

1 We thank all reviewers for their insightful reviews which have helped us improve our work.

2 **To Reviewer #4.** We thank R4 for pointing out our contributions as well as some confusions. This review summarizes  
3 the highlights in our work, including the advantage of independence of dimensionality, analysis of optimal quantization  
4 parameter setting, and the extension to non-linear quantization.

5 *The dimension dependent bound is tight so why does the re-parameterization help?* Our paper is intended to provide  
6 insights for scaling real-life low-precision applications based on theoretical analysis. Because the dimension-dependent  
7 bound is tight under the standard assumptions, this means that if we want to understand when low-precision SGD  
8 may be dimension-independent, we need to look at it under different assumptions. The reparameterization in terms  
9 of  $L_1, \sigma_1$  gives us these different assumptions, allowing us to show that (under our new assumptions) the outcome of  
10 low-precision training on large-scale problems does not necessarily go worse as the dimension scales. And based on  
11 this, we made analyzed non-linear quantization schemes such as logarithmic and floating-point quantization to show  
12 that, given the problem we are faced with, there is a way to adjust our quantization setting to optimize the performance.  
13 We made a comparison between our results and previous results in Table 1. As we stated in the analysis after Theorems  
14 1 and 2, our result only reduces to the error level of prior analyses in the worst case, where parameter  $L_1, \sigma_1$  are unfixed  
15 and can only be bounded by  $L_1 \leq \sqrt{d}L, \sigma_1 \leq \sqrt{d}\sigma$ ; otherwise, our result is an improvement.

16 **To Reviewer #5.** We thank R5 for the encouraging feedback. This review perfectly summarizes our work and points  
17 out the advantages of our results and possible applications.

18 *Presentation.* We presented Table 1 with comparisons of our results and previous results to point out the novelty and  
19 improvement of our work. But as you suggested, we have moved it to appear after the assumptions. Parameter  $\kappa, \kappa_1$  are  
20 introduced between Assumptions 3 and 4 on page four, as condition numbers.  $\sigma$  and  $\sigma_1$  are bounds for loss gradients in  
21 different norms,  $\sigma_0$  is the bound for the gradient variance, which is normally defined as  $\sigma$  in other works, so we added  
22 subscripts to distinguish them.

23 *Analysis of other algorithms.* We have applied our analysis to two other algorithms which use low-precision models:  
24 low-precision SVRG [13] and HALP [8]. We achieved similar dimension-independence conclusion and explored  
25 the application of non-linear quantization schemes for these algorithms—however, since the analysis was essentially  
26 identical and merely repeated our other claims about SGD, we did not include it. Other low-precision algorithms  
27 (e.g. [22,23]) use low-precision arithmetic in different ways (such as to store intermediate values used during gradient  
28 computation) to which our theory does not directly apply: we plan to explore theory for these algorithms in future work.

29 *Parameter values.* The parameter values mentioned by R5 were measured for the MNIST dataset: these are the smallest  
30 values for  $L, L_1, \sigma$  and  $\sigma_1$  for which Assumptions 1–4 hold for multiclass logistic regression on MNIST. If larger (i.e.  
31 loose) values were used here instead of the reported ones, this would just result in a looser theoretical bound.

32 **To Reviewer #6.** This review helps us understand what parts of our work are not explained well enough and may cause  
33 confusion. We thank R6 for the useful constructive feedback.

34 R6 points out that the parameters  $L_1$  and  $\sigma_1$  may depend on  $d$ , and is concerned about our overall bounds being  
35 dependent on  $d$  in this case. We have two responses to this. First, as shown in Fig 1(a) the standard dimension-  
36 dependent bound is in some sense tight, so we should expect to see classes of problems for which the performance  
37 depends strongly on  $d$ . For these classes of problems, our parameters  $L_1$  and  $\sigma_1$  will also increase strongly with  $d$ .  
38 However, there are classes of problems for which this does not happen, and for the class we study in Figure 2(a), the  
39 performance does not depend on  $d$  either, which is what our theory predicts. Second, even in the worst-case scenario  
40 when the parameters  $L_1$  and  $\sigma_1$  do depend strongly on dimension, our results in Table 1 show that, by using non-linear  
41 quantization, we can actually put those terms inside double log and get a  $\mathcal{O}(\log \log d)$  upper bound when it comes to  
42 the number of bits required. This is better than the  $\mathcal{O}(\log d)$  bound from previous work on linear quantization.

43 In the experiment, we choose the number of non-zero entries to be a fixed number  $s = 16$  to guarantee that the  
44 parameters in Assumptions 1–4 are fixed, so that we can validate the dimension-free bound from our theorems. The  
45 results showed no dependence on the dimension  $d$ , as we expected. For denser cases with an increasing  $s$ , which would  
46 result in non-fixed model parameters, we do expect the performance to change, like what happened in Figure 2(b). But  
47 this does not contradict the results presented in our work, and is in concordance with what our theory predicts.

48 **Other issues.** •  $\text{dom}(\delta, b)$  denotes the domain of low-precision numbers that can be represented based on parameter  
49  $\delta$  and  $b$ , i.e.  $\{-2^{b-1}, -2^{b-1} + 1, \dots, -1, 0, 1, 2, \dots, 2^{b-1} - 1\}$  times  $\delta$ . • As we showed in Table 1 and analysis  
50 after Theorem 4, non-linear quantization is better than linear quantization by a  $\sqrt{d}$  factor when it comes to the number  
51 of bits needed. • Proving convergence with a constant step size is actually a stronger result than using decreasing step  
52 size. • We included the different setting to show that our work can be applied to various problem classes in real-life  
53 applications. Though the 8-page length limits our discussion to some extent, we were able to present the highlights of  
54 our results and include detailed analysis in the appendix.