

1 Many thanks to the reviewers for their deep, thoughtful reviews and constructive suggestions.

2 **Main Contributions:** We would firstly explain the theoretical contribution of our work. Our algorithm is the first to
3 provide a provably communication-efficient and distributional robust FL algorithm, and both the design and the analysis
4 of the algorithm is technically non-trivial.

5 **(R4):** Convergence analysis. To solve the minimax problem: $\min_w \max_\lambda F(w, \lambda)$, our goal is to find the saddle point
6 (w^*, λ^*) such that $\forall w, \lambda, F(w^*, \lambda) \leq F(w^*, \lambda^*) \leq F(w, \lambda^*)$. We use the primal-dual gap: $\max_\lambda F(\hat{w}, \lambda) -$
7 $\min_w F(w, \hat{\lambda})$ as convergence measure, a classic convergence measure used in convex-concave minimax optimization.
8 It measures how much our solution $(\hat{w}, \hat{\lambda})$ is far from the saddle point (w^*, λ^*) (if our solution is exactly the saddle
9 point, then the primal-dual gap is zero). Our theoretical novelties are: (1) Our problem setting is novel. We study the
10 minimax optimization problem where primal variable is trained locally, averaged periodically, and the dual variable
11 can only be updated periodically. (2) To mitigate the issue that dual variable cannot get updated every iteration, we
12 propose a novel gradient sampling method (as pointed by **R2** as a “clever algorithm” which we appreciate). By this
13 gradient sampling update scheme, we have a unbiased approximation on the history gradients of dual variable. (3)
14 We give the convergence analysis on convex-linear and nonconvex-linear settings. Nonconvex-concave (linear) is one
15 of the hardest cases in minimax optimization, and in this work we even have more difficulties, due to the irregular
16 updating scheme compared to a single machine case. In addition, we also give the analysis for regularized setting:
17 strongly-convex-strongly-concave, and PL-strongly-concave cases. We totally agree with **(R4)** that we can not have both
18 strong convexity and bounded gradients assumptions at the same time **unless** the problem is defined over a compact set.
19 In our analysis we assume that parameter domain is bounded (Assumption 3), thus having both assumptions is sound.
20 Currently we cannot remove bounded gradient assumption, because we need to characterize the variance of history
21 gradient estimation of λ . **(R3)** commented that “server choosing clients” mode cannot solve low participation issue.
22 The solution can be that server sends request to more clients and waits until K clients responded. The convergence
23 analysis will be trickier though.

24 For the questions about convergence rate, we admit that to achieve communication efficiency, we sacrifice slightly on
25 convergence rate **(R4)**. However, that does NOT mean our algorithm is slower **(R4)**. The communication cost is the
26 major bottleneck that slows down the distributed training. As shown in the experiments, our algorithm converges faster
27 than vanilla Agnostic Federated Learning (AFL). Compared to gradient descent **(R1)**, in [27] they get $O(1/\sqrt{T})$ rate,
28 slightly better than us, but not communication efficient. As for the choice of τ **(R3)**, it is a very good question. We do
29 optimize on η, γ and τ to get the best convergence rate. The convergence rate is presented as polynomial of τ in the
30 Appendix proofs. **(R4)** also had a comment that the same τ for all clients are not nice since clients will have different
31 amounts of data. We respectfully disagree with this opinion, since each client will only compute gradient on a small
32 fixed mini-batch of data, so it does not matter how many data they have in total.

33 For the privacy issue, **(R2)** proposed a very good question: since server asks client to evaluate its loss at \tilde{w} with data
34 point $\xi = (x, y)$, and return the loss $f_i(\tilde{w}, \xi) = a$, does it mean server can infer the data point information equipped
35 with \tilde{w} and $f_i(\tilde{w}, \xi)$? Let us consider the simple linear regression case. If server wants to infer data point $\xi = (x, y)$, it
36 needs to solve the following problem: $\|\tilde{w}^\top x - y\|^2 = a$ to find out x and y . This problem does not admit a unique
37 solution, or namely, server does not have enough information to determine x and y at the same time. However, we
38 should admit, there is still chance that information could be leaked due to some model inversion attack, but it is not just
39 the problem in our work, but the problem in whole FL community. Thus the suggestions from **(R2)** about adding noise
40 or model compression will be good future directions.

41 For experiment part, **(R1)** pointed out that our dataset partition is not practical and we should follow AFL’s setup. We
42 would argue that we exactly followed them: (1) We use the Fashion MNIST dataset, and adult dataset (in Appendix),
43 the same as them; (2) We partition the dataset by giving one client only one class data, the same as them; (3) We use the
44 full 10 classes Fashion MNIST dataset but AFL only use part of dataset, so our experiments are even more convincing;
45 and (4) We compare the model trained in average domain (FedAvg) and trained by distributionally robust domain,
46 the same as what they did with Fashion MNIST dataset, yet one difference is that we do not train the model fully on
47 local data, but it is already widely observed (even in AFL paper) local model will have very poor out-of-distribution
48 performance, let alone the worst distribution performance. The other question is that in Fig.2, the DRFA has the same
49 performance with conventional FL **(R4)**. This is because Fig.2 is showing the average distribution performance, hence,
50 it is reasonable that in average our algorithm performs similar to FedAvg. We should mention that the main point is that
51 our method is better than FedAvg in worst distribution performance, as shown in Fig.3. **(R4)** commented that our work
52 may not be as efficient as asynchronous FL. We are confused in which aspect our work is worse, can asynchronous FL
53 achieve more distributional robustness? or it can solve the DRO with less communication cost? If so, it would be great
54 if providing related works and we are very happy to compare with them in subsequent version.

55 For some minor questions, Alg.1 L4, 12, 13 should be \tilde{w} instead of w , and in L12, 13 local models should not have
56 bars. L166, it should be minimizer of loss functions not gradients. L187 we missed square sign. L218, we mean single
57 machine case. We would also thank **R1** for different FL protocol papers, and we will discuss them.