

1 **Response to Reviewer 1.** Thank you for your compliments on our contributions as well as our presentation!

2 **Response to Reviewer 2.** Thank you for your thoughtful comments and questions!

3 **Q. Significance of relaxing component convexity?** For practical significance, it is crucial for understanding the
4 (local) behavior of shuffling algorithms for nonconvex problems, such as neural network training. As Reviewer 3 also
5 mentioned, in such settings, the function F would behave like convex near a neighborhood around a local minimum,
6 while each component function f_i could be highly nonconvex in the neighborhood.

7 For theoretical significance, this relaxation signifies an innovation in our proof techniques. In fact, the component
8 convexity is heavily exploited in the previous works [10, 13]. At a high level, the previous works use this assumption to
9 ensure that the algorithm makes a sufficient progress during *each* iteration. In contrast, our proof technique proves fast
10 convergence rates without having to show such a per-iterate progress. We also highlight that it is thanks to our proof
11 technique that we obtain the first optimal rate for nonconvex PL function class in Theorem 1.

12 **Q. Constant step algorithms vs. decreasing step size algorithms are not comparable?** We view the decreasing step
13 size case as choosing a different set of hyperparameters for the *same* algorithm. In light of this, our main goal in
14 analyzing decreasing step sizes was to see whether the common limitation in the prior works (large epoch requirement)
15 is inherent to the algorithm. Our main results indeed show that the common limitation is a derivative of the constant step
16 size choice. On a technical note, capturing the desired convergence rate with decreasing step sizes requires a non-trivial
17 analysis (Section 6.3 or Appendix E), and we believe that our work develops a toolkit for analyzing decreasing step
18 sizes for epoch-based algorithms like shuffling SGD.

19 **Q. Guarantees without Assumption 1?** We agree that Assumption 1 is a bit unsatisfactory, but please note that
20 this assumption is also present in many prior results (see Line 144). While our Theorems 1 and 2 *do not require*
21 Assumption 1, Theorems 3 and 4 *do* rely on Assumption 1 because they are built on existing results [10, 13] that
22 make use of Assumption 1. We believe that removing this assumption, as done in [12] or a concurrent work “Random
23 Reshuffling: Simple Analysis with Vast Improvements,” is an important future research direction. We will add a remark
24 on this in our revision.

25 **Q. Comparison to GD?** Our current proof techniques do not show a regime where rates for shuffling SGD can be
26 better than that of GD. Indeed, this limitation is shared among the existing works, as their analyses rely on comparing
27 the epoch progress of shuffling SGD to that of GD. In other words, this question is currently beyond the scope of
28 existing proof techniques, and in particular, finding a regime where shuffling SGD outperforms both SGD and GD
29 would be an interesting future direction to pursue. We will add this discussion in our revision.

30 **Response to Reviewer 3.** Thank you for detailing our contributions into a list. In particular, we agree that the removal
31 of individual convexity has important consequences in practical applications. Also, as per your suggestion, we will
32 explicitly state the definition of the constant G in Theorem 2.

33 **Response to Reviewer 4.** Thank you for appreciating the strengths of our paper! We did not spell out G in Theorem 1
34 because existence of G is implied by the other assumptions. We will explicitly state the constant G in the theorem
35 statement as per your suggestion.

36 On a separate note, we would like to fix one technical mistake in the proof of Theorem 1. After the submission
37 deadline, it came to our attention that the Hoeffding-Serfling (HS) inequality used in the proof of Theorem 1 can be
38 applied only to RANDOMSHUFFLE. In SINGLESHUFFLE, the first iterate x_0^k of the k -th epoch is not independent
39 of the permutation σ as soon as $k > 1$, hence we cannot apply the HS inequality. In light of this, we note that
40 Theorem 1 only holds for RANDOMSHUFFLE. As to our claims on SINGLESHUFFLE, we will add an additional theorem
41 in the revision. The new theorem shows that if F is μ -strongly convex and f_i 's are L -smooth, SINGLESHUFFLE
42 with number of epochs $K \geq 10\kappa^2 \log(n^{1/2}K)$, step size $\eta_i^k = \eta := 2 \log(n^{1/2}K) / \mu n K$, and initialization x_0 satisfies
43 $\mathbb{E}[F(x_0^{K+1})] - F^* \leq 2L \|x_0 - x^*\|^2 / n K^2 + c \log^3(nK) / n K^2$, for some $c = O(\kappa^4)$.

44 Although this change slightly weakens our initial claim on SINGLESHUFFLE, we believe this new theorem is still
45 interesting progress, because (i) it is a tight bound (up to log factors) for SINGLESHUFFLE on strongly convex functions;
46 (ii) it does not require convexity of individual components or bounded iterates assumption (Assumption 1); (iii) it
47 shows that the optimal rates for minimizing strongly convex functions are the same for RANDOMSHUFFLE and SIN-
48 GLESHUFFLE; and (iv) the proof can be easily extended to any algorithms between the spectrum of RANDOMSHUFFLE
49 and SINGLESHUFFLE. The proof of the new theorem involves an end-to-end analysis similar to the one-dimensional
50 quadratic result [16] and a modified version of the approximate matrix AM-GM inequality in eq (6.4). For this theorem,
51 the HS inequality is applied only to the partial sums of the individual gradients at the global minimum x^* , which is
52 always independent of the permutation σ .