

1 We thank the reviewers for their positive and constructive feedback. We address several points in the review below.
 2 **Connecting Section 5 to DP-SGD.** The bias reduction technique in Section 5 is designed for DP-SGD with clipping.
 3 When it is applied to DP-SGD, the update rule is shown below.

$$x_{t+1} = x_t - \alpha \left(\left(\frac{1}{|S_t|} \sum_{i \in S_t} \text{clip}(\nabla f(x_t) + \xi_{t,i} + k\zeta_{t,i}, c) \right) + Z_t \right),$$

4 where $\zeta_{t,i} \sim \mathcal{N}(0, I)$ is the noise added to reduce clipping bias and $Z_t \sim \mathcal{N}(0, \sigma^2 I)$ is the noise added for privacy (see
 5 (8) in the paper for other notations). The privacy guarantee of this algorithm is the same as original DP-SGD since the
 6 noise added before clipping does not scale with the corresponding L_2 sensitivity and does not provide formal privacy
 7 guarantee. We also compared original DP-SGD and the above algorithm with $\sigma = 1$ in the 3 experiment settings in
 8 Section 5, the results are similar to those in the figure in Section 5 (which used $\sigma = 0$ for all algorithms). We will
 9 include the experiments for $\sigma = 1$ in the paper.

10 **Typos:** Thank you for pointing them out, we will correct the typos.

11 **Reviewer 1. Q1: Discussion on the convergence rate of DP-SGD.** This is indeed a very important point to discuss
 12 in the paper. DP-SGD achieves convergence rate of $O(\sqrt{d}/(n\epsilon))$ in the existing literature. As shown in Theorem 5,
 13 with gradient clipping, the rate becomes $O(\sqrt{d}/(n\epsilon)) + \text{clipping bias}$. When gradient distribution is symmetric, the
 14 convergence rate can recover to $O(\sqrt{d}/(n\epsilon))$. We will add a more detailed discussion around Theorem 5.

15 **Q2: Better connect the clipping bias reduction method with DP-SGD.** Please see "Connecting Section 5 to DP-SGD".

16 **Q3: Additional feedback.** Thank you for the suggestions, we will revise accordingly.

17 **Reviewer 2. Q1: Experiments to support the Theorems in Section 2.** We will add such experiments in the paper.
 18 Theorem 1 is a standard convergence bound for optimization; we will verify it on MNIST. Theorem 2 is our new
 19 contribution, we conducted some new experiments for it. We considered a 1-dimensional example and choose a
 20 symmetric noise $\xi_t \sim \mathcal{N}(0, 1)$ and set $c = 1$. Then we compare $\mathbb{E}_{\xi_t \sim \tilde{p}}[\langle \nabla f(x_t), g_t \rangle]$ (estimated by averaging 10^5
 21 samples) with the lower bound in Theorem 2 for different $\|\nabla f(x_t)\|$. The results below verify our lower bound.

$\ \nabla f(x_t)\ $	0.05	0.1	1	2	10	100
$\mathbb{E}_{\xi_t \sim \tilde{p}}[\langle \nabla f(x_t), g_t \rangle]$	1.7e-4	6.6e-3	0.612	1.83	10	100
lower bound	4e-5	2e-3	0.148	0.3	1.48	14.8

23 **Reviewer 3. Q1: Interpretation and verification about \tilde{p} in practice.** We will provide a few examples to illustrate good
 24 choices of \tilde{p} and add more discussion on what is the desired property for \tilde{p} . Essentially, we want to choose \tilde{p} such that
 25 the Wasserstein distance term in Theorem 5 is small and the probability term in Theorem 5 is large. For verification of
 26 symmetricity, we found the random projection method used in Figure 1 and 2 to be quite effective in practice. It is
 27 interesting to investigate better ways to measure symmetricity for very high-dimensional distributions.

28 **Q2: The symmetric distribution may not hold for neural nets such as LSTM.** This is an interesting point. In general, it
 29 will be very interesting to investigate the geometric structure among the gradients for different types of neural network,
 30 even absent privacy or clipping concerns. We will put this in our future work discussion.

31 **Reviewer 4. Q1: Unbounded bias terms make the bounds less useful.** The bias terms in the bounds are bounded
 32 by constants in worst case, and could be 0 for symmetric distributions. For example, in Theorem 5, the clipping
 33 bias is given by the Wasserstein distance term. For iteration t , $W_{\nabla f(x_t), c}(\tilde{p}_t, p_t) \leq 2\|\nabla f(x_t)\|c$ always holds. More
 34 importantly, $W_{\nabla f(x_t), c}(\tilde{p}_t, p_t) = 0$ if $\tilde{p}_t = p_t$ which can happen when p_t (gradient noise distribution) is symmetric.

35 **Q2: Application of the technique in Section 5 to DP-clipping.** Please see Section "Connecting section 5 to DP-SGD".

36 **Q3: How different the conclusions of this paper are from previous work.** As noticed by the reviewer, the use of
 37 symmetric noise distributions in analysis for clipping is new. The key conclusion of our paper is the 0 (small) clipping
 38 bias for symmetric (approximate symmetric) gradient noise distributions. We also proposed a trick to symmetrify
 39 gradient noise distributions to reduce clipping bias. The most related previous works are Pichapati et al. [2019] and
 40 Zhang et al. [2019]. Both works give bounds on clipping bias essentially similar to $c\|\nabla f(x_t)\|\mathbb{P}[\text{gradient is clipped}]$
 41 which could be large when the clipping probability is large. Compare with these works, we give sharper bounds for
 42 (approximate) symmetric gradient noise distributions. Another closely related work is [1], which was uploaded to arxiv
 43 after the submission deadline. The work proves that clipping can lead to constant regret for DP-SGD, like in Example 1
 44 in our paper. We will add a detailed discussion on comparison with these works.

45 **Q4: Smoothness assumption limits the applicability, symmetricity assumption may not hold for RNNs, show least
 46 symmetric projections.** Smooth assumption is commonly used for DP-SGD in literature and it holds for smooth
 47 approximations for relu such as softplus and swish. Extending the analysis of DP-SGD to non-smooth problems is
 48 important and interesting even absent clipping. The gradient distribution for RNNs is also a very interesting point to
 49 explore. We will mention these in our future work discussion. We will add figures for projections with least symmetry.

50 [1] Shuang Song, Om Thakkar, Abhradeep Thakurta. Characterizing private clipped gradient descent on convex
 51 generalized linear problems. *arXiv preprint arXiv:2006.06783*, 2020.