

1 We thank the reviewers for their helpful feedback. The reviews emphasized the significance of our results for the
 2 analysis of finite-width neural networks (**R1**, **R3**) and the strength of our technical contributions (**R1**, **R2**). One primary
 3 concern was that our results are only asymptotic. We address this concern and other questions below.

4 **1. Theory**

5 **R2** & **R4**: *Being asymptotic in both n and t , Theorem 3.3 is not very interesting.* We agree that nonasymptotic bounds
 6 would be highly desirable, but such results are difficult and have only been established in limited settings that are not
 7 applicable to neural network training (e.g. weak interactions in [1]) beyond the *lazy* regime. Despite being asymptotic
 8 in n (just like the classical CLT), our results give insights on the evolution of the deviations from mean-field limit and
 9 set the groundwork for nonasymptotic results. In particular, they reveal a benign dependency on ambient dimension
 10 (see point below). Regarding the dependence on t , existing bounds on the fluctuations (cf. [2]) grow exponentially in t ,
 11 whereas we show a finite bound as $t \rightarrow \infty$. We also discuss the decay of the fluctuations at finite t in the first remark
 12 under Sec. 2 below. Finally, we check the validity of our results in numerical experiments under finite n and t .

13 **R1** & **R3**: *Assumption on \tilde{D} .* We assume it is a *closed manifold* – a compact manifold without boundary – like a sphere.

14 **R1**: *Do the assumptions on activations hold for tanh and erf?* Yes. The only requirements are the Universal
 15 Approximation Theorem and our smoothness assumptions in Sec. 2.1 in the paper.

16 **R2**: *Is Prop. 2.1 tangential?* Not really - it shows by adding the regularization we can control the 2-norm of the loss
 17 minimizer (consistent with Col. 1 & 5, Row 2 in Fig. 1), on which the fluctuation bound in Thm. 3.3 crucially depends.

18 **R2**: *The bound in Theorem 3.3 is not truly dimension free.* Our bound depends on the norm of the target function
 19 in the variation-norm space introduced in Sec. 2.1. As the reviewer points out, the dependency on dimension is
 20 implicit through this norm, which is shown in [3] to depend polynomially in the dimension for functions with hidden
 21 dependency on low-dim structures.

22 **R2**: *Results under time discretization are lacking.* Indeed, we leave this for future work, though our continuous-time
 23 analysis already yields insights. Moreover, we see from Col. 3 of Fig. 1 that the loss evolutions under gradient descent
 24 nicely agree across different n , showing empirical consistency of the discretization scheme.

25 **R3**: *Are $\limsup_{t \rightarrow \infty}$ and $\lim_{n \rightarrow \infty}$ exchangeable in Thm. 3.3? How about changing n to n^α for $\alpha \in (0, 1)$?* Not in our
 26 result (including in the latter case), unfortunately, as the $O(n^{-1})$ scaling of the fluctuations in n at finite time (thanks to
 27 CLT and the continuity of the flow) may not be preserved at the $t \rightarrow \infty$ limit. This is worthy of future investigations.

28 **R3**: *Concrete examples for assumptions (71) and (131).* We show in C.1.1 and C.2.1 that (71) and (131) can be satisfied
 29 if the flow $\Theta_t(\theta)$ converges at a uniform rate of $O(t^{-\alpha})$ with $\alpha > 2$ and $\alpha > \frac{3}{2}$, respectively. Also, an alternative to
 30 (71) is $\int_0^\infty (\mathcal{L}(\mu_t))^{1/2} dt < \infty$, and so a sufficient condition is for the loss value to decrease faster than $O(t^{-2})$.

31 **R3**: *More explanations of \mathbf{T}_t .* We will add: $\forall \theta \in D$, $\mathbf{T}_t(\theta)$ captures the deviation of the flow $\Theta_t(\theta)$ due to the “initial
 32 deviation” ω_0 , i.e. $\mathbf{T}_t = \lim_{n \rightarrow \infty} n^{1/2}(\Theta_t^{(n)} - \Theta_t)$. \mathbf{T}_t satisfies an infinite-dim. linear ODE, (86). To control it, we
 33 show that 1) the asymptotic linear operator is PSD; 2) the source term lies in the range of the linear operator; and 3)
 34 finite time perturbations are controlled.

35 **2. Experimental results**

36 **R1**: *Why do the average fluctuations decay with t ?* While Thm. 3.3 only speaks about the $t \rightarrow \infty$ limit,
 37 we can study the long-term behavior of $\mathbb{E}_0 \|g_t\|_{\hat{\nu}}$ by analyzing the t -asymptotic version of (25) or (86),

38 $\dot{\mathbf{T}}_t = -(\mathcal{A}_\infty^{(K)} + \mathcal{A}_\infty^{(V)})\mathbf{T}_t + \mathbf{b}_\infty$. (Note that this also describes the exact dynamics of the fluctuations if we
 39 set $\mu_0 = \mu_\infty$.) In the unregularized case, we expect that f_∞ interpolates the data, and hence $\mathcal{A}_\infty^{(V)} = 0$. Thus,

40 the solution to above is $\mathbf{T}_t = (1 - e^{-t\mathcal{A}_\infty^{(K)}})(\mathcal{A}_\infty^{(K)})^\dagger \mathbf{b}_\infty$. Also, in the ERM setting, $\mathcal{A}_t^{(K)}$ is a PSD operator with
 41 finitely many nonzero eigenvalues, and hence its nonzero eigenspaces are spanned by eigenfunctions v_1, \dots, v_k

42 associated with eigenvalues $\lambda_1, \dots, \lambda_k > 0$. By Lemma C.3, we can express $\mathbf{b}_\infty = \sum_{i=1}^k c_i v_i$. Using (26), we

43 get $\|g_t\|_{\hat{\nu}}^2 = \|\bar{g}_\infty\|_{\hat{\nu}}^2 - \langle \mathbf{b}_\infty, (I - e^{-2t\mathcal{A}_\infty^{(K)}})(\mathcal{A}_\infty^{(K)})^\dagger \mathbf{b}_\infty \rangle = \|\bar{g}_\infty\|_{\hat{\nu}}^2 - \sum_{i=1}^k \lambda_i^{-1} (1 - e^{-2\lambda_i t}) c_i^2$, which decreases
 44 monotonically in t . This is consistent with the long-term behavior of the fluctuations in Col. 2, Rows 1 & 3 in Fig. 1.

45 **R2**: *Why do the fluctuations decrease faster for smaller n ?* This is due to finite- n effects in the fluctuation dynamics,
 46 which are captured in (63) but not (25), and which decrease as n grows.

47 **R2**: *Why the alignment or lack thereof with teacher network’s TV- and 2-norm in Fig. 1?* When (and only when)
 48 regularized, both the TV- and 2-norm of the student are controlled by those of the teacher, consistent with Prop. 2.1.

49 **R2**: *Why are some of the behaviors plotted in Fig. 1 non-monotonic?* Thm. 3.3 does not guarantee monotonic decay of
 50 the fluctuations during finite time, but only prescribes its behavior as $t \rightarrow \infty$. The calculations above of the fluctuation’s
 51 decay is also an asymptotic analysis. As for the norms, their non-monotonic evolution indicates that the regularization’s
 52 effect become relatively stronger later in training, when the function reconstruction loss is low.

53 **R1**: *In Col. 3 of Fig. 1, why is the loss independent from width?* We first note that Col. 3 plots the training / population
 54 loss, as the study of generalization is beyond the scope of this paper (which is why we don’t distinguish between ν and
 55 $\hat{\nu}$). The good agreement of loss evolution for different n validates empirically the convergence to a mean-fields solution.

56 **R1**: *How do we choose the teacher’s neuron in Fig. 1?* We chose the teacher’s neurons to have $c = 1$ and \mathbf{z} randomly
 57 sampled on the hypersphere. For the plots on Col. 1, we chose one of the neurons as the “marker”.

58

59 [1] Durmus et al. “An Elementary...”. [2] Mei, Misiakiewicz, Montanari, “Mean-field...”. [3] Bach, “Breaking...”.