1 We thank the reviewers for their insightful comments and encouraging feedback. We hope that the concerns raised are
2 addressed adequately below and that our work will be appropriately re-evaluated.

3 **Speedups (R1)**  Reviewer 1 raises two concerns about speedups which we believe to be based on a misunderstanding.
4 Firstly, our large reductions in communication (e.g. $10\times$) lead to smaller reductions in wall-clock time (e.g. $2\times$). We
5 think this is expected, as all mentioned wall-clock times *include forward and backward passes* in addition to gradient
6 compression and communication. A $2\times$ reduction in this metric seems significant. We will clarify this in the paper.

7 Secondly, the reviewer suspects that Tables 6 and 7 show timings for the slower GLOO backend. Let us clarify that
8 all such timings are measured in default conditions: NCCL, all-reduce, 16 GPUs, and end-to-end as described before.
9 The scaling plots in Figure 3 show speedups of 9.3 (PowerSGD) vs 7.1 (SGD) on Cifar over single worker SGD. This
10 is consistent with the 23% savings (9.3 vs. 7.1) reported in Table 6. We include results for an LSTM (Table 7) for
11 completeness. The LSTM's speedups are better due to their higher communication-to-computation ratio.

12 **Failure cases for PowerSGD (R1)**  We agree with Reviewer 1 that an outline of when PowerSGD works and when it
13 breaks would be helpful. To date, we have not observed any failure cases of the 1-step power iteration in the algorithm.
14 To achieve good accuracy in the same number of steps as SGD, a sufficiently high rank (2 or 4 in practice) is required.

15 **Larger models and clusters (R1, R3, R5)**  We are currently running additional experiments on a larger cluster (64
16 GPUs) with larger models (ResNet-50). This should further test the effect of network latency (R5). Extrapolating the
17 scaling plots in Figure 3, we expect PowerSGD to perform favorably in those conditions.

18 **Convergence of Algorithms 1 and 2 (R1, R3)**  While we do not currently include an end-to-end convergence proof
19 for PowerSGD, each of its core components are well studied. Algorithm 2 (EF-SGD with Momentum) adds momentum
20 to the well-studied EF-SGD algorithm (as in Karimireddy et al. 2019). EF-SGD is guaranteed to converge if the
21 compressor $\mathcal{C}$ satisfies $\|X - \mathcal{C}(X)\|_2^2 \leq (1 - \delta)\|X\|_2^2$. This condition is satisfied by PowerSGD with SVD for best
22 rank-$k$ approximation (see Appendices A.1 and A.2). The cheaper 1-step power iteration with warm start is akin to the
23 famous Oja's algorithm (Oja, 1982) and is empirically shown to yield the same performance as a full SVD.

24 **Linearity of PowerSGD (R3)**  We use the term linearity to mean that PowerSGD on a single worker with gradient
25 matrix $\bar{M} := (M_1 + M_2)/2$ is equivalent to PowerSGD with two workers with their own gradients $M_1$ and $M_2$ (for
26 any number of workers and any split of $M$.) To see that this holds, consider that $P$ in line 4 of Algorithm 1 in the
27 two-worker example amounts to $P = \frac{1}{2}(M_1 + M_2)Q = \bar{M}Q$. The matrices $M_1$ and $M_2$ are never multiplied with
28 each other. The same is true for $Q$ in line 7. This makes PowerSGD just a function of the average gradient $\bar{M}$.

29 **PowerSGD without feedback (R3)**  Because PowerSGD's very-low-rank gradient approximations are coarse, it
30 required error feedback to converge in our experiments. We will include the requested comparison in the Appendix.

31 **GradiVeQ (Yu et al. 2018) (R3)**  We thank Reviewer 3 for pointing us to this interesting work. We will include this
32 method in our discussion.

33 **High Cifar-10 accuracy (R3)**  We use a ResNet-18 based on `torchvision`. Compared to the ResNet-20 model used
34 for Cifar-10 in the original paper, the layers have more feature maps (are wider), explaining the superior performance.
35 He et al. use this wider architecture for ImageNet.

36 **Global batch size (R5)**  Reviewer 5 mentions that some related papers scale the number of workers while keeping
37 the global batch size fixed. Most of the work we are aware of instead keep the local batch size fixed (e.g. Goyal et al.
38 (2017)) since it better utilizes the computational power of the workers. Moreover, if we kept a fixed global batch size
39 the computation performed per bit communicated would decrease with the number of workers—only further favoring
40 compressed algorithms such as ours.

41 **End-to-end speedup results (R5)**  Time-to-accuracy results, as requested by Reviewer 5, can currently be found in
42 Appendix C of the submission. We will consider including these plots in the main paper.

43 Goyal, P., et al. "Accurate, large minibatch SGD: Training ImageNet in 1 hour." arXiv 2017.
44 He, K., et al. "Deep residual learning for image recognition." CVPR 2016.
45 Karimireddy, S.P. et al., "Error feedback fixes SignSGD and other gradient compression schemes." ICML 2019.
46 Oja, E. "Simplified neuron model as a principal component analyzer." Journal of Mathematical Biology, 1982.