

1 We thank the reviewers for their careful reading and useful feedback. Below we briefly address each reviewer’s
2 comments; these will be discussed in greater details in the revised version.

3 **Common question on how to choose the problem parameters:** As reviewers pointed out, the requirement to know
4 the problem parameters (strong convexity and smoothness) is endemic to many first order methods, including the
5 original Catalyst framework for minimization. It is worth pointing out that for many applications, the strong convexity
6 comes from the regularization term that is usually known. In practice, using *a rough estimate* of the strong convexity or
7 smoothness constant would also work as long as the choice of regularization parameter τ leads to a more balanced
8 auxiliary problem (which is one of the key insights of our framework). Such a rough estimate can be obtained by
9 approximately computing the spectrum of the Hessian or through backtracking. As an alternative, these constants can
10 be simply treated as hyper-parameters and tuned with the stepsize.

11 **To Reviewer 1:**

- 12 • **Numerical comparison:** Implementations of previous methods such as Minimax-APPA and primal-dual smoothing,
13 *are not available* in the original papers. These algorithms are difficult to implement as they require smoothing
14 and knowing the desired accuracy for each loop, which must be tuned in practice. That is why we pick DIAG as
15 a representative benchmark among these three-loop algorithms for numerical comparison. We will move these
16 comparisons from the appendix to the main paper and provide more numerical evidence.
- 17 • **Dependence on \mathcal{D}_y and \mathcal{D}_x :** All algorithms in Table 2 have the same dependency on \mathcal{D}_y . Minimax-APPA assumes
18 \mathcal{X} to be bounded and has logarithmic dependency on \mathcal{D}_x in the complexity, and DIAG requires $\mathcal{X} = \mathbb{R}^d$. Our
19 algorithm assumes \mathcal{X} is convex and only has logarithmic dependency on $\|x_0 - x^*\|$ in SC-C setting.
- 20 • **Negligible cost for checking stopping criterion:** For stochastic methods such as SVRG and SVRE, the cost for
21 checking the stopping criterion is $\mathcal{O}(n)$ in each epoch and does not increase the overall complexity. In addition, the
22 full gradient is already computed in each epoch, so the cost of checking this criterion is still almost negligible.
- 23 • **Gain in other setting.** As the reviewer pointed out, there is no gain for the convex-concave setting as EG is
24 already optimal. However, there could be gains in other settings, such as SC-SC and NC-SC settings. In fact, we
25 have recently shown that when extending to (μ_x, μ_y) -SC-SC setting, our framework achieves the complexity of
26 $\mathcal{O}(\ell/\sqrt{\mu_x\mu_y} \log(1/\epsilon))$, which again improves over EG and matches the lower bound up to logarithmic factors.
- 27 • **Lower bound:** No matching lower bound is known for the NC-C setting. A valid lower bound from nonconvex
28 minimization is $\mathcal{O}(1/\epsilon^2)$, which may not be tight. For the SC-C setting, we do not know if the current lower bound
29 is missing a log factor or not. We will explicitly discuss the open questions of lower bounds related to our settings.
- 30 • **Typos and appendix:** We will fix these issues and clean the appendix. *Thanks for the careful reading!*

31 **To Reviewer 2:**

- 32 • **Narrow niche:** We should have mentioned that strongly-convex-concave minimax optimization itself has broad
33 applications in game theory, imaging, distributionally robust optimization, etc. The bilinear case of this setting has
34 been studied extensively, but the general case remains unexplored. Besides, the framework we present for SC-C
35 setting can be extended to NC-C setting (as shown in the paper) and potentially other settings such as C-C, SC-SC,
36 and NC-SC settings (see our response to Review 1).
- 37 • **Different learning rates:** We agree with the reviewer that some related work uses different learning rates for the
38 two players for NC-C setting, e.g., GDA [Lin et al., 2019] uses $\mathcal{O}(\epsilon^4)$ stepsize for the min player and $\mathcal{O}(1)$ for the
39 max player. Theoretically, they achieve a much worse $\mathcal{O}(\epsilon^{-6})$ complexity than the $\mathcal{O}(\epsilon^{-3})$ complexity achieved by
40 our Catalyst with GDA under *the same constant stepsize* for both players. We will also add numerical comparisons.
- 41 • **Related work:** We thank the reviewer for the pointers. We will add these references and discuss them in the paper.
42 Competitive gradient descent [Schäfer & Anandkumar, 2019] and some recent work [Xu et al., 2020] add quadratic
43 regularization terms in the forms of $\|x\|^2$ and $\|y\|^2$. This is different from our framework, as we only add the term
44 $\|y - z_t\|^2$, where z_t comes from extrapolation.

45 **To Reviewer 3:**

46 **Pre-defining the number of inner loop or accuracy:** This is a good point. In Theorem 4.1, we choose the desired
47 accuracy for inner loop to be $\bar{\epsilon} = \mathcal{O}(\epsilon^2)$. In fact, this can be replaced by an adaptive accuracy, i.e., $\bar{\epsilon}_t = \frac{g(x_0) - g^*}{t+1}$, which
48 will not increase the overall complexity. Hence, the number of inner loops does not necessarily have to be pre-defined.

49 **To Reviewer 4:**

50 We thank the reviewer for the acknowledgement of our contribution. Although Catalyst framework in minimization
51 provide a strong intuition for our work, the analysis for the outer-loop complexity of our framework is a nontrivial
52 extension of their analysis as it requires evaluating solution performance in terms of primal-dual gap (which is much
53 stronger than dual optimality).