

GENERALIZED MOMENTUM-BASED METHODS: A HAMILTONIAN PERSPECTIVE*

JELENA DIAKONIKOLAS[†] AND MICHAEL I. JORDAN[‡]

Abstract. We take a Hamiltonian-based perspective to generalize Nesterov’s accelerated gradient descent and Polyak’s heavy ball method to a broad class of momentum methods in the setting of (possibly) constrained minimization in Euclidean and non-Euclidean normed vector spaces. Our perspective leads to a generic and unifying nonasymptotic analysis of convergence of these methods in both the function value (in the setting of convex optimization) and in norm of the gradient (in the setting of unconstrained, possibly nonconvex, optimization). Our approach relies upon a time-varying Hamiltonian that produces generalized momentum methods as its equations of motion. The convergence analysis for these methods is intuitive and is based on the conserved quantities of the time-dependent Hamiltonian.

Key words. acceleration, momentum-based methods, stationary points, Hamiltonian dynamics

1. Introduction. Accelerated, momentum-based, methods enjoy optimal iteration complexity for the minimization of smooth convex functions over convex sets, which has led to their broad acceptance as algorithmic primitives in many applications, notably applications in machine learning. Further, decades of empirical experience suggest that momentum methods are capable of exploring multiple local minima (see, e.g., [8, 15, 60]), which gives them advantages over gradient flows. The latter have optimal worst-case complexity for convergence to stationary points, but are strongly attracted to local minima. Moreover, recent theoretical results have established that momentum methods escape saddle points faster than standard gradient descent [34, 51], providing further evidence of their value in nonconvex optimization.

The first (locally) accelerated method for smooth and strongly convex minimization¹ is from the 1960s and is due to Polyak [52]. Working with continuous-time dynamics, Polyak introduced the following (momentum-based) second-order ordinary differential equation (ODE):

$$(HBD) \quad \ddot{\mathbf{x}}_t = \alpha_1 \dot{\mathbf{x}}_t + \alpha_2 \nabla f(\mathbf{x}_t),$$

where f is the function being minimized and α_1, α_2 are constants. He also studied its two-step discretization, which can be written as:

$$(HB) \quad \mathbf{x}_{k+1} = \mathbf{x}_k - \alpha \nabla f(\mathbf{x}_k) + \beta(\mathbf{x}_k - \mathbf{x}_{k-1}),$$

where α, β are constants. In particular, under a suitable choice of α, β , [52] showed that, when initialized “sufficiently close” to the optimal solution \mathbf{x}^* , the method converges at rate $(\frac{1-\sqrt{\kappa}}{1+\sqrt{\kappa}})^k$ where κ denotes the condition number of the objective function f . This convergence rate was later proved to be optimal and globally achievable [46].

*
Funding: Research presented in this paper was partially supported by the NSF grant #CCF-1740855, by the Mathematical Data Science program of the Office of Naval Research under grant number N00014-18-1-2764, and by the Office of the Vice Chancellor for Research and Graduate Education at the University of Wisconsin–Madison with funding from the Wisconsin Alumni Research Foundation. Part of this work was done while the authors were visiting Simons Institute for the Theory of Computing.

[†]Department of Computer Sciences, UW-Madison, Madison, WI (jelena@cs.wisc.edu).

[‡]Department of Statistics and EECS, UC Berkeley, Berkeley, CA (jordan@cs.berkeley.edu).

¹In the special case of quadratic objectives, Polyak’s method attains the globally-accelerated (iteration-complexity optimal) convergence rate.

For the setting of smooth minimization, Nesterov [47] discovered a method that has an optimal convergence rate of $1/k^2$. Nesterov [47] also showed that when this method is applied to strongly convex functions with a known strong convexity parameter and coupled with scheduled restart, it leads to the optimal rate for the class of smooth and strongly convex minimization problems. In later work (see, e.g., the textbook [49] and references therein), Nesterov introduced a separate, more direct method for the minimization of smooth and strongly convex functions that enjoys the same optimal rate as (HB) but for which the optimal rate holds globally while the algorithm has the same continuous-time limit (HBD).

A flurry of research has followed these seminal papers on momentum-based methods [2, 4, 7, 8, 13, 16, 24–26, 33, 38, 41, 42, 44, 45, 53, 55–59, 61], with many of these works [7, 13, 24, 25, 27, 38, 44, 45, 53, 55, 56, 58, 59, 61] seeking to interpret Nesterov acceleration as a discretization of a continuous-time dynamical system. Further, some of these works led to physical interpretations of Polyak’s [4, 8, 27] and Nesterov’s [13, 58] methods in the Lagrangian and Hamiltonian formalism, and the physical interpretation of momentum has even led to new algorithms (e.g., [13, 27, 43]). For the setting of nonconvex optimization, however, and, more broadly, convergence to points with small norm of the gradient, the continuous-time perspective has been much less explored [6, 8, 34, 55].

In this paper, we take a Hamiltonian-based perspective to derive a broad class of momentum methods that yield Nesterov’s and Polyak’s methods as special cases. As a specific example, a class of methods obtained from the introduced Hamiltonian and parametrized by $\lambda \in [0, 1]$ interpolates between Nesterov’s method for smooth minimization [47] (when $\lambda = 1$) and a generalization of the heavy ball method [52] (when $\lambda = 0$). We show that because the methods are obtained as the equations of motion of this Hamiltonian, we can deduce invariants (conserved quantities of the Hamiltonian) that can be used to argue about convergence in function value (for convex optimization) and convergence to stationary points (for possibly nonconvex optimization). The techniques are general and lead to results in general normed vector spaces.

We note that non-Euclidean normed spaces are frequently encountered in applications, particularly those arising in sparsity-oriented machine learning, where it is natural to use the ℓ_1 norm, and in fast algorithms for network flow problems, where the natural norm is ℓ_∞ [35, 40, 54]. Here, by “natural,” we mean that the use of an alternative ℓ_p norm would incur a polynomial dependence on the dimension in the method’s iteration complexity. This is because both the initial distance to the optimum or the diameter of the set *and* the Lipschitz constant of the objective function or its gradient that determine the iteration complexity are defined w.r.t. the norm that is adopted for the underlying vector space (see, e.g., [23] for a similar discussion).

In Section 3 we provide analysis of the methods in terms of convergence in function value in the setting of (possibly) constrained convex optimization in possibly non-Euclidean normed vector spaces. We show that the entire class of methods parametrized by λ (as mentioned above) converges at rate $1/k^2$ as long as λ is bounded away from zero. When $\lambda = 0$, the convergence slows down to $1/k$. This agrees with previously obtained results for the heavy-ball method in the setting of smooth (non-strongly convex) minimization [29] (our case $\lambda = 0$). As a byproduct of this approach, we obtain a generalization of the heavy-ball method to constrained convex optimization in non-Euclidean vector spaces and show that it converges at rate $1/k$. Such a result was previously known only for unconstrained convex optimization in Euclidean spaces [29].

In terms of the convergence to stationary points, we consider the unconstrained case and focus on finding points with small norm of the gradient, in either Euclidean or non-Euclidean spaces. We show that when f is convex, any method from the class satisfies $\min_{0 \leq i \leq k} \|\nabla f(\mathbf{x}_k)\|_*^2 = O(\frac{L(f(\mathbf{x}_0 - \mathbf{x}^*)}{k})$. Note that any of these methods, when run for $k/2$ iterations after running Nesterov’s method for $k/2$ iterations, satisfies $\min_{k/2 \leq i \leq k} \|\nabla f(\mathbf{x}_k)\|_*^2 = O(\frac{L\|\mathbf{x}_0 - \mathbf{x}^*\|^2}{k^3})$. While this is suboptimal for the case of convex functions—the optimal rate is $\min_{0 \leq i \leq k} \|\nabla f(\mathbf{x}_k)\|_*^2 = O(\frac{L(f(\mathbf{x}_0 - \mathbf{x}^*)}{k^2})$ [18, 37, 48] and it is achieved by [37]—we conjecture that it is tight. In particular, [36] demonstrated that the convergence of the form $\min_{k/2 \leq i \leq k} \|\nabla f(\mathbf{x}_k)\|_*^2 = O(\frac{L\|\mathbf{x}_0 - \mathbf{x}^*\|^2}{k^3})$ is tight for Nesterov’s method.

For the case of nonconvex functions, we show that methods that are instantiations of the heavy-ball method converge at the optimal [18] rate of $\min_{0 \leq i \leq k} \|\nabla f(\mathbf{x}_k)\|_*^2 = O(\frac{L(f(\mathbf{x}_0 - \mathbf{x}^*)}{k})$. While a similar result exists for the case of Euclidean spaces [30], we are not aware of any other results for the more general normed vector spaces in which the method does not lose its favorable properties. For example, it is possible to establish similar rates for modifications of Nesterov’s method that turn it into a descent method that makes at least as much progress as gradient descent, as in, e.g., [31, 48, 50]. In this case, the analysis of convergence in norm of the gradient boils down to the analysis of (standard) gradient descent. Unfortunately, because such methods monotonically decrease the function value, they lose the property of utilizing the momentum to escape shallow local minima. Note that, as mentioned earlier, global exploration of local minima is one of the primary reasons for considering momentum-based methods in nonconvex optimization [8].

Finally, we note that the notions of convergence considered in this paper are only the standard notions for first-order smooth (not necessarily strongly convex) optimization: (i) convergence to (neighborhoods of) points with small optimality gap, and (ii) convergence to (neighborhoods of) points with small gradient norm. Our analysis only guarantees that such points can be output by the algorithm and does not deal with the convergence of the trajectories, which is generally a more challenging problem (see, e.g., [17]). We note that weak convergence results of the trajectories have been established by prior work for variants of inertial proximal algorithms that generalize the Nesterov accelerated method to the setting of composite objectives (smooth plus nonsmooth) [10, 20].

1.1. Related Work. In addition to the work already mentioned above, we provide a few more remarks regarding related work. First, for convex minimization (in function value), there are several approaches that apply to constrained minimization and general normed vector spaces [13, 24, 25, 38, 58], with a subset of them being directly motivated by Lagrangian [58] and Hamiltonian [13] mechanics. The latter make use of special Lyapunov functions (which are sometimes also called the energy functions) to characterize convergence rates, and their applicability to convergence to stationary points is unclear. In contrast, our work is not based on separately constructed Lyapunov functions; rather, our analysis of convergence rates stems from the analysis of conserved quantities of the Hamiltonian. These conserved quantities can also be viewed as Lyapunov functions. The main difference compared to the prior work is that we do not need to “guess” these functions—they are directly derived from the same Hamiltonian whose equations of motion are the momentum dynamics we consider.

A significant body of recent work in nonconvex optimization has focused on

convergence to approximate local minima, with many of the methods having (near-)optimal iteration complexities (see, e.g., [1, 18, 34]). The only work that we are aware of that has used a Hamiltonian perspective on convergence to stationary points in the nonconvex setting is [34]. However, the connection to Hamiltonian systems in [34] is limited—it essentially relies on showing that (HBD) dissipates energy of the form $f(\mathbf{x}_t) + \frac{1}{2}\|\dot{\mathbf{x}}_t\|_2^2$, and is specialized to a Euclidean setting. In fact, the main contribution of [34] is in providing a near-optimal method for convergence to approximate local minima and not in providing a Hamiltonian perspective on nonconvex optimization. Also worth mentioning are inertial approaches to nonconvex optimization [11, 19], based on Kurdyka-Łojasiewicz (KL) property of the objective, as introduced in [14]. The KL property ensures that the gradients of the objective do not vanish too quickly around critical points. Such a property is not assumed to hold for the problems considered in this work.

1.2. Preliminaries. While there is nothing in our approach that prevents one from working in general (possibly infinitely-dimensional) Banach spaces, we will limit our attention to finite-dimensional real vector spaces, as they suffice for the motivating applications discussed in the introduction. The primal, n -dimensional real vector space is denoted by E . The space E is normed, endowed with a norm $\|\cdot\|$. Its dual space, consisting of all linear functions on E , is denoted by E^* . For $\mathbf{z} \in E^*$ and $\mathbf{x} \in E$, we denote by $\langle \mathbf{z}, \mathbf{x} \rangle$ the value of \mathbf{z} at \mathbf{x} . The dual norm (associated with space E^*) is defined in the standard way as $\|\mathbf{z}\|_* = \max_{\mathbf{x} \in E} \frac{\langle \mathbf{z}, \mathbf{x} \rangle}{\|\mathbf{x}\|}$. For Euclidean spaces, $\langle \cdot, \cdot \rangle$ is the standard inner product and $\|\cdot\| = \|\cdot\|_* = \|\cdot\|_2$.

We assume that $f : \mathcal{X} \rightarrow \mathbb{R}$ is a (possibly nonconvex) continuously-differentiable function, and $\mathcal{X} \subseteq E$ is closed and convex. $\mathbf{x}^* \in \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$ denotes any fixed minimizer of f . To avoid vacuous statements, we assume that $f(\mathbf{x}^*) > -\infty$.

For all the methods, $\mathbf{x}_t \in \mathcal{X}$ will be the running solution, and $\mathbf{z}_t \in \mathcal{Z} = \mathbf{z}_0 + \operatorname{Lin}\{\nabla f(\mathbf{x}) : \mathbf{x} \in \mathcal{X}\}$ will be the sum of $\mathbf{z}_0 \in E^*$ and some linear combination of the gradients $\nabla f(\mathbf{x}_\tau)$ for $\tau \in [0, t]$. In the Hamiltonian formalism, \mathbf{z}_t will correspond to the conjugate momenta, $f(\mathbf{x})$ will correspond to the potential energy, and $\psi^*(\mathbf{z})$ will correspond to the kinetic energy, where $\psi^* : \mathcal{Z} \rightarrow \mathbb{R}$ is a convex conjugate (defined below) of some strongly convex function $\psi : \mathcal{X} \rightarrow \mathbb{R}$ (e.g., $\psi^*(\mathbf{z}) = \frac{1}{2}\|\mathbf{z}\|_2^2$ if $\|\cdot\| = \|\cdot\|_2$ and $\psi(\mathbf{x}) = \frac{1}{2}\|\mathbf{x}\|_2^2$).

We now outline some useful definitions and facts that are used in the paper.

To carry out the analysis of the cases of convex and nonconvex objectives in a unified way, we use the following notion of weak convexity, similar to e.g., [1, 22].

DEFINITION 1.1. *We say that a continuously-differentiable function $f : \mathcal{X} \rightarrow \mathbb{R}$ is ϵ_H -weakly convex for some $\epsilon_H \in \mathbb{R}_+$ if*

$$(\forall \mathbf{x}, \mathbf{y} \in \mathcal{X}) : \quad f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle - \frac{\epsilon_H}{2} \|\mathbf{y} - \mathbf{x}\|^2.$$

We will mainly be concerned with cases $\epsilon_H = 0$ and $\epsilon_H = L$. Observe that a 0-weakly convex function is convex, by the standard first-order definition of convexity for continuously-differentiable functions that can be stated as

$$(\forall \mathbf{x}, \mathbf{y} \in \mathcal{X}) : \quad f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle.$$

DEFINITION 1.2. *A continuously-differentiable function $f : \mathcal{X} \rightarrow \mathbb{R}$ is L -smooth for $L \in \mathbb{R}_+$, if $\forall \mathbf{x}, \mathbf{y} \in \mathcal{X}$, $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_* \leq L\|\mathbf{x} - \mathbf{y}\|$.*

²Here we use the notation ψ^* to emphasize that \mathbf{x} and \mathbf{z} do not, in general, belong to the same vector space.

Recall that L -smoothness of a function implies that

$$(\forall \mathbf{x}, \mathbf{y} \in \mathcal{X}) : f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|^2.$$

It is not hard to show that an L -smooth function is also L -weakly convex. We will be assuming throughout that there exists $L < \infty$ such that f is L -smooth.

DEFINITION 1.3. *A continuously-differentiable function $f : \mathcal{X} \rightarrow \mathbb{R}$ is μ -strongly convex for $\mu \in \mathbb{R}_+$, if $\forall \mathbf{x}, \mathbf{y} \in \mathcal{X}$, $f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{\mu}{2} \|\mathbf{y} - \mathbf{x}\|^2$.*

DEFINITION 1.4. *The convex conjugate of $\psi : \mathcal{X} \rightarrow \mathbb{R}$ is defined as $\psi^*(\mathbf{z}) = \sup_{\mathbf{x} \in \mathcal{X}} \{\langle \mathbf{z}, \mathbf{x} \rangle - \psi(\mathbf{x})\}$.³*

The following standard fact is a corollary of Danskin's Theorem (see, e.g., [12]).

FACT 1.5. *Let $\psi : \mathcal{X} \rightarrow \mathbb{R}$ be a strongly convex function.⁴ Then ψ^* is continuously differentiable and $\nabla \psi^*(\mathbf{z}) = \operatorname{argmax}_{\mathbf{x} \in \mathcal{X}} \{\langle \mathbf{z}, \mathbf{x} \rangle - \psi(\mathbf{x})\}$.*

Another useful property of convex conjugacy is the duality between smoothness and strong convexity, which can be seen as a strengthening of Fact 1.5.

FACT 1.6. *Let $\psi : \mathcal{X} \rightarrow \mathbb{R}$ be a μ -strongly convex function. Then ψ^* is $\frac{1}{\mu}$ -smooth.*

For a convex function $\psi : \mathcal{X} \rightarrow \mathbb{R}$ that is continuously differentiable on \mathcal{X} , Bregman divergence is defined in a usual way as $D_\psi(\mathbf{y}, \mathbf{x}) = \psi(\mathbf{y}) - \psi(\mathbf{x}) - \langle \nabla \psi(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle$, where $\mathbf{x}, \mathbf{y} \in \mathcal{X}$. Some useful properties of Bregman divergence are stated below.

FACT 1.7. *(Properties of Bregman Divergence.) Let $\psi : \mathcal{X} \rightarrow \mathbb{R}$ be convex and continuously differentiable on \mathcal{X} . Then:*

(i) *For any $\mathbf{u}, \mathbf{v}, \mathbf{w} \in \mathcal{X}$, the following three-point identity holds:*

$$D_\psi(\mathbf{u}, \mathbf{v}) = D_\psi(\mathbf{w}, \mathbf{v}) + \langle \nabla \psi(\mathbf{w}) - \nabla \psi(\mathbf{v}), \mathbf{u} - \mathbf{w} \rangle + D_\psi(\mathbf{u}, \mathbf{w}).$$

(ii) *Let ψ be μ -strongly convex w.r.t. $\|\cdot\|$. Then, $\forall \mathbf{z}, \mathbf{z}'$,*

$$D_{\psi^*}(\mathbf{z}, \mathbf{z}') \geq \frac{\mu}{2} \|\nabla \psi^*(\mathbf{z}) - \nabla \psi^*(\mathbf{z}')\|^2.$$

2. Continuous-Time Methods and their Convergence Analysis. We start this section by providing a brief overview of the well-studied inertial approach that has been widely used in continuous optimization to study the behavior of momentum-based methods. We then provide basic definitions and illustrate the main ideas that underlie our analysis. In doing so, we compare our approach to the inertial one and highlight the connections between them.

2.1. Inertial Approach in Optimization. One of the earliest examples of the inertial approach is the heavy-ball method introduced by Polyak [52] and defined in (HBD), (HB). The method (HBD) applies to *unconstrained Euclidean settings*, and its generalizations to Hilbert spaces and constrained setups has been studied in [4, 5, 8, 9, 39].

³Our definition of a convex conjugate slightly differs from the standard definition, where the convex conjugate is defined w.r.t. the entire vector space E , i.e., $\psi^*(\mathbf{z}) = \sup_{\mathbf{x} \in E} \{\langle \mathbf{z}, \mathbf{x} \rangle - \psi(\mathbf{x})\}$. Thus, our definition corresponds to taking the convex conjugate of $\psi + I_{\mathcal{X}}$ in the standard definition, where $I_{\mathcal{X}}$ is the indicator function of the set \mathcal{X} (equal to zero within the set, and to infinity outside it). This choice incurs no loss of generality and simplifies the notation throughout the paper.

⁴Alternatively, it suffices that ψ is strictly convex and either \mathcal{X} is compact or the gradient of ψ blows up at the boundary of \mathcal{X} .

The basic variant of the second-order differential equation considered as part of the (damped) inertial approach can be stated as:

$$(2.1) \quad \ddot{\mathbf{x}}_t + \alpha_1(t)\dot{\mathbf{x}}_t + \alpha_2(t)\nabla f(\mathbf{x}_t) = 0,$$

where $\ddot{\mathbf{x}}_t$ is the inertial term, $\alpha_1(t)\dot{\mathbf{x}}_t$ is the friction term that dissipates energy, and $\alpha_2(t)\nabla f(\mathbf{x}_t)$ is the term that corresponds to potential forces that drive the motion. When $\alpha_1(t) = \alpha_1$ and $\beta_1(t) = \beta_1$ are constants, Eq. (2.1) reduces to Polyak's heavy-ball method as stated in (HBD). When $\alpha_1(t) = \frac{\alpha}{t}$ for $\alpha > 0$ and $\alpha_2(t) = 1$, Eq. (2.1) reduces to the differential equation introduced by Su et al. [56] for studying Nesterov acceleration and its generalizations. In particular, [56] showed that Nesterov acceleration corresponds to the case $\alpha = 3$, while choices of $\alpha \geq 3$ can generally lead to accelerated $1/t^2$ rates for smooth minimization. The subcritical case $\alpha < 3$ and the degradation of the rates from $1/t^2$ to $1/t^{\frac{2}{3}\alpha}$ was characterized by Attouch et al. [7].

While the inertial approach is powerful and can lead to different qualitative results (e.g., it can be used to study convergence of the trajectories that is not considered in this work), its main limitation is that technical obstacles are encountered once generalizations to non-Euclidean and constrained setups are considered.

In its basic form, the dynamics from Eq. (2.1) cannot be directly applied to non-Euclidean setups, as the points \mathbf{x}_t and the gradients $\nabla f(\mathbf{x}_t)$ do not in general belong to the same vector space. This issue was resolved in the work by Wibisono et al. [58] using the concept of Bregman Lagrangian and the dynamics that are generated from it. However, it is not immediately clear how to generalize the results from [58] to constrained setups.

A specific challenge that is encountered when imposing constraints directly in the second order ODE from Eq. (2.1) is that the direct enforcement of the constraints typically involves the use of maps that are not differentiable and which turn differential equation problems into differential inclusion problems (see, e.g., [5, 9]). More importantly, the introduction of the constraints within this inertial approach is well-known to cause non-elastic shocks, due to the discontinuities in the velocity $\dot{\mathbf{x}}$ encountered at the boundary of the feasible region [5, 9]. This problem is overcome in our work (and similarly in prior work [25, 38]) through the use of constraint regularization.

Finally, the convergence analysis in the inertial approach is carried out using Lyapunov functions that have the interpretation of total energy of the system. They are typically constructed as the sum of the (possibly) scaled optimality gap $f(\mathbf{x}) - f(\mathbf{x}^*)$ and another function that has the interpretation of kinetic energy. While the energy interpretation is intuitive for these Lyapunov functions, it is not always *a priori* clear how to construct them. In our approach, it is the Hamiltonian itself that has the interpretation of the total energy, and both our Lyapunov functions and the continuous-time dynamics are obtained from this same Hamiltonian.

2.2. Background on Hamiltonian Mechanics. In Hamiltonian mechanics, the motion of a particle, or, more broadly, a physical system, is described by canonical coordinates (\mathbf{x}, \mathbf{z}) ⁵, where \mathbf{x} corresponds to the generalized coordinates (position of a particle) and \mathbf{z} are their conjugate momenta. Given a Hamiltonian $\mathcal{H}(\mathbf{x}, \mathbf{z}, t)$, which is typically interpreted as the total energy of the system, the time evolution of the physical system is described by Hamilton's equations:

$$(2.2) \quad \frac{d\mathbf{x}}{dt} = \nabla_{\mathbf{z}} \mathcal{H}(\mathbf{x}, \mathbf{z}, t), \quad \frac{d\mathbf{z}}{dt} = -\nabla_{\mathbf{x}} \mathcal{H}(\mathbf{x}, \mathbf{z}, t).$$

⁵In physics, it is standard to use the notation (\mathbf{p}, \mathbf{q}) instead to denote the generalized coordinates and the conjugate momenta. We use (\mathbf{x}, \mathbf{z}) for consistency with the optimization literature.

An immediate implication of Eq. (2.2) that will be used in deriving the conserved quantities used in the convergence analysis is the following key relationship:

$$(2.3) \quad \begin{aligned} \frac{d\mathcal{H}(\mathbf{x}, \mathbf{z}, t)}{dt} &= \nabla_{\mathbf{x}}\mathcal{H}(\mathbf{x}, \mathbf{z}, t) \cdot \frac{d\mathbf{x}}{dt} + \nabla_{\mathbf{z}}\mathcal{H}(\mathbf{x}, \mathbf{z}, t) \cdot \frac{d\mathbf{z}}{dt} + \frac{\partial\mathcal{H}(\mathbf{x}, \mathbf{z}, t)}{\partial t} \\ &= \frac{\partial\mathcal{H}(\mathbf{x}, \mathbf{z}, t)}{\partial t}. \end{aligned}$$

Perhaps the simplest Hamiltonian that one can formulate is the following separable function: $\mathcal{H}(\mathbf{x}, \mathbf{z}) = f(\mathbf{x}) + \psi^*(\mathbf{z})$. Here, $f(\mathbf{x})$ can be viewed as the potential energy of a particle at position \mathbf{x} , while $\psi^*(\mathbf{z})$ is its kinetic energy. The corresponding continuous-time dynamics is:

$$(HD) \quad \begin{aligned} \dot{\mathbf{x}}_t &= \nabla_{\mathbf{z}_t}\mathcal{H}(\mathbf{x}_t, \mathbf{z}_t) = \nabla\psi^*(\mathbf{z}_t), \\ \dot{\mathbf{z}}_t &= -\nabla_{\mathbf{x}_t}\mathcal{H}(\mathbf{x}_t, \mathbf{z}_t) = -\nabla f(\mathbf{x}_t). \end{aligned}$$

This dynamics is meaningful only in the unconstrained regime, since otherwise we cannot guarantee that $\mathbf{x}_t \in \mathcal{X}$. Hence, we assume here that $\mathcal{X} = \mathbb{R}^n$. As $\mathcal{H}(\mathbf{x}, \mathbf{z})$ does not explicitly depend on time, we have $\frac{d}{dt}\mathcal{H}(\mathbf{x}, \mathbf{z}) = 0$. Equivalently, $\mathcal{H}(\mathbf{x}, \mathbf{z})$ is conserved with time. An immediate implication is that the norm of the averaged gradient decays as $1/t$, as stated in the following proposition.

PROPOSITION 2.1. *Let $\mathbf{x}_t, \mathbf{z}_t$ evolve according to (HD), for $\mathbf{z}_0 = \mathbf{0}$ and arbitrary (but fixed) $\mathbf{x}_0 \in \mathbb{R}^n$, and let $\psi^*(\mathbf{0}) = 0$. If ψ^* is μ -strongly convex, then, $\forall t \geq 0$:*

$$\left\| \frac{1}{t} \int_0^t \nabla f(\mathbf{x}_\tau) d\tau \right\|_* \leq \frac{\sqrt{2(f(\mathbf{x}_0) - f(\mathbf{x}^*))}/\mu}{t}.$$

Proof. Since the Hamiltonian is conserved, $\psi^*(\mathbf{z}_t) = f(\mathbf{x}_0) - f(\mathbf{x}_t) + \psi^*(\mathbf{z}_0) = f(\mathbf{x}_0) - f(\mathbf{x}_t)$. Recall that \mathbf{x}^* minimizes f . As $f(\mathbf{x}_t) \geq f(\mathbf{x}^*)$, it follows that $\psi^*(\mathbf{z}_t) \leq f(\mathbf{x}_0) - f(\mathbf{x}^*)$. By the μ -strong convexity of ψ^* , $\frac{\mu}{2}\|\mathbf{z}_t\|_* \leq f(\mathbf{x}_0) - f(\mathbf{x}^*)$. Finally, integrating the second equation from (HD), we have that $\mathbf{z}_t = \mathbf{z}_0 - \int_0^t \nabla f(\mathbf{x}_\tau) d\tau = -\int_0^t \nabla f(\mathbf{x}_\tau) d\tau$, which, combined with the last inequality (after dividing both sides by $t^2\mu/2$ and taking the square root of both sides), gives the claimed bound. \square

While Proposition 2.1 shows that the *average* of the gradients converges in norm $\|\cdot\|_*$ to zero at a sublinear rate, it does not guarantee that the dynamics converges to or even visits any stationary point of f . Indeed, since the energy (equal to $\mathcal{H}(\mathbf{x}_t, \mathbf{z}_t)$) is conserved with time, the dynamics is well-known to be non-convergent. Hence, the fact that the average gradient converges in the dual norm only implies that the path of the dynamics consists of cycle-like segments over which the gradients cancel out.

2.3. Generalized Momentum Dynamics. The standard Hamiltonian dynamics from the previous subsection is overly aggressive as a function of the history of the gradients (i.e., the momentum \mathbf{z}_t). As a consequence of energy conservation, the energy is exchanged between the potential and kinetic energy, which makes the dynamics exhibit non-convergent behavior. For the dynamics to be attracted to stationary points, it needs to be dampened. The most common approach is to introduce friction into the equations of motion, which leads to second order ODEs such as those described by Eq. (2.1). Here, we take an alternative approach that directly modifies the Hamiltonian. As we will see, this approach allows us to consider constrained optimization problems over general normed spaces, unlike most of the friction-based approaches.

For example, the Nesterov acceleration method for smooth constrained minimization in general normed spaces can be expressed in continuous time as follows [25, 38]:

$$(AD) \quad \begin{aligned} \dot{\mathbf{x}}_t &= \frac{\dot{\alpha}_t(\nabla\psi^*(\mathbf{z}_t) - \mathbf{x}_t)}{\alpha_t}, \\ \dot{\mathbf{z}}_t &= -\dot{\alpha}_t\nabla f(\mathbf{x}_t). \end{aligned}$$

Using Eq. (2.2), this dynamics can be shown to correspond to the following Hamiltonian: ⁶

$$(2.4) \quad \mathcal{H}(\bar{\mathbf{x}}, \mathbf{z}, \tau) = \tau f(\bar{\mathbf{x}}/\tau) + \psi^*(\mathbf{z}), ⁷$$

after a suitable time reparametrization $\tau = \alpha_t$ (see also, e.g., [58]), where α_t is a strictly increasing function of time t and $\bar{\mathbf{x}} = \tau\mathbf{x}$. To see that the dynamics from (AD) corresponds to the equations of motion of the Hamiltonian (2.4), observe that

$$\frac{d}{dt}\bar{\mathbf{x}}_t = \frac{d\tau}{dt}\frac{d}{d\tau}\bar{\mathbf{x}}_t = \dot{\alpha}_t\nabla_{\mathbf{z}}\mathcal{H}(\bar{\mathbf{x}}_t, \mathbf{z}_t, \tau) = \dot{\alpha}_t\nabla\psi^*(\mathbf{z}_t),$$

which, using $\bar{\mathbf{x}}_t = \alpha_t\mathbf{x}_t$, is exactly the first equation from (AD). Similarly, the second equation of motion for the Hamiltonian from Eq. (2.4) $\dot{\mathbf{z}}_t = -\dot{\alpha}_t\nabla_{\mathbf{x}}\mathcal{H}(\bar{\mathbf{x}}_t, \mathbf{z}_t, \alpha_t) = -\dot{\alpha}_t\nabla f(\bar{\mathbf{x}}_t/\alpha_t)$ is exactly the second equation from (AD), as $\bar{\mathbf{x}}_t = \alpha_t\mathbf{x}_t$.

We now show that it is possible to generalize the Hamiltonian from Eq. (2.4) and its resulting equations of motion to capture a much broader class of convergent momentum-based methods that contains a generalization of Polyak's heavy ball method [52]. In particular, consider:

$$(2.5) \quad \mathcal{H}_M(\bar{\mathbf{x}}, \mathbf{z}, \tau) = h(\tau)f(\bar{\mathbf{x}}/\tau) + \psi^*(\mathbf{z}),$$

where, as before, $\bar{\mathbf{x}} = \tau\mathbf{x}$, $h(\tau)$ is a positive function of τ , and we reparametrize time as $\tau = \alpha_t$. We will mainly be considering the case $h(\tau) = \tau^\lambda$ for $\lambda \in [0, 2]$ (see Sections 3 and 4). The resulting equations of motion (after time reparametrization) of this Hamiltonian are:

$$(MoD) \quad \begin{aligned} \dot{\mathbf{x}}_t &= \frac{\dot{\alpha}_t(\nabla\psi^*(\mathbf{z}_t) - \mathbf{x}_t)}{\alpha_t}, \\ \dot{\mathbf{z}}_t &= -h(\alpha_t)\frac{\dot{\alpha}_t}{\alpha_t}\nabla f(\mathbf{x}_t). \end{aligned}$$

Observe that the assumption that the vector space E is non-Euclidean, i.e., that \mathbf{x}_t and \mathbf{z}_t belong to different (mutually dual) spaces, is seamlessly handled through the use of the map $\nabla\psi^* : E^* \rightarrow \mathcal{X} \subseteq E$. Due to Fact 1.6, $\nabla\psi^*$ is a Lipschitz-continuous map whenever ψ is strongly convex, which is a basic assumption we make throughout. Further, it is not hard to see that if $\mathbf{x}_0 \in \mathcal{X}$, then $\mathbf{x}_t \in \mathcal{X}$, $\forall t > 0$. To see this, observe that, after suitably rearranging the terms, we can equivalently write the first equation from (MoD) as

$$\frac{d(\alpha_t\mathbf{x}_t)}{dt} = \dot{\alpha}_t\nabla\psi^*(\mathbf{z}_t).$$

⁶This Hamiltonian was obtained in the discussions between J. Diakonikolas and Lorenzo Orecchia.

⁷To avoid division by zero, the dynamics can be started from time $\tau = 1$. Alternatively, one can replace τ by $\tau + 1$ and start the dynamics from $\tau = 0$. In the sequel, we will always assume that α_0 is bounded away from zero.

Integrating both sides of the last equation, it follows that

$$\mathbf{x}_t = \frac{\alpha_0}{\alpha_t} \mathbf{x}_0 + \frac{1}{\alpha_t} \int_0^t \nabla \psi^*(\mathbf{z}_\tau) \dot{\alpha}_\tau d\tau.$$

Thus, as $\mathbf{x}_0 \in \mathcal{X}$ and $\nabla \psi^*(\mathbf{z}) \in \mathcal{X}$, $\forall \mathbf{z} \in E^*$ (see Fact 1.5), it follows that $\mathbf{x}_t \in \mathcal{X}$. Additionally, if the dynamics is started from the relative interior of the set \mathcal{X} , then \mathbf{x}_t always remains in the relative interior. As $\nabla \psi^*$ is a Lipschitz continuous map, unlike in the inertial approach, no non-elastic shocks can arise at the boundary. We note that the use of strongly convex functions ψ resulting in smooth functions ψ^* can also be viewed as constraint regularization.

Clearly, the Hamiltonian \mathcal{H}_M and its equations of motion (MoD) generalize the accelerated dynamics: (2.4) and (AD) correspond to the case $h(\tau) = \tau$. It is possible to show that the class of methods captured by the equations of motion of (2.5) also contains a generalization of Polyak's heavy ball method, as shown in the following proposition.

PROPOSITION 2.2. *Polyak's heavy ball method is equivalent to (MoD) when $\mathcal{X} = \mathbb{R}^n$, $\|\cdot\| = \|\cdot\|_2$, $\psi^*(\mathbf{z}) = \frac{1}{2\mu} \|\mathbf{z}\|_2^2$, $h(\tau) = \tau^0 = 1$, and $\frac{\dot{\alpha}_t}{\alpha_t} = \eta > 0$.*

Proof. Under the assumptions of the proposition, $\dot{\mathbf{x}}_t = \eta(\frac{1}{\mu} \mathbf{z}_t - \mathbf{x}_t)$ and $\dot{\mathbf{z}}_t = -\eta \nabla f(\mathbf{x}_t)$. Hence, $\ddot{\mathbf{x}}_t = -\eta \dot{\mathbf{x}}_t - \frac{\eta^2}{\mu} \nabla f(\mathbf{x}_t)$. For suitable choices of η , μ , this is equivalent to (HBD) from [52]. \square

More generally, we can relate (MoD) to the inertial dynamics from Eq. (2.1) in the special case of *unconstrained Euclidean setups* as follows. Let $\mathcal{X} \equiv E$, $\|\cdot\| = \|\cdot\|_2$ and $\psi^*(\mathbf{z}) = \frac{1}{2\mu} \|\mathbf{z}\|_2^2$ for some $\mu > 0$. Then $\nabla \psi^*(\mathbf{z}) = \frac{1}{\mu} \mathbf{z}$, and we have:

$$\frac{d(\alpha_t \mathbf{x})}{dt} = \dot{\alpha}_t \mathbf{x}_t + \alpha_t \dot{\mathbf{x}}_t = \dot{\alpha}_t \mathbf{z}_t / \mu.$$

Dividing both sides by $\dot{\alpha}_t > 0$ and differentiating w.r.t. t , we have, using the second equation in (MoD):

$$\left(1 + \frac{\dot{\alpha}_t^2 - \alpha_t \ddot{\alpha}_t}{\dot{\alpha}_t^2}\right) \dot{\mathbf{x}}_t + \frac{\alpha_t}{\dot{\alpha}_t} \ddot{\mathbf{x}}_t = -\frac{1}{\mu} h(\alpha_t) \frac{\dot{\alpha}_t}{\alpha_t} \nabla f(\mathbf{x}_t).$$

Rearranging the last inequality:

$$(2.6) \quad \ddot{\mathbf{x}}_t + \frac{\dot{\alpha}_t}{\alpha_t} \left(1 + \frac{\dot{\alpha}_t^2 - \alpha_t \ddot{\alpha}_t}{\dot{\alpha}_t^2}\right) \dot{\mathbf{x}}_t + \frac{1}{\mu} h(\alpha_t) \left(\frac{\dot{\alpha}_t}{\alpha_t}\right)^2 \nabla f(\mathbf{x}_t) = 0,$$

which is precisely the inertial ODE from Eq. (2.1) with $\alpha_1(t) = \frac{\dot{\alpha}_t}{\alpha_t} \left(1 + \frac{\dot{\alpha}_t - \ddot{\alpha}_t}{\dot{\alpha}_t^2}\right)$ and $\alpha_2(t) = \frac{1}{\mu} h(\alpha_t) \left(\frac{\dot{\alpha}_t}{\alpha_t}\right)^2$. In particular, for $\alpha_t = t^p$, $p > 0$, $\mu = p$, and $h(\alpha_t) = \alpha_t^{2/p}/p$, Eq. (2.6) reduces to

$$\ddot{\mathbf{x}}_t + \frac{p+1}{t} \dot{\mathbf{x}}_t + \nabla f(\mathbf{x}_t) = 0,$$

which is the damped inertial ODE studied by Su et al. [56].

The main usefulness of Hamiltonian \mathcal{H}_M is that it can be used to argue about convergence in both the function value (for convex optimization problems) and convergence to stationary points (for potentially nonconvex problems). In the following lemma we exhibit two different conserved quantities (or invariants) of (2.5) that can be used towards this goal.

LEMMA 2.3. Let $\mathbf{x}_t, \mathbf{z}_t$ evolve according to (MoD) for an arbitrary initial point $\mathbf{x}_0 = \nabla\psi^*(\mathbf{z}_0) \in \mathcal{X}$ and some differentiable $\psi^*(\cdot)$. Denote $\beta_t = h(\alpha_t)\alpha_t$. Then, $\forall t \geq 0$, $\frac{d}{dt}\mathcal{C}_t^f = 0$ and $\frac{d}{dt}\mathcal{C}_t = 0$, where:

(2.7)

$$\mathcal{C}_t^f \stackrel{\text{def}}{=} h(\alpha_t)f(\mathbf{x}_t) - \int_0^t f(\mathbf{x}_\tau) \frac{d(h(\alpha_\tau))}{d\tau} d\tau + \int_0^t h(\alpha_t) \frac{\dot{\alpha}_\tau}{\alpha_\tau} \langle \nabla f(\mathbf{x}_\tau), \mathbf{x}_\tau \rangle d\tau + \psi^*(\mathbf{z}_t),$$

(2.8)

$$\mathcal{C}_t \stackrel{\text{def}}{=} \beta_t f(\mathbf{x}_t) - \beta_0 f(\mathbf{x}_0) - \int_0^t \dot{\beta}_\tau d\tau f(\mathbf{x}_\tau) d\tau + \alpha_0 D_{\psi^*}(\mathbf{z}_t, \mathbf{z}_0) + \int_0^t D_{\psi^*}(\mathbf{z}_t, \mathbf{z}_\sigma) \dot{\alpha}_\sigma d\sigma.$$

The proof of Lemma 2.3 is provided in Appendix A.

Let us now provide some context for how conserved quantities \mathcal{C}_t^f and \mathcal{C}_t lead to convergence in function value and convergence in gradient norm, respectively. First, when $h(\alpha_t) = \alpha_t$ (in which case (MoD) is equivalent to (AD)), conservation of \mathcal{C}_t^f can be shown to be equivalent to the conservation of the scaled approximate duality gap from [24, 25]. More generally, as we show in Section 3, the conservation of \mathcal{C}_t^f can be used to upper bound the optimality gap $f(\hat{\mathbf{x}}_t) - f(\mathbf{x}^*)$ for some $\hat{\mathbf{x}}_t \in \mathcal{X}$ that is constructed as a convex combination of $\{\mathbf{x}_\tau\}_{\tau \in [0, t]}$.

Consider now \mathcal{C}_t . If $\frac{d}{dt}\mathcal{C}_t = 0$, then it is not hard to check that it must be the case that $\mathcal{C}_t = 0, \forall t$. Equivalently, as it also holds that $\frac{\mathcal{C}_t}{h(\alpha_t)\alpha_t} = 0, \forall t$, we have:

$$\begin{aligned} f(\mathbf{x}_t) - \frac{h(\alpha_0)\alpha_0 f(\mathbf{x}_0) + \int_0^t \frac{d(h(\alpha_\tau)\alpha_\tau)}{d\tau} f(\mathbf{x}_\tau) d\tau}{h(\alpha_t)\alpha_t} \\ (2.9) \qquad \qquad \qquad = - \frac{\alpha_0 D_{\psi^*}(\mathbf{z}_t, \mathbf{z}_0) + \int_0^t D_{\psi^*}(\mathbf{z}_t, \mathbf{z}_\sigma) \dot{\alpha}_\sigma d\sigma}{h(\alpha_t)\alpha_t}. \end{aligned}$$

Observe that the right-hand side of (2.9) is always non-positive, as ψ^* is assumed to be convex. Suppose for now that the right-hand side of (2.9) is strictly negative and less than $-\delta$ for some $\delta > 0$. We then have: $f(\mathbf{x}_t) - \frac{h(\alpha_0)\alpha_0}{h(\alpha_t)\alpha_t} f(\mathbf{x}_0) - \frac{1}{h(\alpha_t)\alpha_t} \int_0^t \frac{d(h(\alpha_\tau)\alpha_\tau)}{d\tau} f(\mathbf{x}_\tau) d\tau < -\delta$. In other words, the function value at the last point \mathbf{x}_t is strictly smaller than a weighted average of function values at points \mathbf{x}_τ for $\tau \in [0, t]$. This means that the (weighted) average function value must be strictly decreasing with time. As the function is bounded below, after some finite time it must be that the right-hand side is at least $-\delta$. Observe that this argument can be made for any $\delta > 0$. The main idea in the analysis is to show that the inequality

$$(2.10) \qquad \frac{\alpha_0 D_{\psi^*}(\mathbf{z}_t, \mathbf{z}_0) + \int_0^t D_{\psi^*}(\mathbf{z}_t, \mathbf{z}_\sigma) \dot{\alpha}_\sigma d\sigma}{h(\alpha_t)\alpha_t} \leq \delta$$

implies that the dynamics must visit at least one point \mathbf{x} such that $\|\nabla f(\mathbf{x})\|_* \leq \epsilon(\delta)$.

3. Convergence in Function Value. In this section, we show that the invariants implied by the Hamiltonian that generates the momentum-based methods can be used to argue about convergence in function value. We start by arguing about the continuous-time case, and then show how the same invariant can be used analogously to argue about convergence of discretized versions of (MoD).

All the results will be obtained for the following choice of $h(\alpha_t)$:

$$h(\alpha_t) = \alpha_t^\lambda, \quad \text{where } \lambda \in [0, 1],$$

with the same relationship holding between their corresponding discrete-time counterparts (A_k and H_k). This choice of $h(\alpha_t)$ interpolates between the accelerated method (AD) (for $\lambda = 1$) and the generalized heavy ball method (for $\lambda = 0$).

3.1. Convergence of the Continuous-Time Dynamics. We now show how Lemma 2.3 can be used to argue about the convergence in function value of (MoD).

LEMMA 3.1. *Let $\mathbf{x}_t, \mathbf{z}_t$ evolve according to (MoD) for $h(\alpha_t) = \alpha_t^\lambda$, $\lambda \in [0, 1]$ and $\mathbf{x}_0 = \nabla\psi^*(\mathbf{z}_0) \in \text{relint}(\mathcal{X})$. If $\lambda = 0$, assume that $\frac{\dot{\alpha}_t}{\alpha_t} = \eta > 0$. Denote:*

$$\hat{\mathbf{x}}_t = \begin{cases} \lambda \frac{\alpha_t^\lambda \mathbf{x}_t + (1-\lambda) \int_0^t \dot{\alpha}_\tau \alpha_\tau^{1-\lambda} \mathbf{x}_\tau d\tau + \frac{1-\lambda}{\lambda} \mathbf{x}_0}{\alpha_t^\lambda}, & \text{if } \lambda \in (0, 1], \\ \frac{\mathbf{x}_t + \eta \int_0^t \mathbf{x}_\tau d\tau}{1+\eta t}, & \text{if } \lambda = 0. \end{cases}$$

Then, $\forall t \geq 0$:

$$f(\hat{\mathbf{x}}_t) - f(\mathbf{x}^*) \leq \begin{cases} \lambda \frac{\frac{\alpha_0}{\lambda} (f(\mathbf{x}_0) - f(\mathbf{x}^*)) + D_\psi(\mathbf{x}^*, \mathbf{x}_0)}{\alpha_t^\lambda}, & \text{if } \lambda \in (0, 1], \\ \frac{f(\mathbf{x}_0) - f(\mathbf{x}^*) + D_\psi(\mathbf{x}^*, \mathbf{x}_0)}{1+\eta t}, & \text{if } \lambda = 0. \end{cases}$$

Proof. Lemma 2.3 implies that $\mathcal{C}_t^f = \mathcal{C}_0^f$, $\forall t \geq 0$. Hence, as $h(\alpha_t) = \alpha_t^\lambda$:

$$(3.1) \quad \begin{aligned} \alpha_t^\lambda f(\mathbf{x}_t) - \alpha_0^\lambda f(\mathbf{x}_0) - \lambda \int_0^t \dot{\alpha}_\tau \alpha_\tau^{\lambda-1} f(\mathbf{x}_\tau) d\tau \\ = \psi^*(\mathbf{z}_0) - \psi^*(\mathbf{z}_t) - \int_0^t \dot{\alpha}_\tau \alpha_\tau^{\lambda-1} \langle \nabla f(\mathbf{x}_\tau), \mathbf{x}_\tau \rangle d\tau. \end{aligned}$$

Write $-\int_0^t \dot{\alpha}_\tau \alpha_\tau^{\lambda-1} \langle \nabla f(\mathbf{x}_\tau), \mathbf{x}_\tau \rangle d\tau$ as:

$$(3.2) \quad \begin{aligned} - \int_0^t \dot{\alpha}_\tau \alpha_\tau^{\lambda-1} \langle \nabla f(\mathbf{x}_\tau), \mathbf{x}_\tau \rangle d\tau \\ = \int_0^t \dot{\alpha}_\tau \alpha_\tau^{\lambda-1} \langle \nabla f(\mathbf{x}_\tau), \mathbf{x}^* - \mathbf{x}_\tau \rangle d\tau - \int_0^t \dot{\alpha}_\tau \alpha_\tau^{\lambda-1} \langle \nabla f(\mathbf{x}_\tau), \mathbf{x}^* \rangle d\tau. \end{aligned}$$

Observe that, by convexity of f :

$$(3.3) \quad \int_0^t \dot{\alpha}_\tau \alpha_\tau^{\lambda-1} \langle \nabla f(\mathbf{x}_\tau), \mathbf{x}^* - \mathbf{x}_\tau \rangle d\tau \leq \int_0^t \dot{\alpha}_\tau \alpha_\tau^{\lambda-1} (f(\mathbf{x}^*) - f(\mathbf{x}_\tau)) d\tau.$$

By the definition of \mathbf{z}_t from (MoD),

$$(3.4) \quad - \int_0^t \dot{\alpha}_\tau \alpha_\tau^{\lambda-1} \langle \nabla f(\mathbf{x}_\tau), \mathbf{x}^* \rangle d\tau = \int_0^t \langle \dot{\mathbf{z}}_\tau, \mathbf{x}^* \rangle = \langle \mathbf{z}_t - \mathbf{z}_0, \mathbf{x}^* \rangle.$$

The next step is to combine $\langle \mathbf{z}_t - \mathbf{z}_0, \mathbf{x}^* \rangle$ with $\psi^*(\mathbf{z}_0) - \psi^*(\mathbf{z}_t)$ to write them in the form of Bregman divergences. In particular, define \mathbf{z}^* so that $\nabla\psi^*(\mathbf{z}^*) = \mathbf{x}^*$. Then:

$$\psi^*(\mathbf{z}_0) - \psi^*(\mathbf{z}_t) + \langle \mathbf{z}_t - \mathbf{z}_0, \mathbf{x}^* \rangle = D_{\psi^*}(\mathbf{z}_0, \mathbf{z}^*) - D_{\psi^*}(\mathbf{z}_t, \mathbf{z}^*) \leq D_{\psi^*}(\mathbf{z}_0, \mathbf{z}^*).$$

Using Fact 1.5 (which implies $\psi^*(\mathbf{z}) = \langle \nabla\mathbf{z}, \nabla\psi^*(\mathbf{z}) \rangle - \psi(\nabla\psi^*(\mathbf{z}))$) and $\mathbf{z}_0 = \nabla\psi(\mathbf{x}_0)$ (which follows from the assumption that \mathbf{x}_0 is from the relative interior of \mathcal{X}), it is not

hard to show that: $D_{\psi^*}(\mathbf{z}_0, \mathbf{z}^*) = D_{\psi}(\nabla\psi^*(\mathbf{z}^*), \nabla\psi^*(\mathbf{z}_0)) = D_{\psi}(\mathbf{x}^*, \mathbf{x}_0)$. Combining with (3.1)-(3.4):

$$\begin{aligned} \alpha_t^\lambda f(\mathbf{x}_t) - \alpha_0^\lambda f(\mathbf{x}_0) - \lambda \int_0^t \dot{\alpha}_\tau \alpha_\tau^{\lambda-1} f(\mathbf{x}_\tau) d\tau \\ \leq \int_0^t \dot{\alpha}_\tau \alpha_\tau^{\lambda-1} (f(\mathbf{x}^*) - f(\mathbf{x}_\tau)) d\tau + D_{\psi}(\mathbf{x}^*, \mathbf{x}_0). \end{aligned}$$

Assume first that $\lambda > 0$. Integrating and rearranging the terms in the last inequality:

$$\begin{aligned} \alpha_t^\lambda f(\mathbf{x}_t) + (1 - \lambda) \int_0^t \dot{\alpha}_\tau \alpha_\tau^{\lambda-1} f(\mathbf{x}_\tau) d\tau + \frac{1 - \lambda}{\lambda} f(\mathbf{x}_0) - \frac{\alpha_t^\lambda}{\lambda} f(\mathbf{x}^*) \\ \leq \frac{\alpha_0}{\lambda} (f(\mathbf{x}_0) - f(\mathbf{x}^*)) + D_{\psi}(\mathbf{x}^*, \mathbf{x}_0). \end{aligned}$$

It remains to divide both sides by $\frac{\alpha_t^\lambda}{\lambda}$ and apply Jensen's inequality.

If $\lambda = 0$, then, assuming $\frac{\dot{\alpha}_t}{\alpha_t} = \eta$:

$$f(\mathbf{x}_t) + \eta \int_0^t f(\mathbf{x}_\tau) d\tau - (1 + \eta t) f(\mathbf{x}^*) \leq f(\mathbf{x}_0) - f(\mathbf{x}^*) + D_{\psi}(\mathbf{x}^*, \mathbf{x}_0).$$

It remains to divide both sides by $1 + \eta t$ and apply Jensen's inequality. \square

Observe that when $\lambda = 1$ (that is, when (MoD) is equivalent to (AD)), $\hat{\mathbf{x}}_t = \mathbf{x}_t$, and we recover the standard guarantee on the last iterate of the accelerated dynamics [24, 25, 38]. When $\lambda = 0$, we obtain a $1/t$ convergence rate for the generalization of the heavy ball method. The result applies to constrained optimization and non-Euclidean spaces. We note that a generalization of the heavy ball method to constrained convex optimization was previously considered in [4]. However, the result from [4] applies only to Hilbert spaces and provides weak (asymptotic) convergence results. The second-order ODE considered in [4] seems to correspond to a different continuous-time dynamics than (MoD) with $h(\alpha_t) = 1$, and it is unclear how to compare it to (MoD).

3.2. Discrete-Time Convergence. Define the discrete-time counterpart to the continuous-time conserved quantity $\mathcal{C}_t^f, \mathcal{C}_k^f$, as:

$$\mathcal{C}_k^f = H_k f(\mathbf{y}_k) - \sum_{i=1}^k h_i f(\mathbf{x}_i) + \sum_{i=1}^k H_i \frac{a_i}{A_i} \langle \nabla f(\mathbf{x}_i), \mathbf{x}_i \rangle + \psi^*(\mathbf{z}_k),$$

where $A_k = \sum_{i=0}^k a_i$, $H_k = \sum_{i=0}^k h_i$, and $H_k = A_k^\lambda$.

The discretization of the continuous-time dynamics that we will use is:

$$\begin{aligned} \mathbf{x}_k &= \frac{H_{k-1}/H_k}{H_{k-1}/H_k + a_k/A_k} \mathbf{y}_{k-1} + \frac{a_k/A_k}{H_{k-1}/H_k + a_k/A_k} \nabla\psi^*(\mathbf{z}_{k-1}), \\ (\text{GMD}_{\hat{f}}) \quad \mathbf{z}_k &= \mathbf{z}_{k-1} - H_k \frac{a_k}{A_k} \nabla f(\mathbf{x}_k), \\ \mathbf{y}_k &= \mathbf{x}_k + \frac{a_k}{A_k} (\nabla\psi^*(\mathbf{z}_k) - \nabla\psi^*(\mathbf{z}_{k-1})). \end{aligned}$$

The motivation for this particular choice of the discretization will become clear from Proposition 3.3. (In particular, the discretization was chosen to ensure that $\mathcal{C}_k^f \leq$

\mathcal{C}_{k-1}^f .) Note that when $H_k = A_k$ ($\lambda = 1$), the method is precisely the AGD+ method from [21]. Note that AGD+ is closely related to Nesterov's method [47] and the accelerated proximal method of Güler [32]: both Nesterov's and Güler's methods can be seen as alternative discretizations of (AD) (or (MoD) with $h(\tau) = \tau$), where \mathbf{y}_k is replaced by a correction step, which is the gradient descent step for Nesterov's method and the proximal step for Güler's method. This interpretation of Nesterov's method can also be found in [25].

For (GMD_f) to apply to constrained minimization, we need to show that the iterates \mathbf{y}_k remain in the feasible set. This is established by the following proposition.

PROPOSITION 3.2. *Let $\mathbf{x}_k, \mathbf{y}_k, \mathbf{z}_k$ evolve as in (GMD_f), where the initial point satisfies $\mathbf{y}_0 = \nabla\psi^*(\mathbf{z}_0) \in \text{relint}\mathcal{X}$, $\mu \leq L$, and $a_k, A_k, H_k > 0$ satisfy: $A_k = \sum_{i=0}^k a_i$, $\frac{a_k^2}{A_k^2} = c \frac{\mu}{LH_k}$, for some $c \in (0, 1]$, and $H_k = A_k^\lambda$ for $\lambda \in [0, 1]$. Then $\mathbf{y}_k \in \mathcal{X}$, $\forall k \geq 0$.*

Proof. The claim clearly holds for $k = 0$, by the initialization. To simplify the notation, denote $\theta_k = \frac{a_k}{A_k}$, $\theta'_k = \frac{a_k/A_k}{H_{k-1}/H_k + a_k/A_k} = \frac{\theta_k}{H_{k-1}/H_k + \theta_k}$, $\mathbf{v}_k = \nabla\psi^*(\mathbf{z}_k)$. By the definition of a convex conjugate (Definition 1.4) and Fact 1.5, $\mathbf{v}_k \in \mathcal{X}$, $\forall k$.

Under the assumptions of the proposition, it is not hard to see that $\theta_k \leq \theta_{k-1}$, $\forall k \geq 1$. Namely, this condition is equivalent to $\frac{c\mu}{LH_k} \leq \frac{c\mu}{LH_{k-1}}$, which, by the definition of H_k , is equivalent to $A_{k-1}^\lambda \leq A_k^\lambda$. As A_k is non-increasing with k (as $A_k = \sum_{i=0}^k a_i$ and $a_i \geq 0$, $\forall i$) and $\lambda \geq 0$, we clearly have that $A_{k-1}^\lambda \leq A_k^\lambda$, and, thus $\theta_k \leq \theta_{k-1}$.

To prove the proposition, we will first show that \mathbf{y}_k can be expressed as a non-negative linear combination of $\{\mathbf{v}_i\}_{i=0}^k$. We subsequently show by induction on k that the coefficients of that linear combination must sum to one, which completes the proof.

Using (GMD_f), we can write \mathbf{y}_k in the following recursive form:

$$\mathbf{y}_k = (1 - \theta'_k)\mathbf{y}_{k-1} + \theta'_k\mathbf{v}_{k-1} + \theta_k(\mathbf{v}_k - \mathbf{v}_{k-1}).$$

Applying this definition recursively over $i = 0, \dots, k$ and using $\mathbf{y}_0 = \mathbf{v}_0$, we have

$$\mathbf{y}_k = \sum_{i=0}^k \gamma_{i,k} \mathbf{v}_i,$$

$$\gamma_{i,k} = \begin{cases} \theta_k, & \text{if } i = k; \\ \left[\prod_{j=i+2}^k (1 - \theta'_j) \right] \left[\theta'_{i+1}(1 - \theta_i) + \theta_i - \theta_{i+1} \right], & \text{if } 1 \leq i \leq k-1; \\ \left[\prod_{j=1}^k (1 - \theta'_j) \right], & \text{if } i = 0, \end{cases}$$

where, by convention, we take $\prod_i^j(\cdot) = 1$ whenever $j < i$. Given that for all $i \geq 0$, we have $\theta_i, \theta'_i \in [0, 1]$ and $\theta_{i+1} \leq \theta_i$, it immediately follows that $\gamma_{i,k} \geq 0$, $\forall i \in \{0, 1, \dots, k\}$.

We now show by induction on k that it must be the case that $\sum_{i=0}^k \gamma_{i,k} = 1$. This clearly holds for $k = 0$. Suppose that it holds for some $k-1 \geq 0$. Then $\mathbf{y}_{k-1} \in \mathcal{X}$. As $\mathbf{x}_k = (1 - \theta'_k)\mathbf{y}_{k-1} + \theta'_k\mathbf{v}_{k-1}$, it follows that $\mathbf{x}_k \in \mathcal{X}$, and, moreover, \mathbf{x}_k is a convex combination of $\{\mathbf{v}_i\}_{i=0}^{k-1}$. Now, observe from the definition of \mathbf{x}_k that $\theta'_k\mathbf{v}'_{k-1} = \mathbf{x}_k - (1 - \theta'_k)\mathbf{y}_{k-1}$. Hence, we can express \mathbf{y}_k as:

$$\mathbf{y}_k = (1 - \theta_k/\theta'_k)\mathbf{x}_k + \frac{\theta_k}{\theta'_k}(1 - \theta'_k)\mathbf{y}_{k-1} + \theta_k\mathbf{v}_k.$$

As $1 - \frac{\theta_k}{\theta'_k} + \frac{\theta_k}{\theta'_k}(1 - \theta'_k) + \theta_k = 1$ and each $\mathbf{x}_k, \mathbf{y}_k, \mathbf{z}_k$ are convex combinations of $\{\mathbf{v}_i\}_{i=0}^k$, it follows that $\sum_{i=0}^k \gamma_{i,k} = 1$, which, together with the fact that $\gamma_{i,k} \geq 0, \forall i$, completes the proof. \square

PROPOSITION 3.3. *Let $\mathbf{x}_k, \mathbf{y}_k, \mathbf{z}_k$ evolve according to (GMD_f), where $\psi : \mathcal{X} \rightarrow \mathbb{R}$ is a μ -strongly convex function, $\mathbf{y}_0 = \nabla\psi^*(\mathbf{z}_0) \in \text{relint}\mathcal{X}$ and $a_k, A_k, H_k > 0$ satisfy: $A_k = \sum_{i=0}^k a_i$ and $\frac{a_k^2}{A_k^2} \leq \frac{\mu}{LH_k}$. Then $\mathcal{C}_k^f \leq \mathcal{C}_{k-1}^f, \forall k \geq 1$.*

Proof. By the definition of \mathcal{C}_k^f , we have:

$$(3.5) \quad \begin{aligned} \mathcal{C}_k^f - \mathcal{C}_{k-1}^f &= H_k f(\mathbf{y}_k) - H_{k-1} f(\mathbf{y}_{k-1}) - h_k f(\mathbf{x}_k) + H_k \frac{a_k}{A_k} \langle \nabla f(\mathbf{x}_k), \mathbf{x}_k \rangle \\ &\quad + \psi^*(\mathbf{z}_k) - \psi^*(\mathbf{z}_{k-1}). \end{aligned}$$

Observe first, by smoothness and convexity of f :

$$(3.6) \quad \begin{aligned} H_k f(\mathbf{y}_k) - H_{k-1} f(\mathbf{y}_{k-1}) - h_k f(\mathbf{x}_k) \\ &= H_k (f(\mathbf{y}_k) - f(\mathbf{x}_k)) + H_{k-1} (f(\mathbf{x}_k) - f(\mathbf{y}_{k-1})) \\ &\leq \langle \nabla f(\mathbf{x}_k), H_k \mathbf{y}_k - H_{k-1} \mathbf{y}_{k-1} - h_k \mathbf{x}_k \rangle + H_k \frac{L}{2} \|\mathbf{y}_k - \mathbf{x}_k\|^2. \end{aligned}$$

On the other hand, by the definitions of a Bregman divergence and \mathbf{z}_k :

$$(3.7) \quad \begin{aligned} \psi^*(\mathbf{z}_k) - \psi^*(\mathbf{z}_{k-1}) &= -D_{\psi^*}(\mathbf{z}_{k-1}, \mathbf{z}_k) - \langle \nabla\psi^*(\mathbf{z}_{k-1}), \mathbf{z}_{k-1} - \mathbf{z}_k \rangle \\ &= -D_{\psi^*}(\mathbf{z}_{k-1}, \mathbf{z}_k) - H_k \frac{a_k}{A_k} \langle \nabla f(\mathbf{x}_k), \nabla\psi^*(\mathbf{z}_{k-1}) \rangle. \end{aligned}$$

As ψ is μ -strongly convex, we have (by Fact 1.7): $D_{\psi^*}(\mathbf{z}_{k-1}, \mathbf{z}_k) \geq \frac{\mu}{2} \|\nabla\psi^*(\mathbf{z}_k) - \nabla\psi^*(\mathbf{z}_{k-1})\|^2$. Hence, combining (3.5)-(3.7):

$$\begin{aligned} \mathcal{C}_k^f - \mathcal{C}_{k-1}^f &\leq H_k \frac{L}{2} \|\mathbf{y}_k - \mathbf{x}_k\|^2 - \frac{\mu}{2} \|\nabla\psi^*(\mathbf{z}_k) - \nabla\psi^*(\mathbf{z}_{k-1})\|^2 \\ &\quad + \left\langle \nabla f(\mathbf{x}_k), H_k \mathbf{y}_k - H_{k-1} \mathbf{y}_{k-1} - h_k \mathbf{x}_k + H_k \frac{a_k}{A_k} (\mathbf{x}_k - \nabla\psi^*(\mathbf{z}_k)) \right\rangle. \end{aligned}$$

Note that we want to make the right-hand side of the last inequality non-positive. To do so, we can make the last term equal to zero by setting: $H_k \mathbf{y}_k - H_{k-1} \mathbf{y}_{k-1} - h_k \mathbf{x}_k + H_k \frac{a_k}{A_k} (\mathbf{x}_k - \nabla\psi^*(\mathbf{z}_k)) = 0$. To make the first two terms non-positive, we require $\mathbf{y}_k - \mathbf{x}_k = \frac{a_k}{A_k} (\nabla\psi^*(\mathbf{z}_k) - \nabla\psi^*(\mathbf{z}_{k-1}))$.⁸ Solving these last two equations for \mathbf{x}_k gives (GMD_f). We conclude by using $H_k \frac{a_k^2}{A_k^2} \leq \frac{\mu}{L}$. \square

To obtain a convergence rate for (GMD_f), it remains to show that $\mathcal{C}_k^f \leq \mathcal{C}_0^f$ implies a convergence in function value for (GMD_f), as was done for the continuous-time case in Lemma 3.1.

THEOREM 3.4. *Let $\mathbf{x}_k, \mathbf{y}_k, \mathbf{z}_k$ evolve according to (GMD_f), where $\mathbf{x}_0 = \mathbf{y}_0 = \nabla\psi^*(\mathbf{z}_0) \in \text{relint}\mathcal{X}$, $\psi : \mathcal{X} \rightarrow \mathbb{R}$ is μ -strongly convex, $H_k = A_k^\lambda, \lambda \in [0, 1], A_k = \sum_{i=0}^k a_k, a_0 = 1$, and $\frac{a_k^2}{A_k^2} = c \frac{\mu}{LH_i}$, for $c \in (0, 1], k \geq 1$. Define:*

$$\hat{\mathbf{x}}_k = \frac{H_k \mathbf{y}_k + \sum_{i=1}^k \left(\frac{a_i}{A_i} H_i - h_i\right) \mathbf{x}_i}{H_0 + \sum_{i=1}^k \frac{a_i}{A_i} H_i}.$$

⁸Note that the multiplier $\frac{a_k}{A_k}$ on the right-hand side is necessary here for \mathbf{x}_k to be explicitly defined. Any other factor would make \mathbf{x}_k depend on $\nabla\psi^*(\mathbf{z}_k)$, which is a function of \mathbf{x}_k (as $\mathbf{z}_k = \mathbf{z}_{k-1} - H_k \frac{a_k}{A_k} \nabla f(\mathbf{x}_k)$), and would thus make \mathbf{x}_k be only implicitly defined.

Then, $\forall k \geq 1$, $\hat{\mathbf{x}}_k \in \mathcal{X}$ and:

$$f(\hat{\mathbf{x}}_k) - f(\mathbf{x}^*) \leq \begin{cases} \frac{f(\mathbf{x}_0) - f(\mathbf{x}^*) + D_\psi(\mathbf{x}^*, \mathbf{x}_0)}{1 + \sqrt{c\mu/Lk}}, & \text{if } \lambda = 0, k \geq 1 \\ \Theta(1)\lambda^\lambda \frac{L}{c\mu} \frac{f(\mathbf{x}_0) - f(\mathbf{x}^*) + D_\psi(\mathbf{x}^*, \mathbf{x}_0)}{k^2}, & \text{if } \lambda \in (0, 1], k = \Omega(1/\lambda). \end{cases}$$

Proof. By Proposition 3.2, $\mathbf{y}_k \in \mathcal{X}$, $\forall k$. As $\mathbf{x}_k \in \mathcal{X}$ (as a convex combination of \mathbf{y}_{k-1} , $\nabla\psi^*(\mathbf{z}_{k-1}) \in \mathcal{X}$), we have that $\hat{\mathbf{x}}_k$ is a convex combination of points from the feasible space \mathcal{X} , and, thus, it must be the case that $\hat{\mathbf{x}}_k \in \mathcal{X}$, $\forall k$.

By Proposition 3.3, $\mathcal{C}_k^f \leq \mathcal{C}_0^f$. Hence:

$$(3.8) \quad H_k f(\mathbf{y}_k) - H_0 f(\mathbf{y}_0) - \sum_{i=1}^k h_i f(\mathbf{x}_i) \leq \psi^*(\mathbf{z}_0) - \psi^*(\mathbf{z}_k) - \sum_{i=1}^k H_i \frac{a_i}{A_i} \langle \nabla f(\mathbf{x}_i), \mathbf{x}_i \rangle.$$

As in Lemma 3.1, write $-\sum_{i=1}^k H_i \frac{a_i}{A_i} \langle \nabla f(\mathbf{x}_i), \mathbf{x}_i \rangle$ as:

$$(3.9) \quad -\sum_{i=1}^k H_i \frac{a_i}{A_i} \langle \nabla f(\mathbf{x}_i), \mathbf{x}_i \rangle = \sum_{i=1}^k H_i \frac{a_i}{A_i} \langle \nabla f(\mathbf{x}_i), \mathbf{x}^* - \mathbf{x}_i \rangle - \sum_{i=1}^k H_i \frac{a_i}{A_i} \langle \nabla f(\mathbf{x}_i), \mathbf{x}^* \rangle.$$

By convexity of f :

$$(3.10) \quad \sum_{i=1}^k H_i \frac{a_i}{A_i} \langle \nabla f(\mathbf{x}_i), \mathbf{x}^* - \mathbf{x}_i \rangle \leq \sum_{i=1}^k H_i \frac{a_i}{A_i} (f(\mathbf{x}^*) - f(\mathbf{x}_i)).$$

Let \mathbf{z}^* be such that $\nabla\psi^*(\mathbf{z}^*) = \mathbf{x}^*$. Then, by the same arguments as in the proof of Lemma 3.1:

$$(3.11) \quad \begin{aligned} \psi^*(\mathbf{z}_0) - \psi^*(\mathbf{z}_k) - \sum_{i=1}^k H_i \frac{a_i}{A_i} \langle \nabla f(\mathbf{x}_i), \mathbf{x}^* \rangle &= D_{\psi^*}(\mathbf{z}_0, \mathbf{z}^*) - D_{\psi^*}(\mathbf{z}_k, \mathbf{z}^*) \\ &\leq D_\psi(\mathbf{x}^*, \mathbf{x}_0). \end{aligned}$$

Combining (3.8)-(3.11):

$$H_k f(\mathbf{y}_k) - H_0 f(\mathbf{y}_0) - \sum_{i=1}^k h_i f(\mathbf{x}_i) \leq \sum_{i=1}^k H_i \frac{a_i}{A_i} (f(\mathbf{x}^*) - f(\mathbf{x}_i)) + D_\psi(\mathbf{x}^*, \mathbf{x}_0).$$

To complete the proof, it remains to rearrange the terms in the last equation. Notice that $\sum_{i=1}^k h_i = H_k - H_0$, and, thus, the coefficients multiplying $f(\cdot)$ sum up to zero. Notice also that, as $H_i = A_i^\lambda$, $a_i = A_i - A_{i-1}$, and $h_i = H_i - H_{i-1}$, it must be the case that $H_i \frac{a_i}{A_i} - h_i \geq 0$. We have:

$$\begin{aligned} H_k f(\mathbf{y}_k) + \sum_{i=1}^k \left(H_i \frac{a_i}{A_i} - h_i \right) f(\mathbf{x}_i) - \left(H_0 + \sum_{i=1}^k H_i \frac{a_i}{A_i} \right) f(\mathbf{x}^*) \\ \leq H_0 (f(\mathbf{x}_0) - f(\mathbf{x}^*)) + D_\psi(\mathbf{x}^*, \mathbf{x}_0). \end{aligned}$$

Dividing both sides of the last equation by $\left(H_0 + \sum_{i=1}^k H_i \frac{a_i}{A_i} \right)$ and applying Jensen's inequality:

$$f(\hat{\mathbf{x}}_k) - f(\mathbf{x}^*) \leq \frac{H_0 (f(\mathbf{y}_0) - f(\mathbf{x}^*)) + D_\psi(\mathbf{x}^*, \mathbf{x}_0)}{\left(H_0 + \sum_{i=1}^k H_i \frac{a_i}{A_i} \right)}.$$

Recall that, as $a_0 = 1$, we must have $H_0 = 1$. To bound $H_0 + \sum_{i=1}^k H_i \frac{a_i}{A_i}$, we need to argue about the growth of $\frac{a_i}{A_i} H_i = \frac{a_i}{A_i^{1-\lambda}}$. When $\lambda = 0$, we have $\frac{a_i}{A_i^{1-\lambda}} = \sqrt{\frac{c\mu}{L}}$. Assume now that $\lambda > 0$ and $k = \Omega(\frac{1}{\lambda})$. By assumption of the theorem, $\frac{a_i^2}{A_i^{2-\lambda}} = \frac{c\mu}{L}$. If $a_i \propto (\frac{c\mu\lambda}{2L})^{1/\lambda} i^{2/\lambda-1}$ and $i = \Omega(\frac{1}{\lambda})$, then $A_i \propto \frac{\lambda}{2} (\frac{c\mu\lambda}{2L})^{1/\lambda} i^{2/\lambda}$, and it can be ensured that $\frac{a_i^2}{A_i^{2-\lambda}} = \frac{c\mu}{L}$ holds. Thus, for $i \geq k/2$, $\frac{a_i}{A_i^{1-\lambda}} = \Theta(1)\lambda^\lambda \frac{c\mu}{L} i$. Hence, $\sum_{i=1}^k \frac{a_i}{A_i} H_i$ scales as $\Theta(k)\sqrt{\frac{c\mu}{L}}$ for $\lambda = 0$ and $\Theta(k^2)\lambda^\lambda \frac{c\mu}{L}$ for $\lambda > 0$ and $k = \Omega(\frac{1}{\lambda})$. \square

Observe that a generalization of the heavy ball method (obtained from (GMD_f) for $H_k = A_k^0 = 1$) converges at rate $1/k$ for smooth convex functions. The result applies to constrained minimization and general normed vector spaces. Note that such a (nonasymptotic) result was previously known only for the setting of unconstrained minimization in Euclidean spaces [29].

4. Convergence in Norm of the Gradient. In this section we turn to convergence in terms of the norm of the gradient. We focus on discrete-time methods, based on a discrete-time counterpart to the conserved quantity \mathcal{C}_t . Throughout the section, we assume that $\mathcal{X} \equiv E$. We start the section with an overview of the approach and a structural (algorithm-independent) lemma that is used later in proving the convergence results. We then present the results for Euclidean spaces in Section 4.2, and show how these results can be generalized to non-Euclidean spaces in Section 4.3.

4.1. An Overview of the Approach and a Structural Lemma. We consider a counterpart to \mathcal{C}_t that was derived for general momentum methods and defined in Lemma 2.3. This counterpart is defined as:

$$(4.1) \quad \mathcal{C}_k = B_k f(\mathbf{y}_k) - \sum_{i=0}^k b_i f(\mathbf{y}_i) + \sum_{i=0}^k a_i D_{\psi^*}(\mathbf{z}_k, \mathbf{z}_i),$$

where $A_k = \sum_{i=0}^k a_i$, $B_k = \sum_{i=0}^k b_i$, and $a_i, b_i > 0, \forall i \geq 0$. Going from continuous to the discrete time, $\alpha_t, h(\alpha_t), \alpha_t h(\alpha_t)$ translate into A_k, H_k, B_k , respectively.

To characterize convergence to stationary points, we denote by $E_k \stackrel{\text{def}}{=} \mathcal{C}_k - \mathcal{C}_{k-1}$ the discretization error between the iterations $k-1$ and k (recall that in the continuous time domain, \mathcal{C}_t was conserved). As $\mathcal{C}_0 = 0$, we clearly have $\mathcal{C}_k = \sum_{i=1}^k E_i$. Given specific assumptions about the objective function (e.g., if it is convex or nonconvex, its degree of smoothness, etc.), we will need to argue that the total discretization error $\sum_{i=1}^k E_i$ is “sufficiently small” (possibly zero, or even negative) under those assumptions. In general, the magnitude of the discretization error will be determined by the step sizes a_i, b_i , which will be one of the constraining factors determining the rate of convergence.

Average Function Value Decrease. The following (algorithm-independent) lemma implies that if \mathcal{C}_k is non-increasing with k (namely, if $E_k \leq 0, \forall k$) then the average function value taken at all points constructed by the algorithm is decreasing with the iteration count k . The lemma will be crucial in obtaining the results for the convergence in norm of the gradient.

LEMMA 4.1. *Let $\mathcal{C}_k = B_k f(\mathbf{y}_k) - \sum_{i=0}^k b_i f(\mathbf{y}_i) + \sum_{i=0}^k a_i D_{\psi^*}(\mathbf{z}_k, \mathbf{z}_i)$ and $E_k =$*

$\mathcal{C}_k - \mathcal{C}_{k-1}$, where $\forall i, a_i, b_i > 0$ and $\forall k, A_k = \sum_{i=0}^k a_i, B_k = \sum_{i=0}^k b_i$. Then:

$$\begin{aligned} \frac{1}{B_k} \sum_{i=0}^k b_i f(\mathbf{y}_i) &= f(\mathbf{y}_0) - \sum_{i=1}^k \left(\frac{1}{B_{i-1}} - \frac{1}{B_i} \right) \sum_{j=0}^{i-1} a_j D_{\psi^*}(\mathbf{z}_i, \mathbf{z}_j) \\ &\quad + \sum_{i=1}^k \left(\frac{1}{B_{i-1}} - \frac{1}{B_k} \right) E_i. \end{aligned}$$

Proof. Denote $S_k = \sum_{i=0}^k b_i f(\mathbf{y}_i)$. The proof is by induction on k . The base case is immediate, as, by the definition of S_k and B_k , we have $S_0 = b_0 f(\mathbf{y}_0) = B_0 f(\mathbf{y}_0)$. Now assume that the statement is true for $k \geq 0$. Then, by the definition of S_k ,

$$(4.2) \quad S_{k+1} = S_k + b_{k+1} f(\mathbf{y}_{k+1}).$$

On the other hand, by the definitions of \mathcal{C}_k and E_k , we have $\mathcal{C}_k = \sum_{i=1}^k E_i$, and, thus:

$$(4.3) \quad f(\mathbf{y}_{k+1}) = \frac{1}{B_k} \left(S_k - \sum_{j=0}^k a_j D_{\psi^*}(\mathbf{z}_{k+1}, \mathbf{z}_j) + \sum_{i=1}^{k+1} E_i \right).$$

Combining Equations (4.2) and (4.3):

$$\begin{aligned} S_{k+1} &= \left(1 + \frac{b_{k+1}}{B_k} \right) S_k - \frac{b_{k+1}}{B_k} \sum_{j=0}^k a_j D_{\psi^*}(\mathbf{z}_{k+1}, \mathbf{z}_j) + \frac{b_{k+1}}{B_k} \sum_{i=1}^{k+1} E_i \\ &= \frac{B_{k+1}}{B_k} S_k - B_{k+1} \sum_{j=0}^k \left(\frac{1}{B_k} - \frac{1}{B_{k+1}} \right) a_j D_{\psi^*}(\mathbf{z}_{k+1}, \mathbf{z}_j) \\ &\quad + B_{k+1} \sum_{i=1}^{k+1} \left(\frac{1}{B_k} - \frac{1}{B_{k+1}} \right) E_i, \end{aligned}$$

where we have used $B_{k+1} = B_k + b_{k+1}$. Applying the inductive hypothesis and grouping terms into appropriate summations completes the proof. \square

4.2. Convergence to Stationary Points in Euclidean Spaces. In this section, we take $(E, \|\cdot\|)$ to be a Euclidean space. The following simple claim is useful for passing from Bregman divergences to gradient norms.

CLAIM 4.2. *Let $a, b > 0$. Then $a\|\mathbf{z} + \Delta\mathbf{z}\|^2 + b\|\mathbf{z}\|^2 \geq \frac{ab}{a+b} \|\Delta\mathbf{z}\|^2$.*

Proof. As $\|\cdot\| = \sqrt{\langle \cdot, \cdot \rangle}$, we have that $a\|\mathbf{z} + \Delta\mathbf{z}\|^2 + b\|\mathbf{z}\|^2 = (a+b)\|\mathbf{z}\|^2 + 2a\langle \mathbf{z}, \Delta\mathbf{z} \rangle + a\|\Delta\mathbf{z}\|^2$. By the Cauchy-Schwarz inequality, $\langle \mathbf{z}, \Delta\mathbf{z} \rangle \leq \|\mathbf{z}\| \|\Delta\mathbf{z}\|$. Since the claim holds trivially for $\|\Delta\mathbf{z}\| = 0$, assume $\|\Delta\mathbf{z}\| \neq 0$ and let $c = \frac{\|\mathbf{z}\|}{\|\Delta\mathbf{z}\|}$. Then $a\|\mathbf{z} + \Delta\mathbf{z}\|^2 + b\|\mathbf{z}\|^2 \geq \|\Delta\mathbf{z}\|^2 ((a+b)c^2 - 2ac + a)$. As $(a+b)c^2 - 2ac + a$ is minimized for $c = \frac{a}{a+b}$, the claim follows. \square

To relate the Bregman divergences in the definition of \mathcal{C}_k to norms of the gradients, we will use the following application of Claim 4.2.

PROPOSITION 4.3. *Let $\psi(\mathbf{x}) = \frac{\mu}{2} \|\mathbf{x}\|^2$ (so that $D_{\psi^*}(\mathbf{w}, \mathbf{z}) = \frac{1}{2\mu} \|\mathbf{w} - \mathbf{z}\|^2$). Then:*

$$\sum_{j=0}^{i-1} a_j D_{\psi^*}(\mathbf{z}_i, \mathbf{z}_j) \geq \frac{1}{2\mu} \sum_{j=1}^i \nu_j^i \left\| \nabla f(\mathbf{x}_j) \right\|^2,$$

where $\nu_i^i = \frac{a_i^2}{A_i^2} H_i^2 \left(\frac{a_i}{2} + \frac{a_i a_{i-1}}{a_i + a_{i-1}} \right)$ and $\nu_j^i = \frac{a_{j-1} a_j^3}{a_{j-1} + a_j} \frac{H_j^2}{A_j^2}$ for $1 \leq j \leq i-1$.

Proof. By the choice of function ψ and by $\|\cdot\|$ being induced by an inner product, we have that $D_{\psi^*}(\mathbf{w}, \mathbf{z}) = \frac{1}{2\mu}\|\mathbf{w} - \mathbf{z}\|^2$. By the definition of \mathbf{z}_k , we have that, $\forall i > j \geq 0$, $D_{\psi^*}(\mathbf{z}_i, \mathbf{z}_j) = \frac{1}{2\mu}\|\sum_{k=j+1}^i \frac{a_k}{A_k} H_k \nabla f(\mathbf{x}_k)\|^2$. By Claim 4.2, $\forall i > j + 1 > 0$:

$$(4.4) \quad a_j D_{\psi^*}(\mathbf{z}_i, \mathbf{z}_j) + a_{j+1} D_{\psi^*}(\mathbf{z}_i, \mathbf{z}_{j+1}) \geq \frac{1}{2\mu} \frac{a_j a_{j+1}}{a_j + a_{j+1}} \left\| \frac{a_{j+1}}{A_{j+1}} H_{j+1} \nabla f(\mathbf{x}_{j+1}) \right\|^2.$$

Write $\sum_{j=0}^{i-1} a_j D_{\psi^*}(\mathbf{z}_i, \mathbf{z}_j)$ as:

$$\begin{aligned} \sum_{j=0}^{i-1} a_j D_{\psi^*}(\mathbf{z}_i, \mathbf{z}_j) &= \frac{a_0}{2} D_{\psi^*}(\mathbf{z}_i, \mathbf{z}_0) + \frac{a_{i-1}}{2} D_{\psi^*}(\mathbf{z}_i, \mathbf{z}_{i-1}) \\ &\quad + \frac{1}{2} \sum_{j=0}^{i-2} \left(a_j D_{\psi^*}(\mathbf{z}_i, \mathbf{z}_j) + a_{j+1} D_{\psi^*}(\mathbf{z}_i, \mathbf{z}_{j+1}) \right). \end{aligned}$$

Combine the last equation with Eq. (4.4) and $\frac{a_0}{2} D_{\psi^*}(\mathbf{z}_i, \mathbf{z}_0) \geq 0$. \square

Discrete-Time Methods. The particular two-step discretization (reminiscent of a predictor-corrector method) that we consider for general momentum dynamics (MoD) is:

$$\begin{aligned} \mathbf{x}_k &= \frac{A_{k-1}}{A_k} \mathbf{y}_{k-1} + \frac{a_k}{A_k} \nabla \psi^*(\mathbf{z}_{k-1}), \\ (\text{GMD}) \quad \mathbf{z}_k &= \mathbf{z}_{k-1} - \frac{a_k}{A_k} H_k \nabla f(\mathbf{x}_k), \\ \mathbf{y}_k &= \mathbf{x}_k + \frac{a_k}{A_k} (\nabla \psi^*(\mathbf{z}_k) - \nabla \psi^*(\mathbf{z}_{k-1})). \end{aligned}$$

When $H_k = A_k$, (GMD) is equivalent to (GMD_f), AGD+ [21], and the method of similar triangles [28], which generalize Nesterov's accelerated method [47] and accelerated extra-gradient method [24]. When $H_k = 1$, (GMD) is a slightly different discretization of the generalized heavy-ball dynamics than (GMD_f) with $H_k = 1$, where the difference lies only in the size of the extrapolation step \mathbf{x}_k . Working with the general momentum method (GMD) will allow us to obtain results for AGD+ and the generalized heavy-ball method as special cases, and it will also allow us to understand how the different choices of H_k affect the convergence in norm of the gradient.

Discretization Error. Characterization of the discretization error is what crucially determines the convergence of the methods in norm of the gradient, as well as in function value. Here, we show how the discretization error is affected by the choice of the step size and assumptions about the objective function, such as smoothness and convexity.

LEMMA 4.4. *Let $E_k \stackrel{\text{def}}{=} C_k - C_{k-1}$, where C_k was defined in (4.1), with $a_i, b_i > 0, \forall i$, $A_k = \sum_{i=0}^k a_i$, $B_k = \sum_{i=0}^k b_i$, and $B_{k-1} = H_k A_{k-1}$. Let $\mathbf{x}_k, \mathbf{y}_k, \mathbf{z}_k$ evolve according to (GMD), where $\mathbf{x}_0 = \mathbf{y}_0$ is an arbitrary initial point such that $\nabla \psi^*(\mathbf{z}_0) = \mathbf{x}_0$, $\psi: E \rightarrow \mathbb{R}$ is a μ -strongly convex function w.r.t. $\|\cdot\|$, and f is an L -smooth function w.r.t. $\|\cdot\|$. If f is ϵ_H -weakly convex for $\epsilon_H \in [0, L]$ and $\frac{a_k^2}{A_k^2} = \frac{c\mu}{LH_k}$ for $c \in [0, 1]$, then*

$$E_k \leq -(1-c)A_{k-1}D_{\psi^*}(\mathbf{z}_{k-1}, \mathbf{z}_k) + B_{k-1} \frac{\epsilon_H}{2} \|\mathbf{x}_k - \mathbf{y}_{k-1}\|^2.$$

Proof. By the definitions of E_k and \mathcal{C}_k and $D_{\psi^*}(\mathbf{z}, \mathbf{z}) = 0$, we have:

$$(4.5) \quad E_k = B_{k-1} [f(\mathbf{y}_k) - f(\mathbf{y}_{k-1})] + \sum_{i=0}^{k-1} a_i [D_{\psi^*}(\mathbf{z}_k, \mathbf{z}_i) - D_{\psi^*}(\mathbf{z}_{k-1}, \mathbf{z}_i)].$$

As f is L -smooth and ϵ_H -weakly convex, we have:

$$(4.6) \quad \begin{aligned} f(\mathbf{y}_k) - f(\mathbf{y}_{k-1}) &= f(\mathbf{y}_k) - f(\mathbf{x}_k) + f(\mathbf{x}_k) - f(\mathbf{y}_{k-1}) \\ &\leq \langle \nabla f(\mathbf{x}_k), \mathbf{y}_k - \mathbf{y}_{k-1} \rangle + \frac{L}{2} \|\mathbf{y}_k - \mathbf{x}_k\|^2 + \frac{\epsilon_H}{2} \|\mathbf{x}_k - \mathbf{y}_{k-1}\|^2 \\ &= \langle \nabla f(\mathbf{x}_k), \mathbf{y}_k - \mathbf{y}_{k-1} \rangle + \frac{L}{2} \frac{a_k^2}{A_k^2} \|\nabla \psi^*(\mathbf{z}_k) - \nabla \psi^*(\mathbf{z}_{k-1})\|^2 \\ &\quad + \frac{\epsilon_H}{2} \|\mathbf{x}_k - \mathbf{y}_{k-1}\|^2. \end{aligned}$$

By the first part of Fact 1.7 and the definition of \mathbf{z}_k , $\forall i \leq k-1$:

$$\begin{aligned} D_{\psi^*}(\mathbf{z}_{k-1}, \mathbf{z}_i) &= D_{\psi^*}(\mathbf{z}_k, \mathbf{z}_i) + \langle \nabla \psi^*(\mathbf{z}_k) - \nabla \psi^*(\mathbf{z}_i), \mathbf{z}_{k-1} - \mathbf{z}_k \rangle + D_{\psi^*}(\mathbf{z}_{k-1}, \mathbf{z}_k) \\ &= D_{\psi^*}(\mathbf{z}_k, \mathbf{z}_i) + \frac{a_k}{A_k} H_k \langle \nabla f(\mathbf{x}_k), \nabla \psi^*(\mathbf{z}_k) - \nabla \psi^*(\mathbf{z}_i) \rangle \\ &\quad + D_{\psi^*}(\mathbf{z}_{k-1}, \mathbf{z}_k). \end{aligned}$$

Hence, we have:

$$(4.7) \quad \begin{aligned} &\sum_{i=0}^{k-1} a_i [D_{\psi^*}(\mathbf{z}_k, \mathbf{z}_i) - D_{\psi^*}(\mathbf{z}_{k-1}, \mathbf{z}_i)] \\ &= -A_{k-1} D_{\psi^*}(\mathbf{z}_{k-1}, \mathbf{z}_k) - \frac{a_k}{A_k} H_k \left\langle \nabla f(\mathbf{x}_k), A_{k-1} \nabla \psi^*(\mathbf{z}_k) - \sum_{i=0}^{k-1} \nabla \psi^*(\mathbf{z}_i) \right\rangle \\ &= -A_{k-1} D_{\psi^*}(\mathbf{z}_{k-1}, \mathbf{z}_k) - a_k H_k \left\langle \nabla f(\mathbf{x}_k), \nabla \psi^*(\mathbf{z}_k) - \frac{1}{A_k} \sum_{i=0}^k \nabla \psi^*(\mathbf{z}_i) \right\rangle \\ &= -A_{k-1} D_{\psi^*}(\mathbf{z}_{k-1}, \mathbf{z}_k) - a_k H_k \langle \nabla f(\mathbf{x}_k), \nabla \psi^*(\mathbf{z}_k) - \mathbf{y}_k \rangle, \end{aligned}$$

where the last equality is by $\mathbf{y}_k = \frac{1}{A_k} \sum_{i=0}^k a_i \nabla \psi^*(\mathbf{z}_i)$, which follows by applying the definition of \mathbf{y}_k from (GMD) recursively and using $\mathbf{y}_0 = \nabla \psi^*(\mathbf{z}_0)$.

By Fact 1.7, $D_{\psi^*}(\mathbf{z}_{k-1}, \mathbf{z}_k) \geq \frac{\mu}{2} \|\nabla \psi^*(\mathbf{z}_k) - \nabla \psi^*(\mathbf{z}_{k-1})\|^2$. Hence, combining Eqs. (4.5)–(4.7):

$$\begin{aligned} E_k &\leq B_{k-1} \left\langle \nabla f(\mathbf{x}_k), \mathbf{y}_k - \mathbf{y}_{k-1} - \frac{a_k H_k}{B_{k-1}} (\nabla \psi^*(\mathbf{z}_k) - \mathbf{y}_k) \right\rangle + \frac{\epsilon_H B_{k-1}}{2} \|\mathbf{x}_k - \mathbf{y}_{k-1}\|^2 \\ &\quad + \left(B_{k-1} \frac{L}{2} \frac{a_k^2}{A_k^2} - \frac{c\mu}{2} A_{k-1} \right) \|\nabla \psi^*(\mathbf{z}_k) - \nabla \psi^*(\mathbf{z}_{k-1})\|^2 \\ &\quad - (1-c) A_{k-1} D_{\psi^*}(\mathbf{z}_{k-1}, \mathbf{z}_k) \\ &\leq \frac{\epsilon_H B_{k-1}}{2} \|\mathbf{x}_k - \mathbf{y}_{k-1}\|^2 - (1-c) A_{k-1} D_{\psi^*}(\mathbf{z}_{k-1}, \mathbf{z}_k), \end{aligned}$$

where we have used $B_{k-1} = H_k A_{k-1}$, $\mathbf{y}_k = \frac{A_{k-1}}{A_k} \mathbf{y}_{k-1} + \frac{a_k}{A_{k-1}} \nabla \psi^*(\mathbf{z}_k)$ (which implies $\mathbf{y}_k - \mathbf{y}_{k-1} - \frac{a_k H_k}{B_{k-1}} (\nabla \psi^*(\mathbf{z}_k) - \mathbf{y}_k) = 0$), and $\frac{a_k^2}{A_k^2} = c \cdot \frac{\mu}{L H_k}$. \square

Final Convergence Bound. To be able to bound the non-negative term in the discretization error E_k (which comes from ϵ_H -weak convexity), we will make use of the following proposition.

PROPOSITION 4.5. *Let $\mathbf{x}_k, \mathbf{y}_k, \mathbf{z}_k$ evolve according to (GMD), for $\mathbf{x}_0 = \mathbf{y}_0 = \nabla\psi^*(\mathbf{z}_0)$ and μ -strongly convex ψ . Then:*

$$\frac{1}{2}\|\mathbf{x}_k - \mathbf{y}_{k-1}\|^2 \leq \frac{1}{\mu} \frac{a_k^2}{A_k^2 A_{k-1}} \sum_{i=0}^{k-2} a_i D_{\psi^*}(\mathbf{z}_{k-1}, \mathbf{z}_i).$$

Proof. By applying the definition of \mathbf{y}_k recursively in (GMD), we have that $\mathbf{y}_k = \frac{1}{A_k} \sum_{i=0}^k a_i \nabla\psi^*(\mathbf{z}_i)$, $\forall k \geq 0$. Further, by the definition of \mathbf{x}_k in (GMD), we have $\mathbf{x}_k = \frac{A_{k-1}}{A_k} \mathbf{y}_{k-1} + \frac{a_k}{A_k} \nabla\psi^*(\mathbf{z}_{k-1})$, and it follows that:

$$\begin{aligned} \mathbf{x}_k - \mathbf{y}_{k-1} &= \frac{a_k}{A_k} (\nabla\psi^*(\mathbf{z}_{k-1}) - \mathbf{y}_{k-1}) \\ &= \frac{a_k}{A_k A_{k-1}} \sum_{i=0}^{k-1} a_i (\nabla\psi^*(\mathbf{z}_{k-1}) - \nabla\psi^*(\mathbf{z}_i)). \end{aligned}$$

Applying Jensen's inequality:

$$\|\mathbf{x}_k - \mathbf{y}_{k-1}\|^2 \leq \frac{a_k^2}{A_k^2 A_{k-1}} \sum_{i=0}^{k-1} a_i \|\nabla\psi^*(\mathbf{z}_{k-1}) - \nabla\psi^*(\mathbf{z}_i)\|^2.$$

The rest of the proof follows by using (by Fact 1.7), $D_{\psi^*}(\mathbf{z}, \mathbf{w}) \geq \frac{\mu}{2} \|\nabla\psi^*(\mathbf{z}) - \nabla\psi^*(\mathbf{w})\|^2$, $\forall \mathbf{z}, \mathbf{w}$. \square

Using Lemmas 4.1 and 4.4, we now show how to bound the minimum norm of the gradient, under suitable step sizes, so that Lemma 4.4 applies.

THEOREM 4.6. *Let $\mathbf{x}_k, \mathbf{y}_k, \mathbf{z}_k$ evolve according to (GMD), for arbitrary $\mathbf{x}_0 = \mathbf{y}_0 \in \mathbb{R}^n$ such that $\mathbf{x}_0 = \nabla\psi^*(\mathbf{z}_0)$, where $\psi(\mathbf{x}) = \frac{\mu}{2} \|\mathbf{x}\|^2$. Let $\frac{a_i^2}{A_i^2} = c \frac{\mu}{L H_i}$ for some $c \in [0, 1]$ and $B_{k-1} = A_{k-1} H_k$. If, for some $c' \in [0, 1]$ and all $i \leq k$: $(1 - c') \left(\frac{1}{B_{i-1}} - \frac{1}{B_i} \right) \geq \frac{c \epsilon_H}{L} \left(\frac{1}{B_i} - \frac{1}{B_k} \right)$, then:*

$$\begin{aligned} c' \sum_{i=1}^k \left[\frac{1}{B_{i-1}} - \frac{1}{B_i} \right] \sum_{j=0}^{i-1} a_j D_{\psi^*}(\mathbf{z}_i, \mathbf{z}_j) + \frac{c(1-c)}{2L} \sum_{i=1}^k \left(1 - \frac{B_{i-1}}{B_k} \right) \|\nabla f(\mathbf{x}_i)\|^2 \\ \leq f(\mathbf{x}_0) - f(\mathbf{x}^*). \end{aligned}$$

Proof. Using Lemma 4.1, we have:

$$\begin{aligned} (4.8) \quad \sum_{i=1}^k \left(\frac{1}{B_{i-1}} - \frac{1}{B_i} \right) \sum_{j=0}^{i-1} a_j D_{\psi^*}(\mathbf{z}_i, \mathbf{z}_j) - \sum_{i=1}^k \left(\frac{1}{B_{i-1}} - \frac{1}{B_k} \right) E_i \\ = f(\mathbf{y}_0) - \frac{1}{B_k} \sum_{i=0}^k b_i f(\mathbf{y}_i) \leq f(\mathbf{x}_0) - f(\mathbf{x}^*), \end{aligned}$$

as $\mathbf{y}_0 = \mathbf{x}_0$ and \mathbf{x}^* minimizes f .

To prove the theorem, it suffices to bound from below the left-hand side of Eq. (4.8). Let us first bound the discretization error. Using Proposition 4.5, we have that, $\forall i$:

$$\begin{aligned} E_i &\leq \frac{\epsilon_H B_{i-1}}{\mu} \frac{a_i^2}{A_i^2 A_{i-1}} \sum_{j=0}^{i-2} a_j D_{\psi^*}(\mathbf{z}_{i-1}, \mathbf{z}_j) - (1-c) A_{i-1} D_{\psi^*}(\mathbf{z}_{i-1}, \mathbf{z}_i) \\ &= c \frac{\epsilon_H}{L} \sum_{j=0}^{i-2} a_j D_{\psi^*}(\mathbf{z}_{i-1}, \mathbf{z}_j) - (1-c) A_{i-1} D_{\psi^*}(\mathbf{z}_{i-1}, \mathbf{z}_i). \end{aligned}$$

Therefore:

$$(4.9) \quad \begin{aligned} \sum_{i=1}^k \left(\frac{1}{B_{i-1}} - \frac{1}{B_k} \right) E_i &\leq c \frac{\epsilon_H}{L} \sum_{i=2}^k \left(\frac{1}{B_{i-1}} - \frac{1}{B_k} \right) \sum_{j=0}^{i-2} a_j D_{\psi^*}(\mathbf{z}_{i-1}, \mathbf{z}_j) \\ &\quad - (1-c) \sum_{i=1}^k \left(\frac{1}{B_{i-1}} - \frac{1}{B_k} \right) A_{i-1} D_{\psi^*}(\mathbf{z}_{i-1}, \mathbf{z}_i). \end{aligned}$$

As $D_{\psi^*}(\mathbf{z}_{i-1}, \mathbf{z}_i) = \frac{1}{2\mu} \|\mathbf{z}_i - \mathbf{z}_{i-1}\|^2 = \frac{1}{2\mu} \frac{a_i^2}{A_i^2} H_i^2 \|\nabla f(\mathbf{x}_i)\|^2$, we further have:

$$(4.10) \quad \begin{aligned} \sum_{i=1}^k \left(\frac{1}{B_{i-1}} - \frac{1}{B_k} \right) A_{i-1} D_{\psi^*}(\mathbf{z}_{i-1}, \mathbf{z}_i) &= \sum_{i=1}^k \left(\frac{1}{B_{i-1}} - \frac{1}{B_k} \right) A_{i-1} \frac{1}{2\mu} \frac{a_i^2}{A_i^2} H_i^2 \|\nabla f(\mathbf{x}_i)\|^2 \\ &= \frac{c}{2L} \sum_{i=1}^k \left(1 - \frac{B_{i-1}}{B_k} \right) \|\nabla f(\mathbf{x}_i)\|^2, \end{aligned}$$

where we have used $B_{i-1} = A_{i-1} H_i$ and $\frac{a_i^2}{A_i^2} = \frac{c}{H_i} \frac{\mu}{L}$, both from the statement of the theorem.

Combining Eqs. (4.8)–(4.10), we have:

$$\begin{aligned} f(\mathbf{x}_0) - f(\mathbf{x}^*) &\geq \sum_{i=1}^k \left[\frac{1}{B_{i-1}} - \frac{1}{B_i} - \frac{c\epsilon_H}{L} \left(\frac{1}{B_i} - \frac{1}{B_k} \right) \right] \sum_{j=0}^{i-1} a_j D_{\psi^*}(\mathbf{z}_i, \mathbf{z}_j) \\ &\quad + \frac{c(1-c)}{2L} \sum_{i=1}^k \left(1 - \frac{B_{i-1}}{B_k} \right) \|\nabla f(\mathbf{x}_i)\|^2. \end{aligned}$$

To complete the proof, it remains to use $(1-c') \left(\frac{1}{B_{i-1}} - \frac{1}{B_i} \right) \geq \frac{c\epsilon_H}{L} \left(\frac{1}{B_i} - \frac{1}{B_k} \right)$. \square

To obtain useful convergence bounds, we need to show that it is possible to satisfy the assumptions of Theorem 4.6. We start by providing examples for the case of convex objectives, and then discuss the nonconvex case.

The Convex Case. When f is convex, $\epsilon_H = 0$, and Theorem 4.6 can be applied with $c' = 0$. Further, once $c \in [0, 1]$, μ , and H_i are specified, all other parameters are set, since a_i and A_i can be computed from $\frac{a_i^2}{A_i^2} = \frac{c\mu}{H_i L}$ and $A_k = \sum_{i=0}^k a_i$, and, finally, we have that $B_i = A_{i-1} H_i$. The only restriction that Theorem 4.6 imposes is that B_i is a non-decreasing sequence.

To illustrate the results, we take $H_i = A_i^\lambda$, for $\lambda \in [0, 2]$. As mentioned before, $\lambda = 1$ corresponds to AGD+ [21] and $\lambda = 0$ corresponds to a generalization of the heavy-ball method. For this choice of H_i , we have that $\frac{a_i^2}{A_i^2} = \frac{c\mu}{A_i^\lambda L}$, or, equivalently:

$\frac{a_i^2}{A_i^{2-\lambda}} = c\frac{\mu}{L}$. It is not hard to verify (using the asymptotic formula $\sum_{i=1}^k i^p = \frac{k^{p+1}}{p+1} + \frac{k^p}{2} + O(k^{p-1})$) that a sequence $\{a_i\}_i$ that satisfies this condition will grow as:

$$(4.11) \quad a_i \propto \begin{cases} \left(\frac{c\mu}{L}\right)^{1/\lambda} i^{(2-\lambda)/\lambda}, & \text{if } \lambda > 0 \text{ and } i = \Omega\left(\frac{1}{\lambda}\right), \\ \sqrt{\frac{c\mu}{L}} \left(1 - \sqrt{\frac{c\mu}{L}}\right)^{-(i-1)}, & \text{if } \lambda = 0. \end{cases}$$

The following corollary (of Theorem 4.6) shows that any generalized momentum method (GMD) with $H_i = A_i^\lambda$, $\lambda \in [0, 2]$, converges to a point with small gradient norm at rate $1/k$.

COROLLARY 4.7. *Let $\mathbf{x}_k, \mathbf{y}_k, \mathbf{z}_k$ evolve as in (GMD), for convex f , $\psi^*(\mathbf{z}) = \frac{1}{2\mu}\|\mathbf{z}\|^2$, $0 < \mu < L$, and $\frac{a_i^2}{A_i^{2-\lambda}} = c\frac{\mu}{L}$ for some $\lambda \in [0, 2]$, $c \in (0, 1]$, and $c\frac{\mu}{L} < 1$. Then, $\forall k \geq 1$:*

$$\min_{1 \leq i \leq k} \|\nabla f(\mathbf{x}_i)\|^2 = O\left(\frac{L(f(\mathbf{x}_0) - f(\mathbf{x}^*))}{c(1-c)k + c \min\{\log k/\lambda, k\sqrt{c\mu/L}\}}\right).$$

In particular, for $c = \frac{1}{2}$: $\min_{1 \leq i \leq k} \|\nabla f(\mathbf{x}_i)\|^2 = O\left(\frac{L(f(\mathbf{x}_0) - f(\mathbf{x}^*))}{k}\right)$.

Proof. Applying Theorem 4.6, we have:

$$(4.12) \quad \begin{aligned} f(\mathbf{x}_0) - f(\mathbf{x}^*) &\geq \sum_{i=1}^k \left[\frac{1}{B_{i-1}} - \frac{1}{B_i} \right] \sum_{j=0}^{i-1} a_j D_{\psi^*}(\mathbf{z}_i, \mathbf{z}_j) \\ &+ \frac{c(1-c)}{2L} \sum_{i=1}^k \left(1 - \frac{B_{i-1}}{B_k}\right) \|\nabla f(\mathbf{x}_i)\|^2. \end{aligned}$$

Consider first the case $\lambda \in (0, 2]$. Then, from Eq. (4.11), we have $a_i \propto \left(\frac{c\mu}{L}\right)^{1/\lambda} i^{(2-\lambda)/\lambda}$ and $A_i \propto \frac{\lambda}{2} \left(\frac{c\mu}{L}\right)^{1/\lambda} i^{2/\lambda}$ for $i = \Omega\left(\frac{1}{\lambda}\right)$, in which case also

$$B_i = \Omega\left(\lambda^{1+\lambda} \left(\frac{c\mu}{L}\right)^{(1+\lambda)/\lambda} i^{\frac{2(1+\lambda)}{\lambda}}\right).$$

When $\lambda = 0$, we have: $B_i = A_{i-1} \propto \left(1 - \sqrt{\frac{c\mu}{L}}\right)^{-(i-1)}$. In either case:

$$\frac{a_{j-1}}{A_{j-1}} = \Omega\left(\min\left\{\sqrt{\frac{c\mu}{L}}, \frac{1}{\lambda_j}\right\}\right).$$

As B_i grows as a function of i at least cubically, we have that:

$$\frac{c(1-c)}{2L} \sum_{i=1}^k \left(1 - \frac{B_{i-1}}{B_k}\right) = \frac{c(1-c)}{2L} \Theta(k).$$

It remains to bound $\sum_{i=1}^k \left[\frac{1}{B_{i-1}} - \frac{1}{B_i}\right] \sum_{j=0}^{i-1} a_j D_{\psi^*}(\mathbf{z}_i, \mathbf{z}_j)$. By Proposition 4.3 and $a_j \geq a_{j-1}$,

$$\sum_{j=0}^{i-1} a_j D_{\psi^*}(\mathbf{z}_i, \mathbf{z}_j) \geq \frac{1}{2\mu} \sum_{j=1}^i \frac{a_{j-1}}{2} \frac{a_j^2}{A_j^2} H_j^2 \|\nabla f(\mathbf{x}_j)\|^2 = \frac{c}{4L} \sum_{j=1}^i \frac{a_{j-1}}{A_{j-1}} B_{j-1} \|\nabla f(\mathbf{x}_j)\|^2.$$

As $\frac{a_{j-1}}{A_{j-1}} = \Omega(\min\{\sqrt{\frac{c\mu}{L}}, \frac{1}{\lambda_j}\})$, we have that:

$$\begin{aligned} \sum_{i=1}^k \left[\frac{1}{B_{i-1}} - \frac{1}{B_i} \right] \sum_{j=0}^{i-1} a_j D_{\psi^*}(\mathbf{z}_i, \mathbf{z}_j) \\ = \Omega\left(\frac{c}{L}\right) \min_{1 \leq i \leq k} \|\nabla f(\mathbf{x}_i)\|^2 \sum_{i=1}^k \left[\frac{1}{B_{i-1}} - \frac{1}{B_i} \right] \sum_{j=1}^i \min\left\{ \sqrt{\frac{c\mu}{L}}, \frac{1}{\lambda_j} \right\} B_{j-1}. \end{aligned}$$

Finally, observe that:

$$\begin{aligned} \sum_{i=1}^k \left[\frac{1}{B_{i-1}} - \frac{1}{B_i} \right] \sum_{j=1}^i \min\left\{ \sqrt{\frac{c\mu}{L}}, \frac{1}{\lambda_j} \right\} B_{j-1} \\ = \sum_{j=1}^k \Omega\left(\min\left\{ \sqrt{\frac{c\mu}{L}}, \frac{1}{\lambda_j} \right\}\right) B_{j-1} \sum_{i=j}^k \left[\frac{1}{B_{i-1}} - \frac{1}{B_i} \right] \\ = \sum_{j=1}^k \Omega\left(\min\left\{ \sqrt{\frac{c\mu}{L}}, \frac{1}{\lambda_j} \right\}\right) \left[1 - \frac{B_{j-1}}{B_k} \right] \\ = \Omega\left(\min\left\{ \log(k)/\lambda, k\sqrt{c\mu/L} \right\}\right). \end{aligned}$$

Hence, we obtain:

$$\min_{1 \leq i \leq k} \|\nabla f(\mathbf{x}_i)\|^2 = O\left(\frac{L(f(\mathbf{x}_0) - f(\mathbf{x}^*))}{c(1-c)k + c \min\{\log k/\lambda, k\sqrt{c\mu/L}\}}\right),$$

as claimed. \square

A few remarks are in order here. It is not hard to see that for an *arbitrary* positive sequence of numbers a_i, H_i, B_i that satisfy $\frac{a_i^2}{A_i^2} = c\frac{\mu}{LH_i}$, $A_i = \sum_{j=0}^i a_j$, and $B_{i-1} = A_{i-1}H_i$, it is not possible to get better than a $1/k$ rate for convergence to stationary points, as long as Proposition 4.3 is used. This rate is known to be suboptimal—the optimal rate for smooth convex functions is $1/k^2$ and it is achieved by the OGM-G algorithm from [37]. The rate $1/k$ for the generalized momentum methods is not surprising, and in the case of $\lambda = 1$ ($H_i = A_i$, in which case the method is essentially equivalent to Nesterov’s accelerated method in Euclidean spaces), this rate is known to be tight [36]. We expect that the same is true for an arbitrary (but fixed) value of λ , though this may be possible to show only numerically [36].

The Nonconvex Case. The main restriction for obtaining the results in the non-convex case is ensuring that:

$$(1 - c') \left(\frac{1}{B_{i-1}} - \frac{1}{B_i} \right) \geq \frac{c\epsilon_H}{L} \left(\frac{1}{B_i} - \frac{1}{B_k} \right),$$

for some c, c' . Let us first study what kind of a constraint on the sequence B_i such a condition imposes. Rearranging the terms in the last expression, we have that:

$$(4.13) \quad \frac{B_i}{B_{i-1}} \geq 1 + \frac{c\epsilon_H}{(1-c')L} \left(1 - \frac{B_i}{B_k} \right).$$

Since the last expression needs to be satisfied for all i , when B_i grows polynomially with i , it is not hard to verify that to ensure $c' \geq 0$, we would need to have $\frac{c\epsilon_H}{L} = O(\frac{1}{i})$,

$\forall i$; that is, for a fixed number of iterations k , we would need $\frac{c\epsilon_H}{L} = O(\frac{1}{k})$. This leads to uninformative convergence bounds, unless $\epsilon_H = O(L/k)$ (in which case one can show that $\min_{1 \leq i \leq k} \|\nabla f(\mathbf{x}_i)\|^2 = O(\frac{L(f(\mathbf{x}_0) - f(\mathbf{x}^*))}{c(1-c)k})$).

When $\epsilon_H = \omega(L/k)$, we are unable to show the convergence rate of $1/k$ for polynomially growing sequences a_i (and, consequently, polynomially growing A_i, B_i). Instead, we can only show such a convergence rate for constant H_i . In particular, let us choose μ and H_i so that $\frac{\mu}{LH_i} = 1$. Then $\frac{a_i}{A_i} = \sqrt{c}$, and it follows that $\frac{B_i}{B_{i-1}} = (1 - \sqrt{c})^{-1}$. As $\epsilon_H \leq L$ and $1 - \frac{B_i}{B_k} \leq 1$, to satisfy the condition from Eq. (4.13), it suffices to have $\sqrt{c} \geq \frac{c}{1-c'}$. Equivalently, Eq. (4.13) is satisfied with:

$$c' \leq 1 - \sqrt{c}.$$

By the same arguments as in the proof of Corollary 4.7, we immediately have:

COROLLARY 4.8. *Let $\mathbf{x}_k, \mathbf{y}_k, \mathbf{z}_k$ evolve according to (GMD), for $\psi^*(\mathbf{z}) = \frac{1}{2\mu} \|\mathbf{z}\|^2$, $\mu > 0$, $\frac{\mu}{LH_i} = 1$, $\frac{a_i^2}{A_i^2} = c$, and $c \in (0, 1)$. Then, $\forall k \geq 1$:*

$$\min_{1 \leq i \leq k} \|\nabla f(\mathbf{x}_i)\|^2 = O\left(\frac{L(f(\mathbf{x}_0) - f(\mathbf{x}^*))}{c(1-c)k + \sqrt{c}(1-\sqrt{c})ck}\right).$$

In particular, for $c = \frac{1}{2}$: $\min_{1 \leq i \leq k} \|\nabla f(\mathbf{x}_i)\|^2 = O(\frac{L(f(\mathbf{x}_0) - f(\mathbf{x}^))}{k})$.*

4.3. Convergence to Stationary Points in non-Euclidean Spaces. We now show that it is possible to obtain results for convergence to stationary points even in non-Euclidean spaces $(E, \|\cdot\|)$. We are only able to show such a result, however, for a different discretization of (MoD). To obtain the result, we require only that $\psi(\cdot)$ is μ -strongly convex with respect to the norm $\|\cdot\|$. By Fact 1.6, this implies that ψ^* is $\frac{1}{\mu}$ -smooth with respect to the dual norm $\|\cdot\|_*$.

The alternative discretization uses a gradient descent step for \mathbf{y}_k to ensure a decrease in \mathcal{C}_k that depends on $\|\nabla f(\mathbf{y}_k)\|_*^2$. However, to ensure the right change in \mathcal{C}_k over iterations k , such a choice of \mathbf{y}_k requires changes to the extrapolation step \mathbf{x}_k . In particular, the discrete-time algorithm is:

$$\begin{aligned} \mathbf{x}_k &= \mathbf{y}_{k-1} + \frac{a_k}{A_k} \nabla \psi^*(\mathbf{z}_{k-1}) - \frac{a_k}{A_k A_{k-1}} \sum_{i=0}^{k-1} a_i \nabla \psi^*(\mathbf{z}_i), \\ (\text{GMD}_B) \quad \mathbf{z}_k &= \mathbf{z}_{k-1} - \frac{a_k}{A_k} H_k \nabla f(\mathbf{x}_k), \\ \mathbf{y}_k &= \underset{\mathbf{u}}{\operatorname{argmin}} \left\{ \langle \nabla f(\mathbf{x}_k), \mathbf{u} - \mathbf{x}_k \rangle + \frac{L}{2} \|\mathbf{u} - \mathbf{x}_k\|^2 \right\}. \end{aligned}$$

Note that in Euclidean spaces, when $\psi^*(\mathbf{z}) = \frac{\mu}{2} \|\mathbf{z}\|^2$, (GMD_B) is equivalent to (GMD). Hence, (GMD_B) can be seen as a generalization of (GMD) to non-Euclidean spaces.

Discretization Error. As in the previous section, we start by bounding the discretization error $E_k \stackrel{\text{def}}{=} \mathcal{C}_k - \mathcal{C}_{k-1}$.

LEMMA 4.9. *Let $E_k \stackrel{\text{def}}{=} \mathcal{C}_k - \mathcal{C}_{k-1}$, where \mathcal{C}_k was defined in (4.1), with $a_i, b_i > 0, \forall i$, $A_k = \sum_{i=0}^k a_i$, $B_k = \sum_{i=0}^k b_i$, and $B_{k-1} = H_k A_{k-1}$. Let $\mathbf{x}_k, \mathbf{y}_k, \mathbf{z}_k$ evolve according to (GMD_B), where $\mathbf{x}_0 = \mathbf{y}_0 \in E$ is an arbitrary initial point such that $\nabla \psi^*(\mathbf{z}_0) = \mathbf{x}_0$, $\psi : E \rightarrow \mathbb{R}$ is a μ -strongly convex function w.r.t. $\|\cdot\|$, and f is an L -smooth function w.r.t. $\|\cdot\|$. If f is ϵ_H -weakly convex for $\epsilon_H \in [0, L]$ and $\frac{a_k^2}{A_k^2} = \frac{c\mu}{LH_k}$ for $c \in [0, 1]$, then*

$$E_k \leq -(1-c) \frac{B_{k-1}}{2L} \|\nabla f(\mathbf{x}_k)\|_*^2 + B_{k-1} \frac{\epsilon_H}{2} \|\mathbf{x}_k - \mathbf{y}_{k-1}\|^2.$$

Proof. The proof uses similar arguments as the proof of Lemma 4.4. By the definitions of E_k and \mathcal{C}_k and $D_{\psi^*}(\mathbf{z}, \mathbf{z}) = 0$, we have:

$$(4.14) \quad E_k = B_{k-1} [f(\mathbf{y}_k) - f(\mathbf{y}_{k-1})] + \sum_{i=0}^{k-1} a_i [D_{\psi^*}(\mathbf{z}_k, \mathbf{z}_i) - D_{\psi^*}(\mathbf{z}_{k-1}, \mathbf{z}_i)].$$

As f is L -smooth and ϵ_H -weakly convex, we have:

$$(4.15) \quad \begin{aligned} f(\mathbf{y}_k) - f(\mathbf{y}_{k-1}) &= f(\mathbf{y}_k) - f(\mathbf{x}_k) + f(\mathbf{x}_k) - f(\mathbf{y}_{k-1}) \\ &\leq \langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{y}_{k-1} \rangle - \frac{1}{2L} \|\nabla f(\mathbf{x}_k)\|_*^2 + \frac{\epsilon_H}{2} \|\mathbf{x}_k - \mathbf{y}_{k-1}\|^2. \end{aligned}$$

By the first part of Fact 1.7 and the definition of \mathbf{z}_k , $\forall i \leq k-1$:

$$\begin{aligned} D_{\psi^*}(\mathbf{z}_k, \mathbf{z}_i) &= D_{\psi^*}(\mathbf{z}_{k-1}, \mathbf{z}_i) + \langle \nabla \psi^*(\mathbf{z}_{k-1}) - \nabla \psi^*(\mathbf{z}_i), \mathbf{z}_k - \mathbf{z}_{k-1} \rangle + D_{\psi^*}(\mathbf{z}_k, \mathbf{z}_{k-1}) \\ &= D_{\psi^*}(\mathbf{z}_{k-1}, \mathbf{z}_i) - \frac{a_k}{A_k} H_k \langle \nabla f(\mathbf{x}_k), \nabla \psi^*(\mathbf{z}_{k-1}) - \nabla \psi^*(\mathbf{z}_i) \rangle \\ &\quad + D_{\psi^*}(\mathbf{z}_k, \mathbf{z}_{k-1}). \end{aligned}$$

Hence, we have:

$$(4.16) \quad \begin{aligned} &\sum_{i=0}^{k-1} a_i [D_{\psi^*}(\mathbf{z}_k, \mathbf{z}_i) - D_{\psi^*}(\mathbf{z}_{k-1}, \mathbf{z}_i)] \\ &= A_{k-1} D_{\psi^*}(\mathbf{z}_k, \mathbf{z}_{k-1}) - \frac{a_k}{A_k} H_k \left\langle \nabla f(\mathbf{x}_k), A_{k-1} \nabla \psi^*(\mathbf{z}_{k-1}) - \sum_{i=0}^{k-1} \nabla \psi^*(\mathbf{z}_i) \right\rangle \\ &= A_{k-1} D_{\psi^*}(\mathbf{z}_k, \mathbf{z}_{k-1}) - B_{k-1} \langle \nabla f(\mathbf{x}_k), \mathbf{x}_k - \mathbf{y}_{k-1} \rangle, \end{aligned}$$

where the last line is by $B_k = A_{k-1} H_k$ and the definition of \mathbf{x}_k , which can implies

$$\mathbf{x}_k - \mathbf{y}_{k-1} = \frac{a_k}{A_k A_{k-1}} (A_{k-1} \nabla \psi^*(\mathbf{z}_{k-1}) - \sum_{i=0}^{k-1} a_i \nabla \psi^*(\mathbf{z}_i)).$$

By Fact 1.6, $D_{\psi^*}(\mathbf{z}_k, \mathbf{z}_{k-1}) \leq \frac{1}{2\mu} \|\mathbf{z}_k - \mathbf{z}_{k-1}\|_*^2 = \frac{1}{2\mu} \frac{a_k^2 H_k^2}{A_k^2} \|\nabla f(\mathbf{x}_k)\|_*^2$. Hence, combining Eqs.(4.14)–(4.16):

$$\begin{aligned} E_k &\leq B_{k-1} \left(-\frac{1}{2L} + \frac{a_k^2}{2\mu A_k^2} H_k \right) \|\nabla f(\mathbf{x}_k)\|_*^2 + B_{k-1} \frac{\epsilon_H}{2} \|\mathbf{x}_k - \mathbf{y}_{k-1}\|^2 \\ &= -\frac{(1-c)B_{k-1}}{2L} \|\nabla f(\mathbf{x}_k)\|_*^2 + B_{k-1} \frac{\epsilon_H}{2} \|\mathbf{x}_k - \mathbf{y}_{k-1}\|^2, \end{aligned}$$

where we have used that $B_{k-1} = H_k A_{k-1}$ and $\frac{a_k^2}{A_k^2} = c \cdot \frac{\mu}{L H_k}$. \square

Final Convergence Bound. Since the proof of Proposition 4.5 only required that ψ be μ -strongly convex and that $\mathbf{x}_k - \mathbf{y}_{k-1} = \frac{a_k}{A_k} (\nabla \psi^*(\mathbf{z}_{k-1}) - \frac{1}{A_{k-1}} \sum_{i=0}^{k-1} a_i \nabla \psi^*(\mathbf{z}_i))$, the same claim holds for the iterates of (GMD_B). Thus, we can draw a similar conclusion as for (GMD).

THEOREM 4.10. *Let $\mathbf{x}_k, \mathbf{y}_k, \mathbf{z}_k$ evolve according to (GMD_B), for arbitrary $\mathbf{x}_0 = \mathbf{y}_0 \in E$ such that $\mathbf{x}_0 = \nabla \psi^*(\mathbf{z}_0)$, where $\psi(\mathbf{x})$ is strongly convex w.r.t. $\|\cdot\|$. Let*

$\frac{a_i^2}{A_i^2} = c \frac{\mu}{LH_i}$ for some $c \in [0, 1]$ and $B_{k-1} = A_{k-1}H_k$. If for some $c' \in [0, 1]$ and all $i \leq k$: $(1 - c')\left(\frac{1}{B_{i-1}} - \frac{1}{B_i}\right) \geq \frac{cc\mu}{L}\left(\frac{1}{B_i} - \frac{1}{B_k}\right)$, then:

$$\begin{aligned} c' \sum_{i=1}^k \left[\frac{1}{B_{i-1}} - \frac{1}{B_i} \right] \sum_{j=0}^{i-1} a_j D_{\psi^*}(\mathbf{z}_i, \mathbf{z}_j) + \frac{c(1-c)}{2L} \sum_{i=1}^k \left(1 - \frac{B_{i-1}}{B_k}\right) \|\nabla f(\mathbf{x}_i)\|_*^2 \\ \leq f(\mathbf{x}_0) - f(\mathbf{x}^*). \end{aligned}$$

The proof is the same as the proof of Theorem 4.6 and is thus omitted.

The main difference between (GMD) and (GMD_B) in terms of the conclusions about the convergence to stationary points is that, because we are no longer assuming strong convexity of ψ^* , we can no longer bound $\sum_{j=0}^{i-1} a_j D_{\psi^*}(\mathbf{z}_i, \mathbf{z}_j)$ below as a function of the norms of the gradients. However, the term $\frac{c(1-c)}{2L} \sum_{i=1}^k \left(1 - \frac{B_{i-1}}{B_k}\right) \|\nabla f(\mathbf{x}_i)\|_*^2$ from the theorem statement is still sufficient for obtaining $1/k$ asymptotic convergence, as shown in the following corollary.

COROLLARY 4.11. *Let $\mathbf{x}_k, \mathbf{y}_k, \mathbf{z}_k$ evolve as in (GMD_B), for some μ -strongly convex ψ , where $\mu > 0$, and where $\frac{a_i^2}{A_i^2} = c \frac{\mu}{LH_i}$, $c \in (0, 1)$.*

(i) *If $\frac{\mu}{LH_i} = 1$, then $\forall k \geq 1$: $\min_{1 \leq i \leq k} \|\nabla f(\mathbf{x}_i)\|_*^2 = O\left(\frac{L(f(\mathbf{x}_0) - f(\mathbf{x}^*))}{c(1-c)k}\right)$.*

In particular, for $c = \frac{1}{2}$: $\min_{1 \leq i \leq k} \|\nabla f(\mathbf{x}_i)\|_^2 = O\left(\frac{L(f(\mathbf{x}_0) - f(\mathbf{x}^*))}{k}\right)$.*

(ii) *If f is convex and $H_i = A_i^\lambda$, then: $\min_{1 \leq i \leq k} \|\nabla f(\mathbf{x}_i)\|_*^2 = O\left(\frac{L(f(\mathbf{x}_0) - f(\mathbf{x}^*))}{c(1-c)k}\right)$.*

In particular, for $c = \frac{1}{2}$: $\min_{1 \leq i \leq k} \|\nabla f(\mathbf{x}_i)\|_^2 = O\left(\frac{L(f(\mathbf{x}_0) - f(\mathbf{x}^*))}{k}\right)$.*

Proof. The first part of the corollary follows because under the assumption that $\frac{\mu L}{H_i} = 1$, the condition of Theorem 4.10 can be satisfied with $c' = 1 - \sqrt{c} \geq 0$, as discussed in the previous subsection. Further, in this case B_i grows exponentially fast and $\sum_{i=1}^k \left(1 - \frac{B_{i-1}}{B_k}\right) = \Omega(k)$. As B_i is increasing and Bregman divergences are non-negative, we have:

$$\frac{c(1-c)}{2L} \sum_{i=1}^k \left(1 - \frac{B_{i-1}}{B_k}\right) \|\nabla f(\mathbf{x}_i)\|_*^2 \leq f(\mathbf{x}_0) - f(\mathbf{x}^*),$$

which implies the claimed statement.

For the second part of the corollary, convexity of f implies that the condition from Theorem 4.10 can be satisfied with $c' = 1$. As discussed in the proof of Corollary 4.7, for $H_i = A_i^\lambda$, B_i grows at least cubically in i , which, again, implies that $\sum_{i=1}^k \left(1 - \frac{B_{i-1}}{B_k}\right) = \Omega(k)$, and leads to the same conclusion as in the first part of the proof. \square

5. Conclusion. We presented a generic Hamiltonian-based framework for the analysis of general momentum methods in non-Euclidean spaces and in the settings of both convex and nonconvex optimization. Several questions that merit further investigation remain. For example, while convergence to stationary points in Euclidean spaces is well-understood [18, 37], much less is known in terms of both upper and lower bounds in general non-Euclidean spaces. Another interesting direction for future research is a rigorous characterization of the use of momentum to escape shallow local minima. Finally, it is worth noting that, even though it has led to intuitive interpretations and fruitful results in classical and more recent optimization literature, the continuous-time perspective in optimization is not a panacea and certain phenomena are crucially discrete (see [3] for a stimulating discussion on this topic). Thus, it is interesting to also understand the limitations of the continuous-time perspective.

Acknowledgements. We thank Guilherme França, Michael Muehlebach, and Uri Ascher for useful comments and suggestions.

REFERENCES

- [1] Z. ALLEN-ZHU, *Natasha: Faster non-convex stochastic optimization via strongly non-convex parameter*, in Proc. ICML'17, 2017.
- [2] Z. ALLEN-ZHU AND L. ORECCHIA, *Linear coupling: An ultimate unification of gradient and mirror descent*, in Proc. ITCS'17, 2017.
- [3] U. M. ASCHER, *Discrete processes and their continuous limits*, arXiv preprint arXiv:1910.02098, (2019).
- [4] H. ATTOUCH AND F. ALVAREZ, *The heavy ball with friction dynamical system for convex constrained minimization problems*, in Optimization, Springer, 2000, pp. 25–35.
- [5] H. ATTOUCH, A. CABOT, AND P. REDONT, *The dynamics of elastic shocks via epigraphical regularization of a differential inclusion. Barrier and penalty approximations*, Advances in Mathematical Sciences and Applications, 12 (2002), pp. 273–306.
- [6] H. ATTOUCH, Z. CHBANI, J. FADILI, AND H. RIAHI, *First-order optimization algorithms via inertial systems with Hessian driven damping*, arXiv preprint arXiv:1907.10536, (2019).
- [7] H. ATTOUCH, Z. CHBANI, AND H. RIAHI, *Rate of convergence of the Nesterov accelerated gradient method in the subcritical case $\alpha \leq 3$* , ESAIM Contr. Optim. Ca., 25 (2019), p. 2.
- [8] 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 (2000), pp. 1–34.
- [9] H. ATTOUCH, P.-E. MAINGÉ, AND P. REDONT, *A second-order differential system with Hessian-driven damping; Application to non-elastic shock laws*, Differential Equations and Applications, 4 (2012), pp. 27–65.
- [10] H. ATTOUCH AND J. PEYPOUQUET, *Convergence rate of proximal inertial algorithms associated with Moreau envelopes of convex functions*, in Splitting Algorithms, Modern Operator Theory, and Applications, Springer, 2019, pp. 1–44.
- [11] P. BÉGOUT, J. BOLTE, AND M. A. JENDOUBI, *On damped second-order gradient systems*, Journal of Differential Equations, 259 (2015), pp. 3115–3143.
- [12] D. P. BERTSEKAS, *Control of uncertain systems with a set-membership description of the uncertainty.*, PhD thesis, MIT, 1971.
- [13] M. BETANCOURT, M. I. JORDAN, AND A. C. WILSON, *On symplectic optimization*, arXiv preprint arXiv:1802.03653, (2018).
- [14] J. BOLTE, A. DANILIDIS, O. LEY, AND L. MAZET, *Characterizations of Łojasiewicz inequalities: Subgradient flows, talweg, convexity*, Transactions of the American Mathematical Society, 362 (2010), pp. 3319–3363.
- [15] E. BOSTAN, M. SOLTANOLKOTABI, D. REN, AND L. WALLER, *Accelerated Wirtinger flow for multiplexed Fourier ptychographic microscopy*, in Proc. IEEE ICIP'18, 2018.
- [16] S. BUBECK, Y. T. LEE, AND M. SINGH, *A geometric alternative to Nesterov's accelerated gradient descent*, arXiv preprint, arXiv:1506.08187, (2015).
- [17] A. CABOT, H. ENGLER, AND S. GADAT, *On the long time behavior of second order differential equations with asymptotically small dissipation*, Transactions of the American Mathematical Society, 361 (2009), pp. 5983–6017.
- [18] Y. CARMON, J. C. DUCHI, O. HINDER, AND A. SIDFORD, *Lower bounds for finding stationary points i*, Mathematical Programming, (2019).
- [19] C. CASTERA, J. BOLTE, C. FÉVOTTE, AND E. PAUWELS, *An inertial Newton algorithm for deep learning*, arXiv preprint arXiv:1905.12278, (2019).
- [20] A. CHAMBOLE AND C. DOSSAL, *On the convergence of the iterates of the “fast iterative shrinkage/thresholding algorithm”*, Journal of Optimization theory and Applications, 166 (2015), pp. 968–982.
- [21] M. B. COHEN, J. DIAKONIKOLAS, AND L. ORECCHIA, *On acceleration with noise-corrupted gradients*, in Proc. ICML'18, 2018.
- [22] D. DAVIS AND B. GRIMMER, *Proximally guided stochastic subgradient method for nonsmooth, nonconvex problems*, SIAM Journal on Optimization, 29 (2019), pp. 1908–1930.
- [23] J. DIAKONIKOLAS AND C. GUZMÁN, *Lower bounds for parallel and randomized convex optimization.*, Journal of Machine Learning Research, 21 (2020), pp. 1–31.
- [24] J. DIAKONIKOLAS AND L. ORECCHIA, *Accelerated extra-gradient descent: A novel, accelerated first-order method*, in Proc. ITCS'18, 2018.

- [25] J. DIAKONIKOLAS AND L. ORECCHIA, *The approximate duality gap technique: A unified theory of first-order methods*, SIAM J. Optimiz., 29 (2019), pp. 660–689.
- [26] D. DRUSVYATSKIY, M. FAZEL, AND S. ROY, *An optimal first order method based on optimal quadratic averaging*, SIAM J. Optimiz., 28 (2018), pp. 251–271.
- [27] G. FRANÇA, J. SULAM, D. P. ROBINSON, AND R. VIDAL, *Conformal symplectic and relativistic optimization*, arXiv preprint arXiv:1903.04100, (2019).
- [28] A. V. GASNIKOV AND Y. E. NESTEROV, *Universal method for stochastic composite optimization problems*, Comput. Math. & Math. Phys., 58 (2018), pp. 48–64.
- [29] E. GHADIMI, H. R. FEYZMAHDAVIAN, AND M. JOHANSSON, *Global convergence of the heavy-ball method for convex optimization*, in Proc. IEEE ECCV'15, 2015.
- [30] S. GHADIMI AND G. LAN, *Accelerated gradient methods for nonconvex nonlinear and stochastic programming*, Math. Program., 156 (2016), pp. 59–99.
- [31] S. GHADIMI, G. LAN, AND H. ZHANG, *Generalized uniformly optimal methods for nonlinear programming*, Journal of Scientific Computing, 79 (2019), pp. 1854–1881.
- [32] O. GÜLER, *New proximal point algorithms for convex minimization*, SIAM Journal on Optimization, 2 (1992), pp. 649–664.
- [33] B. HU AND L. LESSARD, *Control interpretations for first-order optimization methods*, in Proc. IEEE ACC'17, 2017.
- [34] C. JIN, P. NETRAPALLI, AND M. I. JORDAN, *Accelerated gradient descent escapes saddle points faster than gradient descent*, in Proc. COLT'18, 2018.
- [35] J. A. KELNER, Y. T. LEE, L. ORECCHIA, AND A. SIDFORD, *An almost-linear-time algorithm for approximate max flow in undirected graphs, and its multicommodity generalizations*, in Proc. ACM-SIAM SODA'14, 2014.
- [36] D. KIM AND J. A. FESSLER, *Generalizing the optimized gradient method for smooth convex minimization*, SIAM J. Optimiz., 28 (2018), pp. 1920–1950.
- [37] D. KIM AND J. A. FESSLER, *Optimizing the efficiency of first-order methods for decreasing the gradient of smooth convex functions*, arXiv preprint, arXiv:1803.06600, (2018).
- [38] W. KRICHENE, A. BAYEN, AND P. L. BARTLETT, *Accelerated mirror descent in continuous and discrete time*, in Proc. NIPS'15, 2015.
- [39] R. LARAKI AND P. MERTIKOPOULOS, *Inertial game dynamics and applications to constrained optimization*, SIAM Journal on Control and Optimization, 53 (2015), pp. 3141–3170.
- [40] Y. T. LEE, S. RAO, AND N. SRIVASTAVA, *A new approach to computing maximum flows using electrical flows*, in Proc. ACM STOC'13, 2013.
- [41] L. LESSARD AND P. SEILER, *Direct synthesis of iterative algorithms with bounds on achievable worst-case convergence rate*, arXiv preprint arXiv:1904.09046, (2019).
- [42] H. LIN, J. MAIRAL, AND Z. HARCHAOUI, *A universal catalyst for first-order optimization*, in Proc. NIPS'15, 2015.
- [43] C. J. MADDISON, D. PAULIN, Y. W. TEH, B. O'DONOGHUE, AND A. DOUCET, *Hamiltonian descent methods*, arXiv preprint arXiv:1809.05042, (2018).
- [44] M. MUEHLEBACH AND M. JORDAN, *A dynamical systems perspective on Nesterov acceleration*, in Proc. ICML'19, 2019.
- [45] A. S. NEMIROVSKI AND Y. E. NESTEROV, *Optimal methods of smooth convex minimization*, Zh. Vychisl. Mat. i Mat. Fiz., 25 (1985), pp. 356–369.
- [46] A. NEMIROVSKI AND D. B. YUDIN, *Problem complexity and method efficiency in optimization*, Wiley, 1983.
- [47] Y. NESTEROV, *A method of solving a convex programming problem with convergence rate $O(1/k^2)$* , in Doklady AN SSSR (translated as Soviet Mathematics Doklady), vol. 269, 1983, pp. 543–547.
- [48] Y. NESTEROV, *How to make the gradients small*, Optima, 88 (2012), pp. 10–11.
- [49] Y. NESTEROV, *Lectures on Convex Optimization*, Springer, 2018.
- [50] Y. NESTEROV, A. GASNIKOV, S. GUMINOV, AND P. DVURECHENSKY, *Primal-dual accelerated gradient descent with line search for convex and nonconvex optimization problems*, arXiv preprint, arXiv:1809.05895, (2018).
- [51] M. O'NEILL AND S. J. WRIGHT, *Behavior of accelerated gradient methods near critical points of nonconvex functions*, Math. Program., (2017), pp. 1–25.
- [52] B. T. POLYAK, *Some methods of speeding up the convergence of iteration methods*, USSR Comput. Math. & Math. Phys., 4 (1964), pp. 1–17.
- [53] D. SCIEUR, V. ROULET, F. BACH, AND A. D'ASPREMONT, *Integration methods and accelerated optimization algorithms*, in Proc. NIPS'17, 2017.
- [54] J. SHERMAN, *Area-convexity, l_∞ regularization, and undirected multicommodity flow*, in Proc. ACM STOC'17, 2017.
- [55] B. SHI, S. S. DU, M. I. JORDAN, AND W. J. SU, *Understanding the acceleration phenomenon*

- via high-resolution differential equations, arXiv preprint, arXiv:1810.08907, (2018).
- [56] W. SU, S. BOYD, AND E. J. CANDÉS, *A differential equation for modeling Nesterov’s accelerated gradient method: Theory and insights*, J. Mach. Learn. Res., 17 (2016), pp. 1–43.
- [57] P. TSENG, *On accelerated proximal gradient methods for convex-concave optimization*, 2008.
- [58] A. WIBISONO, A. C. WILSON, AND M. I. JORDAN, *A variational perspective on accelerated methods in optimization*, in Proc. Natl. Acad. Sci. U.S.A., 2016.
- [59] A. C. WILSON, B. RECHT, AND M. I. JORDAN, *A Lyapunov analysis of momentum methods in optimization*, arXiv preprint, arXiv:1611.02635, (2016).
- [60] R. XU, M. SOLTANOLKOTABI, J. P. HALDAR, W. UNGLAUB, J. ZUSMAN, A. F. LEVI, AND R. M. LEAHY, *Accelerated Wirtinger flow: A fast algorithm for ptychography*, arXiv preprint arXiv:1806.05546, (2018).
- [61] J. ZHANG, A. MOKHTARI, S. SRA, AND A. JADBABAIE, *Direct Runge-Kutta discretization achieves acceleration*, in Proc. NeurIPS’18, 2018.

Appendix A. Proof of Lemma 2.3. The simplest way of proving the lemma is by directly computing $\frac{d}{dt}\mathcal{C}_t^f$ and $\frac{d}{dt}\mathcal{C}_t$, and showing that (MoD) implies that both are equal to zero. Here, we provide a longer, but more constructive proof that highlights how \mathcal{C}_t^f , \mathcal{C}_t arise as invariants of (2.5).

As $\bar{\mathbf{x}}_t$, \mathbf{z}_t evolve according to the equations of motion of (2.5), we have that $\frac{d}{dt}\mathcal{H}_M(\bar{\mathbf{x}}_t, \mathbf{z}_t, \alpha_t) = \frac{\partial}{\partial t}\mathcal{H}_M(\bar{\mathbf{x}}_t, \mathbf{z}_t, \alpha_t) = \frac{d}{d\alpha_t}\mathcal{H}_M(\bar{\mathbf{x}}_t, \mathbf{z}_t, \alpha_t) \cdot \dot{\alpha}_t$. Observe that:

$$\begin{aligned} \frac{d}{d\alpha_t}\mathcal{H}_M(\bar{\mathbf{x}}_t, \mathbf{z}_t, \alpha_t) &= h'(\alpha_t)f(\bar{\mathbf{x}}_t/\alpha_t) - h(\alpha_t) \left\langle \nabla f(\bar{\mathbf{x}}_t/\alpha_t), \frac{\bar{\mathbf{x}}_t}{\alpha_t^2} \right\rangle \\ &= h'(\alpha_t)f(\mathbf{x}_t) - \frac{h(\alpha_t)}{\alpha_t} \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t \rangle. \end{aligned}$$

Hence, we have that:

$$\frac{d}{dt}(h(\alpha_t)f(\mathbf{x}_t) + \psi^*(\mathbf{z}_t)) = \frac{dh(\alpha_t)}{dt}f(\mathbf{x}_t) - h(\alpha_t)\frac{\dot{\alpha}_t}{\alpha_t} \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t \rangle.$$

Equivalently, using the product rule of differentiation:

$$(A.1) \quad h(\alpha_t)\frac{d}{dt}f(\mathbf{x}_t) + \frac{d}{dt}\psi^*(\mathbf{z}_t) = -h(\alpha_t)\frac{\dot{\alpha}_t}{\alpha_t} \langle \nabla f(\mathbf{x}_t), \mathbf{x}_t \rangle.$$

Integrating both sides of (A.1) from 0 to t and using integration by parts leads to \mathcal{C}_t^f .

To obtain \mathcal{C}_t , observe (from (MoD)) that $\dot{\mathbf{z}}_t = -h(\alpha_t)\frac{\dot{\alpha}_t}{\alpha_t}\nabla f(\mathbf{x}_t)$. Multiplying both sides of (A.1) by α_t , we thus have:

$$(A.2) \quad \alpha_t h(\alpha_t)\frac{d}{dt}f(\mathbf{x}_t) + \alpha_t \frac{d}{dt}\psi^*(\mathbf{z}_t) = \alpha_t \langle \dot{\mathbf{z}}_t, \mathbf{x}_t \rangle.$$

As for \mathcal{C}_t^f , to obtain \mathcal{C}_t , we integrate both sides of the last equation from 0 to t . Integrating the left-hand side and applying integration by parts gives:

$$(A.3) \quad \begin{aligned} \int_0^t \left(\alpha_\tau h(\alpha_\tau)\frac{d}{d\tau}f(\mathbf{x}_\tau) + \alpha_\tau \frac{d}{d\tau}\psi^*(\mathbf{z}_\tau) \right) d\tau \\ = h(\alpha_t)\alpha_t f(\mathbf{x}_t) - h(\alpha_0)\alpha_0 f(\mathbf{x}_0) - \int_0^t \frac{d(h(\alpha_\tau)\alpha_\tau)}{d\tau} f(\mathbf{x}_\tau) d\tau \\ + \alpha_t \psi^*(\mathbf{z}_t) - \alpha_0 \psi^*(\mathbf{z}_0) - \int_0^t \dot{\alpha}_\tau \psi^*(\mathbf{z}_\tau) d\tau. \end{aligned}$$

On the other hand, by the definition of \mathbf{x}_t in (MoD), $\mathbf{x}_t = \frac{\alpha_0}{\alpha_t} \mathbf{x}_0 + \frac{1}{\alpha_t} \int_0^t \dot{\alpha}_\sigma \nabla \psi^*(\mathbf{z}_\sigma) d\sigma$. Thus, integrating the right-hand side of (A.2), we have:

$$(A.4) \quad \int_0^t \alpha_\tau \langle \dot{\mathbf{z}}_\tau, \mathbf{x}_\tau \rangle d\tau = \int_0^t \left\langle \dot{\mathbf{z}}_\tau, \alpha_0 \mathbf{x}_0 + \int_0^\tau \nabla \psi^*(\mathbf{z}_\sigma) \dot{\alpha}_\sigma d\sigma \right\rangle d\tau \\ = \alpha_0 \langle \mathbf{z}_t - \mathbf{z}_0, \nabla \psi^*(\mathbf{z}_0) \rangle + \int_0^t \int_0^\tau \langle \dot{\mathbf{z}}_\tau, \nabla \psi^*(\mathbf{z}_\sigma) \rangle \dot{\alpha}_\sigma d\sigma d\tau,$$

where we have used $\mathbf{x}_0 = \nabla \psi^*(\mathbf{z}_0)$. By elementary calculus, it is possible to exchange the order of integration on the right-hand side of (A.4), which leads to:

$$(A.5) \quad \int_0^t \alpha_\tau \langle \dot{\mathbf{z}}_\tau, \mathbf{x}_\tau \rangle d\tau = \int_0^t \left\langle \dot{\mathbf{z}}_\tau, \alpha_0 \mathbf{x}_0 + \int_0^\tau \nabla \psi^*(\mathbf{z}_\sigma) \dot{\alpha}_\sigma d\sigma \right\rangle d\tau \\ = \alpha_0 \langle \mathbf{z}_t - \mathbf{z}_0, \nabla \psi^*(\mathbf{z}_0) \rangle + \int_0^t \langle \mathbf{z}_t - \mathbf{z}_\sigma, \nabla \psi^*(\mathbf{z}_\sigma) \rangle \dot{\alpha}_\sigma d\sigma.$$

By the definition of Bregman divergence, $D_{\psi^*}(\mathbf{z}, \mathbf{w}) = \psi^*(\mathbf{z}) - \psi^*(\mathbf{w}) - \nabla \psi^*(\mathbf{w}, \mathbf{z} - \mathbf{w})$. Thus, combining Eqs. (A.3) and (A.5) leads to $\mathcal{C}_t = 0$. As this holds for an arbitrary t , the proof is complete.