



Figure I: (A) MLP α fitting. (B + C) Training loss curves at early and late stages for optimal primal/dual ResNext models (orange and green) and baseline (blue), trained on CIFAR100 (for details see section 4.2.1 in the main text)

1 **R1+R2+R3+R4:** We would like to thank the reviewers for their thoughtful comments, suggestions and constructive
 2 feedback. Typos and missing citations as pointed out by all reviewers will be addressed in future versions. The main
 3 point of contention articulated by all 4 reviewers is the validity of $\text{var}(\mathbf{K})$ as a good proxy for generalization and
 4 the correctness of Proposition 1 in the general case, and so we will address these in detail. **var(K) as a proxy for**
 5 **generalization:** We agree with the reviewers that this is a strong assumption, and our main hypothesis is indeed to
 6 use $\text{var}(\mathbf{K})$ as a proxy for *trainability*, as stated in the abstract (line 12) and introduction. As a direct evidence, we
 7 additionally show in Fig. I above that a lower training loss is achieved with the primal formulation, and the loss is
 8 roughly maintained in the dual formulation. In the ablation study and section 2, we empirically show how lower $v(\mathbf{k})$
 9 translates to better performance in the tested models under practical scenarios. We will eliminate the claims regarding
 10 $v(\mathbf{k})$ and generalization in the main text, and instead leave the precise characterization of it as future work. Meanwhile,
 11 we mention [A] where a rigorous connection is made between $v(\mathbf{k})$ and generalization (see section 3), advocating
 12 for ensembling as variance reducing alternatives. **Validity of Proposition 1:** Additional experiments on its validity
 13 (including Fig. I(A) above) will strengthen the paper, and will be added to the appendix. Note that in the context of
 14 section 2, its correctness in general is less important than its usefulness. In fact, any formula $v(\mathbf{k}) = f(n)$ that can be nicely
 15 fitted to Monte Carlo simulations in algorithm 1 can be used to derive optimal ensembles. We find that Proposition 1
 16 holds nicely in general. **R1:Undefined notations.** K_∞ is indeed the infinite width kernel. We will make this explicit in
 17 future versions. **R1:Gradient flow assumption.** We agree with the reviewer. However with some additional work the
 18 $\mathcal{O}((nm)^{-1})$ results can be derived in the infinite width limit for gradient descent with a small learning rate. We leave
 19 this extension to future work. **R2:Ablation study.** The ablation study in the paper is conducted on mnist, however
 20 figures 6,7 clearly demonstrate the effect of variance and capacity on larger datasets with modern architectures. Note
 21 that we do not make the claim that capacity in itself is harmful in the general case. **R2:Full Imagenet results.** We were
 22 unable to produce the full Imagenet results in time for the deadline of the main text, and so they were included in the
 23 appendix. In future versions they will be included in the main text. **R3: Connection of var(k) to generalization.** We
 24 do not see our results being at odds with previous work since we are not advocating for eliminating variance all together
 25 or using small learning rates. (thus eliminating the prospects of representation learning). Please see our above response
 26 for further details. **R3:Proposition 1 as an assumption.** We agree with the reviewer and will change its definition to
 27 an assumption in future versions. **R3:Lemma 1.** The proof of Lem 1 in the appendix mentions the finiteness of first
 28 moments (line 47), and we will add more details on that matter in future versions. **R3+R4:On Lemma 1 and Theorem**
 29 **1.** The discussion of NTKs in the context of collegial ensemble requires us to expand the notion of the "kernel" regime,
 30 since width can be interchanged with multiplicity (number of models in the ensemble), to which section 1 is dedicated.
 31 The $\mathcal{O}(1/m)$ result in The 1 is non-trivial, implying that large collegial ensembles can approach the kernel regime
 32 with linear, instead of quadratic dependence on the number of weights even after training, where the independence
 33 assumption of the ensemble models is broken. The novelty here is both technical and in the final result, and we feel will
 34 be appreciated by the community. **R4:On α in proposition 1.** For a given architecture, α represents the rate in which
 35 $v(\mathbf{k})$ changes with width. Computing its theoretical value for fully connected networks is a technical challenge that is
 36 beyond the scope of the paper (see [B]). Empirically, any architecture parameterized by its width can be fitted with α as
 37 described in Algorithm 1. **R1+R2+R3+R4:Additional experiments that reviewers may have missed** Our proposed
 38 method can be adapted to efficiency metrics other than parameter count. In section E of the appendix we adopt a FLOPS
 39 metric instead and show how for a given FLOPS budget, optimally performant ensembles can be derived.

40 REFERENCES

41 [A] M. Geiger et al: Scaling description of generalization with number of parameters in deep learning. (Journal of
 42 Statistical Mechanics Theory and Experiment 2020). [B] B. Hanin et al: Finite Depth and Width Corrections to the
 43 Neural Tangent Kernel. (ICLR 2019).