

---

# On the Algorithmic Stability of SGD in Deep Learning

## Abstract

1 Many existing deep learning generalization bounds do not seem to be informative  
2 and can even increase with the sample size, which has further motivated the study  
3 of algorithmic stability as a possible approach for overcoming these limitations. In  
4 this work, we present empirical evidence that uniform stability might not appear in  
5 practical deep learning settings with sufficient strength to explain generalization  
6 and, further, that a key requirement of existing arguments is not satisfied: for  
7 two datasets differing by a single point, the distance between the final learned  
8 parameters does not decrease with dataset size. However, deeper investigation  
9 reveals that these failures might not be as bleak as they appear: despite separation  
10 by a large distance, these parameters can still sometimes end up in the same basin  
11 of attraction. We use our insights to suggest promising directions for algorithmic  
12 stability as a tool for explaining generalization in deep learning.

## 13 1 Introduction

14 Despite the impressive empirical success of deep learning models, their ability to generalize well  
15 (on a significant set of data distributions) despite overparameterization has thus far largely eluded  
16 the research community [25, 23]. Various flavors of generalization bounds have been applied to  
17 neural networks, including various norm- and margin-based bounds [2, 22, 7, 17, 19], PAC-Bayes  
18 bounds [6, 24], and VC-dimension-based bounds [3]. However, many such bounds have been  
19 shown to be insufficiently-correlated with generalization as various model components are varied  
20 (e.g., number of parameters) [13]. Recently, Nagarajan and Kolter [20] demonstrated that some of  
21 these bounds can even increase with sample size in certain settings, underscoring the importance  
22 of *empirically* evaluating proposed bounds' behavior as a function of dataset size. Furthermore,  
23 Nagarajan and Kolter [20] suggest that all uniform-convergence-based approaches might inherently  
24 be unable to explain deep learning's generalization performance, even *after* uniform convergence  
25 is restricted to the smallest possible set of models determined by the implicit bias of the learning  
26 algorithm. If this is true, then what tools for proving deep learning generalization bounds remain?  
27 One such tool, as acknowledged by Nagarajan and Kolter [20], is *algorithmic stability*.

28 **Algorithmic stability.** Algorithmic stability typically refers to a sensitivity analysis of the algorithm  
29 itself; specifically, how much can swapping (or removing) one point in an  $m$ -item training set  
30  $S$  change the output of an algorithm  $\mathcal{A}(S)$ ? Bousquet and Elisseeff [4] formalized and proved  
31 generalization bounds under various different flavors of algorithmic stability; since then, additional  
32 variants of algorithmic stability have been developed [1, 9, 15, 18]. However, to this day, the main  
33 variant for obtaining bounds that hold with high probability over the random draw of the training set  
34 is *uniform stability*, the strictest of the requirements. Specifically, a learning algorithm  $\mathcal{A}$  is called  $\beta$ -  
35 uniformly stable with respect to loss  $\ell$  if:

$$\forall S \in \mathcal{Z}^m, \forall i \in \{1, \dots, m\}, \forall z \in \mathcal{Z} : |\ell(\mathcal{A}(S), z) - \ell(\mathcal{A}(S^{\setminus i}), z)| \leq \beta,$$

36 where  $S^{\setminus i}$  is  $S$  with element  $i$  removed. Often, uniform stability is expressed with respect to the  
37 *swapping* of one point, instead of the *removal*:

$$\forall S, S' \in \mathcal{Z}^m, \forall z \in \mathcal{Z} : |\ell(\mathcal{A}(S), z) - \ell(\mathcal{A}(S'), z)| \leq \beta,$$

38 where  $S$  and  $S'$  only differ at one index.

39 **Algorithmic stability of stochastic gradient descent (SGD).** Various works thus far have studied  
40 whether the framework of algorithmic stability can be applied to the analysis of stochastic gradient  
41 descent (typically including at least some extension to nonconvex loss landscapes) [11, 8, 16].  
42 However, each of these results has some *subset* of the following weaknesses when applied to practical  
43 deep learning:

- 44 • The bound is only in *expectation* with respect to the draw of the sample  $S$ . In general, we  
45 ultimately seek bounds that will hold with high probability over the draw of the sample  
46  $S \sim \mathcal{D}^m$ , although such bounds are generally more difficult to prove theoretically.
- 47 • The stability parameter  $\beta$  relies on smoothness parameters of the loss landscape that might  
48 not be particularly favorable for neural networks.
- 49 • The result heavily relies on a learning rate of  $\mathcal{O}(1/t)$ , where  $t$  is the parameter update (vs.  
50 epoch). This ensures that, in expectation over the algorithm’s randomness, the learning  
51 rate has decayed more for larger samples by the time the swapped point is encountered.  
52 In contrast, in deep learning, the learning rate typically stays constant for at least the first  
53 epoch.
- 54 • The proof relies on controlling the (expected) distance between  $\mathcal{A}(S)$  and  $\mathcal{A}(S')$  in param-  
55 eter space, which seems unlikely to decrease sufficiently with dataset size in practice (without  
56 the aforementioned  $1/t$  learning rate schedule). We explore this in more detail in Section 4.

57 **Our work.** Inspired by the growing literature empirically analyzing the shortcomings of current  
58 deep learning generalization bounds and the anticipated algorithmic stability weaknesses discussed  
59 above, in this work we initiate a study of the following question: Does SGD empirically satisfy  
60 uniform stability in *practical* deep learning settings, in a manner sufficient to yield generalization  
61 bounds that hold with high probability (over the draw of the dataset and the algorithm’s randomness)?  
62 Unfortunately, analyzing uniform stability *empirically* is incredibly challenging due to the many  
63 suprema in the definition (i.e.,  $\forall S, S', z$ ), and we thus do not claim that any empirical analysis  
64 can *definitively* answer whether or not SGD in deep learning is uniformly stable. However, to our  
65 knowledge, this is the most *extensive* empirical examination of uniform stability in deep learning to  
66 date. Our contributions are as follows:

- 67 • Discussion of challenges in the empirical evaluation of uniform stability, with suggested  
68 methodology for overcoming them. Crucially, we validate our methodology in the simpler  
69 setting of logistic regression.
- 70 • Evidence that uniform stability (with respect to the cross-entropy loss) does not decrease  
71 sufficiently with dataset size to fully explain deep learning’s generalization.
- 72 • Evidence that  $\|\mathcal{A}(S) - \mathcal{A}(S')\|_2$  (when the output of  $\mathcal{A}$  is treated as a single vector of  
73 concatenated parameters) does not sufficiently decrease with dataset size in practical deep  
74 learning settings; in some cases, it can even increase despite strong generalization. We  
75 suggest that, if there is a form of algorithmic stability at play in deep learning, it does not  
76 stem from parameter closeness. We argue that future theoretical attempts to prove stability  
77 of SGD in deep learning should proceed through a different key path.
- 78 • Discovery of settings with insufficient cross-entropy uniform stability to explain generaliza-  
79 tion but for which  $\mathcal{A}(S)$  and  $\mathcal{A}(S')$  are in the same basin of attraction (see Section 4 for a  
80 precise definition), suggesting that convex settings with large basins of attraction could also  
81 share these same failure modes and thus pave the way for more tractable analyses.

## 82 2 Methodology

83 Here, we describe the key aspects of our methodology for empirically evaluating uniform stability,  
84 with additional details in Appendix A. We first applied our methodology to logistic regression, which  
85 we used to help validate our methodology. We then applied our methodology to deep neural networks.

86 Throughout the paper, we use  $\mathcal{A}(S)$  to denote to the output of the algorithm on dataset  $S$ . Although  
87 this object is really a function, we slightly abuse notation and treat it as a vector, i.e., with all of the  
88 model’s parameters concatenated into a single vector. We occasionally use  $W_S$  instead to denote the  
89 parameters output by  $\mathcal{A}$  on  $S$ , concatenated into a single vector.

90 **Random seeds.** Since we are seeking bounds that hold with high probability over the randomness  
91 of the algorithm, each plot we produce examines a *single* setting of the seed controlling initialization  
92 and the seed controlling SGD order. Thus, for each dataset/hyperparameter configuration, we present  
93 a single setting of the seeds in the main paper and defer our plots for other seeds to Appendix B.

94 **Datasets:  $S$  and  $S'$ .** We used MNIST for logistic regression (divided into two classes for binary  
95 classification: labels 0-4 and labels 5-9), and we used CIFAR-10 [14] and SVHN (Street View  
96 House Numbers) [21] for neural network training (10-class classification). In order to thoroughly  
97 study behavior (e.g., test/train error, various stability metrics, etc.) as a function of dataset size, we  
98 examined the following dataset sizes:  $\{800, 1600, 3200, 6400, 12800\}$  for logistic regression and  
99  $\{15000, 20000, 25000, 30000, 35000, 40000, 45000, 50000\}$  for neural networks. See Appendix A  
100 for more details regarding this choice of dataset sizes.

101 We emphasize that we intentionally do *not* use data augmentation; we want to precisely measure  
102 behavior as a function of dataset size, and to our knowledge, there is no widely-accepted approach  
103 for calculating the *effective* dataset size with data augmentation.

104 **Multiple trials per dataset size.** To study the quantifier  $\forall S, S' \in \mathcal{Z}^m$  *empirically* as thoroughly  
105 as possible, we sampled multiple  $(S, S')$  pairs per dataset size  $m$ . Each  $(S, S')$  pair was sampled as  
106 follows: we first randomly drew a subset of size  $m$  from the relevant training dataset (uniformly at  
107 random without replacement) to form  $S$ , we then uniformly sampled a single element of the relevant  
108 test set (call this element  $z$ ), and finally we uniformly sampled an index  $i \in \{1, \dots, m\}$  of  $S$  in  
109 which to swap in  $z$ , thus forming  $S'$ . We then trained two models in parallel, one on training dataset  
110  $S$  and one on  $S'$ . This procedure was repeated 90 times per dataset size for logistic regression and 40  
111 times per dataset size for each neural network configuration (due to the higher cost of each run).

112 Crucially, the only difference between training on  $S$  and  $S'$  was the appearance of  $z$  in  $S'$  in a single  
113 batch per epoch. All other data points were the same and were visited in the same order. Furthermore,  
114 we explicitly disabled all sources of GPU nondeterminism to ensure that we were fully isolating the  
115 effect of swapping in  $z$ .

116 **Models and training.** The logistic regression model is a 784-dimensional linear classifier plus a  
117 bias term, and the neural networks are residual networks, specifically ResNet-20 [12].

118 The logistic regression models were trained via stochastic gradient descent (SGD) with learning rate  
119 0.1, batch size 128, and no momentum.

120 On SVHN, we trained a ResNet-20 via SGD with learning rate 0.01, batch size 32, and no momentum.  
121 On CIFAR-10, we explored two different hyperparameter configurations: one without momentum  
122 and one with momentum 0.9. The other hyperparameters were the same across both configurations: a  
123 decaying learning rate schedule (starting at 0.1 and dividing by 10 at iterations 32,000 and 48,000)  
124 and batch size 128 [12].

125 **Stopping criterion.** We train each model for 100,000 iterations (i.e., parameter updates). See  
126 Appendix A for a more detailed discussion of stopping criteria.

127 **Uniform stability with respect to the cross-entropy loss.** In the uniform stability definition,  
128 instead of a supremum over the domain, we calculate a max over the test set. A priori, it might not be  
129 clear how effective this would be, and we thus validate our methodology via logistic regression in  
130 Section 3 before proceeding to deep learning.

131 **Plots and curve fitting.** Many of the quantities examined in this paper take the form of  $g(m) =$   
132  $\sup_{S \in \mathcal{Z}^m} f(S)$  or  $g(m) = \sup_{S, S' \in \mathcal{Z}^m} f(S, S')$  for some function  $f$ , and we expect  $g(m)$  to have  
133 the form  $g(m) = am^b$  for some constants  $a, b$ . Thus, for these quantities, we use the following  
134 plotting motif: all trials per dataset size are displayed as blue dots, the maximum value per dataset

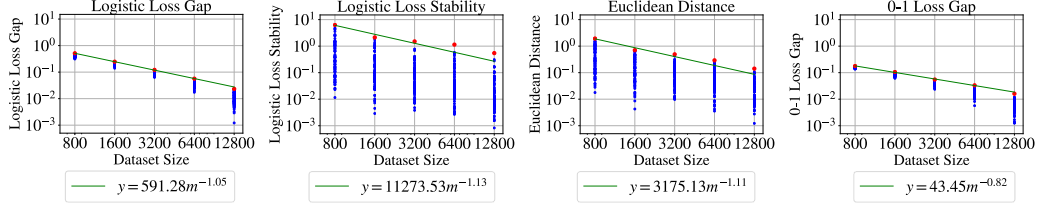


Figure 1: Generalization and stability curves. 90 samples per dataset size  $m$ . Each sample involves independently drawing  $S \sim \mathcal{D}_{\text{train}}^m$ ,  $z \sim \mathcal{D}_{\text{test}}$ ,  $i \sim U([m])$ .  $S' := S^{i \leftarrow z}$ .

135 size is a red dot, and a green curve of the form  $g(m) = am^b$  is fit to the *red* dots (see Appendix A for  
 136 curve fitting details). To emphasize  $b$ , the rate of decrease (or occasionally increase) with  $m$ , these  
 137 plots display both the  $x$ - and  $y$ -axes in log scale.

### 138 3 Uniform Stability and Generalization

139 For many years, obtaining useful generalization bounds via uniform stability required  $\beta = \mathcal{O}(1/m)$ ,  
 140 but [8] (followed by [5]) recently derived tighter bounds of the form: with probability at least  $1 - \delta$   
 141 over the choice of  $S \sim \mathcal{D}^m$ ,

$$R_{\mathcal{D}}(\mathcal{A}(S)) \leq \widehat{R}_S(\mathcal{A}(S)) + c \left( \beta \log(m) \log(m/\delta) + \frac{\sqrt{\log(1/\delta)}}{\sqrt{m}} \right) \quad (1)$$

142 for some constant  $c$ . Here,  $R_{\mathcal{D}}(\mathcal{A}(S))$  is the expected loss over the true distribution, and  
 143  $\widehat{R}_S(\mathcal{A}(S))$  is the empirical loss evaluated on  $S$ . This bound suggests that  $R_{\mathcal{D}}(\mathcal{A}(S)) - \widehat{R}_S(\mathcal{A}(S))$   
 144 is bounded by  $\tilde{\mathcal{O}}(\max\{\beta, 1/\sqrt{m}\})$ , hiding logarithmic dependencies inside the  $\tilde{\mathcal{O}}$ . Thus, if  
 145  $R_{\mathcal{D}}(\mathcal{A}(S)) - \widehat{R}_S(\mathcal{A}(S))$  empirically decays more slowly than  $1/\sqrt{m}$ , providing empirical evi-  
 146 dence that  $R_{\mathcal{D}}(\mathcal{A}(S)) - \widehat{R}_S(\mathcal{A}(S))$  and  $\beta$  decay *similarly* with  $m$  would suggest that uniform  
 147 stability has sufficient strength to explain generalization.

148 In this section, we present the results of our uniform stability experiments for both logistic regression  
 149 and neural networks. In both sections, we also carefully estimate  $R(\mathcal{A}(S)) - \widehat{R}_S(\mathcal{A}(S))$  as a function  
 150 of dataset size, under both the 0-1 loss and the logistic or cross-entropy loss, to understand to what  
 151 degree our uniform stability results are able to capture the strength of generalization. For convenience,  
 152 we use the phrase “generalization gap” or “loss gap” to denote this difference in test and train loss.

#### 153 3.1 Logistic regression

154 As there are obvious challenges in the empirical investigation of uniform stability with respect to  
 155 the cross-entropy loss, we began by analyzing logistic regression, which presents many of the same  
 156 challenges (e.g., the logistic loss, how to analyze the suprema over the domain, etc.) but provides a  
 157 much simpler and better-understood testbed in which to explore our methodology.

158 **Results.** In Figure 1, we plot the logistic loss generalization gap, our empirical estimate of the  
 159 logistic loss uniform stability, the Euclidean distance between the final parameters of  $\mathcal{A}(S)$  and  
 160  $\mathcal{A}(S')$ , and the 0-1 loss generalization gap. We fit a curve to the maximum value per dataset size,  
 161 as described in detail in Section 2 and Appendix A, and we compare the dependence on  $m$  of our  
 162 curves. Among the first three plots, we see a very similar dependence on  $m$ , ranging from  $m^{-1.05}$   
 163 to  $m^{-1.13}$ . The dependence on  $m$  in the 0-1 loss generalization gap plot is a bit weaker ( $m^{-0.82}$ ),  
 164 but we include this primarily for completeness and as a frame of reference; we are more interested  
 165 in whether *logistic loss stability* can explain the strength (with respect to  $m$ ) of generalization with  
 166 respect to the *logistic loss*.

167 **Conclusions.** These plots demonstrate the potential of our methodology to capture, via uniform  
 168 stability with a *finite* maximum over the test set, the dependence on  $m$  of the Euclidean distance  
 169 between parameters and, most importantly, the logistic loss generalization gap. Thus, although there

ID	Model	Dataset	Learning rate	Batch size	Momentum
1a	ResNet-20	SVHN	0.01 (constant)	32	0.0
2a	ResNet-20	CIFAR-10	0.1, 0.01 at 32k, 0.001 at 48k	128	0.0
2b	ResNet-20	CIFAR-10	0.1, 0.01 at 32k, 0.001 at 48k	128	0.9

Table 1: Deep neural network settings studied.

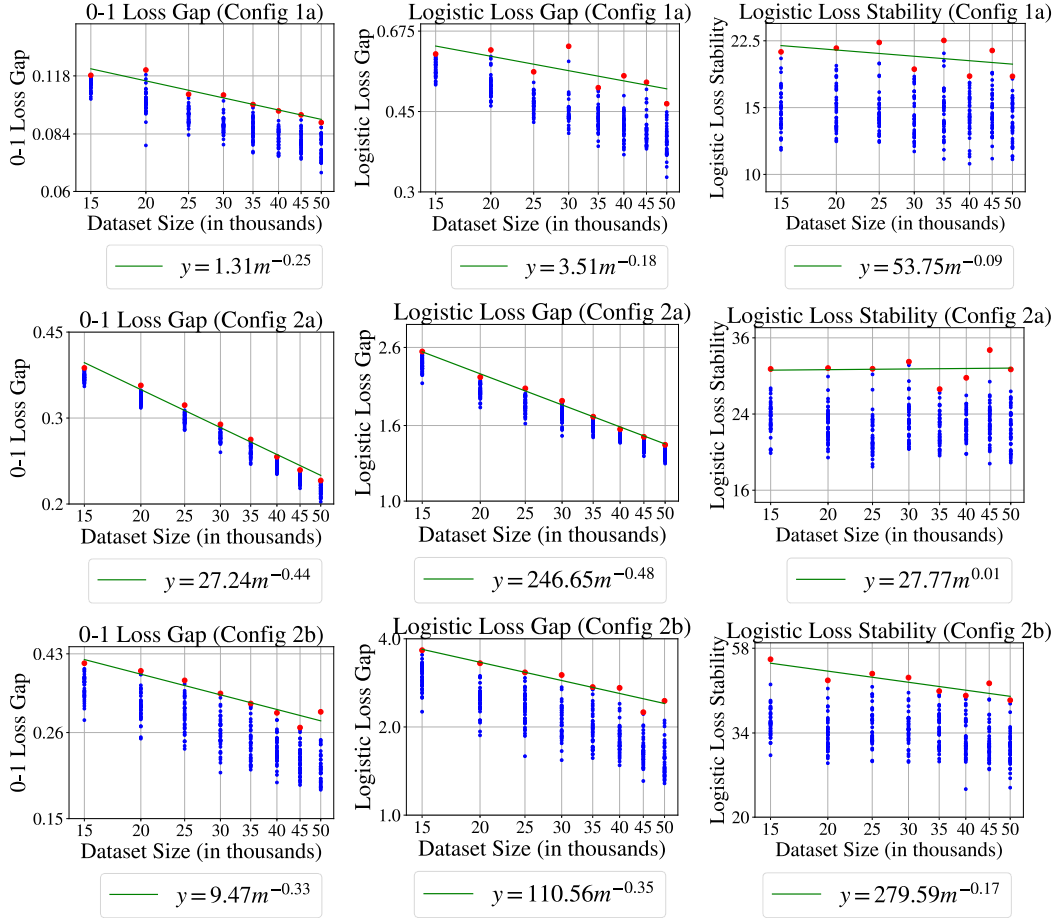


Figure 2: Generalization and stability curves. 40 samples per dataset size  $m$ . Each sample involves independently drawing  $S \sim \mathcal{D}_{\text{train}}^m$ ,  $z \sim \mathcal{D}_{\text{test}}$ ,  $i \sim U([m])$ .  $S' := S^{i \leftarrow z}$ . Note: We use “logistic loss” and “cross-entropy” loss interchangeably here; all models in this figure were trained and evaluated with the cross-entropy loss.

170 are obvious differences between the suprema in the definition of uniform stability and our empirical  
 171 evaluation with finite maxima, our results suggest that there is nevertheless some promise of obtaining  
 172 informative empirical results.

### 173 3.2 Deep learning

174 After validating our methodology in the simpler setting of logistic regression, we now extend our  
 175 methodology to the three deep learning configurations described in Section 2.

176 **Results.** Figure 2 displays the generalization and stability results for our three neural network  
 177 settings. In contrast with logistic regression, we postpone examining the parameters themselves  
 178 until Section 4, in which we conduct an analysis more targeted to deep learning’s nonconvex loss  
 179 landscape.

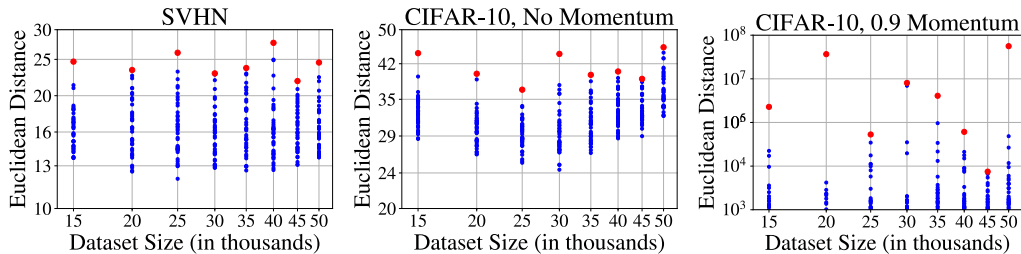


Figure 3:  $\|\mathcal{A}(S) - \mathcal{A}(S')\|_2$ , at  $t = 100,000$ , with 40 samples per dataset size  $m$ . Each sample involves independently drawing  $S \sim \mathcal{D}_{\text{train}}^m$ ,  $z \sim \mathcal{D}_{\text{test}}$ ,  $i \sim U([m])$ .  $S' := S^{i \leftarrow z}$ .

180 Most significantly, we compare the cross-entropy loss generalization gap to the uniform stability  
 181 curve. For the ResNet-20 on SVHN, the stability curve displays a mild decrease with  $m$  (specifically,  
 182  $m^{-0.09}$ ), compared to  $m^{-0.18}$  for the cross-entropy loss generalization gap. For the ResNet-20 on  
 183 CIFAR-10 *without* momentum, the stability curve does *not* decrease with  $m$ , despite the cross-entropy  
 184 loss generalization gap having a dependence of  $m^{-0.48}$ . For the ResNet-20 on CIFAR-10 *with*  
 185 momentum, the stability curve displays a mild decrease with  $m$  (specifically,  $m^{-0.17}$ ), compared to  
 186  $m^{-0.35}$  for the cross-entropy loss generalization gap.

187 To provide a frame of reference, we also compare the cross-entropy loss generalization gap to the  
 188 0-1 loss generalization gap and note that, at least for these particular configurations, attempting to  
 189 explain the rate of decrease with  $m$  of the cross-entropy loss generalization gap does not leave us too  
 190 far from the 0-1 generalization gap either.

191 Appendix B includes the same experiments repeated with more seeds (for initialization and SGD data  
 192 order) and includes plots at other stopping points (other than 100,000 iterations).

193 **Conclusions.** Overall, in our deep learning experiments, uniform stability with respect to the cross-  
 194 entropy loss does not appear with sufficient strength to explain observed generalization with respect  
 195 to the cross-entropy loss.

## 196 4 Behavior of Parameters

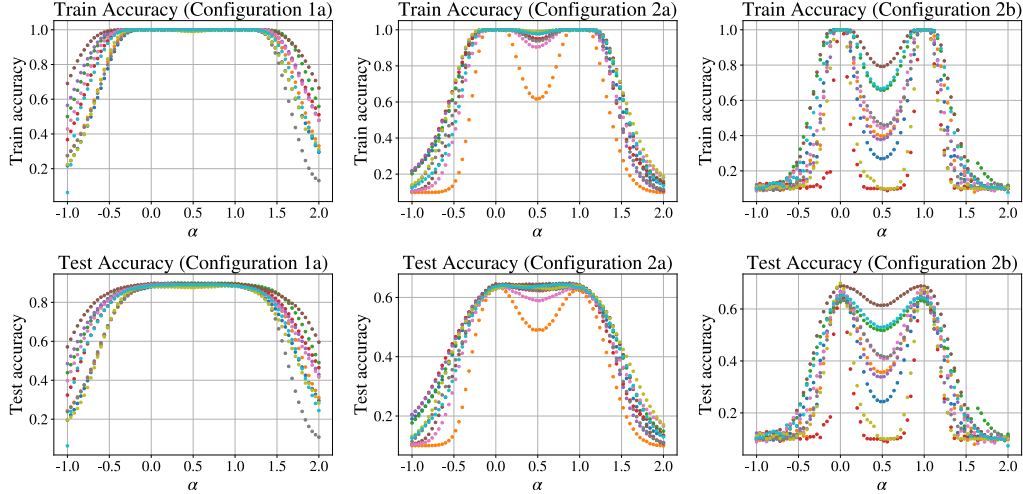
197 In this section, we analyze the behavior of the underlying parameters to try to disentangle the effect  
 198 of the cross-entropy loss and the supremum over the domain (estimated via the max over the test set)  
 199 from the learned models themselves in parameter space.

### 200 4.1 Euclidean Distance

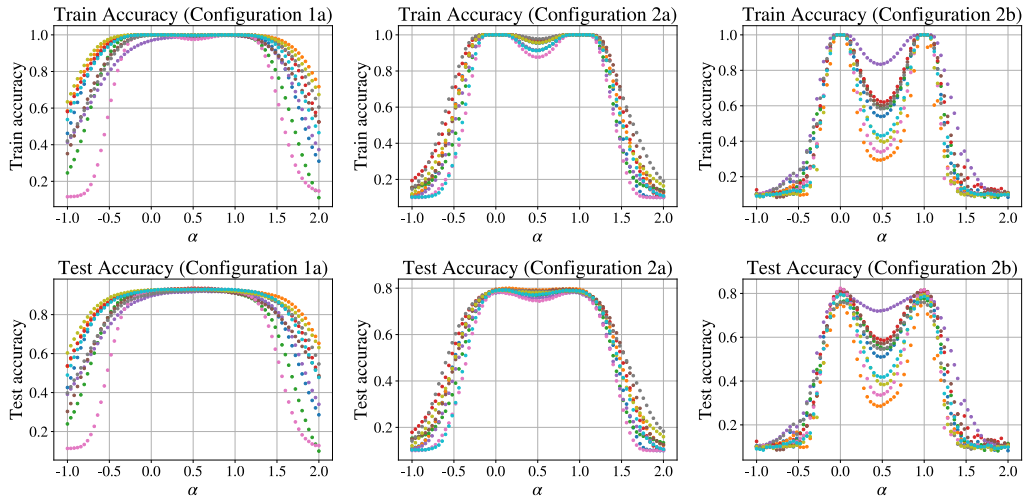
201 As mentioned in Section 1, we are further interested in studying the Euclidean distance between  
 202 the final learned parameters to help understand whether the key proof strategy introduced by Hardt  
 203 et al. [11] extends to practical deep learning settings. Since this paper, most proofs of the stability  
 204 of SGD (even in nonconvex settings) proceed by bounding the Euclidean distance in parameter  
 205 space between  $\mathcal{A}(S)$  and  $\mathcal{A}(S')$  and then appealing to the Lipschitzness of the loss. However, if  
 206 the Euclidean distance between  $\mathcal{A}(S)$  and  $\mathcal{A}(S')$  does not decrease with dataset size in our trained  
 207 models, this suggests that this proof strategy might not be sufficient for obtaining generalization  
 208 bounds in practical deep learning settings that hold with high probability (over the random draw of  
 209 the dataset and the random initialization and SGD data order of the algorithm).

210 **Results.** Figure 3 presents  $\|\mathcal{A}(S) - \mathcal{A}(S')\|_2$  for our three neural network configurations. We see  
 211 that, from  $m = 15k$  to  $m = 50k$ , the distances do not decrease with dataset size at a sufficient rate to  
 212 explain generalization and actually even increase in some dataset size ranges.

213 **Conclusions.** These results suggest that a decrease in Euclidean distance of the parameters with  
 214 dataset size is likely not a viable path through which to prove stability in *practical* deep learning  
 215 settings.



(a)  $m = 15,000$ .



(b)  $m = 50,000$ .

Figure 4: Interpolation at  $t = 100,000$ .

216 One might ask whether the nondecreasing Euclidean distance we observe here is caused by the norms  
 217 in parameter space *themselves* growing with dataset size. We first emphasize that this question does  
 218 not impact our conclusions, as the proofs to which we have referred invoke the raw Euclidean distance  
 219 between the parameters. However, for completeness, we refer the interested reader to Appendix C for  
 220 an extensive analysis of norms and normalized Euclidean distances.

## 221 4.2 Linear Mode Connectivity

222 We now ask the question: Is nonconvexity causing optimization on  $S$  and  $S'$  to diverge to different  
 223 basins of attraction, thus thwarting efforts to extend analyses from convex settings to deep learning's  
 224 nonconvex setting (as is done by [11] and follow-up works)? Here, we use *basin of attraction* to  
 225 mean a convex set of solutions (in parameter space) all with comparable training and/or test loss.

226 To make this more precise, we invoke the *linear mode connectivity* framework of Frankle et al. [10]  
 227 to study this question. Specifically, linear mode connectivity asks whether, at all networks along the  
 228 linear path between two candidate networks (in parameter space), the training and/or test error does  
 229 not increase. In our setting,  $W_S$  and  $W_{S'}$  qualify as linearly connected modes if, for all  $\alpha \in [0, 1]$ ,

230 the (test or train) accuracy of the model with parameters  $\alpha W_S + (1 - \alpha)W_{S'}$  is not significantly  
231 below that of  $W_S$  or  $W_{S'}$  (roughly 2%, per [10]).

232 **Results.** We plot the train accuracy (on  $S$ ) and test accuracy at  $\alpha W_S + (1 - \alpha)W_{S'}$  at 76 equally-  
233 spaced values of  $\alpha \in [-1, 2]$ . The learned parameters are at  $\alpha = 0$  and  $\alpha = 1$ , but we include  
234 additional values of  $\alpha$  on either end as a frame of reference. We randomly select 10 trials among  
235 the 40 trials described in Section 2 and, for each trial, we plot the train and test accuracy for each  
236 value of  $\alpha$ . Figure 4a has results for  $m = 15,000$  and Figure 4b has results for  $m = 50,000$ .  
237 The ResNet-20 on SVHN has nondecreasing accuracy when linearly interpolating between  $W_S$   
238 and  $W_{S'}$ , the ResNet-20 on CIFAR-10 without momentum has slightly decreasing accuracy when  
239 linearly interpolating between  $W_S$  and  $W_{S'}$ , and the ResNet-20 on CIFAR-10 with momentum has  
240 significantly decreasing accuracy when interpolating.

241 **Conclusions.** A priori, it is not obvious what one should expect when linearly interpolating, and it  
242 is thus perhaps surprising that our three configurations largely span the space of possibilities. Thus,  
243 in order to further study the weaknesses of uniform stability in practical deep learning settings, we  
244 suggest that moving to a regime such as a ResNet-20 on CIFAR-10 *with* momentum, in which the  
245 solutions are *not* connected by a path of nondecreasing accuracy, might not be immediately necessary  
246 from a scientific standpoint. Perhaps the limitations of uniform stability can be explored and better  
247 understood with a configuration such as our ResNet-20 on SVHN (without momentum). Notably, the  
248 SVHN interpolation results suggest that nonconvexity might not be necessary at all to investigate the  
249 particular weaknesses of algorithmic stability experienced by deep learning; rather, *convex* settings  
250 with large enough basins of attraction (defined for our purposes as convex sets of parameters yielding  
251 approximately equal training and/or test loss) to host a reasonable degree of functional diversity  
252 might actually be subject to these same weaknesses. Thus, our findings might open the door to the  
253 study of more tractable, convex settings in which one can study the same limitations of algorithmic  
254 stability that appear in deep learning.

## 255 5 Conclusions and Future Work

256 In this work, we have initiated the challenging endeavor of empirically studying the uniform stability  
257 of deep learning. Although we freely admit that no reasonable empirical results could *definitively*  
258 rule out uniform stability (due to its formulation as several maxima over the domain), we believe that  
259 our results present compelling evidence that (a) uniform stability (with respect to the cross-entropy  
260 loss) might not be present in practical deep learning with sufficient strength to explain generalization,  
261 and (b) that typical theoretical approaches based on parameter distance decreasing with dataset size  
262 are likely not the driving force behind any form of algorithmic stability that nevertheless might exist  
263 in deep learning. Ultimately, if some form of algorithmic stability (perhaps weaker than uniform  
264 stability) is at play in deep learning, we suspect that it will stem from a function-space view that  
265 appropriately handles divergence to different basins of attraction after swapping one data point (as  
266 seen in Section 4.2, Configuration 2b). However, in the meantime, we present compelling evidence  
267 that many of the weaknesses of uniform stability can already be seen empirically in simpler, perhaps  
268 even convex, settings. We believe that formalizing and further investigating these more tractable  
269 settings presents an exciting direction for future work.



## 270 **Broader Impact**

271 Overall, we believe that this work has relatively minimal societal impact. Ultimately, in the long term,  
272 we do hope that this work will contribute to our collective understanding of how deep learning works,  
273 which is increasingly critical as deep learning is deployed in an ever-growing range of real-world  
274 applications. We see this as a potential positive benefit of our work. On the negative side, we do  
275 not envision any major *harm* from this work, other than the environmental cost of running these  
276 experiments. However, our hope is that, as we gain a better understanding of deep learning as a  
277 society, the need for so many large-scale scientific experiments will eventually subside, and we will  
278 be better equipped to predict the behavior of deep learning through theory (at least more than we are  
279 able to do at present).

## 280 **References**

- 281 [1] K. Abou-Moustafa and C. Szepesvári. An exponential efron-stein inequality for  $l_q$  stable  
282 learning rules. In *Proceedings of the 30th International Conference on Algorithmic Learning*  
283 *Theory*, 2019.
- 284 [2] P. Bartlett, D. J. Foster, and M. Telgarsky. Spectrally-normalized margin bounds for neural  
285 networks. *arXiv preprint arXiv:1706.08498*, 2017.
- 286 [3] P. Bartlett, N. Harvey, C. Liaw, and A. Mehrabian. Nearly-tight vc-dimension and pseudodi-  
287 mension bounds for piecewise linear neural networks. *Journal of Machine Learning Research*,  
288 20:1–17, 2019.
- 289 [4] O. Bousquet and A. Elisseeff. Stability and generalization. *Journal of Machine Learning*, 2:  
290 499–526, 2002.
- 291 [5] O. Bousquet, Y. Klochkov, and N. Zhivotovskiy. Sharper bounds for uniformly stable algorithms.  
292 *arXiv preprint arXiv:1910.07833*, 2019.
- 293 [6] G. K. Dziugaite and D. M. Roy. Computing nonvacuous generalization bounds for deep  
294 (stochastic) neural networks with many more parameters than training data. *arXiv preprint*  
295 *arXiv:1703.11008*, 2017.
- 296 [7] G. F. Elsayed, D. Krishnan, H. Mobahi, K. Regan, and S. Bengio. Large margin deep networks  
297 for classification. In *Proceedings of NeurIPS*, 2018.
- 298 [8] V. Feldman and J. Vondrak. High probability generalization bounds for uniformly stable  
299 algorithms with nearly optimal rate. In *Proceedings of COLT*, 2019.
- 300 [9] D. J. Foster, S. Greenberg, S. Kale, H. Luo, M. Mohri, and K. Sridharan. Hypothesis set stability  
301 and generalization. In *Proceedings of NeurIPS 2019*, pages 6726–6736, 2019.
- 302 [10] J. Frankle, G. K. Dziugaite, D. M. Roy, and M. Carbin. Linear mode connectivity and the lottery  
303 ticket hypothesis. In *Proceedings of the 37th International Conference on Machine Learning*,  
304 2020.
- 305 [11] M. Hardt, B. Recht, and Y. Singer. Train faster, generalize better: Stability of stochastic gradient  
306 descent. In *Proceedings of the 33rd International Conference on Machine Learning*, 2016.
- 307 [12] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In  
308 *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016.
- 309 [13] Y. Jiang, B. Neyshabur, H. Mobahi, D. Krishnan, and S. Bengio. Fantastic generalization  
310 measures and where to find them. In *Proceedings of ICLR*, 2020.
- 311 [14] A. Krizhevsky. Learning multiple layers of features from tiny images. 2009.
- 312 [15] S. Kutin and P. Niyogi. Almost-everywhere algorithmic stability and generalization error. In  
313 *Proceedings of UAI*, pages 275–282, 2002.
- 314 [16] I. Kuzborskij and C. H. Lampert. Data-dependent stability of stochastic gradient descent. In  
315 *Proceedings of the 35th International Conference on Machine Learning*, 2018.
- 316 [17] T. Liang, T. Poggio, A. Rakhlin, and J. Stokes. Fisher-rao metric, geometry, and complexity of  
317 neural networks. In *Proceedings of the 22nd International Conference on Artificial Intelligence*  
318 *and Statistics*, 2019.

- 319 [18] T. Liu, G. Lugosi, G. Neu, and D. Tao. Algorithmic stability and hypothesis complexity. In  
320 *Proceedings of the 34th International Conference on Machine Learning*, 2017.
- 321 [19] P. M. Long and H. Sedghi. Generalization bounds for deep convolutional neural networks. In  
322 *Proceedings of ICLR*, 2020.
- 323 [20] V. Nagarajan and J. Z. Kolter. Uniform convergence may be unable to explain generalization in  
324 deep learning. In *Proceedings of NeurIPS*, 2019.
- 325 [21] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng. Reading digits in nat-  
326 ural images with unsupervised feature learning. In *NIPS Workshop on Deep Learning and*  
327 *Unsupervised Feature Learning.*, 2011.
- 328 [22] B. Neyshabur, R. Tomioka, and N. Srebro. Norm-based capacity control in neural networks. In  
329 *Conference on Learning Theory*, pages 1376–1401, 2015.
- 330 [23] B. Neyshabur, S. Bhojanapalli, D. McAllester, and N. Srebro. Exploring generalization in deep  
331 learning. In *Advances in Neural Information Processing Systems*, pages 5947–5956, 2017.
- 332 [24] B. Neyshabur, S. Bhojanapalli, and N. Srebro. A PAC-Bayesian approach to spectrally-  
333 normalized margin bounds for neural networks. In *Proceedings of ICLR*, 2018.
- 334 [25] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals. Understanding deep learning requires  
335 rethinking generalization. In *Proceedings of ICLR*, 2017.

## 336 A Further Methodology Details

337 **Dataset splits.** MNIST: From the 60,000 training examples, we randomly sampled subsets as  
338 specified in Section 2 for training. We used the full 10,000-element test set for evaluation (including  
339 computation of uniform stability, as specified in Section 2).

340 CIFAR-10: From the 50,000 training examples, we randomly sampled subsets as specified in Section  
341 2 for training. We used the full 10,000-element test set for evaluation (including computation of  
342 uniform stability, as specified in Section 2).

343 SVHN: From the 73,257 training examples, we randomly sampled subsets as specified in Section  
344 2 for training. To maintain consistency with MNIST and CIFAR-10, we randomly sampled 10,000  
345 elements from the 26,032-element test set for evaluation. All datasets were normalized in the same  
346 manner, by dividing each coordinate by 255.

347 **Selection of hyperparameters.** The hyperparameters for ResNet-20 on CIFAR-10 are derived  
348 from [12]. The hyperparameters for logistic regression were chosen similarly, intentionally without  
349 momentum (since our primary goal was to study SGD) and without a decaying learning rate. The  
350 hyperparameters for Configuration 1a were intentionally chosen to vary from Configurations 2a  
351 and 2b, in order to create more diversity in our hyperparameter settings; specifically, we deemed  
352 it valuable to investigate a smaller batch size (i.e., 32) without momentum, and the corresponding  
353 learning rate of 0.01 worked fairly well with this batch size.

354 **Stopping criterion details.** We considered three different possible stopping criteria: parameter  
355 updates, epoch number, and average training loss. We performed preliminary analyses with all  
356 stopping criteria, but after careful consideration, we ultimately chose to focus our analysis on  
357 *parameter updates* for the following reasons: (1) parameter updates align with theoretical analyses of  
358 uniform stability, such as [11, 8], in which the stability parameter is expressed as a function of the  
359 number of parameter updates; and (2) parameter updates appear to give uniform stability the *best*  
360 chance at succeeding in explaining generalization, thus making our *negative* results more significant.  
361 Specifically, if instead we were to hold the number of epochs fixed across dataset sizes, this means  
362 that *larger* datasets would take more steps. This is true for average training loss as a stopping criterion  
363 as well, as it typically takes more steps for larger datasets to reach the same average training loss as  
364 smaller datasets. Thus, although these stopping criteria are perhaps truer to practice, we believe that  
365 they make it even *easier* for uniform stability to fail to explain the strength of generalization.

366 **Dataset size range.** We chose to limit our analysis to the ranges specified in Section 2 for the fol-  
367 lowing reason. In order to ask the question *Can the strength of decrease with  $m$  in our generalization*  
368 *gap be explained by uniform stability?*, we wanted a rate of decrease with  $m$  that would be roughly  
369 constant in our dataset size range. Figure 5 shows a plot and curve fit on a normal-scale plot, followed  
370 by a log-log plot. Although the curve fit displays some room for improvement in the original plot,  
371 the log-log plot reveals different regions of decrease with  $m$ . Through this plot and additional such  
372 investigations, we noticed that the dataset size range 15,000-50,000 yielded the largest window with  
373 a roughly consistent rate of decrease with  $m$ . Thus, we chose to focus our analysis on this window  
374 in order to draw more meaningful conclusions. As deep learning models are typically trained in  
375 large-data regimes, this decision aligns with practical considerations as well.

376 **Curve fitting details.** We used `scipy`'s `optimize` package, specifically the `curve_fit` function.  
377 We fit parameters  $a$  and  $b$  in  $y = am^b$ , where  $m$  is the dataset size and  $y$  is the metric of interest.

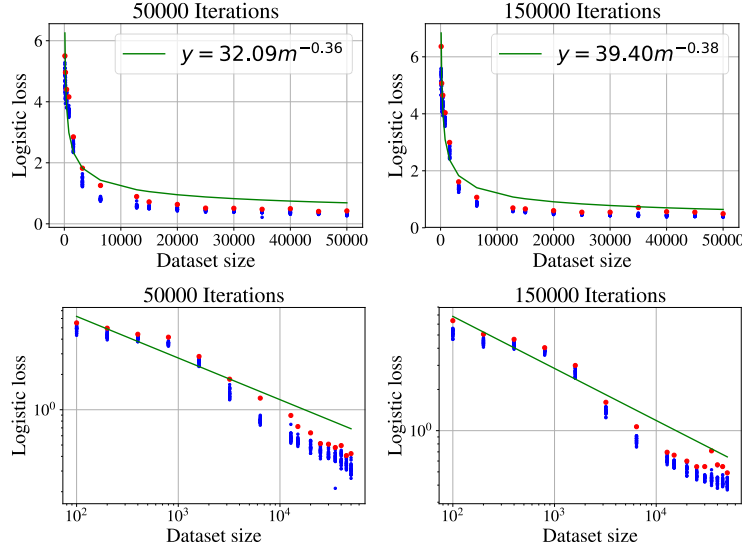


Figure 5: Configuration 1a, with 20 samples per dataset size  $m$ . Each sample involves independently drawing  $S \sim \mathcal{D}_{\text{train}}^m$ ,  $z \sim \mathcal{D}_{\text{test}}$ ,  $i \sim U([m])$ .  $S' := S^{i \leftarrow z}$ .

### 378 B Further Experiments for Section 3

379 In this section, we present stability and generalization results for two additional seeds (Trials 2 and  
 380 3) and compare them to the original seed presented in the main paper (Trial 1). Our figures are as  
 381 follows:

- 382 • Figure 6 has all trials for Configuration 1a (SVHN).
- 383 • Figure 7 has all trials for Configuration 2a (CIFAR-10, no momentum).
- 384 • Figure 8 has all trials for Configuration 2b (CIFAR-10, 0.9 mometum).
- 385 • Figure 9 has generalization (cross-entropy only) and stability results for Iteration 50,000 for  
 386 Configuration 1a.
- 387 • Figure 10 has generalization (cross-entropy only) and stability results for Iteration 50,000  
 388 for Configuration 2a.
- 389 • Figure 11 has generalization (cross-entropy only) and stability results for Iteration 50,000  
 390 for Configuration 2b.
- 391 • Figure 12 has generalization (cross-entropy only) and stability results for Iteration 150,000  
 392 for Configuration 1a.
- 393 • Figure 13 has generalization (cross-entropy only) and stability results for Iteration 150,000  
 394 for Configuration 2a.
- 395 • Figure 14 has generalization (cross-entropy only) and stability results for Iteration 150,000  
 396 for Configuration 2b.

397 The additional trials are roughly consistent with the trial highlighted in the main paper.

### 398 C Further Experiments for Section 4

399 In this section, we present further experiments regarding regarding the Euclidean distance between  
 400  $\mathcal{A}(S)$  and  $\mathcal{A}(S')$ , parameter norms, and normalized Euclidean distances. Our figures are as follows:

- 401 • Figure 15 presents additional trials for  $\|\mathcal{A}(S) - \mathcal{A}(S')\|_2$  at Iteration 100,000.
- 402 • Figure 16 has  $\|\mathcal{A}(S)\|_2$  at Iteration 100,000.
- 403 • Figure 17 has normalized Euclidean distances, with further details in the caption.

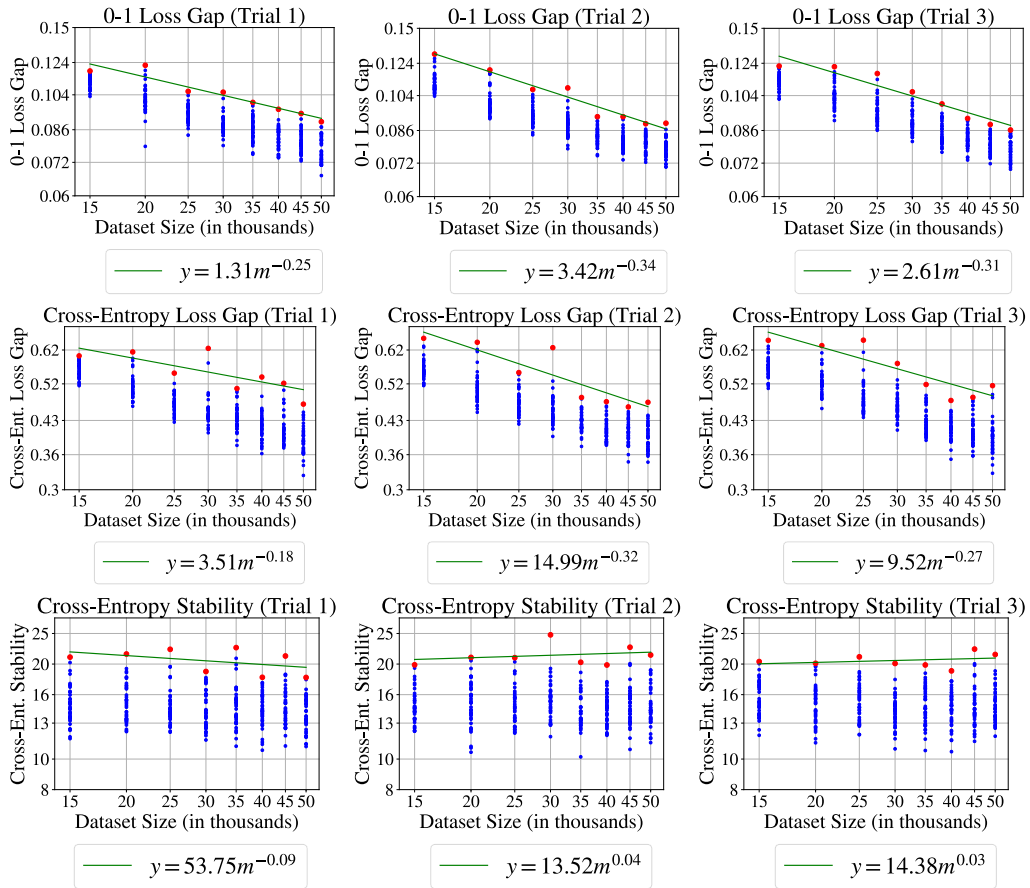


Figure 6: All trials for Configuration 1a (SVHN).

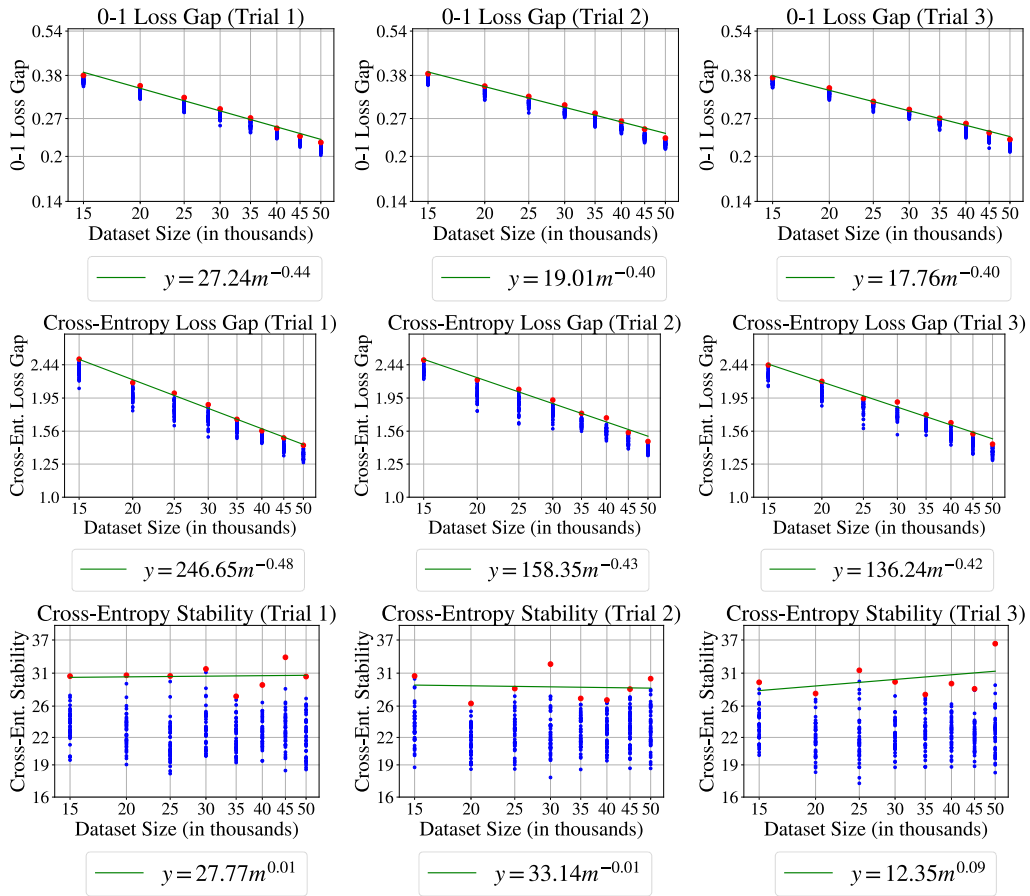


Figure 7: All trials for Configuration 2a (CIFAR-10, no momentum).

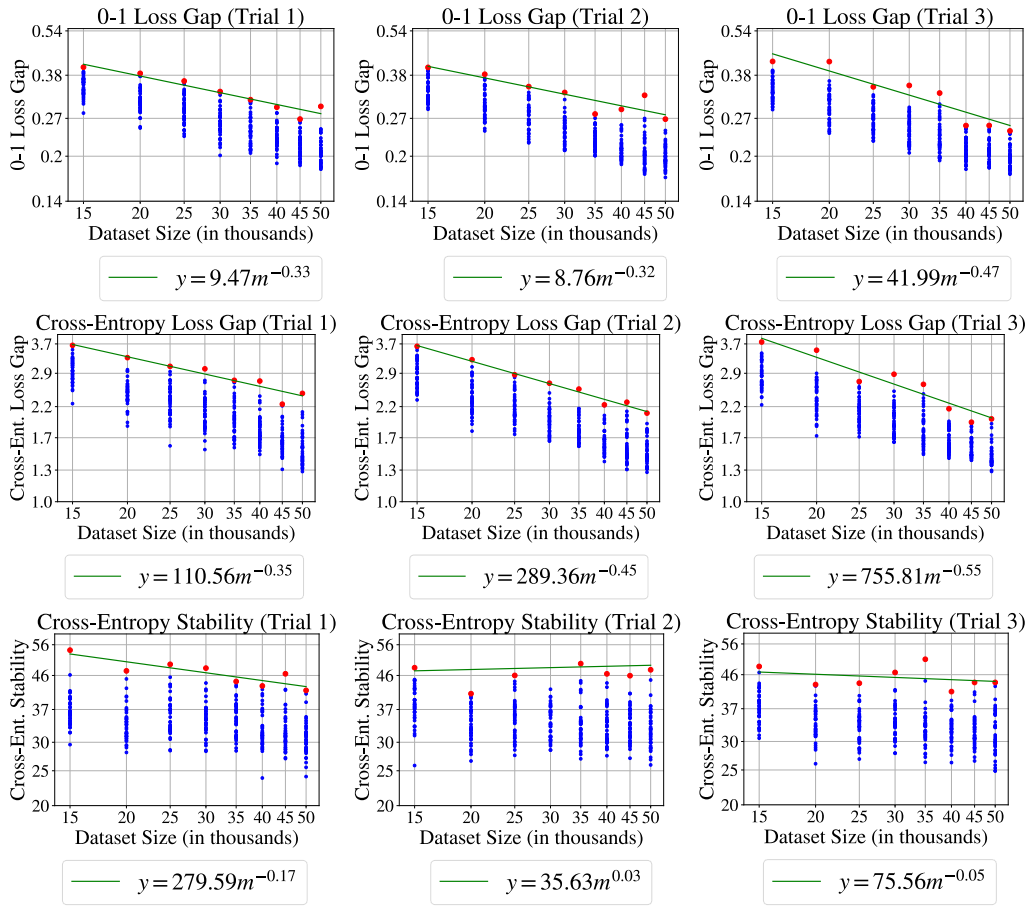


Figure 8: All trials for Configuration 2b (CIFAR-10, 0.9 momentum).

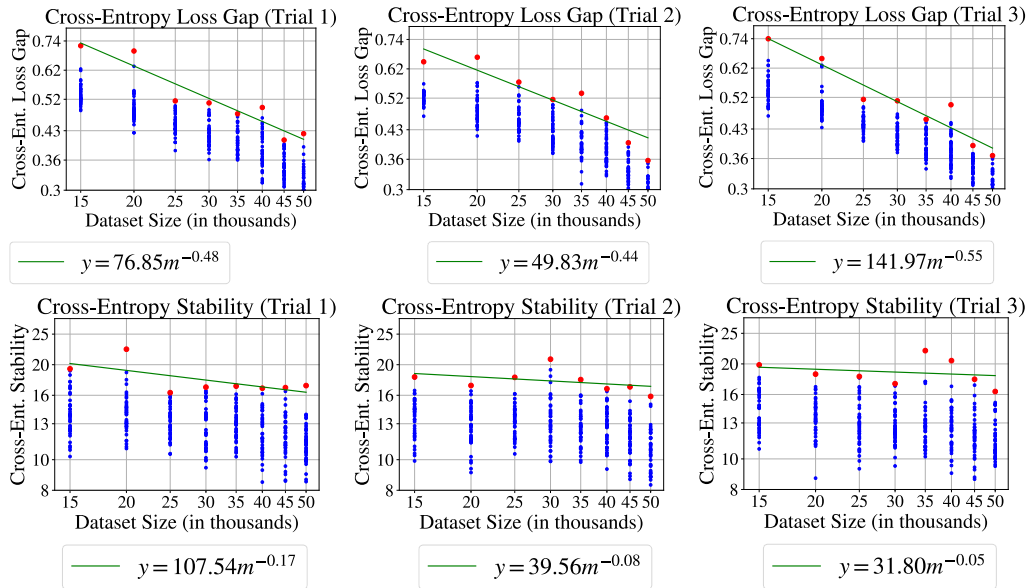


Figure 9: Iteration 50,000 for Configuration 1a.

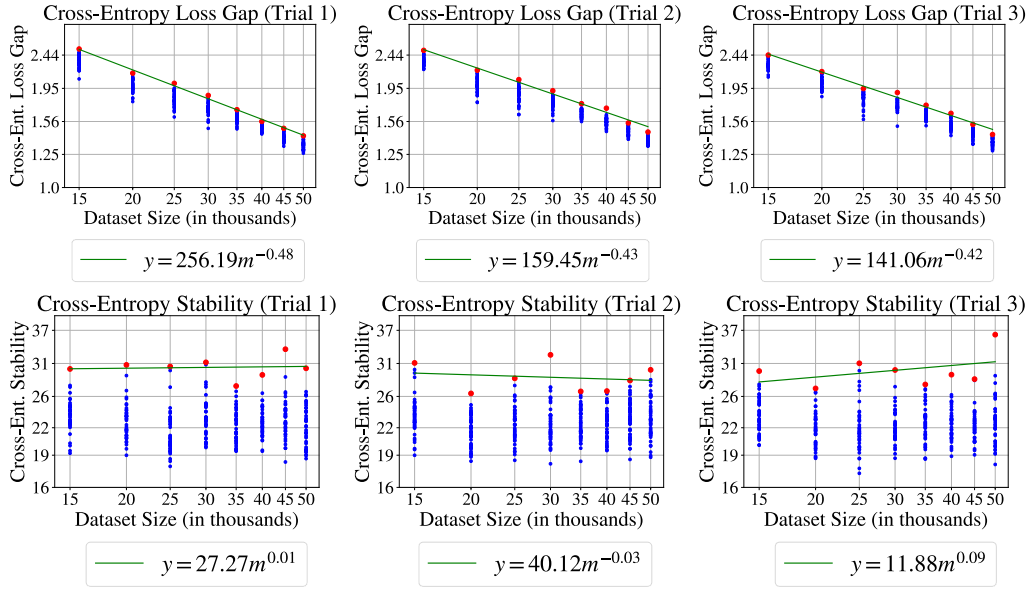


Figure 10: Iteration 50,000 for Configuration 2a.

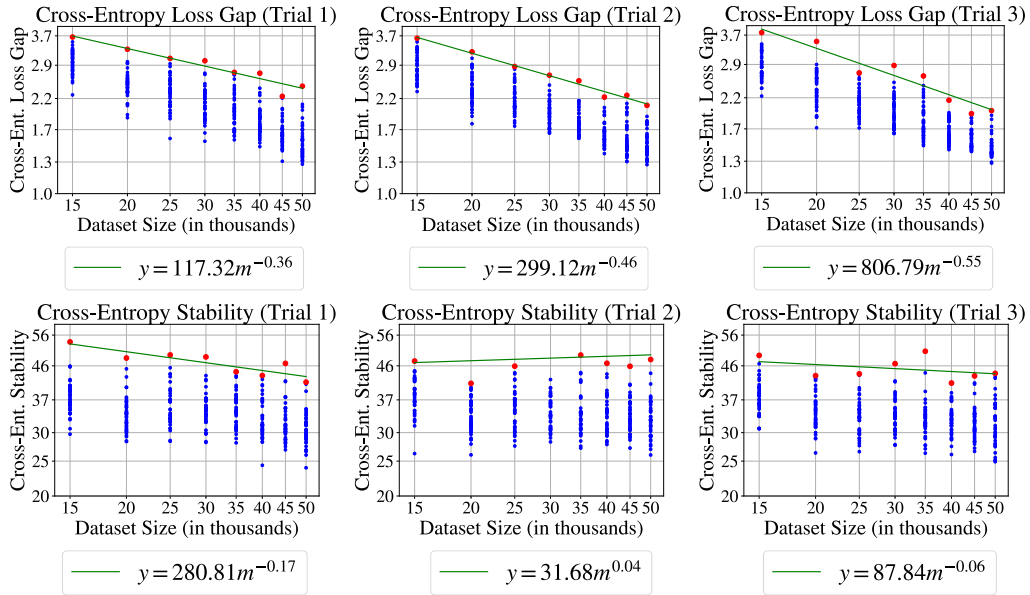


Figure 11: Iteration 50,000 for Configuration 2b.



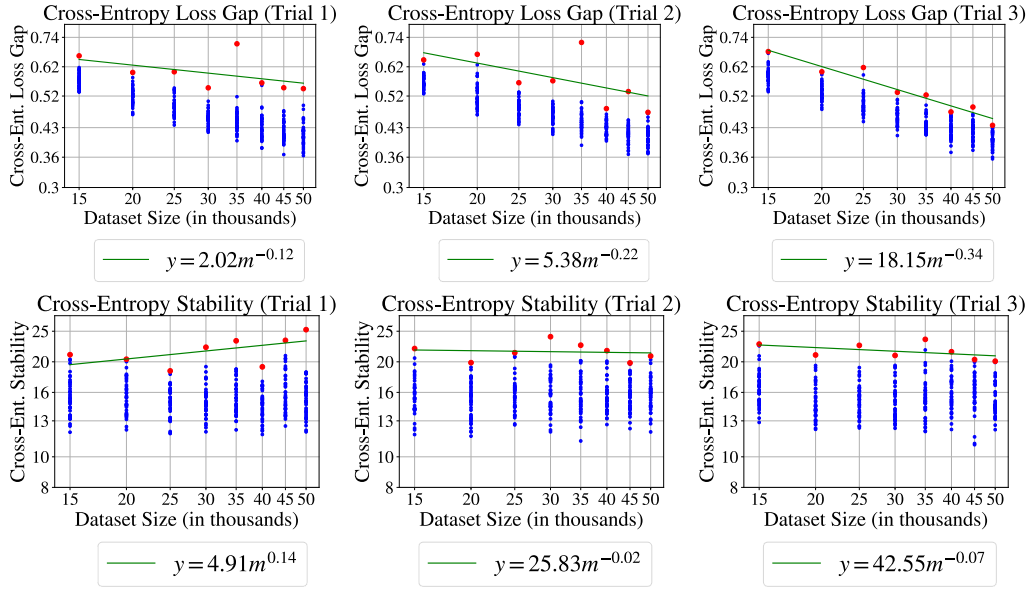


Figure 12: Iteration 150,000 for Configuration 1a.

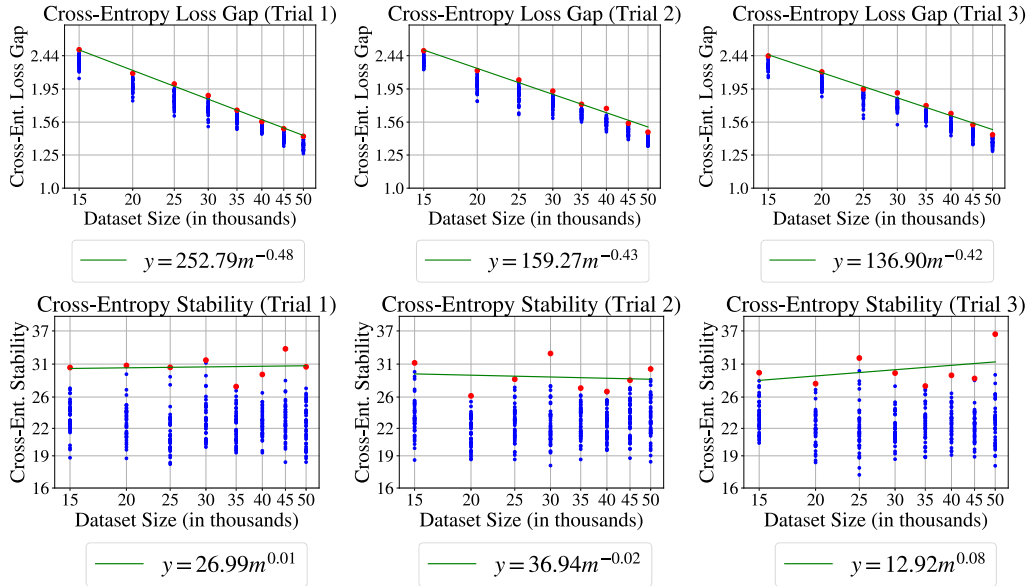


Figure 13: Iteration 150,000 for Configuration 2a.

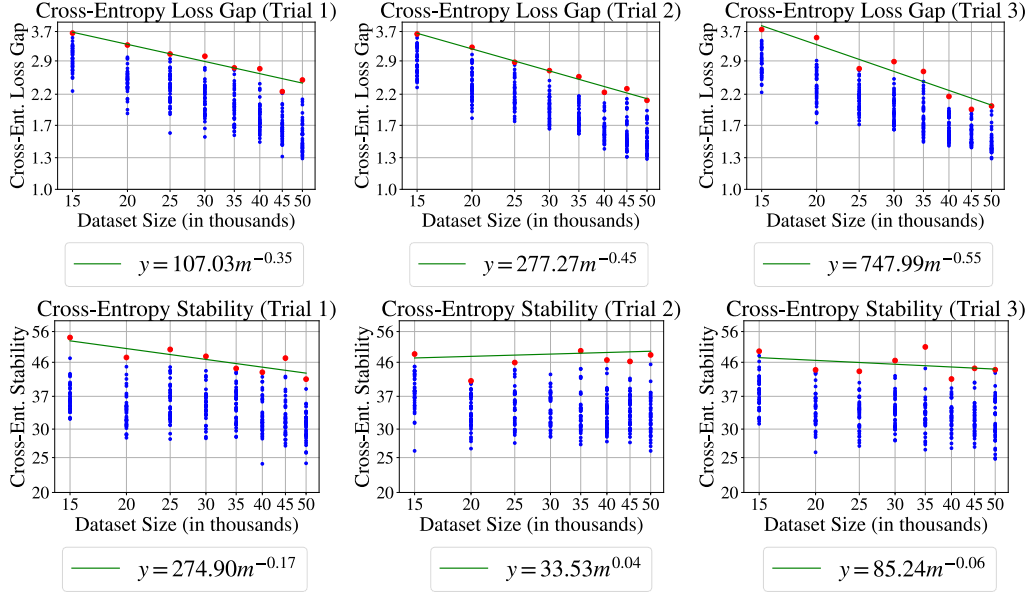


Figure 14: Iteration 150,000 for Configuration 2b.

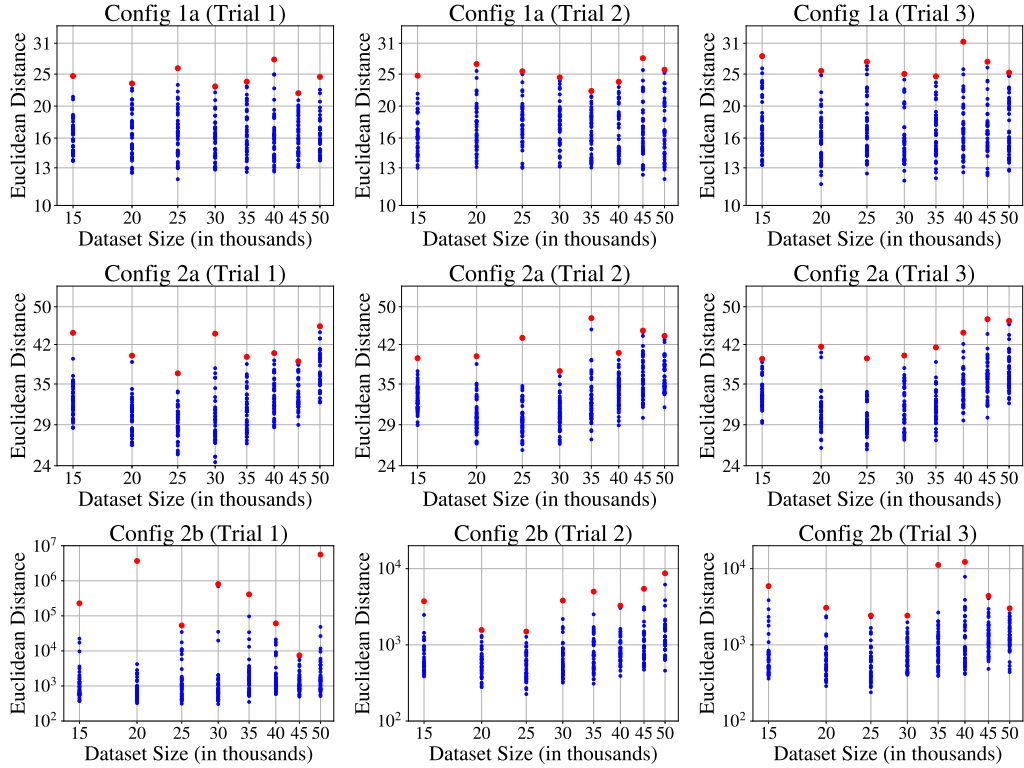


Figure 15: All trials for  $\|\mathcal{A}(S) - \mathcal{A}(S')\|_2$  at  $t = 100,000$ .

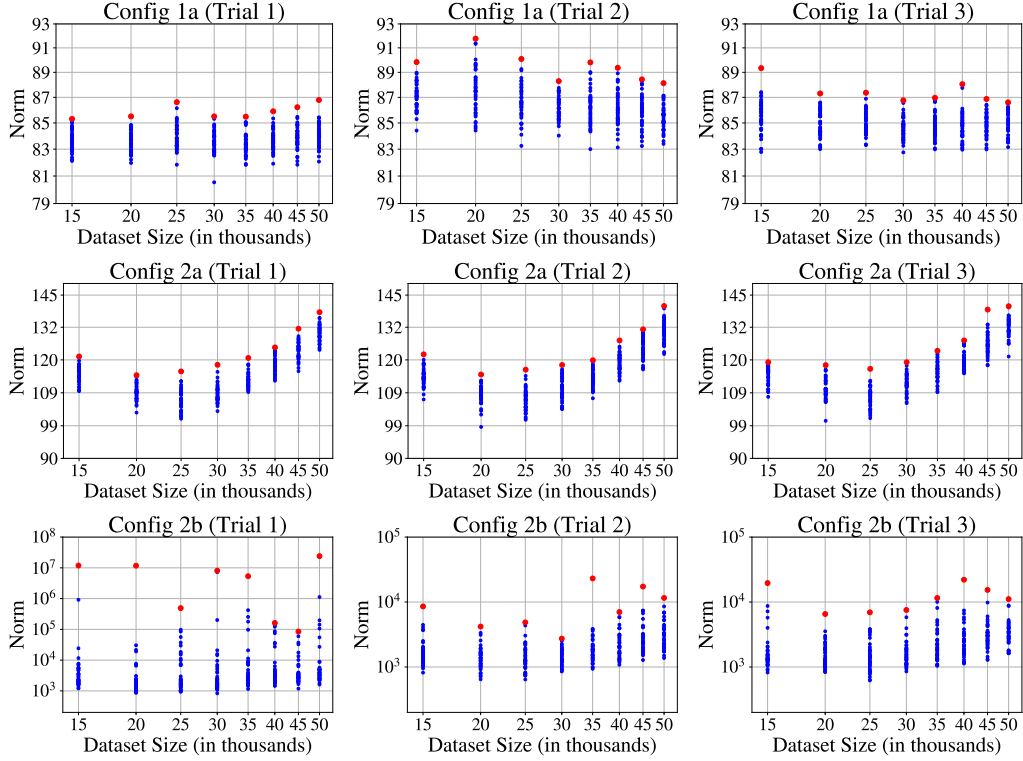


Figure 16:  $\|A(S)\|_2$  at  $t = 100,000$ .

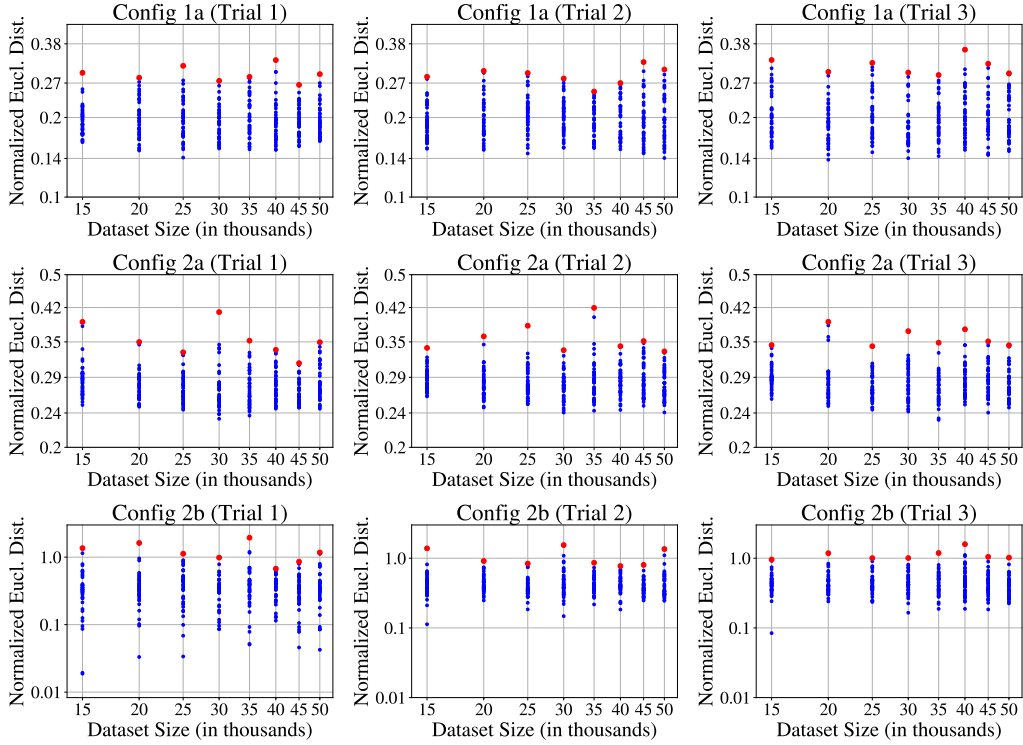


Figure 17: Normalized Euclidean distance. For each Euclidean distance  $\|A(S) - A(S')\|$ , we divide by  $(\|A(S)\| + \|A(S')\|)/2$ . The results suggest that normalizing largely mitigates the growth in Euclidean distance with dataset size; however, this does not appear to yield a significant *decrease* in Euclidean distance with dataset size.

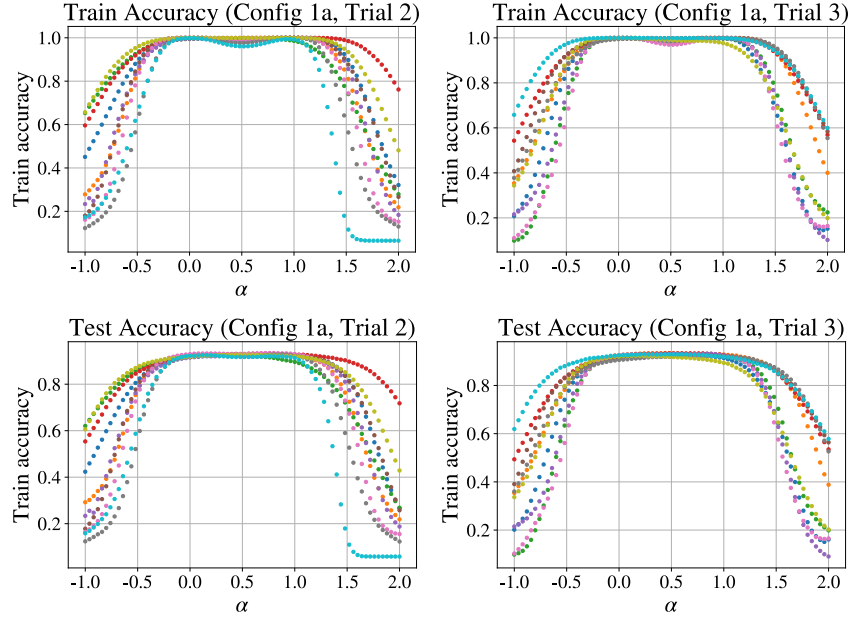


Figure 18: Additional interpolation trials: Configuration 1a (SVHN).

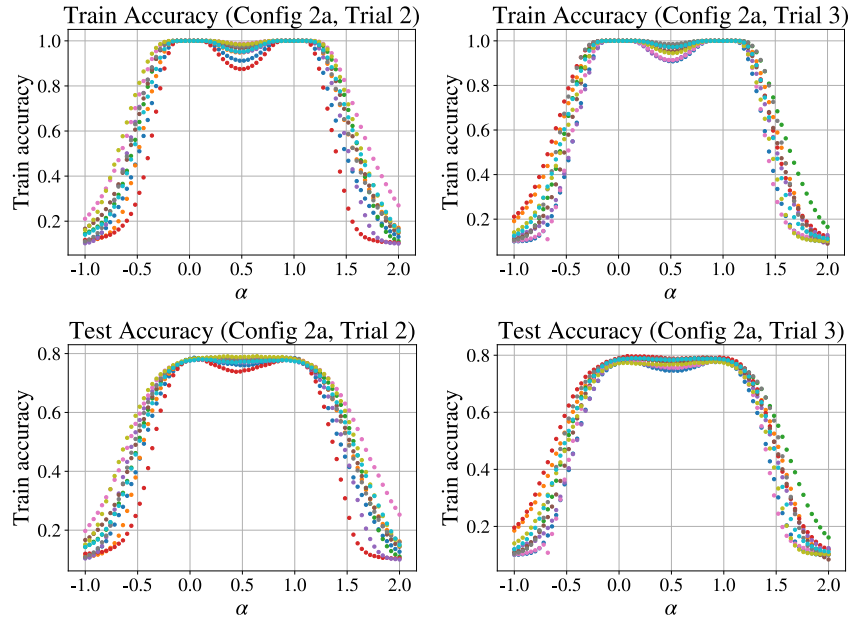


Figure 19: Additional interpolation trials: Configuration 2a (CIFAR-10, no momentum). Note: We omit Configuration 2b from additional trials because the lack of connectivity seen in the body of the paper is not our focus in these additional trials; rather, we are simply interested in confirming cases of linear mode connectivity.