

1 We want to thank the reviewers for their comments.

2 **Reviewer 2:**

3 1. (Experiments) (Challenge in Theoretical Analysis)

4 2.a. “There are eight assumptions used in the analysis which might be too restrictive.”

5 In Tables 1 and 2, in column “Limitations”, we compare how the assumptions between the papers are different. The  
6 main purpose of “Limitations” is to show that our assumptions are not stronger than the assumptions of any papers from  
7 Tables 1 and 2, so we do not agree that they are restrictive. Some papers present these eight assumptions as two-three  
8 assumptions. For instance, we define 3 different smoothness constants, while other papers can define only 1 constant,  
9 which is just the maximum of our 3 smoothness constants.

10 2.b. “In particular, could you please provide more comments on Assumption 6?” This is the mean-squared smooth-  
11 ness property that is used in **all papers with variance reduction**. See (Lower Bounds for Non-Convex Stochastic  
12 Optimization Arjevani et al.; Momentum-Based Variance Reduction in Non-Convex SGD Cutkosky et al., SPIDER:  
13 Near-optimal non-convex optimization via stochastic path integrated differential estimator Fang et al.)

14 3. This part is indeed can be confusing. We mean that the compressors are *statistically* independent. In other words,  
15  $\mathcal{C}_1(x_1), \dots, \mathcal{C}_n(x_n)$  are independent random vectors for all  $x_1, \dots, x_n \in \mathbb{R}^d$ .

16 4. In order to get the computational and communication complexities, we have to substitute an explicit formula of  
17  $\omega$  from Definition 1. It may differ for different compressors (see On Biased Compression for Distributed Learning  
18 Beznosikov et al). So it is not possible to provide a nice corollary that will work for any compressor. The good news is  
19 that RandK is the simplest compressor. By showing an improvement for RandK, we can expect that more advanced  
20 compressors will have even better theoretical and practical guarantees.

21 5. We also discuss it after Corollaries 1, 3, and 4.

22 **Reviewer 4:**

23 1. (Challenge in Theoretical Analysis)

24 2.a If you read papers from Tables 1 and 2, you will find that virtually all of them extend the theory from some previous  
25 papers. The research on the optimization methods is incremental. It is not very obvious why it is a weakness of our  
26 paper and not others.

27 2.b Our assumptions are general and not stronger than the assumptions of methods from Tables 1 and 2. Our analysis is  
28 independent and requires additional mathematical techniques that we provide in our proofs.

29 3. (Experiments)

30 **Reviewer 5:**

31 1. (Experiments)

32 2. We are aware of the work by A. Defazio L. Bottou. But it does not mean that VR methods are hopeless for neural  
33 network optimization. See, for instance, a recent work (Momentum-Based Variance Reduction in Non-Convex SGD  
34 Ashok Cutkosky, Francesco Orabona), where they provided theoretical and practical improvement. The development  
35 and understanding of VR methods are still going.