

1 **Response for paper ID 8808**

2 We thank all reviewers for their thoughtful feedback. Please find detailed responses to your comments below.

3 **Rev1:**

4 Thank you very much for carefully reading our paper and your supportive comments.

- 5 • **Unenglish sentences that could be cleaned up:** We will make the manuscript proofread by a professional English
6 editing service, which we believe resolves grammatical issues.

7 **Rev2:**

8 Thank you very much for carefully reading our paper and your supportive comments.

- 9 • **Theorem 2 is a lot to unpack. Assumptions 1, 2 should be made easier:** The analysis tends to be complicated
10 because the infinite dimensional Langevin dynamics requires a few involved conditions such as smoothness on the
11 objective functions. We are doing our best to make the results presented as intuitively as possible. We will add more
12 digested expositions to the assumptions and convergence rate analysis.

- 13 • **Application to non-residual deep network. It is important to discuss this:** Thank you very much for pointing
14 our an important point. Indeed, our approach *can* be applied to non-residual deep network. The reason why we
15 presented ResNet is just due to space limitation and the fact that ResNet has a continuous-depth representation and
16 such a continuous depth representation is also an interesting application of our analysis. Since we had this application
17 in our mind, we have presented only ResNet in the main text. We will add some more comments about this point in
18 the final version.

- 19 • **I would encourage the authors to improve the exposition:** Thank you for your suggestive comment. We tried to
20 keep technical details as much as possible so that there does not occur confusion and misunderstanding. However, as
21 you pointed out, we would like to use more spaces for intuitive expositions and move some details to the appendix.

22 **Rev3:**

- 23 • **Is it possible to derive a high probability bound on the training loss of one training trajectory?:** Yes, the most
24 direct way is to apply the Markov’s inequality. Moreover, the expectation with respect to the training sample
25 observation is derived from a exponential tail probability bound and thus we can derive a high probability bound with
26 respect to sample observation. As for the training trajectory, the mixing time of the dynamics is fast (exponential
27 with respect to the iteration) and thus as the iteration number increases, the probability in which the trajectory does
28 not contain a “nice” solution satisfying the risk bound decreases exponentially to 0. Since the high probability bound
29 makes the statements complicated (the technical contents are already a bit involved), we have shown the expectation
30 bound for simplicity. We will add more comments on the high probability bound in the final version.

- 31 • **How is ρ_0 reflected in your results on excessive risk?:** As the algorithm progresses, the solution “forgets” the
32 initial solution exponentially fast. In that sense, it does not affect so much on the excess risk. On the other hand, the
33 concentration function is characterized by the relative location between the optimal solution and \mathcal{H}_K , the geometry
34 of $L_2(\rho_0)$ indirectly affects the excess risk through the shape of \mathcal{H}_K . However, it is highly problem dependent.

- 35 • **Assumption 1 (ii), (iii) :** The two layer neural network model presented in the excess risk bound satisfies these
36 conditions (Eq.(8) with bounded input $\|x\| \leq D$ (a.s.) and smooth loss). More specifically, under the setting of
37 Theorem 2, Assumption 1 is satisfied.

- 38 • **ResNet in line 188-192 is strange:** We would like to remark that this is a standard definition where the resid-
39 ual blocks are two layer neural networks. Each layer ℓ receives an output from the previous layer as x_ℓ and
40 it outputs $x_\ell + g_\ell(x_\ell) = (I + g_\ell(\cdot))x_\ell$ to the next layer where g_ℓ is a two layer neural network given as
41 $g_\ell(x) = \int a_{w,\ell} \sigma(W(w, \ell)^\top x) d\rho_0(w)$. This formulation is standard in theoretical analyses of ResNet (e.g., [1, 2]).

42 **Rev4:**

- 43 • **Eigen decay may be a strong assumption since it directly helps to erase the dependence of the dimension for**
44 **the bound. With regularizer, it is always convenient to obtain certain tight generalization bounds:** Indeed, as
45 you pointed out, regularization is the most essential ingredient to obtain a width free generalization bound. Conversely,
46 we can not expect nice generalization without any regularization. Although a global optimal solution for non-convex
47 loss does not directly indicate good generalization, our analysis connects generalization and algorithmic convergence,
48 which we believe is an interesting point. In a real deep learning, we consider that such a regularization is imposed
49 through several explicit/implicit regularizations.

50 **References**

- 51 [1] Y. Lu, C. Ma, Y. Lu, J. Lu, and L. Ying. A mean-field analysis of deep ResNet and beyond: Towards provable
52 optimization via overparameterization from depth. In *Proceedings of ICML2020*, pages 137–147, 2020.
- 53 [2] E. Weinan, J. Han, and Q. Li. A mean-field optimal control formulation of deep learning. *Research in the*
54 *Mathematical Sciences*, 6(1):10, 2019.