1  We thank all the reviewers for their valuable comments. We first highlight our technical contributions.

2  **Main technical contributions:** We emphasize that our main contribution is the first result showing that PSGD and
3  SCRN escape saddle-points and converge to local minimizers faster under Strong Growth Condition (SGC) (which
4  is a consequence of interpolation phenomenon). Prior works (e.g., [VBS18]) considered only convergence to critical
5  points under SGC in the first-order setting. We provide our results in both the zeroth and higher order settings. In
6  the zeroth-order setting, we prove a novel concentration inequality for the zeroth-order gradient estimator which is
7  non-trivial and was not know before. We emphasize that this concentration result does not assume the function is
8  bounded; see also Remark 3. For completeness, we also analyzed the complexity for zeroth-order PSGD without the
9  SGC assumption for unbounded functions, which was not done before in the literature. Furthermore, the analysis of
10 SCRN is also significantly involved under SGC (especially in zeroth-order setup); see also Remark 6 and 7.

11 **Rev #1: NTK regime:** NTK viewpoint provides an alternative explanation on efficiency of optimizing algorithms
12 for DNN training. However, a majority of the results based on NTK approach are for polynomially (in depth and
13 sample-size) large-width networks ([1] mention that their polynomial degrees are impractical). Landscape analysis of
14 DNN training (which, roughly speaking is: all local min are global min and so escaping saddle-points and converging
15 to local-min is needed) provides an alternative view for finite-width multilayer neural networks (see e.g., **P2** and its
16 references). Our contributions in this paper are geared towards the later angle – as DNNs are also interpolators in
17 practice, we show that under SGC we can converge faster to local-min. **Other examples:** Another concrete example
18 satisfying SGC condition is online matrix completion (see **P1** for details). We will add this example in detail in our
19 revision. **Contributions:** Please see Lines 2-10 above. **[AZL18]:** At a high-level, the suggested approach would also
20 work. However, the method in [AL18] is a theoretical computer science style reduction approach. It involves a wrapper
21 algorithm on top of PSGD which increases overall runtime. Our result directly analyzes PSGD iterates, more in line
22 with optimization and machine learning type results. Furthermore, a main drawback of the approach in [AZL18] is that
23 it is not directly applicable for the 0th-order setting due to their bounded variance (of stochastic gradient) assumption.
24 Please also see the discussion in Remark 5 for more details. **Variance reduction (VR):** A motivation for the SGC
25 assumption is that it provides *automatic VR*. As in Table 1, for stochastic setting under SGC, PSGD already achieves the
26 corresponding complexity of its deterministic counterpart (without acceleration). It is interesting future work to examine
27 if similar results hold for finite-sum setting. However, most VR methods invariably involve double-loop algorithms and
28 their empirical performance in deep learning has come under close scrutiny in the recent past; see [DB19] and [Sch20].

29 **Rev #2: Intuitions:** SGC at a high-level could intuitively be described as an assumption satisfied in interpolation
30 models, providing *automatic variance reduction*. Please also see Lemma 3.1 and Remark 1 for more intuition. We
31 will clarify this more in our revision. **Detailed distribution of stochastic grad:** We clarify that *we do not make any*
32 *distributional assumption* on the stochastic gradient. SGC is only a variance/moment based assumption. **Finite-sum**
33 **opt:** Finite-sum opt is a special case of stochastic setting and hence the assumptions are satisfied; see also [MBB18,
34 VBS18]. **Proof Sketch:** Thanks for this suggestion. We will add a proof sketch in our revision. **High-order method:**
35 Making higher-order methods practical is an interesting future work. We will clarify this point and reorganize.

36 **Rev #6: Theoretical contributions:** Please see Lines 2-10 above. **Approximate SGC:** Thanks for raising this
37 extremely interesting question. Considering 1st-order setting, departures from SGC can be modeled, for example, by:
38 $E(\|\nabla F(x_t, \xi)\|^2) \le \rho \|\nabla f(x_t)\|^2 + e_t^2$ where $e_t^2 > 0$ is an additive non-vanishing iteration-dependent noise variance;
39 see also [VBS18]. This assumption could be used for example to model certain specific types of label-noise. When
40 $e_t^2 = \sigma^2, \forall t$, the oracle complexity is $\tilde{O}(\epsilon^{-4})$ since this case is essentially equivalent to the standard stochastic gradient
41 setting and SGC has no effect. But one could obtain rates approaching $\tilde{O}(\epsilon^{-2})$ depending on the decay of the term
42 $R_t = \sum_{i=1}^{t} e_i^2$. If we assume $R_t \approx t^\alpha$ for $\alpha \in (-\infty, 1]$, then we can show that the complexity is $\tilde{O}(\epsilon^{-\max(2, 3.5 + 0.5\alpha)})$.
43 Hence, when $\alpha = -2$, the complexity is $\tilde{O}(\epsilon^{-2.5})$ and when $\alpha \le -3$, it is $\tilde{O}(\epsilon^{-2})$. In this setup, it may be possible to
44 connect label noise (for specific noise models) to $\alpha$. This would provide a possible way of incorporating label noise
45 in this setup. It is intriguing to examine this problem rigorously as future work. We will clarify this in our revision.
46 **Experiments:** We will be happy to add simulations and real-world experiments on DNN and Online matrix completion
47 in our revision. We clarify that our main goal in this work is to provide a plausible explanation for the question *why*
48 *do optimization algorithms for deep-learning models work efficiently in practice despite the associated non-convexity*
49 *?*. As an attempt, [MBB18], [VBS18] and related works proposed the SGC condition which is satisfied due to the
50 interpolating nature of DNNs. However, prior works fell short of providing a complete explanation as they only analyze
51 convergence to critical points. Recently, it has been shown in several works that for DNNs (in the finite-width regime),
52 all approximate local minima are also global minima (see reference **P2**) due to which converging to local-min and
53 escaping saddle-points are important. Hence, based on these motivations, we show that PSGD and SCRN converge to
54 local-min faster with SGC condition. We emphasize that this line of work is only one plausible explanation for the
55 above question and there are other directions (e.g., NTK) attempting alternative explanations in the literature.

56 **References: P1** - *Provable Efficient Online Matrix Completion via Non-convex Stochastic Gradient Descent*, Jin et al.,
57 NeurIPS 2016. **P2** - Elimination of All Bad Local Minima in Deep Learning, Kawaguchi et al., AISTATS, 2020.