

1 We thank the reviewers for their positive feedback and comments. We will incorporate their suggestions into the next  
2 version of the paper. We address points made by each individual reviewer in turn.

### 3 Reviewer 1

4 • In the revision, we will present simplifications of Theorem 2 and Corollary 4. For instance, if we fix  $n, \delta$ , and  $\epsilon$  and  
5 choose  $\eta$  so that  $N_\eta = 1$  (i.e., at least one model is  $\eta$ -similar to all others), then Theorem 2 bounds the number of  
6 testable models as

$$k \leq \frac{c_1}{(1 - \eta)^{c_2}},$$

7 where  $c_1, c_2 \geq 0$  are constants that depend on  $n, \delta$ , and  $\epsilon$ . Then, as the similarity  $\eta$  of the model collection grows,  
8 the number of testable models  $k$  grows as well.

9 • In Corollary 4,  $\alpha$  quantifies the strength of the similarity assumption. For  $\alpha = 0$ , there is no similarity requirement,  
10 and the confidence interval is wide. However, as  $\alpha$  grows, the similarity requirement becomes restrictive while the  
11 confidence interval becomes increasingly tight. We will make this dependence clear and simplify the stated bound.

12 • In the adaptive query case, the  $(n + 1)^{k-1}$  bound on the number of queries is standard in the adaptive data analysis  
13 literature. The bound is indeed conservative and reflects the worst-case behavior of the analyst. Our analysis helps  
14 illustrate why this worst-case behavior does not attain in practice. If the queries are similar, then the effective  
15 number of queries can be much smaller than  $(n + 1)^{k-1}$ .

### 16 Reviewer 2

17 We begin with the main point of Reviewer 2. The reviewer is concerned that the paper’s conclusions could be derived  
18 solely from simulation work rather than mathematical analysis. Simply put, **the  $y$ -axes of our plots heavily rely upon  
19 our theoretical contributions.** As a result, we are not aware of a way to derive our paper’s conclusions purely from  
20 simulation. Since the review is not specific about what types of simulations would yield similar insights into overfitting,  
21 we would appreciate a clarification in the updated review.

22 To expand on this further: the bulk of the theoretical content is careful numerical calculations to obtain sharp  
23 generalization bounds where the improvement due to similar models is apparent on data from ImageNet and CIFAR-  
24 10. The standard union bound does not take advantage of model similarity; the refinement in equation (3) is the  
25 basis for our calculations that highlight the beneficial effect of model similarity. These sharp bounds allow us to  
26 empirically demonstrate that model similarity offers protection against overfitting on heavily used benchmark datasets.  
27 Consequently, the mathematical analysis is not merely an accompaniment to our empirical results, but an important part  
28 in their development.

29 Furthermore, generalization bounds of the type we present in Theorem 2 and Corollary 4 are a core focus in statistical  
30 learning theory. Our bounds demonstrate a link between similarity and protection against overfitting in both the  
31 non-adaptive and adaptive case. The settings considered are purposefully simple to highlight the main ideas, and, as  
32 suggested by Reviewer 1, we will further simplify the bounds to improve their pedagogic value. Finally, generalization  
33 bounds can often be vacuous when instantiated with concrete values from practical settings. A key contribution of our  
34 work is generalization bounds that provide meaningful numerical guarantees on both ImageNet and CIFAR-10 data.

35 We now address the remaining comments of Reviewer 2:

36 • We agree an important future direction is to expand the scope of benchmarks considered. We plan to conduct the  
37 same similarity analysis on both Kaggle competitions and tasks in natural language processing.

38 • Our analysis explicitly considers the case where models cluster into smaller sets with large  $\eta$ . The definition of an  
39  $\eta$ -similarity cover is a clustering such that each cluster has similarity at least  $\eta$ . Figure 3 shows one such example:  
40 for  $\eta \approx 0.8$ , the models can be partitioned into roughly 200 clusters.

41 • Thank you for pointing out the typo for  $p_w$ . In lines 238 and 244,  $p_w$  should be replaced with  $1 - p_w$ . Otherwise,  
42 the definitions and subsequent discussion are unchanged.

### 43 Reviewer 3

44 • We agree an important direction for future work is exploring the extent to which our findings transfer to other  
45 settings. A concrete next step is evaluating model similarity on data from Kaggle competitions, which includes  
46 a diverse set of sample sizes, model classes, and data types. Nevertheless, image classification on ImageNet and  
47 CIFAR-10 is a natural starting place for this line of inquiry. Recent research has shown that both benchmarks do not  
48 suffer from adaptive overfitting despite almost a decade of intense activity, and model similarity offers an important  
49 piece of the explanation for this surprising phenomenon.

50 • We thank the reviewer for pointing out the lack of citations in the introduction. We will update the introduction with  
51 appropriate citations. We would appreciate any pointers to additional related work in the updated review.