

1 We thank all the reviewers for their detailed reviews and valuable comments. We very much appreciate the positive
 2 feedback and are very happy that the reviewers generally like our contribution. In what follows, we focus on the
 3 most significant comments raised by the reviewers; we will carefully implement all minor suggestions as well as
 4 improve the exposition in Section 5 as pointed out by R2 and R3. Code for experiments and proofs can be found in the
 5 supplementary material for reproducibility (R2). Further, we will add a clarification in the related work regarding first
 6 linear convergence results by Garber and Hazan 2014.

7 **Q1. Framework to unify and derive a subset of prior CG/FW results (Reviewer 3, 4).** We highlighted some
 8 connections to prior CG/FW results in Sections 3, 4, Table 1 (appendix), however we will make these more explicit. We
 9 consider the shadow of the gradient, i.e., its directional derivative, as a descent direction. We show that (i) the continuous
 10 time dynamics of moving along the shadow at any point is equivalent to the continuous time dynamics of projected
 11 gradient descent (PGD). (ii) In PGD, a constrained movement along the gradient (e.g., $-\frac{1}{L}\nabla f(\mathbf{x}_t)$) is projected back
 12 to the polytope. We show that the limit of infinite movement along the gradient (i.e., $\lim_{w \rightarrow \infty} -w\nabla f(\mathbf{x})$) followed
 13 by a projection recovers the Frank-Wolfe (FW) vertex (R4: at a high level, we called this “maximal wrap” around the
 14 polytope), thereby establishing a novel interpretation of FW steps [Frank & Wolfe 1956]. (iii) We show the PGD iterate
 15 (i.e., $\Pi_P(\mathbf{x}_t - \frac{1}{L}\nabla f(\mathbf{x}_t))$) lies on the projections curve starting from \mathbf{x}_t . (iv) Next, we show that the shadow direction is
 16 the best normalized (unit norm) local direction for descent in terms of the dot product with the gradient. This implies it
 17 is the best direction in the convex hull of all possible away directions, i.e., the set considered by [Bashiri et. al 2017]. (v)
 18 Pairwise-steps considered by [Jaggi & Lacoste-Julien (2015)] are the sum of FW steps and away directions, therefore,
 19 the best pairwise-steps could similarly be obtained using our framework. (vi) [Bashiri et. al 2017] and [Garber & Hazan
 20 2016] both compute the best away vertex in the minimal face containing the current iterate, whereas the shadow step
 21 recovers the best convex combination of such vertices aligned with the negative gradient. Therefore, these previously
 22 mentioned CG methods are *both* approximations of SHADOW-CG. (vii) Moreover, [Garber & Hazan 2014] restrict the
 23 FW vertex to a ball around the current iterate, thereby normalizing the norm of the descent direction. Their algorithm
 24 can be therefore understood as an approximation of SHADOW-WALK.

25 **Q2. Theoretical differences between SHADOW-CG and SHADOW-WALK (Reviewers 1, 3, 4).** SHADOW-WALK
 26 moves along the shadow direction and traces the projection curve until a “non-boundary step” is taken. SHADOW-CG,
 27 on the other hand, either moves along the FW or shadow direction depending on the dual gap test (choose FW direction
 28 whenever $\langle -\nabla f(\mathbf{x}_t), \mathbf{d}_t^{\text{FW}} \rangle \geq \langle -\nabla f(\mathbf{x}_t), \mathbf{d}_{\mathbf{x}_t}^{\Pi} / \|\mathbf{d}_{\mathbf{x}_t}^{\Pi}\| \rangle = \|\mathbf{d}_{\mathbf{x}_t}^{\Pi}\|$, and shadow direction otherwise). Theoretical bound
 29 on iteration complexity for a given fixed accuracy is better for SHADOW-WALK, since it is closer to PGD. However,
 30 the computational complexity for SHADOW-CG is better compared to SHADOW-WALK since FW steps are much
 31 cheaper to compute compared to the shadow direction and we can avoid the potentially expensive computation via the
 32 TRACE-routine. This is also observed in the experiments (see Q4).

33 **Q3. Independence of the convergence rates on geometric constants (Reviewers 1, 4).** As the reviewers note, we
 34 prove an *algebraic* upper bound of $O(2^m)$ breakpoints in the projections curve, for polytopes with m facets in n
 35 dimensions; however computationally we see much fewer oracle calls (e.g., rarely exceeding 10 in 100 dimensions).
 36 Similar in nature to the PAIRWISE-FW (which is arguably one of the fastest CG variants in practice, yet lacks a more
 37 sophisticated analysis), we unfortunately have to resort to a highly pessimistic worst-case bound. We have talked to
 38 various discrete geometers and there was no consensus on what the right answer for the order of number of breakpoints
 39 should even be. We believe that once the projection curve leaves a facet it can never reenter, and hence we conjecture
 40 $O(mn)$ breakpoints in lines 311-313 of the paper, which would significantly improve upon the $O(2^m)$ bound. Although
 41 our linear convergence rate depends on the the number of facet inequalities m and in fact the combinatorial structure of
 42 the face-lattice of the polytope, *it is invariant under any deformations of the actual geometry of the polytope preserving*
 43 *the face-lattice* (in contrast to vertex-facet distance and pyramidal width), e.g., Figure 4’s discussion shows pyramidal
 44 width can become arbitrarily small while the number of facets is invariant. We will make this distinction more precise.

45 **Q4. Computational differences in SHADOW-CG and SHADOW-WALK (Reviewers 1, 2, 3, 4).** We provide two
 46 ways to compute the shadow: (i) using equation (8), i.e., $\mathbf{d}_{\mathbf{x}}^{\Pi} = (\Pi_P(\mathbf{x} - \epsilon\nabla f(\mathbf{x})) - \mathbf{x})/\epsilon$ for ϵ sufficiently small; (ii) as
 47 stated in line 155 $\mathbf{d}_{\mathbf{x}}^{\Pi} = \arg \min_{\mathbf{d}} \{\|-\nabla f(\mathbf{x}) - \mathbf{d}\|^2 : A_{I(\mathbf{x})}\mathbf{d} \leq \mathbf{0}\}$. Hence, either way, computing the shadow reduces
 48 to a convex quadratic program. Note, that the cost of the shadow computation has to be non-trivial, as it was shown in
 49 [Garber 2020, Diakonikolas et al. 2019] that we cannot eliminate the dependence on the dimension in convergence rates
 50 of CG variants *that solely rely on linear optimization as oracle*, whereas our rates for SHADOW-WALK are independent
 51 of the dimension (i.e., $(1 - \mu/L)$ contraction). That being said, even though computing the shadow is expensive in the
 52 worst case, it can be cheap for specific setups and there are many existing algorithms that provide cheap approximation
 53 to this oracle, e.g., DICG variants [Garber and Hazan 2016, Bashiri et. al 2017]. Our computations confirm that
 54 SHADOW-CG interpolates between CG variants and PGD. In particular, Figure 2 for the Video Co-Localization problem
 55 shows SHADOW-CG has lower iteration count than CG variants (slightly higher than PGD), while also improving on
 56 wall-clock time compared to PGD and SHADOW-WALK being almost as fast as CG without assuming shadow oracle
 57 access. Addition of FW steps in SHADOW-CG significantly reduces the total number of shadow oracle calls and number
 58 of times we enter TRACE-routine as demonstrated for the Lasso regression problem in Figures 10 and 13 (appendix).