

1 We would like to thank the reviewers for their helpful feedback, including the minor concerns that will be clarified but  
 2 which we do not have space to discuss here.

3 **Significance of the graph reformulation:** The augmented graph allows to cast the stochastic decentralized problem  
 4 as a batch decentralized problem with a more complex structure, captured by the matrix  $A$ . In our case, the augmented  
 5 graph contains all the intuition behind the dual formulation (what should the consensus constraints be). More generally,  
 6 we believe that it is a powerful way of reasoning to obtain finite sum algorithms from decentralized algorithms that  
 7 work with subgraphs of communication. Algorithms and rates can then directly be obtained by studying the Laplacian  
 8 of the augmented graph, as done in Appendix D. Yet, it is only needed to obtain the dual formulation.

9 **Extension of past work:** We mainly refer to the extension of the APCG algorithm to arbitrary sampling and to strong  
 10 convexity in subspaces only. This extension is very important for our work because it allows us to choose different  
 11 probabilities for communications and for local computations.

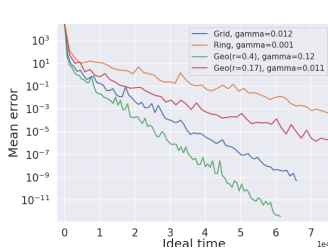
12 **How do we obtain  $W_{k\ell}\Sigma^{-1}y_t$ :** This comes from the coordinate update  $A\nabla_{k\ell}q_A(y_t) = Ae_{k\ell}e_{k\ell}^T A\Sigma^{-1}y_t =$   
 13  $\mu_{k\ell}^2 W_{k\ell}\Sigma^{-1}y_t$ . The  $\mu_{k\ell}$  are indeed related to the communication matrix since the Laplacian matrix of the com-  
 14 munication graph is  $L = \sum_{k\ell} \mu_{k\ell}^2 W_{k\ell}$ . There is a typo and the  $\mu_{k\ell}^2$  of line 188 should not appear.

15 **Hyper parameters:** The only degrees of freedom of ADFS are  $p_{k\ell}$  and  $\mu_{k\ell}$ . The other parameters (such as  $R_{k\ell}$ ,  $\rho$ ,  $\eta$ ,  
 16  $\sigma_A$ ) directly depend on them. For **communication edges**, choosing values of  $\mu$  amounts to choosing a gossip matrix,  
 17 and choosing  $p_{k\ell}$  amounts to tuning how frequently edge  $(k, \ell)$  is sampled. Therefore, choosing uniform  $p$  and  $\mu$  (as in  
 18 the experiments) is efficient as long as the graph is not too heterogeneous. For graphs that have non-regular topologies,  
 19 choosing  $\mu$  better than uniform can be challenging but this is also the case for DSBA and MSDA which consider a  
 20 fixed gossip matrix. Besides, MSDA and DSBA are synchronous algorithms that enforce uniform probabilities for  
 21 all edges (all edges are activated at each step). Note that we would like to choose  $\mu_{k\ell}^2 = p_{k\ell}^2 / (R_{k\ell}[\sigma_k^{-1} + \sigma_\ell^{-1}])$  but  
 22 this is not possible because  $R_{k\ell}$  depends on  $\mu_{k\ell}$ . For **computation edges**, the theory tells us how to set  $p_{k\ell}$  and  $\mu_{k\ell}$   
 23 using an importance sampling scheme. This only requires knowing  $\sigma_k$ , the strong convexity of the problem,  $L_{k\ell}$ , the  
 24 smoothness of individual training examples, and  $\gamma$ , the eigengap of the gossip matrix, which are standard constants  
 25 of decentralized optimization problems. In the end, ADFS can be used with parameters given by theory (as given by  
 26 Theorem 3), **without any extra tuning**. ADFS may perform better with other choices of  $p$  and  $\mu$  in some situations but  
 27 its main competitors are either forced to use uniform values or require the same tuning.

28 **Probability  $p_{k\ell}$ :** The choice of  $p_{k\ell}$  impacts both the convergence rate (precisely captured by Inequality (6)) and the  
 29 average time per iteration (cruder bound of Theorem 2). Contrary to MSDA, the communication  $p_{k\ell}$  can be tuned to  
 30 adapt to the topology of the graph, delays and local condition numbers. Although the optimal choice is only clear for  
 31 regular graphs, heuristics can be designed using Inequality (6). For example, the communication probabilities can be  
 32 chosen as inversely proportional to the edge delay or to the degree of their incident nodes, but this choice is highly  
 33 dependent on the setting. Computing probabilities are chosen optimally in the homogeneous setting but could be set  
 34 differently as well to reduce waiting time, e.g., to reduce the computational burden of a busy node.

35 **Datasets used for experiments:** We favored large networks and datasets with many samples over high-dimensional  
 36 datasets for our experiments. Yet, we believe that the scale of the experiments is already quite significant. As a  
 37 comparison, experiments in the MSDA paper use synthetic datasets of dimension 10 and experiments in the DSBA  
 38 paper use a network of size 10. We use 100 computing nodes and  $10^6$  samples in total, which is significantly larger.  
 39 Finally, serial APCG performs well on RCV1 and ADFS is a decentralized version of APCG, so we expect it to have  
 40 good performances on RCV1 as well. We will add a full set of experiments on RCV1 to a revised version of the paper.

41 **Comparison with Point-SAGA:** As written in the introduction, Point-SAGA is a serial algorithm indeed and we  
 42 apologize if this was not clear enough. Yet, we thought that it was interesting to have it as a competitor because it beats  
 43 many distributed algorithms when the number of machines is not too big, as shown in Figure 3 (a). Our point was that  
 44 ADFS compares nicely with optimal single machine algorithms on one machine (which is not the case of many other  
 45 algorithms) while being able to take advantage of using more machines (Point-SAGA is only designed for one machine,  
 46 as Reviewer 2 pointed out). We will make the distinction clearer by adding “serial” to the legend.



(a) ADFS on several graphs

**Influence of the network topology:** As long as all nodes have roughly the same  
 probability of starting an update (to avoid waiting times), the topology of the network  
 mainly influences the iteration complexity of ADFS through the constant  $\gamma$ . Figure 1a  
 shows experiments with 100 nodes and 500 samples per node on the covtype dataset.  
 For the geographic graph, each node is randomly placed in the unit square and nodes that  
 are at distance less than  $r$  are neighbours. As expected, the convergence rate improves  
 with  $\gamma$ . Yet, ADFS is slower on  $Geo(0.17)$  than on the grid because nodes have uneven  
 degrees so waiting time increases (because all edges have equal probabilities). A  
 non-uniform probability distribution could speed-up convergence for  $Geo(0.17)$ .