

1 We thank all reviewers for the insightful and encouraging comments. Below we provide a point by point response to
2 Reviewers 1,2,3 (R1, R2, R3).

3 **R1 + R2:** [Difference from “lazy training”]. One key contribution of our work is that we identify constancy of tangent
4 kernel, first observed in Jacot et al. (2018), as related to the scaling of the Hessian norm. Notice that the constancy of
5 the tangent kernel cannot be explained from the point of view of “lazy training”: when the last layer is non-linear, the
6 change of parameter from the initialization is of the same order as for the linear case, but the tangent kernel is no longer
7 constant along the optimization path, as the Hessian norm is no longer small (see Section 4).

8 **R1:** *Theorem 3.2 and results in Appendix G has been proved previously (e.g., [1]). ... I tend to think of this paper as a*
9 *summarization of the previous line of NTK papers. ([1] Sanjeev Arora, et al. On Exact Computation...)*

10 Our Hessian analysis results, including Theorem 3.2 and Theorem 3.1, are new. Note that previous works, including
11 [1], only analyze the tangent kernel matrix, which is first order. In contrast, we analyze the Hessian, a second order
12 differential operator. We note that in some related works, including [1], the notation H stands for the tangent kernel,
13 while we use it to denote the Hessian matrix. This can perhaps cause confusion. Our novel contributions include
14 identifying the underlying reasons for constancy of NTK (small Hessian norm, as opposed to “lazy training”), and the
15 finding that NTK is **not constant** when the last layer is non-linear, even in the infinite width limit (Section 4).

16 **R2:** *The paper mostly focuses on squared loss while widely applied NNs use softmax-cross entropy loss. I would*
17 *encourage putting discussion on ... optimization of those networks in the context of this paper’s result.*

18 Thank you for the suggestion. The main focus in this submission is to uncover the underlying reasons for the constancy
19 of NTK (which depends only on the model, rather than the loss function). Still, it is an important issue, we will add a
20 discussion.

21 **R2:** *It (supplementary B) does have significant overlap with current submission ... may be subject to dual submission ...*
22 *or just cite a separate paper distinguishing contribution.*

23 There is no dual submission issue as the supplementary B has not been submitted to NeurIPS or any other confer-
24 ence/journal. Given the space constraints, it does not seem feasible to have a full discussion of the optimization-related
25 issues. For the final version, we are planning to cite the optimization results as a separate document.

26 **R2:** *What is the important point of emphasizing Euclidean norm change is $O(1)$? ... Should I understand the point to be*
27 *while literature casually talks about “small weight change”, one should be aware that in Euclidean norm it could be*
28 *$O(1)$ due to large dimensions?*

29 Yes. The measurement of the change of parameters from initialization depends on the norm. In dimension m , the
30 Euclidean norm and infinity-norm can be different by a factor of \sqrt{m} . When dimension increases, the infinity norm of
31 the difference from the initialization to the solution converges to zero. However, the Euclidean norm of the difference is
32 always $O(1)$. Importantly, the remainder term of the Taylor expansion (and hence the constancy of TK) is controlled by
33 the Euclidean norm of the difference, not the infinity norm. In contrast, “lazy training” suggests that the optimization
34 path stays close to the initialization point. We will clarify this in the paper.

35 **R2:** *Is the condition (b) in Theorem 3.1 violated for these non-linearities (softmax, maxout) and does not have a*
36 *constancy guarantee? Do authors believe these networks would not have constant tangent kernels*

37 This is an open question so far. These non-linearities do not fit in our current analysis. If softmax is in the output layer,
38 our analysis in Section 4 shows that the model does not have a constant tangent kernel.

39 **R2:** *L251: logically violation of conditions in Theorem 3.1 does not necessarily lead to breakdown of linearity since*
40 *Theorem 3.1 is not if and only if statement, correct?*

41 Yes. The conditions in Thm 3.1 are sufficient but not necessary. But note that the neural networks that are shown to
42 have constant NTK satisfy these conditions. Hence, Thm 3.1 is enough to explain the phenomenon of constant NTK.

43 **R3:** Thank you for the positive comments and the helpful suggestions.

44 **R3:** *Could the authors comment on the applicability of their results assuming a hinge loss? It seems that the hinge loss*
45 *ought to satisfy conditions in Eq(13) and Eq(14), given that its second derivative is zero almost everywhere.*

46 First, note that the Hessian is defined for the model, not the loss function. For hinge-loss-like activation functions, the
47 Hessian is zero at most locations, but it is infinite at the hinge point. We note that the Hessian affects the linearity of
48 the model (i.e., constancy of tangent kernel) through Taylor expansion, see Proposition 2.1 and its proof. Hence, the
49 tangent kernel, which is the first-derivative, depends on the integral of the Hessian, which is second-derivative. This
50 infinite Hessian for hinge-like activation functions is likely to have a non-trivial contribution to the tangent kernel after
51 integration, implying non-constant tangent kernel.