

1 We thank the reviewers for their work and for the positive evaluation of our paper. Every reviewer noted that it is
2 well-written and reproducible. There seems to be no disagreement that the paper is of practical value (noted by R1, R2,
3 R4), relevant to the conference (R2, R4), and is generally on an interesting topic (R1, R2).

4 **To all reviewers.** We realized that we downplayed our contribution: Theorems 2 and 3 apply to Shuffle-Once without a
5 single change in the proofs, just as in Theorem 1. Thus, we also provided guarantees for SO without strong convexity.

6 **Reviewer 1.** 1. “*this work uses the fact that all component functions f_i are μ strongly convex*” This assumption is not
7 strong because strong convexity typically comes from ℓ_2 regularization or weight decay, which is normally used in each
8 function. Adding a small amount of regularization is also a common practice for numerical stability.

9 2. “*why removing some of the assumptions like bounded variance and bounded gradients is an important contribution*”
10 This assumption does not hold for any strongly convex function, which includes any convex objective with ℓ_2 regular-
11 ization. It also doesn’t hold for least squares, matrix factorization and neural networks (with smooth activations). These
12 problems are at the core of modern machine learning, and there are a lot of other similar objectives.

13 3. “*The quantity σ_* being finite also implies that all the gradients are finite*” This is not true, σ_* is always finite just
14 because $\nabla f_i(x_*)$ is a finite vector for any i , but globally the gradient is often not bounded (see the previous item).

15 **Reviewer 2.** We appreciate your support of our paper. Since accepted papers will be allowed to have an extra page, if
16 ours gets accepted, we will be happy to add discussion on the shuffling variance, which we agree will be educational.

17 **Reviewer 3.** 1. “*the main theorems (Theorems 1 and 2) need a small step size, similar to previous works. In fact*
18 *Safran and Shamir (2020) show that convergence is only possible for step size $O(1/n)$ ” Firstly, we disagree about
19 Theorem 1—even with step size $\frac{1}{L}$ it guarantees convergence to a neighborhood. In practice, achieving machine
20 precision is sometimes not important, and a neighborhood would suffice. Furthermore, a substantial part of our work is
21 devoted to explaining why the neighborhood’s size is small. Secondly, the argument of Safran and Shamir (2020) does
22 not show convergence is not possible with large stepsizes, indeed the proofs of Propositions 1 (p. 10, 2nd bullet point)
23 and Thm. 2 (p. 16, 4th bullet point) in their work both show that convergence proportional to n and T may be achieved
24 with a large stepsize. We will add these clarifications.*

25 2. “*the dependence on μ has worsened. In particular, Nagaraj et al. (2019) give an error rate of κ^2/μ ” This
26 comparison is not fair because Nagaraj et al. (2019) bound functional gap and our bound is for the distances. If we
27 apply $\|x - x_*\|^2 \leq \frac{1}{\mu}(f(x) - f(x_*))$, we see that the bound of Nagaraj et al. gives rate $O(\kappa^2/\mu^2)$.*

28 3. “*The result of Theorem 3 says that RR beats SGD only after $\Omega(n)$ epochs, which seems really large. This can*
29 *probably be improved.*” This claim is somewhat speculative: we are unaware of any lower bounds in this setting. We
30 devoted a lot of time to tightening this bound and the complexity that we proved is better than any other available in the
31 literature for the same setting, so it is unclear why a better bound should be possible.

32 4. “*Theorem 1 assumes that the individual functions are strongly convex. This is a really strong assumption which, in*
33 *my opinion, the previous papers have deliberately avoided.*” We disagree that it is such a strong assumption (see our
34 response to Reviewer 1). “*I also believe that this assumption makes the analysis much more easier because now the*
35 *sum of every subset of the individual functions is also strongly convex.*” Our proof does not use that any sum of a subset
36 of the individual functions is strongly convex, so we do not see how this remark is relevant.

37 5. “*The literature review is somewhat poor.*” We wrote a longer review but had to make it shorter because of the page
38 limit. If the paper gets accepted, we would be happy to use part of the allowed additional page to put more discussion.

39 6. “*Also, all the relevant upper bounds should be compared in detail with the lower bounds given by Safran and Shamir*
40 *(2019) and Rajput et al. (2020).*” We already compared the bounds with both papers, check page 6, in lines 177-185.

41 **Reviewer 4.** 1. “*Does the polynomial dependence on κ also match the lower bounds?*” To our knowledge, there is no
42 lower bound with an explicit dependence on κ . As a sanity check, if $n = 1$, then $\sigma_* = 0$ and RR reduces to gradient
43 descent and we obtain the standard dependence on κ , so it is tight at least in this sense.

44 2. “*Can the authors discuss if their results for the non strongly case can be improved or are there known lower bounds?*”
45 We are not aware of any such lower bound. For SGD, recent work (Sebbouh et al. "On the convergence of the Stochastic
46 Heavy Ball Method") shows that momentum makes SGD faster when there is no strong convexity, despite SGD being
47 optimal in the strongly convex case. We think that momentum may improve the rate of RR as well.

48 3. “*Can the authors please comment on if a similar analysis might be possible for SO?*” We double-checked these
49 results and they turned out to apply to SO without any change in the proof.

50 4. “*it seems like $\sigma_{\text{shuffle}}^2$ is much bigger than σ_*^2 unless the step size is very small*” If $\kappa = 1$ and $\gamma = \frac{1}{L}$, then indeed the
51 shuffling variance can be in theory n times bigger than σ_*^2 but we didn’t observe this in our experiments.

52 5. “*I think it would be useful to write the SGD convergence rate as in equation 4*” We can add it: the first term will be
53 the same as in (4) and the second term would be $\frac{\sigma_*^2}{\mu^2 n T}$, where nT is the total number of steps.

54 6. “*For figure 2, how are the values $\sigma_{\text{shuffle}}^2$ and σ_* estimated? And, what are the values of L and μ that are used?*” We
55 first obtained x_* by running Nesterov’s acceleration until machine precision, and then we estimated expectations by
56 randomly sampling permutations, and σ_* was computed exactly.