

1 **Reviewer 1** Thank you for the comments on how to improve the notation, and the specified typos. We will incorporate  
2 these changes. Our responses to your main questions are below.

3 - We will restrict our discussion to strongly convex functions rather than strictly convex functions in the differential  
4 geometry section, as we don't currently handle the extra complexity from strict convexity in our discussion. Thank you  
5 for pointing this out.

6 - The view of the gradient map as a diffeomorphism is a useful way of visualizing what's going on, we will mention this  
7 in the updated draft. We will provide the reference regarding geodesics that you request as well.

8 - We do apply the Bregman divergence machinery to both  $\phi$  and  $f$ .  $\phi$  is used in the introductory sections rather and  
9  $f$  throughout because we technically use the conjugate of  $f$ , not  $f$  its self, in the main change of geometry, and we  
10 wished to avoid notation involving conjugate-of-conjugate operations. We will update the wording to make this clearer.

11 - At line 154, the flatness ensures the solution has a simple form, you are correct that a solution will always exist even  
12 without flatness.

13 **Reviewer 2** - The research avenue of understanding acceleration by discretizing continuous time dynamics is certainly  
14 one of the leading approaches currently being pursued. We have followed this literature carefully, as and far as we are  
15 aware all the existing approaches suffer from an "ad-hoc"-ary problem, where going from the continuous to discrete  
16 form requires a non-obvious or non-standard discretization method. The Approximate duality gap technique suffers  
17 from this problem, where as they state "... we can introduce an additional gradient step whose role is to cancel out the  
18 discretization error by reducing the upper bound.". Although certainly a promising approach, we don't consider this a  
19 satisfying explanation for the performance of acceleration, since it would be difficult to arrive at Nesterov's method via  
20 this path without already being aware of its functional form (i.e. working backwards). In contrast, our technique just  
21 relies on replacing the Bregman divergence used in the proximal point method with another more easily computable  
22 divergence, which is a natural step.

23 "As [1] showed, the modern formulation of estimation sequence has strong power to explain acceleration." - We don't  
24 believe that the mentioned paper refutes our statement about estimate sequences, as they only derive composite-dual  
25 averaging via their discretization technique, not the simpler Beck and Teboulle method, which is the explicit claim we  
26 make. Likewise they do not apply their technique to variance reduced objectives, which is problematic to address using  
27 continuous time dynamics. The accelerated form of Lan that we build upon does apply to variance reduced problems.

28 Addressing point (2), our technique does lead to a ODE whose forward Euler discretization is naturally a discrete  
29 accelerated method, unlike the approach in [1]. We believe the explanation that our method provides, i.e. that Nesterov's  
30 method is the proximal point method in disguise, provides substantial insight that other explanations lack.

31 We hope that Reviewer 2 will revise his judgement based on our comments, although we acknowledge that the utility of  
32 new explanations of existing algorithms is subjective.

33 **Reviewer 3** "To me, the weakest aspect of the paper is that the new interpretation of the accelerated gradient method  
34 still does not explain why it achieves acceleration" This is certainly the fundamental question. The high-level approach  
35 we attempt in this work is to show that it is just a form of the proximal point method. We believe the proximal point  
36 method is "intuitively" fast, but of course this is subjective.

37 - The mentioned Bubeck et. al. paper is perhaps the most unique approach to acceleration developed in recent years. We  
38 don't discuss it in detail in our work as the approach of finding points in the intersection of two balls is quite distinct  
39 from the proximal approach we look at, and it is not equivalent to Nesterov's method.

40 - The linear coupling paper by Allen-Zhu and Orecchia provides an interpretation of acceleration as linearly interpolating  
41 between a primal and a mirror step. This interpretation uses an equational form of acceleration similar to the one  
42 we build upon. It's certainly another valid interpretation, but we believe it lacks explanatory power, as it is unclear  
43 *why* such an interpolation would give an accelerated rate. As they state, it requires a substantial amount of analysis  
44 to derive the interpolation constant, which is not the case in our version, where the  $\tau$  factor is directly given by the  
45 change-of-geometry.

46 - The heavy ball section could potentially be cut. We kept it in as we were often receiving questions about the relation to  
47 the heavy ball method whenever we discussed our work with others. We will move it to the appendix.

48 - You are correct that the biorthogonal condition is just a change of variables for the tangent spaces. It's notable just  
49 because of it's particular simplicity for these two coordinate systems.

50 Thank you for the remaining detailed comments. We do not have space here to address them one-by-one but we will  
51 update the manuscript to fix all the mentioned issues and provide clarification where requested.