1  To **Reviewers**, we will make all suggested minor corrections in the final version and address main concerns below.

2  **R1** (and **R3, R5:**) Thanks for your constructive comments! The idea of integrating several key optimization techniques,
3  dual averaging + variance reduction + acceleration, is novel and nontrivial for three reasons. **First**, for the general convex
4  setting, the proposed SVR-ADA gives the rate $O\big(n\log\log n+\frac{\sqrt{nL}}{\sqrt{\epsilon}}\big)$, which improves the SOTA rate $O\big(n\log\frac{1}{\epsilon}+\frac{\sqrt{nL}}{\sqrt{\epsilon}}\big)$
5  of Katyusha, not substantially improved for over $4$ years; for the well-conditioned strongly convex setting (the case
6  of the number of samples $n$ far greater the condition number, *i.e.*, $n \gg \kappa$), SVR-ADA has the rate $O\big(n\log\log n +$
7  $\frac{n}{\log(n/\kappa)}\log\frac{1}{n\epsilon}\big)$ which improves the SOTA rate $O\big(n\log\frac{1}{\epsilon}\big)$ of SAG, unchanged for 8 years; for the ill-conditioned
8  strongly convex setting ($n \leq O(\kappa)$), SVR-ADA matches the lower bound in [25]. **Second**, besides the improved or
9  optimal convergence results, SVR-ADA shares the simplicity of MiG [29] and the unification of Varag (see comments
10  of **R2**) *simultaneously*. **Third**, this work is the first to show that in the finite sum setting, combining accelerated dual
11  averaging (ADA) with variance reduction (VR) gives better convergence rates than that of combining accelerated mirror
12  descent (AMD) with VR as adopted in Katyusha, MiG and Varag. This provides new perspectives to acceleration.

13  In terms of experiments, SVR-ADA is compared with SOTA finite sum solvers. The results send one consistent
14  message: *SVR-ADA performs well in all the three settings, namely general convex, ill-conditioned or well-conditioned*
15  *strongly convex.* In comparison, the existing accelerated algorithms Katyusha and MiG do not perform so well for
16  the well-conditioned strongly convex setting; the non-accelerated algorithm SVRG does not perform well for the
17  general convex and ill-conditioned strongly convex settings. In terms of the choice of regularization, we use different
18  two-norm regularizations to represent the phenomenon of the three settings compactly by Figure 1 and Figure 2.
19  If we use one-norm, then it can only represent the general convex setting. For completeness, we will show the
20  numerical results of one-norm in the final version. The different regularization parameters of the two figures are used to
21  better represent the obvious change in experimental phenomena. In the final version, we will represent the results of
22  $\{0, 10^{-8}, 10^{-7}, 10^{-6}, 10^{-5}, 10^{-4}, 10^{-3}\}$ for both experiments at least in the supplementary material.

23  Thanks for your suggestions about writing! In the final version, we will rewrite the abstract to make it more clear. Per
24  your suggestions, we will use "general convex" to replace "nonstrongly convex" and make the statements in Line 30 and
25  32 more precise. We will add a paragraph to make a simple introduction to stochastic algorithms that do not use VR.

26  **R2:** We are very glad that you obviously recognize significance of our contributions! We were also aware of Varag
27  after our submission! In the final version, we will make systematical comparisons with Varag from three perspectives:
28  1) **Efficiency.** SVR-ADA improves the SOTA rates for general convex and well-conditioned strongly convex settings,
29  while Varag does not improve any SOTA rates. 2) **Simplicity.** SVR-ADA only uses two-point coupling in the inner
30  iteration, while Varag uses three-point coupling. SVR-ADA uses a uniform average in the outer iteration, while Varag
31  uses a weighted one. 3) **Unification.** SVR-ADA unifies the general convex and strongly convex settings with much
32  simplified parameter settings, which can be simply specified in algorithm description. However, to adapt to both settings,
33  the parameter settings of Varag are very complicated and not explicitly and clearly stated in the algorithmic description.

34  In terms of techniques, **first**, we use the combination of VR and ADA, while they use the combination of VR and AMD;
35  **second**, we use completely different and concise convergence analysis by estimation sequence; **third**, we only use
36  negative momentum, while they use both Nesterov's momentum and negative momentum.

37  In terms of applicability in the constrained setting, we can always reformulate a convex constrained convex problem as
38  an unconstrained convex problem by the indicator function of the constrained set. Thus if the indicator function admits
39  an efficient proximal operator, then we can apply SVR-ADA in the constrained setting. Meanwhile, in the final version,
40  we will give a complete proof for the bound of $A_k$. We will cite the latest reference Joulani et al. (ICML 2020)!

41  **R3:** Thanks for your constructive comments! We fully understand your concern with the applicability for nonconvex
42  problems such as deep neural networks. However, it is an open problem for all control variate finite sum solvers, not
43  only for ours. Meanwhile, the *primary area* of our contribution is in convex optimization, which is valuable for many
44  problems in machine learning, signal processing, operational research. Thus we cannot undervalue the widely studied
45  control variate approaches as they significantly reduce the overall complexity for finite sum convex optimization. It is
46  fair to say broad applicability is the main reason why the control variate finite sum solvers are so widely studied!

47  In terms of **theoretical contributions**, as highlighted in our response to **R1**, we have improved rates in two regimes
48  that remain unimproved for many years, in a very actively studied area. In terms of **empirical performance**, as in our
49  response to **R1**, SVR-ADA is the first algorithm that performs well in all the three settings. We will conduct more
50  experimental evaluation as per the request by you and **R1** and report in the final version.

51  **R5:** Thanks for your constructive comments! For the general convex setting, the best-known result $O(n\log\frac{1}{\epsilon}+\frac{\sqrt{nL}}{\sqrt{\epsilon}})$ is
52  optimal up to a $\log$ factor. In convex optimization, a main endeavor is to shave off $\log$ factors to match the corresponding
53  lower bounds. We have made a substantial step forward by reducing $\log n$ to $\log\log n$. Meanwhile, as in our response
54  to **R1**, we also improve the rate for the well-conditioned strongly convex setting. As **R2** commented, the simplicity and
55  unification merits of SVR-ADA are remarkable! We believe our work is of significant value to this area.