

1 Thank you to all the reviewers for the feedback! R2 and R5 are positive about the paper and gave very helpful  
 2 suggestions. We believe that R4 may have a misunderstanding about the algorithm and analysis, and we have tried to  
 3 clarify it below. We hope that R4 will reconsider and increase their score to recommend acceptance. Thanks again!

4 **[R4] Does FedNova decrease the stepsize of each client and slow down the**  
 5 **convergence speed?** FedNova does not degrade the convergence speed of  
 6 FedAvg, as we justify next. Firstly, when clients perform local updates, they use  
 7 the same step size  $\eta$  as FedAvg and not a scaled one  $\tau_{\text{eff}}\eta/\tau_i$ . Secondly, the dif-  
 8 ference between FedNova and FedAvg is the aggregation weights  $w_i$ , which only  
 9 control the direction, and not the magnitude of the accumulated global update.  
 10 The magnitude is determined by  $\tau_{\text{eff}}$ , which is the same in FedAvg and FedNova.  
 11 For example, observe in Figure 1 that the solid green vector ( $\mathbf{x}_{\text{FedNova}}^{(t+1,0)} - \mathbf{x}^{(t,0)}$ )  
 12 has roughly the same magnitude as the FedAvg update ( $\mathbf{x}_{\text{FedAvg}}^{(t+1,0)} - \mathbf{x}^{(t,0)}$ ).

13 **[R4 & R2] What if we force all clients to run the same local steps (e.g., the**  
 14 **minimum of local steps across clients)?** It is true that forcing all clients to  
 15 perform  $\tau = \min_i \tau_i$  local steps (let us call this algorithm FedAvg-min) can also  
 16 ensure objective consistency. However, its convergence rate is *provably worse*  
 17 than FedNova. This is because, in each round, FedAvg-min will go over less  
 18 data samples than FedNova ( $mb\tau_{\text{min}}$  versus  $b\sum_{i=1}^m \tau_i$  where  $b$  is the mini-batch  
 19 size). Using theorem 2, one can show that the convergence rate of FedAvg-min  
 20 is  $1/\sqrt{mT \min_i \tau_i}$ , which is slower than the rate of FedNova  $1/\sqrt{T \sum_{i=1}^m \tau_i}$ .  
 21 Empirically, we evaluate the performance of FedAvg-min and FedNova on the  
 22 synthetic dataset in Figure 2. Observe that FedNova achieves lower loss value  
 23 than FedAvg-min at any round. Another drawback of a fixed  $\tau$  algorithm like  
 24 FedAvg-min is that faster nodes would remain idle in each round while waiting  
 25 for slower nodes. FedNova avoids such straggling delays by allowing nodes to  
 26 make different numbers of local updates.

27 **[R4] Do we need to know  $\tilde{\tau}$  in order to choose a suitable  $\eta$  in Theorem 3?**  
 28 We do not need to know  $\tilde{\tau}$  or  $\tau_i$  beforehand. It is worth noting that  $m\tilde{\tau}T =$   
 29  $\sum_{i=1}^m \sum_{t=0}^{T-1} \tau_i^{(t)}$  is actually the total number of processed mini-batches across  
 30 all clients after  $T$  rounds. Once we have a budget on the total mini-matches  $K =$   
 31  $m\tilde{\tau}T$  to be processed, we can set the learning rate as  $\eta = \sqrt{m/\tilde{\tau}T} = \sqrt{m^2/K}$ ,  
 32 then the optimization error is guaranteed to be bounded by  $\mathcal{O}(1/\sqrt{K})$ . As for the  
 33 upper bound on learning rate, it is only used for theoretical analysis. In practice,  
 34 one always needs to tune the learning rate.

35 **[R4] Is the general analysis framework a marginal contribution?** We believe that the analytical framework proposed  
 36 in Section 4 is an important and impactful contribution, perhaps even more critical than Section 5. This is because:  
 37 1) we identify the objective inconsistency problem in FedAvg by showing that performing the same number of local  
 38 epochs at clients with heterogeneous sizes datasets optimizes a mismatched objective and 2) we provide the first (to  
 39 the best of our knowledge) rigorous understanding of the objective inconsistency problem in federated learning by  
 40 quantifying the non-vanishing gap caused by incorrect weighted aggregation of heterogeneously updated models.

41 **[R4] Extending our theorems to strongly convex case.** We focus on the non-convex case since it is the most practical  
 42 and challenging setting. It is straightforward to extend our analysis to convex or strongly convex cases. For instance,  
 43 one can directly apply Polyak-Łojasiewicz condition to Eqn. (89) in the appendix and obtain an improved rate.

44 **[R2] Is the bias-correction necessary when using local momentum?** Yes, it is necessary because without bias-  
 45 correction, the algorithm will converge to a stationary point of a mismatched objective, the analytical form of which can  
 46 be derived using our framework. We will add some experiments in appendix to further validate this.

47 **[R5] Clarifications on FedProx: hyper-parameters and differences to FedNova.** On the synthetic dataset, we use  
 48 the same model and hyper-parameters as the FedProx paper. We set  $\mu = 1$  because this is the best value reported in that  
 49 paper. On the CIFAR-10 dataset, we tuned the value of  $\mu$  from  $\{0.0005, 0.001, 0.005, 0.01\}$  as stated in the Appendix.  
 50 FedNova with proximal updates is same as FedProx in terms of the local updates, but the aggregation weights  $w_i$  and  
 51 effective steps  $\tau_{\text{eff}}$  are set differently. In our framework (4), the weights and  $\tau_{\text{eff}}$  used in FedProx are given by Eqn. (6)  
 52 while FedNova uses  $w_i = p_i$  and  $\tau_{\text{eff}} = \sum_{i=1}^m p_i \tau_i$ .

53 **[R5] Avoiding confusions on the algorithm name.** Thanks for the suggestion! We will avoid using the term  
 54 ‘normalized gradient’ and clearly state the meaning of normalization in our paper, or even use another term.

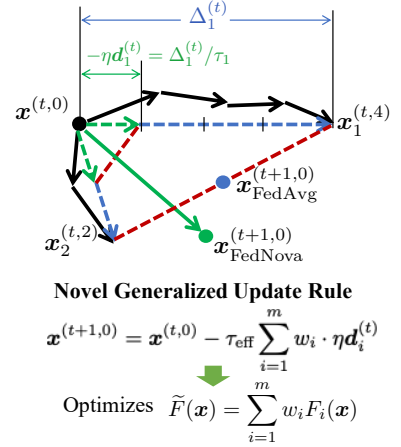


Figure 1: The difference between FedAvg and FedNova is the aggregation weights  $w_i$ , which only controls the direction of the solid green arrow.

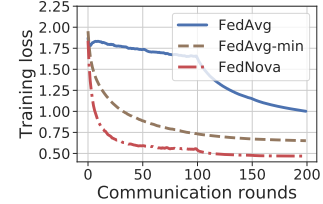


Figure 2: FedAvg-min.