is given by  $(v_N, w_N)_{H_N^\alpha} = (v_N, w_N)_N + ((-\Delta_N)^{\frac{\alpha}{2}} v_N, (-\Delta_N)^{\frac{\alpha}{2}} w_N)_N$ , and  $\|w_N\|_{H_N^\alpha} = \sqrt{(w_N, w_N)_{H_N^\alpha}}$ .

We comment that there are, at least, three different natural choices for the underlying inner product for  $\mathbb{H}_N$ : the  $L_N^2$ -inner product, the  $H_N^\alpha$ -inner product, and the  $\mathcal{L}_N$ -inner product, (6.14). We will choose the first option, i.e., the  $L_N^2$ -inner product for several reasons. To begin with, this is the way numerical experiments are usually done if no special distinction is made between representers of the residual with respect to multiple inner products. In addition, this illustrates the effect of preconditioning more vividly. For example, if we adopt  $H_N^\alpha$ -inner product, this leads to a preconditioned scheme in disguise: finding a representer of the residual with respect to this inner product is equivalent to using  $\mathcal{L}_N$ -inner product (6.14) with  $\nu_N = 1$ . On the other hand,  $L_N^2$ -inner product leads to truly *non-preconditioned* schemes such as GD or AGD.

After having introduced all this notation, we can write our discrete problem as: given  $f_N \in \mathbb{H}_N$ , find  $u_N \in \mathbb{H}_N$  such that

$$(6.12) \quad (-\Delta_N)^\alpha u_N + |u_N|^{p-2} u_N + t u_N = f_N.$$

In this problem,  $f_N \in \mathbb{H}_N$  is some approximation of the problem data  $f$ . For example, if  $f$  is continuous,  $f_N(\mathbf{x}_\mathbf{m}) = f(\mathbf{x}_\mathbf{m})$  is a natural option, and if  $f$  is only an  $L^2(\Omega)$ -function, then the sampling at the nodes of the  $L^2(\Omega)$ -projection of  $f$  onto  $\mathcal{P}_K$ , the trigonometric polynomial of degree at most  $K$ , is natural although these two may not agree even if one starts with the same continuous function  $f$ . In fact, the difference between these two possibilities is very small if  $f$  is smooth and its derivatives are periodic (see [10, pp. 44–45]).

Our discrete problem has a similar energy structure to the continuous problem. It is the Euler-Lagrange equation of the following functional

$$(6.13) \quad G_N(v_N) = \frac{1}{2} \|(-\Delta_N)^{\frac{\alpha}{2}} v_N\|_N^2 + \frac{1}{p} \|v_N\|_{N,p}^2 + \frac{t}{2} \|v_N\|_N^2 - (f_N, v_N)_N.$$

We introduce a (discrete) preconditioner

$$(6.14) \quad \mathcal{L}_N = (-\Delta_N)^\alpha + \nu_N \text{id}_N,$$

where  $\nu_N > 0$  and  $\text{id}_N : \mathbb{H}_N \rightarrow \mathbb{H}_N$  is the identity map. The parameter  $\nu_N > 0$  will be determined later. This preconditioner induces an inner product on  $\mathbb{H}_N$  given by

$$(6.15) \quad (v_N, w_N)_{\mathcal{L}_N} = \nu_N (v_N, w_N)_N + ((-\Delta_N)^{\frac{\alpha}{2}} v_N, (-\Delta_N)^{\frac{\alpha}{2}} w_N)_N,$$

and an associated norm  $\|v_N\|_{\mathcal{L}_N} = \sqrt{(v_N, v_N)_{\mathcal{L}_N}}$ . It is desirable that the convergence of our scheme does not deteriorate as we refine the grid points. We can

ensure this under a certain restriction on  $p$ . The following proposition provides an important tool for that purpose.

**Proposition 6.5** (discrete Sobolev embedding). *Let  $\alpha \in (0, 1]$ . For all  $p \in [2, p^*]$  with  $p^* = \frac{2}{1-\alpha}$  if  $\alpha < 1$  or for all  $p \in [2, \infty)$  if  $\alpha = 1$ , there exists a constant  $C_{p,\alpha} > 0$  such that, for all  $v_N \in \mathbb{H}_N$ ,*

$$(6.16) \quad \|v_N\|_{N,p} \leq C_{p,\alpha} \|v_N\|_{H_N^\alpha}.$$

$C_{p,\alpha}$  is independent of  $v_N$  and  $N$ .

*Proof.* Note that  $\|v_N\|_{N,p} \leq C \|v\|_{L^p(\Omega)}$  for all  $v_N \in \mathbb{H}_N$ ; see, for instance, [18, Lemma 2.48], where  $v$  is the unique trigonometric polynomial of degree less than or equal to  $K$  interpolating  $v_N$  and  $C > 0$  depends only on the dimension of  $\Omega$ . Also, the following Parseval's identity holds in the fractional setting  $\|v_N\|_{H_N^\alpha} = \|v\|_{H_{\text{per}}^\alpha(\Omega)}$ . In conjunction with the Sobolev embedding at the continuous level, (6.3), we have, for any  $v_N \in \mathbb{H}_N$ ,

$$\|v_N\|_{N,p} \leq C \|v\|_{L^p(\Omega)} \leq CC_{p,\alpha} \|v\|_{H_{\text{per}}^\alpha(\Omega)} = CC_{p,\alpha} \|v_N\|_{H_N^\alpha}. \quad \square$$

The following result addresses dimension-independence of the (inverse) condition number as well as the structure of  $G_N$  that is needed to apply the theory we have developed in Section 5. Note that, in the following statement, the sublevel sets of  $G_N$  provide a compatible way to describe bounded, convex sets when we consider multiple resolutions since, strictly speaking, for different values of  $N$ , functions in  $\mathbb{H}_N$  may not be directly comparable.

**Theorem 6.6** (properties of  $G_N$ ). *Let the space of grid functions  $\mathbb{H}_N$  be given by (6.10) and the preconditioner  $\mathcal{L}_N$  by (6.14). Then, the energy functional  $G_N : \mathbb{H}_N \rightarrow \mathbb{R}$  defined by (6.13) is strongly convex and locally Lipschitz smooth with respect to the  $\mathcal{L}_N$ -norm. Moreover, the strong convexity constant  $\hat{\mu}_N$  is independent of  $N$ . Suppose, in addition, that  $p$  satisfies the conditions of Proposition 6.5 and that  $f_N$  is defined in a stable manner when we pose the discrete problem (6.12), i.e., there exists  $C > 0$  independent of  $N$  such that*

$$\|f_N\|_{N,2} \leq C \|f\|_{L^2(\Omega)}.$$

*Then, the local Lipschitz smoothness constant  $\hat{L}_N$  is also independent of  $N$  in the sense that, for each  $M \in \mathbb{R}$ ,  $G_N$  is  $\hat{L}_N$ -Lipschitz smooth on the sublevel set*

$$\{v_N \in \mathbb{H}_N \mid G(v_N) \leq M\}$$

*with  $\hat{L}_N$  independent of  $N$ . Consequently, the (inverse) condition number  $\hat{\mu}_N/\hat{L}_N$  with respect to the  $\mathcal{L}_N$ -norm stays away from 0 as  $N \rightarrow \infty$ .*

*Proof.* The proof of the strong convexity is parallel to that of Proposition 6.3. The proof of local Lipschitz smoothness is also parallel, but we need the assumed stability of  $\|f_N\|_{N,2}$  to proceed from (6.8) to (6.9). Finally, to complete the proof, we simply replace the embedding constant  $C_{emb}$  with its discrete counterpart  $C_{p,\alpha}$  as given in Proposition 6.5.  $\square$

Since the (inverse) condition number  $\hat{\mu}_N/\hat{L}_N$  governs the rate of convergence of PAGD (Corollary 5.9), the previous theorem guarantees that one can achieve the same rate of convergence even if we refine the number of grid points  $N \rightarrow \infty$ .

However, this is true in terms of the number of iterations, but the wall clock time will take longer as the refinement is conducted.

Being in finite dimensions,  $G_N$  is also strongly convex, and locally Lipschitz smooth with respect to *any norm*, for instance  $\|\cdot\|_N$ . The constants in this case, however, are different and depend on the dimension of  $\mathbb{H}_N$ , which obviously depends on the number of grid points, and thus on  $N$ . We label them  $\mu_N$  and  $L_N$  to distinguish them from the dimension-independent constants  $\hat{\mu}_N$  and  $\hat{L}_N$  respectively.

**6.3. A problem with a manufactured solution.** In this first experiment, we solve (6.12) by minimizing the energy (6.13). To compute the errors and energies, the following manufactured solution is used

$$u_N(\mathbf{x}_m) = \exp\left(\sin 2\pi\left(x_{m_1} - \frac{1}{4}\right) + \sin 4\pi\left(y_{m_2} - \frac{3}{8}\right)\right).$$

We set  $\alpha = 0.5$ ,  $p = 4$ ,  $t = 1$ ,  $N = 64$ , and found, experimentally, that the values  $\nu_N = 1.2$ ,  $\mu_N = 1$  are optimal, while we set  $\hat{\mu}_N = 5/6 = \min\{1, t/\nu_N\}$  in view of (6.7). To specify step sizes, recall the step size rules that theoretically guarantee convergence (see Section 3):  $s = 2/(L_N + \mu_N)$  for GD,  $s = 1/L_N$  for AGD,  $s = 2/(\hat{L}_N + \hat{\mu}_N)$  for PGD, and  $s = 1/\hat{L}_N$  for PAGD. Step sizes are set by these relations with  $L_N = 500$  and  $\hat{L}_N = 20$ , which are also experimentally proved to be optimal. However, it must be noted that this is just a way of setting step sizes. We do not really know neither whether the values for  $L_N$  or  $\hat{L}_N$  are the Lipschitz constants of the corresponding energy functionals nor whether the aforementioned step size rules give the optimal results even if we knew the Lipschitz constants. In fact, our last experiment suggests that larger step sizes than what is theoretically proven seem to work.

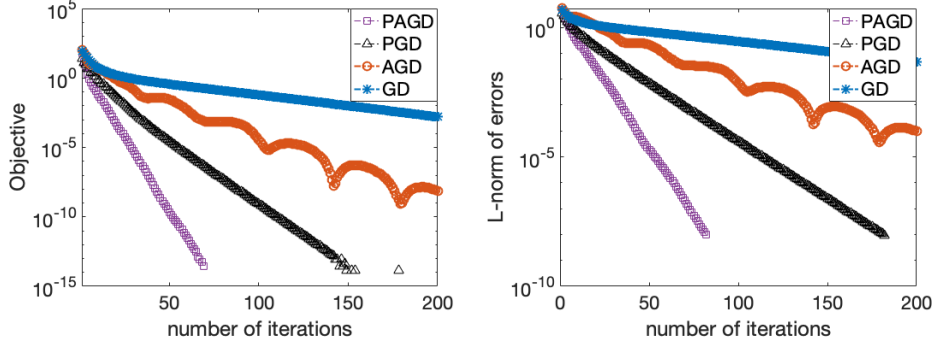
Figure 2 shows the performance of GD, AGD, PGD, and PAGD when used to solve (6.12) by minimizing (6.13), where the data is as described above. The stopping criteria take the following parameters: the tolerance is  $10^{-8}$ , the upper tolerance is  $10^{10}$ , and the maximum number of iterations is 200.

Figure 2 (A) shows the decay of the objective,  $G_N(x_k) - G_N(u_N)$  which is, up to a constant, the same as the decay of the potential energy, for all four schemes. Here  $k$  is the number of iterations. Figure 2 (B) shows the decay of the  $\mathcal{L}_N$ -norm of the errors. Notice that PAGD performs significantly better than all the other methods.

Figure 2 (C) shows the performance of GD and AGD. Since these schemes do not involve a preconditioner, the corresponding total energy is defined by

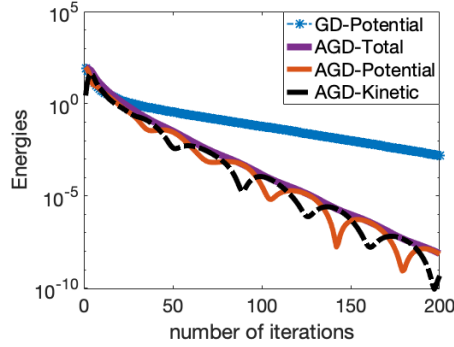
$$(6.17) \quad E_N(x_k, v_k) = \frac{1}{\eta_N} (G_N(x_k) - G_N(u_N)) + \frac{\eta_N}{2} \|v_k - u_N\|_N^2,$$

where  $k$  is the number of iterations and  $\eta_N = \sqrt{\mu_N}$ . The first and second terms can be understood as potential and kinetic energy respectively. Figure 2 (C) shows the decay of various energies for nonpreconditioned schemes. This figure better illustrates our analysis of the previous section than the preconditioned ones since they converge slower. As expected, AGD performs substantially better than GD. The total energy of AGD decreases steadily and exponentially fast. Notice that the vertical axis is in logarithmic scale. This matches what is predicted by the theory in Theorem 5.8. Observe also that the potential and kinetic energies of AGD, by themselves, oscillate; see Remark 5.10. The physical analogy for AGD described



(A) Objective,  $G_N(x_k) - G_N(u_N)$ , plot of GD, AGD, PGD, and PAGD.

(B)  $\mathcal{L}_N$ -norm of errors that are generated by GD, AGD, PGD, and PAGD.



(C) Potential, kinetic, and total energy plot of GD and AGD.

FIGURE 2. Objective, error, and energy decay plots for GD, AGD, PGD, and PAGD. They are implemented to solve (6.12) by minimizing (6.13) ( $\alpha = 0.5$ ,  $p = 4$ ,  $t = 1$ ,  $N = 64$ ,  $\nu_N = 1.2$ ,  $L_N = 500$ ,  $\mu_N = 1$ ,  $\hat{L}_N = 20$ ,  $\hat{\mu}_N = 5/6 = \min\{1, t/\nu_N\}$ , and step sizes are set via  $s = 2/(L_N + \mu_N)$  for GD,  $s = 1/L_N$  for AGD,  $s = 2/(\hat{L}_N + \hat{\mu}_N)$  for PGD, and  $s = 1/\hat{L}_N$  for PAGD). The vertical axes (logarithmic scale) show the value of the objective,  $\mathcal{L}_N$ -norm of errors, or various energies while the horizontal axis (linear scale) shows the number of iterations.

in Remark 4.1 is clear from this picture. A fraction of the potential energy is converted to kinetic energy and they fluctuate as the mechanical system converges to equilibrium.

**6.4. A problem where the solution is unknown.** In this second experiment we, again, solve (6.12) by minimizing the energy (6.13). The discrete right hand

side  $f_N$  is given by

$$(6.18) \quad f_N(\mathbf{x}_m) = \exp(\sin 2\pi(x_{m_1} - 0.25) + \sin 2\pi(y_{m_2} - 0.25)).$$

The parameters of the PDE are set to  $\alpha = 0.5$ ,  $p = 10$ , and  $t = 1$ . Observe that for these values of  $\alpha$  and  $p$  we do not have that  $H_{\text{per}}^\alpha(\Omega) \hookrightarrow L^p(\Omega)$ . We found, experimentally, that the choice  $\nu_N = 0.9$  is optimal for the preconditioner. We also set  $\mu_N = 1$  and  $\hat{\mu}_N = 1 = \min\{1, t/\nu_N\}$  in view of (6.7) as before. Step sizes are set in the same way as in the previous experiment with  $L_N = 300$  or  $3000$ , and  $\hat{L}_N = 9$ . The values of  $\mu_N$ ,  $L_N$ , and  $\hat{L}_N$  were experimentally found to be optimal except for  $L_N = 3000$ . That is, they yield the best convergence rate with all other parameters being fixed. A more detailed explanation about  $L_N = 3000$  will follow. Two different degrees of resolution are used to show the dimension dependence of nonpreconditioned schemes. The stopping criterion parameters are as before.

Figure 3 shows the  $\infty$ -norm of the search directions for GD, AGD, PGD, and PAGD with varying degrees of resolution and with two different step sizes for GD and AGD, which are determined by the same step size rules as in the previous experiment with  $L_N \in \{300, 3000\}$ . In Figure 3 (A), we observe a similar performance as in Figure 2. Recall that we do not have Sobolev embedding. Thus, one can expect the Lipschitz constant  $L_N$ , hence the step size, to depend on the number of grid points. In fact, theory predicts that even  $\hat{L}_N$  depends on it. However, for  $\hat{L}_N$ , such dependence is not observed within the range of  $N$  that we have chosen. We see that the step size for convergence indeed depends on  $N$  in Figure 3 (B). As we increase the resolution of the grid from  $N = 64$  to  $N = 512$ , nonpreconditioned schemes become unstable. Figure 3 (C) shows that the stability of GD and AGD is recovered after  $L_N$  is increased from 300 to 3000, which amounts to decreasing the step size to roughly a tenth of the old one. ( $L_N = 3000$  is not optimally chosen).

Figure 4 shows the dependence of  $L_N$ , hence the step size, on the number of grid points with the same experiment. However, here we use different tolerances and a different maximum number of iterations to best illustrate the dependence. For  $N \in \{16, 32, 64, 128, 256, 512\}$ , Figure 4 records the number of iterations for  $\infty$ -norm of the search direction generated by each scheme to reach a tolerance  $10^{-3}$  (“convergence”) or the maximum number of iterations, which is set to be 1000, if it does not reach the tolerance (“no convergence”). If the  $\infty$ -norm of the search direction reaches an upper tolerance  $10^8$ , the algorithm records the number of iteration taken as 1100, which indicates “blowing up.” Figure 4 (A) shows when the step sizes of the nonpreconditioned schemes correspond to  $L_N = 300$  and those of the preconditioned ones correspond to  $\hat{L}_N = 9$ . GD and AGD converge until  $N = 64$ . However, they become unstable for  $N \geq 128$ . Figure 4 (B) shows the same experiment with smaller step sizes, which correspond to  $L_N = 3000$ . In this case, we recover the stability of GD and AGD.

**6.5. A comparison between PGD and PAGD.** In this final collection of experiments, we aim at comparing the performance of PGD and PAGD in different scenarios. To do so, we solve the discrete problem (6.12) by minimizing the energy (6.13) with the right hand side given by (6.18) as before. The problem parameters are set as  $\alpha \in \{0.1j \mid j = 1, 2, 3, \dots, 10\} \cup \{1.5, 2.0, 2.5, 3.0\} \subset (0, 3]$ ,  $p = 6$ , and  $t = 1$ . We set  $N = 64$  and  $\hat{\mu}_N = \min\{1, t/\nu_N\}$ . Then, for each value of  $\alpha$  (column 1 of Table 1), PGD and PAGD are applied with  $\nu_N \in \{0.1j \mid j = 1, 2, 3, \dots, 100\} \subset (0, 10]$  and the step size  $s \in \{0.01j \mid j = 1, 2, 3, \dots, 200\} \subset (0, 2]$ . Observe that

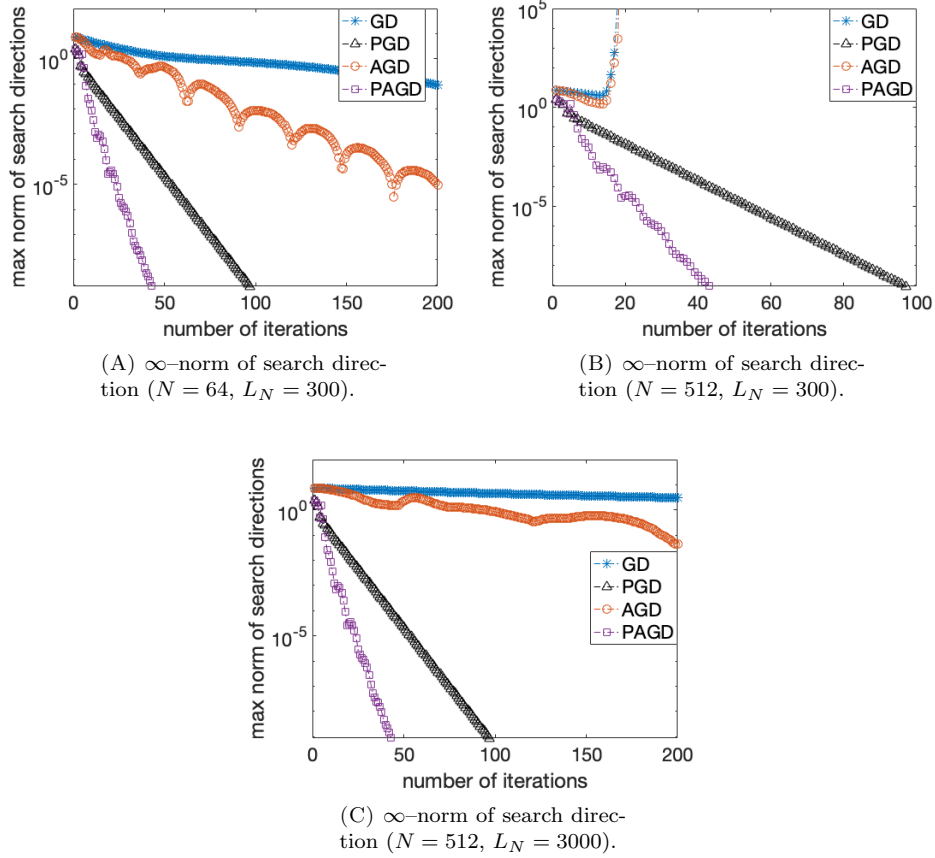


FIGURE 3.  $\infty$ -norm plots of the search directions for GD, AGD, PGD, and PAGD. They are implemented to solve (6.12) by minimizing (6.13) with varying resolutions  $N \in \{64, 512\}$  and varying step sizes for GD and AGD;  $s = 2/(L_N + \mu_N)$  for GD and  $s = 1/L_N$  for AGD with  $L_N \in \{300, 3000\}$ . The other parameters are set to  $\alpha = 0.5, p = 10, t = 1, \nu_N = 0.9, \mu_N = 1, \hat{\mu}_N = 1 = \min\{1, t/\nu_N\}, s = 2/(\hat{L}_N + \hat{\mu}_N)$  for PGD, and  $s = 1/\hat{L}_N$  for PAGD with  $\hat{L}_N = 9$ . The horizontal axis (linear scale) represents the number of iterations. The vertical axis (logarithmic scale) represents  $\infty$ -norm of the search directions.

neither Algorithm 1 nor Algorithm 2 require knowledge of  $\hat{L}_N$  a priori. Instead, we directly set the step size in this last experiment. Among these 20,000 possible values of  $\nu_N$  and  $s$ , the minimal number of iterations for the  $\infty$ -norm of the search direction generated by PGD and PAGD to reach a tolerance of  $10^{-9}$  (convergence) is recorded (column 2 and column 5 of Table 1, respectively). A pair of values,  $\nu_N$  and  $s$ , that led to the minimal number of iterations is also recorded (columns 3 and 4 of Table 1 for PGD and columns 6 and 7 of Table 1 for PAGD). There can be

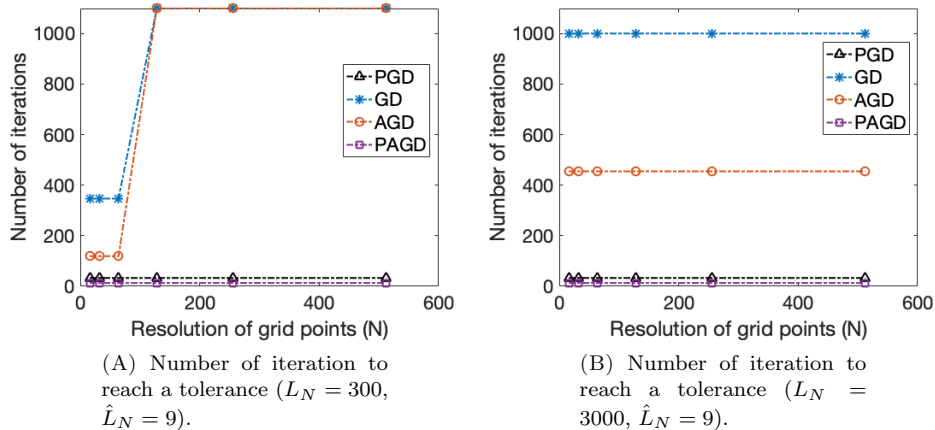


FIGURE 4. Number of iterations for  $\infty$ -norm of the search directions to reach the tolerance  $10^{-3}$  for GD, AGD, PGD, and PAGD. They are implemented to solve (6.12) by minimizing (6.13) with varying resolutions  $N = 16, 32, 64, 128, 256, 512$  and varying  $L_N = 300, 3000$  ( $\alpha = 0.5$ ,  $p = 10$ ,  $t = 1$ ,  $\nu_N = 0.9$ ,  $\mu_N = 1$ ,  $L_N$  as indicated in the subfigures,  $\hat{\mu}_N = 1 = \min\{1, t/\nu_N\}$ ,  $\hat{L}_N = 9$ ). The horizontal axis represents the degrees of resolution,  $N$ . The vertical axis represents the minimum of the number of iterations for the  $\infty$ -norm of the search directions to reach the tolerance  $10^{-3}$  (convergence) or 1000 iterations. The number of iterations being 1100 means that the  $\infty$ -norm of the search directions have reached the upper tolerance  $10^8$  (blow up).

multiple such pairs. If this is the case, the pair  $(\nu_N, s)$  that comes the first in the lexicographic order is recorded.

As we can see from Table 1, for the nonlocal PDE (6.12) with small  $\alpha$  ( $\alpha = 0.1, 0.2, 0.3, 0.4$ ), PAGD performs better than PGD when they are implemented with their own best pair of parameters  $\nu_N$  and  $s$  among those pairs that were considered. In particular, in the cases of  $\alpha = 0.1, 0.2, 0.3$ , PAGD outperforms PGD while the best values of  $\nu_N$  for the two schemes are similar, hence directly comparing their performances roughly make sense. An interesting thing, however, is that one cannot say that PAGD is *always* better than PGD. In fact, for the remaining values of  $\alpha$ , PGD takes fewer iterations to converge in the aforementioned sense than PAGD provided they are equipped with their “best” parameters for each method. It must be noted that this result does not contradict our theory. The theory only tells us some upper bounds about the rate of convergence of the two schemes within a certain range of step size when they involve the same preconditioner. It does not explain what happens outside of that. The result provided here perhaps illustrates the latter case. In any case, we can see an improvement in the convergence of PAGD compared to PGD for “harder” problems ((6.12) with small  $\alpha$ ), where a stronger nonlocality is involved.

$\alpha$	PGD			PAGD		
	# iterations	$\nu_N$	step size	# iterations	$\nu_N$	step size
0.1	64	1.0	0.20	38	0.9	0.14
0.2	50	1.1	0.25	32	1.0	0.18
0.3	39	1.2	0.31	29	1.1	0.22
0.4	29	2.6	0.57	26	1.2	0.26
0.5	22	2.8	0.66	24	1.3	0.30
0.6	16	4.1	0.97	20	5.5	0.83
0.7	13	3.4	0.90	17	5.2	0.91
0.8	11	4.6	1.04	15	4.2	0.88
0.9	12	3.8	0.89	12	5.0	0.96
1.0	10	4.0	0.95	12	4.3	0.92
1.5	9	4.5	0.97	11	4.5	0.97
2.0	8	4.8	1.03	10	4.5	0.96
2.5	8	4.1	0.88	9	4.2	0.90
3.0	8	4.1	0.88	9	4.2	0.90

TABLE 1. The minimal number of iterations needed for the  $\infty$ -norm of the search direction of PGD and PAGD to reach a tolerance of  $10^{-9}$  and the values of  $\nu_N$  and  $s$  (step size) that led to the minimum iterations for a range of values of  $\alpha$ . They are implemented to solve (6.12) by minimizing the energy (6.13).  $N = 64$ ,  $\alpha \in \{0.1j \mid j = 1, 2, 3, \dots, 10\} \cup \{1.5, 2.0, 2.5, 3.0\} \subset (0, 3]$ ,  $p = 6$ ,  $t = 1$ ,  $\hat{\mu}_N = \min\{1, t/\nu_N\}$ , and  $f_N$  is given by (6.18). For each value of  $\alpha$ , we consider  $\nu_N \in \{0.1j \mid j = 1, 2, 3, \dots, 100\} \subset (0, 10]$  and  $s \in \{0.01j \mid j = 1, 2, 3, \dots, 200\} \subset (0, 2]$ . Among the possible 20,000 possible combinations of  $\nu_N$  and  $s$ , we display the values that give the minimal number of iterations.

#### ACKNOWLEDGMENTS

SMW acknowledges partial financial support from NSF-DMS 1719854. AJS has been partially supported by NSF-DMS 1720123.

#### REFERENCES

- [1] R. A. Adams and J. J. F. Fournier. *Sobolev spaces*, volume 140 of *Pure and Applied Mathematics (Amsterdam)*. Elsevier/Academic Press, Amsterdam, second edition, 2003. ISBN 0-12-044143-8.
- [2] Z. Allen-Zhu and L. Orecchia. Linear coupling: an ultimate unification of gradient and mirror descent. In *8th Innovations in Theoretical Computer Science Conference*, volume 67 of *LIPICs. Leibniz Int. Proc. Inform.*, pages Art. No. 3, 22. Schloss Dagstuhl. Leibniz-Zent. Inform., Wadern, 2017.
- [3] H. Attouch, X. Goudou, and P. Redont. The heavy ball with friction method. I. The continuous dynamical system: global exploration of the local minima of a real-valued function by asymptotic analysis of a dissipative dynamical system. *Commun. Contemp. Math.*, 2(1):1–34, 2000. ISSN 0219-1997. doi: 10.1142/S0219199700000025. URL <https://doi.org/10.1142/S0219199700000025>.



- [4] J. W. Barrett and W. B. Liu. Finite element approximation of the  $p$ -Laplacian. *Math. Comp.*, 61(204):523–537, 1993. ISSN 0025-5718. doi: 10.2307/2153239. URL <https://doi.org/10.2307/2153239>.
- [5] A. Beck. *First-order methods in optimization*, volume 25 of *MOS-SIAM Series on Optimization*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA; Mathematical Optimization Society, Philadelphia, PA, 2017. ISBN 978-1-611974-98-0. doi: 10.1137/1.9781611974997.ch1. URL <https://doi.org/10.1137/1.9781611974997.ch1>.
- [6] M. Benyamin, J. Calder, G. Sundaramoorthi, and A. Yezzi. Accelerated variational PDEs for efficient solution of regularized inversion problems. *J. Math. Imaging Vision*, 62(1):10–36, 2020. ISSN 0924-9907. doi: 10.1007/s10851-019-00910-2. URL <https://doi.org/10.1007/s10851-019-00910-2>.
- [7] D. P. Bertsekas. *Nonlinear programming*. Athena Scientific Optimization and Computation Series. Athena Scientific, Belmont, MA, second edition, 1999. ISBN 1-886529-00-0.
- [8] S. Boyd and L. Vandenberghe. *Convex optimization*. Cambridge University Press, Cambridge, 2004. ISBN 0-521-83378-7. doi: 10.1017/CBO9780511804441. URL <https://doi.org/10.1017/CBO9780511804441>.
- [9] J. Calder and A. Yezzi. PDE acceleration: a convergence rate analysis and applications to obstacle problems. *Res. Math. Sci.*, 6(4):Paper No. 35, 30, 2019. ISSN 2522-0144. doi: 10.1007/s40687-019-0197-x. URL <https://doi.org/10.1007/s40687-019-0197-x>.
- [10] C. Canuto, M. Y. Hussaini, A. Quarteroni, and T. A. Zang. *Spectral methods*. Scientific Computation. Springer-Verlag, Berlin, 2006. ISBN 978-3-540-30725-9; 3-540-30725-7. Fundamentals in single domains.
- [11] L. Chen, X. Hu, and S. M. Wise. Convergence analysis of the fast subspace descent method for convex optimization problems. *Math. Comp.*, 89(325):2249–2282, 2020. ISSN 0025-5718. doi: 10.1090/mcom/3526. URL <https://doi.org/10.1090/mcom/3526>.
- [12] P. G. Ciarlet. *Introduction to numerical linear algebra and optimisation*. Cambridge Texts in Applied Mathematics. Cambridge University Press, Cambridge, 1989. ISBN 0-521-32788-1; 0-521-33984-7. With the assistance of Bernadette Miara and Jean-Marie Thomas, Translated from the French by A. Buttigieg.
- [13] P. G. Ciarlet. *Linear and nonlinear functional analysis with applications*. Society for Industrial and Applied Mathematics, Philadelphia, PA, 2013. ISBN 978-1-611972-58-0.
- [14] L. C. Evans. *Partial differential equations*, volume 19 of *Graduate Studies in Mathematics*. American Mathematical Society, Providence, RI, second edition, 2010. ISBN 978-0-8218-4974-3. doi: 10.1090/gsm/019. URL <https://doi.org/10.1090/gsm/019>.
- [15] W. Feng, A. J. Salgado, C. Wang, and S. M. Wise. Preconditioned steepest descent methods for some nonlinear elliptic equations involving  $p$ -Laplacian terms. *J. Comput. Phys.*, 334:45–67, 2017. ISSN 0021-9991. doi: 10.1016/j.jcp.2016.12.046. URL <https://doi.org/10.1016/j.jcp.2016.12.046>.
- [16] W. Feng, Z. Guan, J. Lowengrub, C. Wang, S. M. Wise, and Y. Chen. A uniquely solvable, energy stable numerical scheme for the functionalized Cahn-Hilliard equation and its convergence analysis. *J. Sci. Comput.*, 76(3):1938–1967, 2018. ISSN 0885-7474. doi: 10.1007/s10915-018-0690-1. URL <https://doi.org/10.1007/s10915-018-0690-1>.

- [//doi.org/10.1007/s10915-018-0690-1](https://doi.org/10.1007/s10915-018-0690-1).
- [17] X. Goudou and J. Munier. The gradient and heavy ball with friction dynamical systems: the quasiconvex case. *Math. Program.*, 116(1-2, Ser. B):173–191, 2009. ISSN 0025-5610. doi: 10.1007/s10107-007-0109-5. URL <https://doi.org/10.1007/s10107-007-0109-5>.
- [18] B. S. Jovanović and E. Süli. *Analysis of finite difference schemes*, volume 46 of *Springer Series in Computational Mathematics*. Springer, London, 2014. ISBN 978-1-4471-5459-4; 978-1-4471-5460-0. doi: 10.1007/978-1-4471-5460-0. URL <https://doi.org/10.1007/978-1-4471-5460-0>. For linear partial differential equations with generalized solutions.
- [19] M. Laborde and A. Oberman. A Lyapunov analysis for accelerated gradient methods: from deterministic to stochastic case. In Silvia Chiappa and Roberto Calandra, editors, *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108 of *Proceedings of Machine Learning Research*, pages 602–612, Online, 26–28 Aug 2020. PMLR. URL <http://proceedings.mlr.press/v108/laborde20a.html>.
- [20] H. Luo and L. Chen. From differential equation solvers to accelerated first-order methods for convex optimization, 2020. arXiv:1909.03145.
- [21] Yu. E. Nesterov. A method for solving the convex programming problem with convergence rate  $O(1/k^2)$ . *Dokl. Akad. Nauk SSSR*, 269(3):543–547, 1983. ISSN 0002-3264.
- [22] Yu. E. Nesterov. *Introductory lectures on convex optimization*, volume 87 of *Applied Optimization*. Kluwer Academic Publishers, Boston, MA, 2004. ISBN 1-4020-7553-7. doi: 10.1007/978-1-4419-8853-9. URL <https://doi.org/10.1007/978-1-4419-8853-9>. A basic course.
- [23] R. R. Phelps. *Convex functions, monotone operators and differentiability*, volume 1364 of *Lecture Notes in Mathematics*. Springer-Verlag, Berlin, second edition, 1993. ISBN 3-540-56715-1.
- [24] B. T. Poljak. Some methods of speeding up the convergence of iterative methods. *Ž. Vychisl. Mat i Mat. Fiz.*, 4:791–803, 1964. ISSN 0044-4669.
- [25] H. Schaeffer and T. Y. Hou. An accelerated method for nonlinear elliptic PDE. *J. Sci. Comput.*, 69(2):556–580, 2016. ISSN 0885-7474. doi: 10.1007/s10915-016-0215-8. URL <https://doi.org/10.1007/s10915-016-0215-8>.
- [26] J. Shen, T. Tang, and L.-L. Wang. *Spectral methods*, volume 41 of *Springer Series in Computational Mathematics*. Springer, Heidelberg, 2011. ISBN 978-3-540-71040-0. doi: 10.1007/978-3-540-71041-7. URL <https://doi.org/10.1007/978-3-540-71041-7>. Algorithms, analysis and applications.
- [27] B. Shi, S. S. Du, M. I. Jordan, and W. J. Su. Understanding the acceleration phenomenon via high-resolution differential equations, 2018. arXiv:1810.08907.
- [28] J. W. Siegel. Accelerated first-order methods: Differential equations and Lyapunov functions, 2019. arXiv:1903.05671.
- [29] W. Su, S. Boyd, and E. Candes. A differential equation for modeling Nesterov’s accelerated gradient method: Theory and insights. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 2510–2518. Curran Associates, Inc., 2014.
- [30] A. Wibisono, A. C. Wilson, and M. I. Jordan. A variational perspective on accelerated methods in optimization. *Proc. Natl. Acad. Sci. USA*, 113(47):

E7351–E7358, 2016. ISSN 0027-8424. doi: 10.1073/pnas.1614734113. URL <https://doi.org/10.1073/pnas.1614734113>.

[31] A. C. Wilson, B. Recht, and M. I. Jordan. A Lyapunov analysis of momentum methods in optimization, 2018. arXiv:611.02635.

## APPENDIX A. AN IVP AS THE LIMIT OF THE PAGD METHOD

**A.1. Derivation of the ODE.** Let us start with the same approach as in [29]. We assume, as an *ansatz*, that PAGD is a discretization of an ODE, which has a solution  $X : [0, \infty) \rightarrow \mathbb{H}$ , which we often call a *trajectory*. We also assume that  $X$  is smooth enough, e.g., twice continuously differentiable in time. For a fixed  $t \in (0, \infty)$ , the assumed smoothness on  $X$ , together with the identification  $t = \sqrt{s}k$  and Taylor’s formula in a normed vector space (e.g., [13, Theorem 7.9-1]) implies:

$$\begin{aligned} \frac{x_{k+1} - x_k}{\sqrt{s}} &= \dot{X}(t) + \frac{1}{2}\ddot{X}(t)\sqrt{s} + o(\sqrt{s}) \quad \text{as } s \rightarrow 0, \\ \frac{x_k - x_{k-1}}{\sqrt{s}} &= \dot{X}(t) - \frac{1}{2}\ddot{X}(t)\sqrt{s} + o(\sqrt{s}) \quad \text{as } s \rightarrow 0, \\ \text{(A.1)} \quad \sqrt{s}\mathcal{L}^{-1}G'(y_k) &= \sqrt{s}\mathcal{L}^{-1}G'(X(t)) + o(\sqrt{s}) \quad \text{as } s \rightarrow 0. \end{aligned}$$

The last identity follows from the continuity of  $G'$ , that of  $\mathcal{L}^{-1}$ , and (3.2), from which we can deduce  $y_k \rightarrow X(t)$  as  $s \rightarrow 0$ . Plugging (3.2) into (3.3) and dividing by  $\sqrt{s}$ , we have  $\frac{x_{k+1} - x_k}{\sqrt{s}} - \lambda \frac{x_k - x_{k-1}}{\sqrt{s}} + \sqrt{s}\mathcal{L}^{-1}G'(y_k) = 0$ . Substituting the above Taylor expansions, and then rearranging, we arrive at

$$\text{(A.2)} \quad \frac{1}{2}(1 + \lambda)\ddot{X}(t) + \frac{1 - \lambda}{\sqrt{s}}\dot{X}(t) + \mathcal{L}^{-1}G'(X(t)) + o(1) = 0 \quad \text{as } s \rightarrow 0.$$

To make this estimate consistent, interpret  $\lambda$  as a function of  $s$  and further assume that  $(1 - \lambda)/\sqrt{s} \rightarrow 2\eta$  as  $\sqrt{s} \rightarrow 0$  for some  $\eta \in (0, \infty)$ , which yields

$$\text{(A.3)} \quad \ddot{X}(t) + 2\eta\dot{X}(t) + \mathcal{L}^{-1}G'(X(t)) = 0.$$

**A.2. Derivation of the initial conditions.** The initialization  $y_0 = x_0$  and (3.3) with  $k = 0$  imply

$$\frac{x_1 - x_0}{\sqrt{s}} = \sqrt{s}\mathcal{L}^{-1}G'(x_0).$$

Take the limit  $s \rightarrow 0$  and conclude  $\dot{X}(0) = 0$  since  $G'$  and  $\dot{X}$  are assumed to be continuous. Therefore, we arrive at the desired IVP (4.1).

**Remark A.1** (momentum method). A similar procedure can be carried out far more easily for the so-called *momentum method* (MM). To see this, we recall that

$$\ddot{X}(t) \approx \frac{x_{k+1} - 2x_k + x_{k-1}}{s}, \quad \dot{X}(t) \approx \frac{x_k - x_{k-1}}{\sqrt{s}}, \quad G'(X(t)) \approx G'(x_k).$$

Then, the discrete version of the ODE (4.1) becomes

$$x_{k+1} = x_k - sG'(x_k) + (1 - 2\eta\sqrt{s})(x_k - x_{k-1}),$$

which is MM with the weight  $1 - 2\eta\sqrt{s}$ ; see [24, p. 12 (9)]. This weight is close to  $\lambda$ :

$$\lambda = \frac{1 - \eta\sqrt{s}}{1 + \eta\sqrt{s}} = 1 - \frac{2\eta\sqrt{s}}{1 + \eta\sqrt{s}} \approx 1 - 2\eta\sqrt{s}.$$

In this sense, MM seems more natural and amenable for analysis than AGD. ■

The limiting behavior of MM can also be explained by the IVP (4.1). Observe that the only essential difference between MM and PAGD is where  $G'$  is evaluated, that is,  $x_k$  and  $y_k$  respectively. And in the limit  $s \rightarrow 0$ ,  $x_k$  and  $y_k$  are not distinguishable in this setting. However, PAGD exhibits less oscillation than MM since evaluating  $G'$  at  $y_k$  serves as “foreseeing” the uphill of the objective functional, if exists, along the trajectory and “steering” to avoid unnecessary oscillating behaviors. Recently, a higher order Taylor expansion turns out to help differentiate their performances (see [27]).

#### APPENDIX B. PAGD AS A DISCRETIZATION OF THE IVP

Let us label the step size  $\sqrt{s}$ , rather than  $s$ , in order to make the setting more in line with the PAGD algorithm. Again, it is helpful to have in mind the correspondence: time  $t \longleftrightarrow k\sqrt{s}$  ( $k = 0, 1, 2, \dots$ ) and position  $X(t) \longleftrightarrow x_k$ . First, we will see  $y_k$  corresponds to a “drifted” position without the potential landscape over  $[t, t + \sqrt{s}]$ . This can be modeled by  $\ddot{X}(t) + 2\eta\dot{X}(t) = 0$ , which leads to another energy law  $\frac{1}{2}\|\dot{X}(t + \sqrt{s})\|_{\mathcal{L}}^2 = \frac{1}{2}\|\dot{X}(t)\|_{\mathcal{L}}^2 - 2\eta \int_t^{t+\sqrt{s}} \|\dot{X}(\tau)\|_{\mathcal{L}}^2 d\tau$ . Approximate the speed in the integrand by the average  $\frac{1}{2}(\|\dot{X}(t + \sqrt{s})\|_{\mathcal{L}} + \|\dot{X}(t)\|_{\mathcal{L}})$ , then after a short calculation, one obtains  $\|\dot{X}(t + \sqrt{s})\|_{\mathcal{L}} = \lambda\|\dot{X}(t)\|_{\mathcal{L}}$ . Since the dynamics takes place in a single direction, this implies  $\dot{X}(t + \sqrt{s}) = \lambda\dot{X}(t)$ . The approximations  $\dot{X}(t) \approx \frac{x_k - x_{k-1}}{\sqrt{s}}$  and  $\dot{X}(t + \sqrt{s}) \approx \frac{y_k - x_k}{\sqrt{s}}$  lead us to (3.2).

Next, we discretize the vector  $V(t)$ . Since we do not know the minimizer in practice, we remove it from the definition of  $v_k$  and discretize  $V(t) + x^* = X(t) + \frac{1}{\eta}\dot{X}(t)$ . The approximations  $X(t) \approx y_k$  and  $\dot{X}(t) \approx \frac{y_k - x_k}{\sqrt{s}}$  suggest

$$(B.1) \quad v_k = y_k + \frac{1}{\theta}(y_k - x_k),$$

which leads to the definition of  $\{v_k\}_{k \geq 1}$  (3.4) upon combining with the definition of  $\{y_k\}$ .

Finally, to get the main iterates,  $\{x_k\}_{k \geq 1}$ , we discretize (4.3) using the approximations  $\dot{V}(t) \approx \frac{v_{k+1} - v_k}{\sqrt{s}}$ ,  $\dot{X}(t) \approx \frac{y_k - x_k}{\sqrt{s}}$ , and the evaluation of  $G'$  at  $y_k$ , then it follows  $\eta \frac{v_{k+1} - v_k}{\sqrt{s}} + \eta \frac{y_k - x_k}{\sqrt{s}} + \mathcal{L}^{-1}G'(y_k) = 0$ . Plugging in (3.4) and (B.1), one obtains (3.3), the definition of  $\{x_k\}_{k \geq 1}$ .

#### APPENDIX C. LITERATURE COMPARISON

We summarize our discussion on the existing literature works, and contrast them with our contributions, in Table 2.

*Email address*, J.-H. Park: [jpark79@vols.utk.edu](mailto:jpark79@vols.utk.edu)

*Email address*, A.J. Salgado: [asalgad1@utk.edu](mailto:asalgad1@utk.edu)

*Email address*, S.M. Wise: [swise1@utk.edu](mailto:swise1@utk.edu)

DEPARTMENT OF MATHEMATICS, THE UNIVERSITY OF TENNESSEE, KNOXVILLE, TN 37996, USA

Ref.	$s = 0$	$s > 0$	$L$	$\mathfrak{B}$	$\mathcal{L}$	Numerics	Numerical Analysis
[31]	Opt.	Opt.	Glob.	×	Bregman	×	×
[3]	×	×	Loc.	×	×	✓	×
[19]	Opt.	Opt.	Glob.	×	×	×	×
[25]	×	×	Glob.	×	×	✓	×
[17]	×	×	Loc.	×	×	×	×
[9]	Sub.	×	Glob.	×	×	✓	×
[20]	Opt.	Opt.	Glob.	×	×	×	×
[28]	Opt.	Opt.	Glob.	×	×	×	×
[27]	Sub.	Sub.	Glob.	×	×	✓	×
Ours	Opt.	Opt.	Loc.	✓	✓	✓	✓

TABLE 2. A comparison of recent works from a numerical PDE point of view. All works that provide convergence rates either in the continuous (column  $s = 0$ ) or discrete (column  $s > 0$ ) case assume the global Lipschitz condition (column  $L$ ). No work addresses invariant sets (column  $\mathfrak{B}$ ), incorporates preconditioning explicitly (column  $\mathcal{L}$ ), nor it explains how concrete numerical examples fit the abstract theory (column Numerical Analysis).