

We sincerely thank the four reviewers for their valuable feedback. We would like to address the following concern common to Reviewers **R4**, **R5** and **R6**: about the nature and originality of our contribution:

Purpose of FED-LAMB: In the context of Federated Learning, in particular in the cross-device settings with a large volume of devices, training deep neural networks is a burning challenge. Given the the potentially cumbersome amount of data present in each device, being able to learn high dimensional and nonconvex models per device is of utmost importance. – The nature of our contribution is thus *to improve the local optimization method for each device* so that a better local model is learned in fewer iterations, leading to *a natural improvement of the communication efficiency* of the federated method, hence requiring less rounds of communication to reach similar accuracy.

R1: We thank the reviewer for valuable comments. A proofreading is being done as we clarify:

Notations and Precisions: \bar{L} denotes the sum of the smoothness constants and is stated in the supplementary material. The function $\phi(\cdot)$ is set to the identity function in our runs. The typo in Corollary 1 has been fixed. The usage of d is replaced by i . λ is a weight decaying parameter similar to the original LAMB method. It is tuned on a grid-search for our experiments. This is added to the paper.

Bounds on Theorem 1: The bounds do depend on the number of devices. We precise that our theoretical results hold when **all** are selected. Thus the total number of devices is n , as defined in (1) and appears in our bound, similar dependence is observed in related works.

Bounds on Corollary 1: The bound does depend on L which is the total number of layers. It also depends on the total smoothness which we included in the \mathcal{O} notation. The dependence on then number of devices n is in the denominator of the RHS, which is in accordance with the bound of local AMS in Chen et. al. 2020.

R3: We thank the reviewer for valuable comments. Our point-to-point response is as follows:

Partial selection of devices: The partial selection of devices has practical virtue which we respected in the numerical experiments. Though, as far as convergence bounds, it is common in the literature to consider the total number of workers participating in each round. Either for simplicity or to avoid cumbersome notations, deriving the result for the general case is not an obstacle for the understanding of the convergence behaviour.

Theorem 1 for multiple local updates: We agree that the assumption that $T = 1$ is rather simplistic. We managed to derive the result for *multiple local updates* in time for the supplementary deadline. Please refer to Theorem 3 in the supplementary for the desired result.

R4: We thank the reviewer for the thorough analysis. Our remarks are listed below:

Various remarks: We agree that strictly speaking, the LAMB technique we introduce in our federated method is not a modification of Adam as in the original paper. Yet, Line 56, we explicitly state that we develop a variant of local AMSGrad using the same layerwise adaptivity technique. We would argue that Local AMSGrad is used as a backbone and that periodic averaging is used as the most efficient way to compute a global model from several local ones.

Comparison with "Adaptive Federated Optimization": We thank the reviewer for the reference. This work presents an extension of ADAM for Federated Learning. While the communication cost is lower than ours, due to the extra communication of v (we do not require a broadcast of m to the server), our method matches their convergence bound. A comparison of their method in numerical runs will be done to reinforce our statement.

R5: We thank the reviewer for valuable comments and typos. Our response is as follows:

Notations: We have added some clarification on several notations in our revised paper. $p_{r,i}^t$ is the ratio computed at round r , local iteration t and for device i . $p_{r,i}^{\ell,t}$ denotes its component at layer ℓ

Comparison with FedBN: We thank the reviewer for this reference that we did not consider. After careful reading of the contribution, we argue that our method is purely on the optimization algorithm aspect of things while FedBN is supposedly a new modeling consideration. We argue that our federated learning method can be used with any model, a large variety for numerical runs and certain model satisfying our assumptions for the theoretical part. Batch Normalization could even be an option on top of Fed-LAMB. Though, it would be interesting to include a comparison with FedBN coupled with a vanilla FedSGD in our numerical runs to compare how better or worse such layerwise adaptivity is when performed on the algorithmic level rather than modeling.

R6: We thank the reviewer for valuable comments. We clarify the following points:

Assumption H5: Specializing the upperbound of the estimation of the second order moment is doable and would lead to similar result. For the sake of simplicity, we assumed a constant upperbound.

Convergence of Fed-LAMB: The purpose of this paper is to improve the communication efficiency of the overall learning method in the federated settings. In other words, our method is better than baseline in the sense that for fewer rounds of communications, it reaches a similar accuracy (in terms of stationary point and objective loss function) than other methods. Fair comparison are thus given when the number of local updates are fixed and equal for each method. Indeed, the difference vanished when T goes to infinity but we claim that for a small number of local iterations, which