

1 **Reviewer #1:**

2 **“limited tweak relative to previous work”** The contributions of CSER put it beyond merely a “tweak”. Our two key
 3 contributions are: (1) A novel mechanism “error reset” that uses arbitrary compressors in a way different from error
 4 feedback. Our theoretical analysis (Theorem 1, Lemma 2, Remark 1,2) shows that error reset achieves better error
 5 bounds than error feedback. Empirical results show better convergence of error reset, especially for high compression
 6 ratios. (2) New combination of partial gradient and model synchronization—in particular, by carefully distributing the
 7 communication budgets between two synchronizations, we can further improve the convergence. As shown in Table 4
 8 in Appendix P30, CSER combining 2 synchronizations works better than only using one of them (CSEA) in most cases.
 9 **“results not very surprising”** Our work is the first to push the compression ratio to 1024 for both worker-to-server and
 10 server-to-worker communication, where CSER shows much better convergence than previous work. Particularly, at
 11 points when EF-SGD and QSparse diverge, CSER converges well. So far, the experiments with largest compression ratio
 12 were proposed by Deep Gradient Compression [Lin et al., 2018], which only considers worker-to-server communication
 13 with the compression ratio ≈ 600 , and it had no theoretical analysis. Our GRBS compressor (Section 3.3) reduces
 14 bidirectional communication with the desired high compression ratio (≥ 256), which is not explored in previous work.
 15 **“exact accuracy results for the Imagenet”** See the table below.

16 **“No accuracy loss”** is conditional. For both CIFAR
 17 and ImageNet, the accuracy is slightly better than full-
 18 precision SGD when $R_C \leq 16$. With larger compression
 19 ratios, CSER has some accuracy loss compared to full-
 20 precision SGD, but performs much better than EF-SGD
 21 and QSparse. We will revise the claim by adding context.

$R_C/\text{Optimizer}$	1	16	32	256	1024
SGD	76.41 ± 0.03				
EF-SGD		76.34 ± 0.06	76.19 ± 0.07	69.73 ± 0.66	diverge
QSparse		76.40 ± 0.05	73.89 ± 0.09	diverge	diverge
CSER		76.53 ± 0.05	76.33 ± 0.06	75.94 ± 0.09	74.93 ± 0.11

22 **“is it using NCCL”** All the algorithms use the same communication library: Horovod with NCCL.
 23 **“How is the experimental setup chosen”** Multiple nodes connected by 10Gb/s Ethernet is a typical setup used in
 24 previous works such as [Lin et al., 2018], signSGD [Bernstein et al., 2019] and EF-SGD [32]. When using single node
 25 with multiple GPUs connected by NVLink, the communication will be extremely fast and compression is less necessary.
 26 In this work, we aim to show that we can significantly reduce the heavy inter-node communication. Indeed, one may
 27 increase number of GPUs per node to do large batch training, but this becomes prohibitively expensive due to GPU cost.
 28 **“WMT/Transformer”** typically uses SGD variants with adaptive learning rates such as ADAM. In this paper we focus
 29 on SGD with momentum without adaptive learning rates. Applying error reset to ADAM is future work.
 30 **“speedups relative to QSparse”** CSER, EF-SGD and QSparse use exactly the same amount of communication, thus
 31 theoretically having the same training time with the same overall R_C . The advantage over QSparse is that CSER
 32 converges much better with the same low amount of communication. Figure 1(e), 2(d) show slightly shorter training
 33 time of CSER because of less memory copy in computation, which is irrelevant to communication overhead.

34 **Reviewer #2:**

35 **“High compression ratios”** are useful when network bandwidth is very low and model sizes are very big.
 36 **“other compressors”** Yes. Our theoretical analysis works for arbitrary compressors.
 37 **“non-i.i.d.”** Yes. Our theoretical analysis already applies to non-iid case. In Assumption 2, we do not assume identical
 38 workers. In our proof (line 404 in Appendix), we only need independence to obtain V_1/n variance.
 39 **how to choose two efficient compressors** Theorem 1 shows how the configurations of compressors affect convergence.
 40 With fixed overall R_C , we can enumerate possible configurations (as shown in Table 3 in Appendix P29) to get relatively
 41 smaller error bounds. To find the best configuration in practice, we do grid search.
 42 **“choose a good β ”** is possible, but irrelevant to communication compression. So we just use the common value 0.9.
 43 **“how to choose the learning rate”** We use grid search to tune the learning rate. Details can be found in Section 5.1.
 44 **“detail advantage of the error reset”** We will highlight the advantages of error reset in the revision.

45 **Reviewer #3:**

46 **only reported on ResNet** Prior work such as Deep Gradient Compression [Lin et al., 2018], EF-SGD [32] and QSparse
 47 [3], all used CIFAR-10/100 and ImageNet + ResNet in the experiments. We choose the same to directly contrast results.
 48 **“relationship between training loss and the configuration of H, RC1, and RC2”** Theorem 1 shows the relationship
 49 between squared gradient norm and configurations ($R_{C_1} = 1/\delta_1, R_{C_2} = 1/\delta_2$). Though we cannot translate the
 50 convergence rate of gradient norm into the one of training loss due to non-convexity, in practice better convergence on
 51 the gradient norm implies faster convergence on the training loss. With fixed overall $R_C = 1/[1/(R_{C_1} \times H) + 1/R_{C_2}]$,
 52 we can enumerate possible configurations to get relatively smaller error bounds, but the optimal configuration is
 53 unknown. To find the best configuration in practice, we do grid search.
 54 **“influence from the number of machines”** This is identical behavior as full-precision SGD, in that changing the
 55 number of machines affects the global batch sizes, thus affects the testing accuracy and requires different learning rates.

56 **Reviewer #4:**

57 **“H=1 as a baseline”** CSER with $H = 1$ is a special case called “CSEA”, which is also novel. The results are in
 58 Appendix P30, Table 4 and subsequent figures. CSEA uses only **one compressor**. In Table 4, we can see that CSER
 59 with 2 compressors outperforms CSEA in most cases.
 60 By **“state-of-the-art”** we mean the best latest work combining both local SGD and compression. We will clarify it,
 61 and add decentralized SGD and PowerSGD to the related work.