

1 We appreciate the positive feedbacks from all the reviewers and provide a detailed response as follows.

2 **R1: “Room for improvement in presenting numerical experiments ... it is better to include test error results”**

3 Our initial intention was to conduct experiments reflecting the convergence rate in the theoretical analysis. We  
4 appreciate your suggestions, and will include the test error in the supplementary material.

5 **R1/R2: “It would be interesting to include results using a variance reduction technique”**

6 We thank both reviewers for mentioning variance reduction techniques, this is indeed very related. To apply  
7 variance reduction methods, we need to further assume that the loss function on each node can be decomposed  
8 into a finite-sum structure, say each loss is a sum of  $m$  functions. In this case, combining IDEAL+SVRG will  
9 provide a complexity of order  $\tilde{O}\left((m + \kappa_f)\sqrt{\kappa_W} \log(1/\epsilon)\right)$ , where acceleration is achieved with respect to the  
10 graph condition number  $\kappa_W$  but not with respect to objective’s condition number  $\kappa_f$  (see also ‘Two-fold acceleration’  
11 on page 6). It is conceivable that there might exist a better design for finite sum structure, e.g., by communicating  
12 the full gradient evaluated in SVRG, etc. This is definitely an interesting direction which we leave to future work.

13 **R2: “I wonder whether the communication cost is lower than that of SSDA in practice. I suggest the authors  
14 to plot the communication cost and computation cost separately.”**

15 When the computation/communication  $\tau$  is large, the dominant term is the communication cost. Hence the third  
16 column ( $\tau = 10$ ) in Figure 1 roughly reflects the communication cost. Thanks for the suggestion, we will include a  
17 separate plot in the supplementary material.

18 **R2: “I wonder whether SSDA+AGD+warm start can achieve the same computation cost by using the proof  
19 technique proposed in this paper.”**

20 We appreciate the reviewer raising this question, this is indeed one of the core messages we would like to convey:  
21 IDEAL improves upon SSDA in a non-trivial way. In short, we have applied the same warm-start strategy as  
22 IDEAL+AGD in the complexity analysis of SSDA+AGD, the higher computation cost of SSDA is due to an  
23 intrinsic weakness of the method. The high level intuition is that the regularization parameter  $\rho$  improves the  
24 condition number of the Moreau-envelope, which reduces the number of subproblems. More explicitly, the number  
25 of subproblems to be solved by IDEAL/SSDA is given by the formula

$$K = O\left(\sqrt{\frac{L_\rho}{\mu_\rho}} \log\left(\frac{C_\rho \Delta_{dual}}{\epsilon}\right)\right), \text{ where } L_\rho = \frac{\lambda_{\max}(W)}{\mu + \rho \lambda_{\max}(W)}, \mu_\rho = \frac{\lambda_{\min}^+(W)}{L + \rho \lambda_{\min}^+(W)}. \quad (\text{Eq. 5 on page 5})$$

26 Ignoring the log factor, this quantity is proportional to the regularized condition number  $\sqrt{\frac{L_\rho}{\mu_\rho}}$ .

- 27 • For SSDA (which is equivalent to  $\rho = 0$ ), we have  $\sqrt{L_\rho/\mu_\rho} = \sqrt{\kappa_f \kappa_W}$ ;
- 28 • for IDEAL, by choosing  $\rho = \frac{L}{\lambda_{\max}(W)}$ , we have  $\sqrt{L_\rho/\mu_\rho} \leq \sqrt{2\kappa_W}$ .

29 Hence, IDEAL saves a factor of order  $\sqrt{\kappa_f}$  compared to SSDA in the number of subproblems. Moreover, with the  
30 proposed choice of  $\rho$ , the cost of inexactly solving the subproblems is essentially the same for IDEAL and SSDA.  
31 Therefore we obtain the improvement in computation cost.

32 **R2: “How to determine the regularization parameter, why choose  $\rho = L/\lambda_{max}$ ”**

33 The global complexity of IDEAL is given by the product of  $K$  (number of subproblems) and  $T$  (number of iterations  
34 for each subproblem), both of them are functions of  $\rho$ . The way we determine the regularization parameter is to  
35 minimize the global complexity with respect to the parameter  $\rho$ , then simplify it in an asymptotic manner.

36 **R2: “References”**

37 Thanks for pointing out these references, we will include them in the revision and some other recent works as well.

38 **R3: “The contribution of this article is strong, but my main issue is I don’t believe this is the best format for  
39 presenting these results... it would be very difficult for me to derive and verify the formulas in this table 2”**

40 We are grateful to the reviewer’s suggestion and will provide a more detailed explanation of this derivation in the  
41 revision. Please kindly check our response to R2’s question starting in line 18 for a high level justification.

42 **R3: “The main text of the article is not sufficient to understand the setup for these experiments”**

43 All the experiments we conducted are for binary logistic regression. For the MNIST experiment, we directly use the  
44 normalized image as input feature. For the CIFAR experiment, the dataset is much more complicated, so a linear  
45 model is not rich enough. Hence, we first apply an unsupervised learning model (convolutional kernel network) to  
46 extract the feature and then apply the logistic regression on top of it. This can be approximately viewed as training  
47 the last layer of a conventional neural network by freezing the (well-trained) first layers. We will improve our  
48 presentation on experimental settings with more details in the revision.

49 **R4: “Experiments are only conducted on one task. More tasks could be more compelling...”**

50 Due to limited space, we have only presented the MNIST experiment in the main paper, however the CIFAR  
51 experiment and other ablation studies can be found in the supplementary material.