1 **Rebuttal for "Communication-Efficient Distributed Learning via Lazily Aggregated Quantized Gradients"**

2 **Reviewer 1.** My main concern on the theoretical analysis. Provide a more fine-grained analysis. The linear rate constant
3 can be tightened/simplified, and will be presented it in the final version. With this tighter analysis, the bound explicitly
4 establishes the dependence of convergence rate on the condition number $\kappa$, namely $\sigma_2 = (1 - \frac{a(\xi)+b(\tau)}{\kappa} + \tau^2 c)^{1/\bar{t}}$,
5 where constants $a(\xi)$ and $b(\tau)$ increase as $\xi$ and $\tau$ decrease, and all $a(\xi), b(\tau), c$ do not depend on $\kappa, L, \mu$. Clearly,
6 as quantized values become precise enough ($\tau^2 \to 0$), LAQ approaches the convergence rate of LAG [6], namely,
7 $\sigma_2 = (1 - \frac{a(\xi)}{\kappa})^{1/\bar{t}}$. If there is no quantization or skipping of communication rounds, $\tau^2 \to 0$, $\xi = 0, \bar{t} = 1$, LAQ
8 converges with $\sigma_2 = (1 - \frac{a(0)}{\kappa})$, same as GD. We hope the reviewer will appreciate the merits of this analysis.

9 **Reviewer 2.** 1. The relation of Lyapunov function to the simple risk. Analysis does not apply to neural networks.
10 Linear convergence of the Lyapunov function also implies that $f(\boldsymbol{\theta}^k) - f(\boldsymbol{\theta}^*)$, $\|\nabla f(\boldsymbol{\theta}^k)\|_2^2$, and $\|\boldsymbol{\theta}^k - \boldsymbol{\theta}^*\|_2^2$, all converge
11 with a linear rate. From the definition of a Lyapunov function, it is clear that $f(\boldsymbol{\theta}^k) - f(\boldsymbol{\theta}^*) \le \mathbb{V}(\boldsymbol{\theta}^k) = f(\boldsymbol{\theta}^k) - f(\boldsymbol{\theta}^*) +$
12 $\sum_{d=1}^{D} \sum_{j=d}^{D} \frac{\xi_j}{\alpha} \|\boldsymbol{\theta}^{k+1-d} - \boldsymbol{\theta}^{k-d}\|_2^2 \le \sigma_2^k \mathbb{V}^0$, meaning the risk error $f(\boldsymbol{\theta}^k) - f(\boldsymbol{\theta}^*)$ converges linearly. The $L$-smoothness
13 results in $\|\nabla f(\boldsymbol{\theta}^k)\|_2^2 \le 2L[f(\boldsymbol{\theta}^k - f(\boldsymbol{\theta}^*)] \le 2L\sigma_2^k \mathbb{V}^0$; hence, the gradient norm $\|\nabla f(\boldsymbol{\theta}^k)\|_2^2$ also converges linearly.
14 Similarly, the $\mu$-strong convexity implies $\|\boldsymbol{\theta}^k - \boldsymbol{\theta}^*\|_2^2 \le \frac{2}{\mu}[f(\boldsymbol{\theta}^k - f(\boldsymbol{\theta}^*)] \le \frac{2}{\mu}\sigma_2^k \mathbb{V}^0 - \|\boldsymbol{\theta}^k - \boldsymbol{\theta}^*\|_2^2$ also converges linearly.
15 Convergence analysis of nonconvex nonsmooth objectives is important, and is included in our future research agenda.
16 2. My only concern is that the experiments have been only conducted on MNIST. If accepted, the final version will
17 report experiments on SUSY, IJCNN1 and COVTYPE, for which the results are also promising.
18 3. Connection with [2], [3], [4]. Compared with the innovation-agnostic random selection [2], LAQ explicitly
19 leverages the gradient innovation in both worker selection and gradient quantization, which will result in more effective
20 communication reduction. LAQ also differs from the dynamic averaging [3] in their design principles, as [3] skips
21 communication when the local model does not differ too much from the global model, while LAQ skips communication
22 when the fresh gradient does not differ too much from the stale one. Lazy aggregation and dynamic averaging can be
23 jointly leveraged to further reduce communication. Developing lock-free LAQ or asynchronous quantized method based
24 on [4] is interesting. Due to asynchrony, the resultant algorithm may require re-deriving the communication conditions
25 (based on a counterpart of Lemma 2), which likely will complicate analysis of the new Lyapunov function. Since it
26 needs careful investigation, we will tackle it in our future work. A discussion with these references will be added.
27 4. The proof of them 1, ..., Lyapunov-functions ... not intuitive. The design of Lyapunov function $\mathbb{V}(\boldsymbol{\theta})$ is coupled with
28 the communication rule (7a) that contains a parameter difference term. Intuitively, if no communication is skipped at
29 the current iteration, LAQ behaves as GD that decreases the objective residual in $\mathbb{V}(\boldsymbol{\theta})$; if certain uploads are skipped,
30 LAQ's rule (7a) guarantees that the error of using stale gradients is comparable to the parameter difference in $\mathbb{V}(\boldsymbol{\theta})$ to
31 ensure its descent. Thus, Lyapunov function always decreases. We will add more intuition in the proof. We incorrectly
32 marked the reproducibility, but will make the code publicly available. Thank you for the favorable recommendation.

33 **Reviewer 3.** 1. The scheme requires additional memory. The extra memory is low. The server stores the last aggregated
34 gradient (dimension $p$), and each worker stores the last gradient (dimension $p$) and $D$ model change norms ($D$ scalars).
35 2. Many critical tuning parameters, such as the step-size $\alpha$ and the Lyapunov parameters $\xi$ are quite complex to tune.
36 With the smoothness $L = 19$ in the simulation setting, our parameters used in simulations slightly violate (17). During
37 the rebuttal period, we conducted a simple experiment with $\xi_d = \xi = 1/160, D = 10$ and $\alpha = 0.01(\rho = 0.01, \rho_2 = 0.5)$
38 satisfying (17). It turns out that the result is comparable with that presented in our submission: test accuracy 0.9082,
39 iteration # 2530, communication # 530, and bit # $1.66 \times 10^7$. To assess the sensitivity of parameters, we tested under
40 variable $\alpha$ (with $D = 10$) and variable $D$ (with $\alpha = 0.02$) values, as summarized below. Here, $\xi = 0.8/D$ and $\epsilon = 10^{-6}$.

| | $\alpha = 0.01$ | $\alpha = 0.015$ | $\alpha = 0.02$ | $\alpha = 0.04$ | $D = 2$ | $D = 5$ | $D = 10$ | $D = 15$ |
|---|---|---|---|---|---|---|---|---|
| **Iter #** | 5219 | 3459 | 2663 | 1410 | 2448 | 2503 | 2663 | 2749 |
| **Comm #** | 825 | 626 | 618 | 1908 | 534 | 512 | 618 | 678 |
| **Bit # ($\times 10^7$)** | 2.59 | 1.96 | 1.94 | 5.99 | 2.67 | 1.61 | 1.94 | 2.13 |
| **Accuracy** | 0.9082 | 0.9082 | 0.9082 | 0.9082 | 0.9082 | 0.9082 | 0.9082 | 0.9082 |

41 3. The increase in iteration counts vs the decrease in communication load. With the updated finer-grained analysis, we
42 can obtain the linear rate constant $\sigma_2 = (1 - \frac{a(\xi)+b(\tau)}{\kappa} + \tau^2 c)^{1/\bar{t}}$, where constants $a(\xi)$ and $b(\tau)$ increase as $\xi$ and $\tau$ decrease,
43 and all $a(\xi), b(\tau), c$ do not depend on $\kappa, L, \mu$. Thus, $\bar{t}\frac{\kappa}{a(\xi)+b(\tau)-\kappa\tau^2 c} \log(1/\epsilon)$ iterations, or $bp\bar{t}\frac{\kappa}{a(\xi)+b(\tau)-\kappa\tau^2 c} \log(1/\epsilon)$ bits,
44 are needed to reach $\epsilon$-accuracy. The compressed GD by Khirirat ("Distributed learning with compressed gradients") in a
45 centralized setup requires the same order of iterations $c_1 \frac{(\mu+\bar{L})^2}{4\mu L} \log(1/\epsilon)$, or $(\log_2 p + bp)c_2 c_1 \frac{(\mu+\bar{L})^2}{4\mu L} \log(1/\epsilon)$ bits; in the
46 distributed setup, this approach yields a near-optimal solution. Uncompressed methods, e.g., GD, entail $\frac{\kappa+1}{2} \log(1/\epsilon)$
47 iterations (fewer than LAQ), but more bits (usually encode a float using 32 bits while $b$ bits in LAQ).
48 4. In Them 1, V decays, what about f or gradient norm? Non-convex case? See our reply to Comment 1 of Reviewer 2.
49 5. Error compensated schemes. The error compensation schemes skip communicating certain *entries* of the gradient,
50 but communicate with *all workers*. LAQ skips communicating with certain *workers*, but communicates *all (quantized)*
51 *entries*. The two are not mutually exclusive, and can be used jointly. We will try to empirically compare with [Alistarh
52 etal'18] in the final version. Thanks for recognizing the novelty of our work.