

1 We thank all the reviewers for their insightful feedback which do help us improve the quality of our paper. We explain  
2 how we address your concerns and revise our paper based on your comments. Based on **R1** and **R3** concerns about the  
3 code, we are more than happy to share it. If they wish to see the code right away, we can share the code through AC/PC.

#### 4 **Reviewer 1:**

- 5 – “*There are relatively few stable facts. This paper does not necessarily reduce the entropy.*” The reviewer raised a very  
6 important point. We agree that there are tremendous amount of papers on this topic with few stable facts. We expect our  
7 work to bring new insights to this field, especially in understanding the generalization via the lens of differential privacy.
- 8 – “*Figure 2*” We thank the reviewer for pointing this out. We will improve the quality of the plots in the revision.
- 9 – “*broader impact*” This paragraph will be improved to reflect the idea of avoiding over-fitting (see plot (b) Figure 2).

#### 10 **Reviewer 2:**

- 11 – “*I would have liked to see more thorough and rigorous experiments.*” We mainly follow the method in [Wilson et al.,  
12 2017] to tune the step size, since they highlight that the initial step size and the scheme of decaying have a considerable  
13 impact. We agree with the reviewer that the mini-batch size and hyper-parameter tuning would also play an important  
14 role in the performance. Still, we think that our experiments provide an extensive experimental evaluation of variants of  
15 training algorithms for various tasks such as image classification and language modeling. We believe our experiments  
16 offer a fair comparison since the same effort was done to tune the hyper-parameters for each baseline.
- 17 – “*Does RMSProp offer any particular advantage...*” We agree that DPG-LAG/DPG-SPARSE can be used with any first  
18 order optimization algorithm. The RMSProp can be viewed as SGD when  $\beta_2 = 1$ . We plan to provide a generic stable  
19 adaptive algorithm that encapsulates many popular adaptive and *non-adaptive* methods in the Appendix.
- 20 – “*How do the high probability bounds change when using mini-batches of size  $m$ ?*” The high probability bounds on the  
21 gradient mainly follow the generalization guarantee of differential privacy with conditions on the privacy parameters  
22  $(\epsilon, \delta)$  and sample complexity. In the case of mini-batch, the value of privacy parameters  $(\epsilon, \delta)$  and the condition on  
23 sample complexity get modified. We have provided details in the proof of Theorem 5, see Section B.2 of the Appendix.
- 24 – “*Is data augmentation used in the experiments?*” We used data augmentation for MNIST and CIFAR-10. For MNIST,  
25 we normalize the value of each feature to  $[0,1]$ . For CIFAR-10, we normalize, randomly crop and rotate the images.

#### 26 **Reviewer 3:**

- 27 – “*It is unclear how guaranteeing stationary points that have small gradient norms translates to good generalization*”  
28 Our main theoretical results provide the convergence to the ‘*population stationary point*’. Note that Theorems 2, 4 and  
29 5 show the convergence of the norm of the *population gradient* instead of the empirical gradient. Also, one will be able  
30 to use our results to establish the generalization error of the loss function based on arguments such as the PL condition.
- 31 – “*The Hoeffding’s bound holds true as long as the samples are drawn independently*”. Yes, Hoeffding’s bound holds as  
32 long as the samples are drawn independently. However, in the setting of *sample reuse* (setting in this paper) such as  
33 SGD with multi-pass, the reused samples are not independent anymore, since the posterior distributions of samples  
34 change after training on the reused samples.
- 35 – “*The bounds in Theorem 1 have a dependence on  $d$* ”. The reviewer raised a very interesting question! Yes, the  
36 dependence on  $d$  is a known result for differential privacy (DP) and is hard to avoid (see ref. [1]). Some works on DP try  
37 to improve this dependence on  $d$  by leveraging special structures of the gradients. This will be considered in the future.
- 38 – “*do not depend on the initialization  $\mathbf{w}_0$  but on  $\mathbf{w}_1$* .” We thank the reviewer for this typo: should be  $\mathbf{w}_0$  instead of  $\mathbf{w}_1$ .
- 39 – “*For Penn-Tree bank,... algorithms are not stable w.r.t. train perplexity.*” With respect to train perplexity, all methods  
40 stabilize around a target value (which is of course different given the highly nonconvex loss). We note that the test  
41 perplexity increases after several epochs for most baselines while our method keeps a low and steady one.

#### 42 **Reviewer 4:**

- 43 – “*experiment design mainly follows [Wilson et al., 2017]*” The design is different from [Wilson et al., 2017] (except  
44 for the stepsize tuning, see **Reviewer 2**). Indeed, we study the *generalization* performance of each algorithm with an  
45 *increasing* training sample size  $n$  (see Fig. 1, x-axis is  $n$ ). This is consistent with our theoretical results which show the  
46 convergence of SAGD in terms of  $n$ . However, [Wilson et al., 2017] mainly plotted the training/test accuracy against the  
47 number of epochs. We agree that it would be interesting to add experiments to compare SGD with differential privacy.
- 48 – “*SGD with gradient corrupted by Gaussian noise performs well or not*” Excellent question and nice reference! Actually,  
49 one can also use Gaussian noise to design a differentially private algorithm (namely Gaussian Mechanism [7]). Also,  
50 there are papers showing the connection between SGLD (Stochastic Gradient Langevin Dynamics) and differential  
51 privacy. Yet, the existing generalization bound of SGLD is established by the techniques of algorithmic stability [23,  
52 26], which scales with  $(\sqrt{T})$ . We believe it is of great interest to show how Gaussian noise works in our setting. We  
53 will add a discussion in the paper. We consider the theoretical details and experimental results as a future work.
- 54 – “*whether the proposed method works well for small datasets in terms of generalization*” Figure 1 shows that SAGD  
55 has a slightly better test accuracy than other algorithms when the training sample size  $n$  is small (x-axis).