1  We thank the reviewers for their comments and suggestions. We address their comments separately.

2  **Reviewer #1:**

3  *1.1 It is unclear how the bias is handled and the relation to inexact gradient method.*

4  Reply: First, in this paper we are handling the mean-squared error of estimators, which can be decomposed as
5  $\mathbf{E}[\|\tilde{\nabla}F(x_i^t) - \nabla F(x_i^t)\|^2] = \|\nabla F(x_i^t) - \mathbf{E}[\tilde{\nabla}F(x_i^t)]\|^2 + \mathbf{E}[\|\tilde{\nabla}F(x_i^t) - \mathbf{E}[\tilde{\nabla}F(x_i^t)]\|^2]$. The first term is the squared
6  norm of bias (nonzero in our case) and the second term is the variance. Therefore, our procedure is controlling both of
7  them (see proof of Lemma 1 in Appendix). We shall emphasize this in revision. Second, our method can be considered
8  as an inexact gradient method in the general sense. However, the additive inexactness are carefully controlled by the
9  amount of descent to yield desirable complexity, which is closely related to the property of the estimator. Thus we
10 cannot directly use the standard analysis of inexact gradient method, which will lead to worse complexity.

11 *1.2 The step size of the algorithm needs to be small in typical non-convex optimization, usually depends on $\epsilon$.*

12 Reply: Even for non-convex optimization, constant step size can be used if the objective function is smooth. Small
13 or diminishing step sizes are mostly required in stochastic optimization to combat noise in the stochastic gradients,
14 for both convex and non-convex optimization. We are able to use constant step size in the stochastic optimization
15 setting because of much stronger assumption: we assume each realization of $g_\xi(\cdot)$ is a smooth function, not merely its
16 expectation $E_\xi[g_\xi(\cdot)]$ as in classical stochastic optimization. This assumption is key to effective variance reduction in
17 stochastic optimization, and is satisfied in most machine learning problems. We will elaborate on it in the revision.

18 **Reviewer #3:**

19 *3.1 Significance of improvement and reference to [34].*

20 The problems we consider are the same class of problems addressed in [34], thus lead to similar motivation and
21 introduction. We will revise the introduction to be more concise and rely on reference to [34]. While [34] first showed
22 that variance reduction techniques (SAGA/SVRG) can be used to handle biased gradient estimators in composite
23 stochastic optimization problems, this paper improves the dependence on $n$ from $n^{2/3}$ to $n^{1/2}$ using a different estimator
24 SARAH/SPIDER. From theoretical perspective, the improvement is significant for large $n$. More importantly, it reaches
25 the lower bound for the considered problem class in the inner finite-sum case (see next point).

26 *3.2 Discussion of lower bound in more depth than line 99.*

27 Reference [8] showed that for nonconvex smooth finite-sum optimization problems (without the outer composition), the
28 lower bound on sample complexity of the component gradients is $O(n^{1/2}\epsilon^{-1})$. That is a special case of the composite
29 problem we consider, where the inner smooth mappings are functions (range dimension is 1) and the outer composition
30 is the identity function. From this perspective, our result also reaches the lower bound for the more general nonlinear
31 composite problems. We will clarify further in the revision, especially with the additional page allowed if accepted.

32 *3.2 Questions on the experiments.* We agree with the reviewer that more exhaustive search and tuning of parameters
33 may be needed to compare the bests of different algorithms. However, the main purpose of our experiments are to
34 demonstrate some basic behaviors of the algorithms on a few simple examples. Remember that the complexities
35 cited for different algorithms are their theoretical upper bounds. It is not clear if SVRG or SAGA based estimators
36 can achieve the same complexity as CIVR. At least in the convex finite-sum case, their complexities are the same.
37 Nevertheless, we will do some additional experiments to gain more understanding.

38 **Reviewer #4:**

39 *4.1 On preventing restart.* First, our restart scheme in Section 4 does not use different $\epsilon$ at different periods (outer loop).
40 We use the desired $\epsilon$ to set all algorithmic parameters once and do not change later. The only reason we use restart is
41 that the output of Algorithm 1 is chosen randomly from all past iterates within one period, which we need to use to start
42 the next period for analysis (see Proof of Theorem 5 in Appendix C.1.) If we can use the last iterate of each period, then
43 there is no need to perform "restart". Nevertheless, it can be avoided by using pre-generated stopping times. Note that
44 we can predetermine period length $T$ and epoch length $\tau$ ($\tau_1 = \cdots = \tau_T$). The output is drawn uniformly from the
45 $T\tau$ iterates. Therefore, we can first uniformly randomly generate the "stop time" $(\bar{i}, \bar{t})$. Then as soon as the algorithm
46 arrives this time, we "restart" seamlessly. We shall add this remark to the revision.

47 *4.2 The sensitivity on $\eta$.* Our bounds on $\eta$ in the theorems are to guarantee convergence and order of complexity.
48 Smaller $\eta$ can be used, but will encounter a performance penalty, which is explicit in Theorem 1, 2 and 3 on the bounds
49 in (22), (23) and (24) respectively. For theorems in Section 4, small $\eta$ will make $T \propto 1/\eta$ larger (also explicit in the
50 theorems). This will translate into total iteration cost, which we shall clarify in the revision.

51 *4.3 The "loopless" version (for inner loop, as the outerloop for restart is addressed in 4.1).*

52 We thank the reviewer for pointing out this interesting literature. According to our analysis, this technique can indeed
53 by applied, and we also get the same complexity. However, different from the convex problem in the "loopless" paper,
54 in the non-convex case the "loopless" proof is a bit more complicated then the current analysis. We will point out this
55 potentially interesting variant of the our method in the revised version.