
Distributional Scaling Laws for Emergent Capabilities

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In this paper, we explore the nature of sudden breakthroughs in language model
2 performance at scale, which stands in contrast to smooth improvements governed
3 by scaling laws. While advocates of “emergence” argue that abrupt performance
4 gains arise from acquiring new capabilities at specific scales, recent work has
5 suggested that these are illusions caused by thresholding effects. We propose
6 an alternative explanation: that breakthroughs are driven by random variation,
7 particularly multimodal performance distributions across random seeds. Using a
8 length generalization task as a case study, we show that different random seeds
9 lead to both highly linear or emergent behavior. We further demonstrate that the
10 probability of a model acquiring a breakthrough capability increases continuously
11 with scale, despite apparent discontinuities in performance. Additionally, we
12 find that scaling models in width versus depth has distinct effects: depth impacts
13 the likelihood of sampling from a successful distribution, while width improves
14 the average performance of successful models. These insights suggest a need to
15 consider the role of random variation in scaling and emergent capabilities in LMs.

16 1 Introduction

17 On most benchmarks, language model (LM) performance is determined by a scaling law [12, 8] that
18 responds smoothly to parameter size and overall training compute. There are, however, a number of
19 celebrated exceptions in which performance abruptly improves on specific benchmarks [15].

20 Sudden breakthroughs at scale provide the backdrop of one of the most heated debates in modern
21 machine learning. On one side, advocates of *emergence* claim that performance abruptly improves at
22 particular scales because those scales allow the LM to acquire specific concepts that permit out-of-
23 distribution generalization [17]. On the other side, skeptics argue that these sudden improvements
24 are a *mirage* driven by thresholding effects and alleviated by more appropriate continuous metrics—
25 though a few **breakthrough capabilities** remain stubbornly emergent [13]. Here, we argue that such
26 discontinuities are driven by predictable changes in the *probability* of a breakthrough at each scale.

27 We posit that a breakthrough capability is distinguished not by direct responses to scale, but by
28 *multimodal* random variation. Such effects are currently undetected because scaling laws are generally
29 measured across independent training runs, and resources are rarely committed to correct for random
30 variation by reusing the same hyperparameter settings with different random seeds. Although random
31 variation may be insignificant when model performance is measured in-distribution [6], previous
32 work already suggests that settings which require compositional reasoning may be prone to having
33 performance vary widely across random seeds [19, 20] which can also be seen at larger scales [10].
34 Since training numerous seeds becomes prohibitively expensive at large scales, we study the ability
35 for models to *length generalize* on compositional tasks. Even in these synthetic settings, previous
36 works train on only at most tens of seeds and report summary statistics of these runs. In this work, we
37 seek to characterize the resulting distribution— particularly demonstrating that multimodal variation
38 is present for the studied task. Our contributions can be summarized as follows:

- 39 • **Breakthrough scaling curves can result from bimodal performance distributions.** Using
40 length generalization as a case study, we demonstrate that different random seeds can exhibit
41 either highly linear or highly emergent scale behavior. We connect these differences to
42 the bimodal distribution of this compositional skill across random seeds, a property that
43 materializes at many parameter scales. At these scales, “emergence” is a stochastic property.
- 44 • **Although a given scale curve can exhibit discontinuity (emergence), the probability of a
45 model learning a skill actually responds *continuously* with respect to scale.** Modeling
46 the bimodal distribution as a mix of failure and success distributions, we illustrate that
47 improvements with scale can come from changes in the probability of success or in the
48 average performance of a successful model.
- 49 • **The random variation distribution changes differently when scaling up models with
50 respect to width vs. depth.** Although the scaling laws literature generally treats parameter
51 scale as a single property, we find that scaling in depth changes the probability of sampling
52 from the successful distribution—thereby changing the probability of emergence—whereas
53 scaling in width only changes the average performance of the successful distribution. With
54 respect to emergent capabilities, these are significant differences.

55 2 Methodology

56 Breakthrough capabilities often require compositional reasoning [15, 9, 1]; one such category is the
57 propensity for length generalization [18]. Below we outline our experimental setup, with further
58 details given in Appendix A.

59 **Architecture:** In all of our experiments, we train decoder-only transformer models from scratch,
60 using rotary position embeddings (RoPE) [16] at each layer. To observe the random variation
61 distribution at a variety of scales, we train our models on 250 seeds at a range of 5 scales across three
62 variations of scaling— scaling up width, scaling up depth, and scaling with a fixed parameter count.

63 **Task:** We primarily consider an algorithmic task previously studied in Zhou et al. [20]: learning to
64 count. Given two numbers in increasing order, the model is tasked with outputting a sequence which
65 counts consecutively from the first number to the second number. Examples are given in the form "5,
66 9 >, 5, 6, 7, 8, 9", with the training length representing the length of the counting sequence.
67 It has been previously shown that models can length generalize on this counting task to more than
68 twice their training length; however, we will see in the next section that inspecting the distribution of
69 performance across random variation gives a more nuanced picture about the model’s capacity to
70 length generalize.

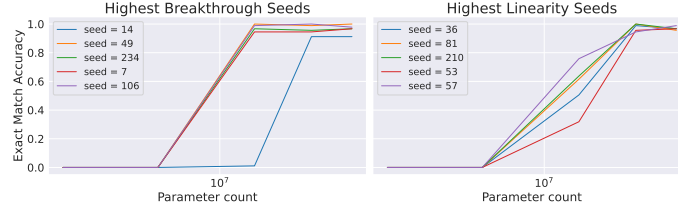
71 **Dataset:** During training, we sample sequences i.i.d from the train set and fill the context with
72 examples, following previous work [5, 20]. The length of examples are sampled uniformly from 1 to
73 the maximum training length, which we fix at 30.

74 3 Results

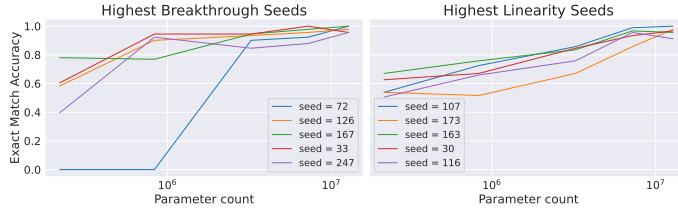
75 3.1 Bimodality leads to emergence

76 Following Srivastava et al. [15], we take the vector of model performances for a fixed seed at test
77 length 60 across different scales and calculate their *breakthroughness* and *linearity* metric (for their
78 definitions, refer to Appendix B). We plot the performance across scale for the top 5 seeds with the
79 highest breakthrough metric and highest linearity metric in Figure 1. Here we fix the random seed
80 which has become the standard practice for reporting benchmark performance for LLMs; although the
81 same random seed has no relation across different scales, this is often a fixed parameter overlooked
82 in practice. In Figure 2, we provide violin plots of the resulting EM accuracy distribution across our
83 three variations of scaling models for test lengths from 30 to 100, with additional figures given in
84 Appendix C.

85 As is clear in Figure 2, many parameter scales exhibit specifically bimodal distributions of length
86 generalization capabilities. The impact of this variability is depicted in Figure 1, which illustrates
87 that we can easily find fixed seeds that show varying levels of emergence and linearity, due to random
88 variation in breakthroughs. Although work on emergence often makes claims as to the specific model



(a) Fixing width to 512 hidden dimension and scaling depth.



(b) Fixing depth to 4 layers and scaling width.

Figure 1: Scaling curves that exhibit emergence (according to the breakthroughness metric in [15]) and those that don't, all sampled from the same task and hyperparameters, fixing random seed for each scale.

89 size required for particular tasks, we show that one can find different levels of emergence and different
 90 breakthrough scales depending on random experimental conditions—particularly in our fixed depth
 91 setting (Figure 1b), where curves range from clear breakthroughs to nearly linear.

92 3.2 Scaling depth versus width

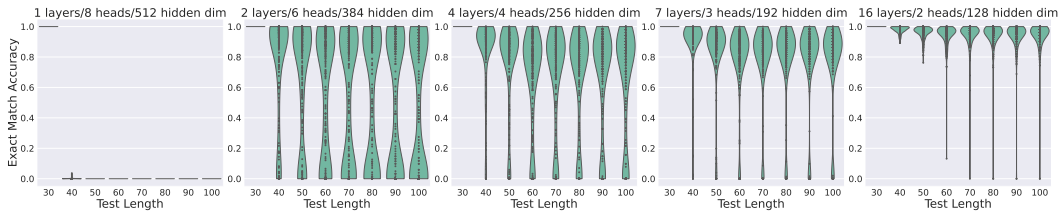
93 Classically, scaling laws treat parameter size as a single scalar value without differentiating between an
 94 increase in width or depth. This difference, though of limited consequence for unimodal capabilities,
 95 becomes crucial when considering breakthrough capabilities. In Figure 2a, we see that, as depth
 96 increases, a fixed parameter count goes from all model runs failing to almost all models scoring over
 97 80% accuracy even at the longest test lengths. When focusing on parameter growth through depth
 98 (Figure 2c), we see mass move from the failed runs to the successful runs; meanwhile, if parameters
 99 grow through hidden width (Figure 2b), the probability of a run failing completely does not decrease,
 100 but the average performance of a successful run increases.

101 The specifics of how these distributions shift is shown in Figure 3. Here, we plot the fraction of runs
 102 which achieve EM accuracy below 10% and above 50% across model scales, as well as the mean of
 103 the runs which achieve above 50% accuracy across scales. First, note that the average accuracy for a
 104 “successful” run (any run with above 50% EM accuracy) increases persistently and monotonically
 105 with width, but not depth. Second, as we increase depth, the proportion of successful runs increases
 106 and the proportion of failed runs (below 10% match accuracy) falls—but as we increase width, neither
 107 trend appears. We conclude that both depth and width can improve performance, but depth improves
 108 performance primarily by shifting the likelihood of a breakthrough, whereas width only improves the
 109 expected performance of a model with the breakthrough capability.

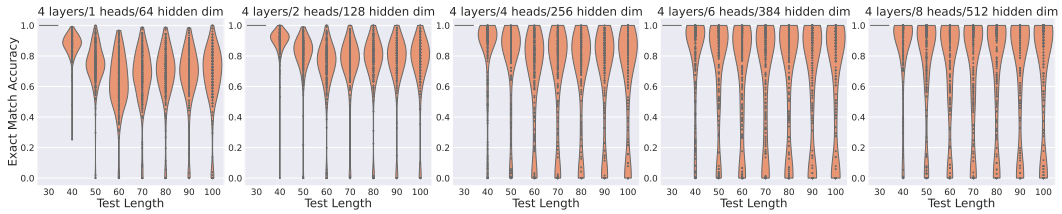
110 4 Discussion and conclusions

111 Zhou et al. [20] first documented the variability of length generalization across random seeds, which
 112 we take advantage of in our work. In general, out-of-distribution behavior like compositional rules
 113 [11] or associative biases [14] often exhibit extreme variation compared to in-distribution performance.
 114 We are also not the first to note bimodal distributions like these, which are also present in text classifier
 115 performance metrics [7] and even in the timing of generalization breakthroughs during training [4].

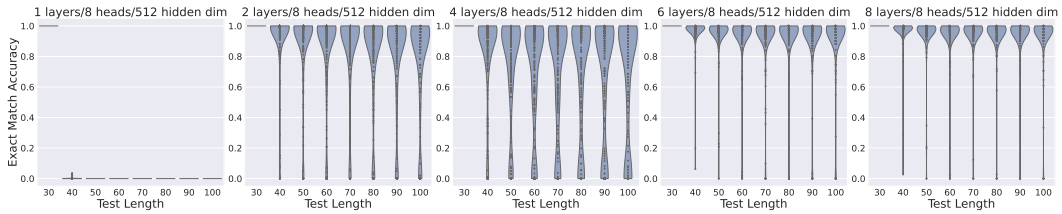
116 Our findings are highly suggestive, but they leave much to future work. Studying other emergent
 117 synthetic tasks, especially other compositional or length generalization settings, would confirm how
 118 general these findings are. Most crucially, we intend to study whether breakthrough capabilities in
 119 language models are associated with bimodal or outlier behavior.



(a) Fixing parameter count to 3.2m parameters.



(b) Fixing network depth to 4 layers.



(c) Fixing network width to 512 hidden dimension.

Figure 2: Violin plots illustrating how the EM accuracy distributions across test lengths vary across a) increasing depth while fixing parameter count, b) increasing width while fixing depth, and c) increasing depth while fixing width. Although the number of parameters remains the same, the distribution in deeper networks exhibits significantly less bimodality. This is further exemplified where at fixed depths, bimodality remains present even at greater widths.

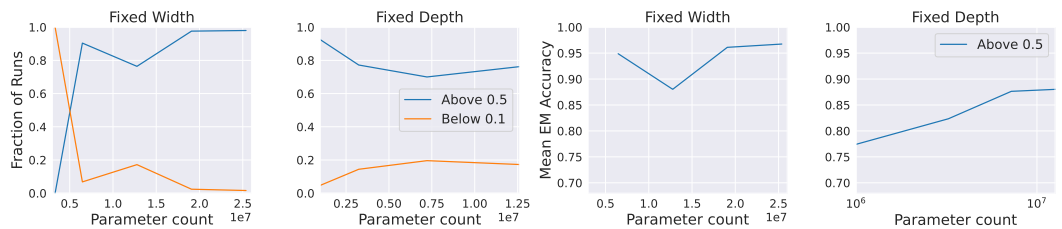


Figure 3: Across fixed width and fixed depth, on the two leftmost plots we show the fraction of runs which achieve above 50% and below 10% EM accuracy, and on the two rightmost plots we show the mean performance of the runs which achieve EM accuracy above 50%.

References

- [1] A. Chen, R. Shwartz-Ziv, K. Cho, M. L. Leavitt, and N. Saphra. Sudden Drops in the Loss: Syntax Acquisition, Phase Transitions, and Simplicity Bias in MLMs. 2024. URL <https://openreview.net/forum?id=M05PikHELW>.
- [2] D. Ganguli, D. Hernandez, L. Lovitt, A. Askell, Y. Bai, A. Chen, T. Conerly, N. Dassarma, D. Drain, N. Elhage, et al. Predictability and surprise in large generative models. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*, pages 1747–1764, 2022.
- [3] J. Hoffmann, S. Borgeaud, A. Mensch, E. Buchatskaya, T. Cai, E. Rutherford, D. d. L. Casas, L. A. Hendricks, J. Welbl, A. Clark, et al. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.
- [4] M. Y. Hu, A. Chen, N. Saphra, and K. Cho. Latent State Models of Training Dynamics. *Transactions on Machine Learning Research*, Aug. 2023. ISSN 2835-8856. URL <https://openreview.net/forum?id=NE2xW0OLF>.
- [5] S. Jelassi, D. Brandfonbrener, S. M. Kakade, et al. Repeat after me: Transformers are better than state space models at copying. In *Forty-first International Conference on Machine Learning*.
- [6] K. Jordan. On the variance of neural network training with respect to test sets and distributions. In *The Twelfth International Conference on Learning Representations*, 2024.
- [7] J. Juneja, R. Bansal, K. Cho, J. Sedoc, and N. Saphra. Linear Connectivity Reveals Generalization Strategies. 2023. URL <https://openreview.net/forum?id=hY6M0JH13uL>.
- [8] J. Kaplan, S. McCandlish, T. Henighan, T. B. Brown, B. Chess, R. Child, S. Gray, A. Radford, J. Wu, and D. Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- [9] A. T. Löwe, L. Touzo, P. S. Muhle-Karbe, A. M. Saxe, C. Summerfield, and N. W. Schuck. Abrupt and spontaneous strategy switches emerge in simple regularised neural networks, 2024. URL <https://arxiv.org/abs/2302.11351>.
- [10] L. Madaan, A. K. Singh, R. Schaeffer, A. Poulton, S. Koyejo, P. Stenertorp, S. Narang, and D. Hupkes. Quantifying variance in evaluation benchmarks. *arXiv preprint arXiv:2406.10229*, 2024.
- [11] R. T. McCoy, J. Min, and T. Linzen. BERTs of a feather do not generalize together: Large variability in generalization across models with similar test set performance. *arXiv:1911.02969 [cs]*, Nov. 2019. URL <http://arxiv.org/abs/1911.02969>. arXiv: 1911.02969.
- [12] J. S. Rosenfeld, A. Rosenfeld, Y. Belinkov, and N. Shavit. A constructive prediction of the generalization error across scales, 2019. URL <https://arxiv.org/abs/1909.12673>.
- [13] R. Schaeffer, B. Miranda, and S. Koyejo. Are Emergent Abilities of Large Language Models a Mirage?, May 2023. URL <http://arxiv.org/abs/2304.15004>. arXiv:2304.15004 [cs].
- [14] T. Sellam, S. Yadlowsky, J. Wei, N. Saphra, A. D’Amour, T. Linzen, J. Bastings, I. Turc, J. Eisenstein, D. Das, I. Tenney, and E. Pavlick. The MultiBERTs: BERT Reproductions for Robustness Analysis. *arXiv:2106.16163 [cs]*, Mar. 2022. URL <http://arxiv.org/abs/2106.16163>. arXiv: 2106.16163.
- [15] A. Srivastava, A. Rastogi, A. Rao, A. A. M. Shueb, A. Abid, A. Fisch, A. R. Brown, A. Santoro, A. Gupta, A. Garriga-Alonso, et al. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *Transactions on Machine Learning Research*, 2023.
- [16] J. Su, M. Ahmed, Y. Lu, S. Pan, W. Bo, and Y. Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.
- [17] J. Wei, Y. Tay, R. Bommasani, C. Raffel, B. Zoph, S. Borgeaud, D. Yogatama, M. Bosma, D. Zhou, D. Metzler, E. H. Chi, T. Hashimoto, O. Vinyals, P. Liang, J. Dean, and W. Fedus. Emergent Abilities of Large Language Models, June 2022. URL <http://arxiv.org/abs/2206.07682>. Number: arXiv:2206.07682 arXiv:2206.07682 [cs].
- [18] H. Zhou, A. Bradley, E. Littwin, N. Razin, O. Saremi, J. Susskind, S. Bengio, and P. Nakkiran. What algorithms can transformers learn? a study in length generalization, 2023. URL <https://arxiv.org/abs/2310.16028>.

- 172 [19] H. Zhou, A. Bradley, E. Littwin, N. Razin, O. Saremi, J. M. Susskind, S. Bengio, and P. Nakkiran.
173 What algorithms can transformers learn? a study in length generalization. In *The Twelfth*
174 *International Conference on Learning Representations, 2024*.
- 175 [20] Y. Zhou, U. Alon, X. Chen, X. Wang, R. Agarwal, and D. Zhou. Transformers can achieve
176 length generalization but not robustly. In *ICLR 2024 Workshop on Mathematical and Empirical*
177 *Understanding of Foundation Models, 2024*.

178 A Additional Experimental Details

179 **Count task:** For all of our training runs, we fix the vocabulary size to 150. For evaluation, we
180 compute the exact match (EM) accuracy across all consecutive subsequences of the test length.

181 **Model scales:** As mentioned in Section 2, we scale up our models in three ways: fixing width and
182 scaling depth, fixing depth and scaling width, and fixing parameter count and scaling depth. The
183 precise parameters for each variation are as follows:

184 • **Fixed depth:** We fix the network depth to 4 layers and vary width by taking hidden
185 dimensions $\{64, 128, 256, 384, 512\}$. The head dimension is fixed to 64.

186 • **Fixed width:** We fix the hidden dimension to be 512 and vary the depth from $\{1, 2, 4, 6, 8\}$
187 layers.

188 • **Fixed parameter count:** We fix the parameter count to 3.2m and vary the depth from
189 $\{1, 2, 4, 7, 16\}$ layers, with the appropriate hidden dimension.

190 **Hyperparameters:** We train all of our models to convergence on the train distribution and use a
191 learning rate of $1e - 3$ with a cosine decay scheduler and weight decay 0.1. We set the maximum
192 training duration to be 10000 steps, with batch size 128 and context length 256.

193 B Breakthroughness and Linearity

194 Srivastava et al. [15] introduced *breakthroughness* and *linearity* metrics to capture model performance
195 improving suddenly or reliably with scale. Given a model’s performances y_i at model scales x_i sorted
196 by ascending model scale, the linearity metric L and breakthroughness metric B are respectively
197 calculated as

$$L = \frac{I(y)}{\text{RootMeanSquare}(\{y_{i+1} - y_i\}_i)},$$
$$B = \frac{I(y)}{\text{RootMedianSquare}(\{y_{i+1} - y_i\}_i)}$$

198 where $I(y) = \text{sign}(\arg \max_i y_i - \arg \min_i y_i)(\max_i y_i - \min_i y_i)$.

199 C Additional Figures

200 In Figures 4, 5, and 6 we provide raw histograms for EM accuracy across test lengths in
201 $\{50, 60, 70, 80, 90, 100\}$ and across model scales fixing parameter count, fixing depth, and fixing
202 width respectively. As highlighted in Section 3.2, for a fixed parameter count we see a transition
203 from all model runs failing to almost all models scoring over 80% accuracy even at the longest test
204 lengths. When fixing depth, there remains a fraction of runs failing even at larger widths, but the
205 average performance of a successful run increases. Finally, when fixing width, mass seems to shift
206 from the failed runs to successful runs, again due to scaling depth.

Exact Match Accuracy - Fixed Parameter Count

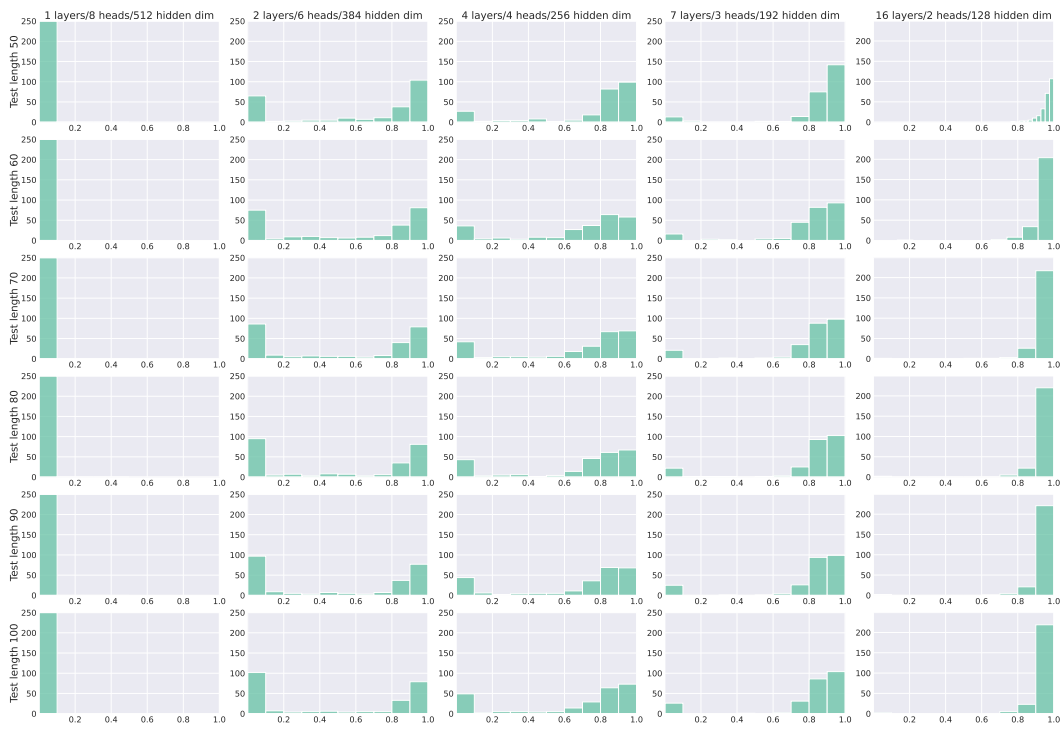


Figure 4: Histogram across model scales fixing parameter count at different test lengths.

Exact Match Accuracy - Fixed Depth



Figure 5: Histogram across model scales fixing network depth at different test lengths.

Exact Match Accuracy - Fixed Width



Figure 6: Histogram across model scales fixing network width at different test lengths.

207 **D Debunking Challenge Submission**

208 **D.1 What commonly-held position or belief are you challenging?**

209 It is now widely recognized that increasing the scale of language models (e.g., in terms of training
210 computation and model parameters) can result in improved performance and sample efficiency across
211 various downstream NLP tasks. However, there are cases where performance can be predicted
212 smoothly in terms of scale via scaling laws [3, 8] as well as cases where performance cannot be
213 predicted as scale increases due to discontinuous improvements [2, 15, 17]. Our work simultaneously
214 addresses both common positions on emergent capabilities. The first position is that continuous
215 metrics can always make them discrete and smooth continuous scaling laws always hold, even
216 in apparent breakthroughs. The second position is that they represent exceptions to scaling laws
217 determined by capacity at scale. Prior work has called into question the empirical validity of
218 emergence; for instance, Schaeffer et al. [13] argued that more than 92% of proposed emergent
219 skills were caused by thresholding effects in accuracy metrics. This does however, still leave a
220 number of potential emergent abilities. Thus, our work addresses both perspectives by considering
221 the distribution across performance due to random variation, where both emergence and linear
222 improvement can be present for the same task and training hyperparameters.

223 **D.2 How are your results in tension with this commonly-held position?**

224 We propose continuous scaling laws in the *distribution* of capabilities rather than the sampled scalar
225 value associated with that capability for a single trained model. By addressing distributions rather than
226 point samples, we explain emergent behavior using scaling laws that are entirely continuous—but
227 *bimodal*. In other words, the notion of a pointwise scaling “law” leads inevitably to discontinuities, but
228 these discontinuities are stochastic in nature and determined by a continuous underlying distribution.
229 We show the potential issue with only considering point samples in Figure 1, where even on the same
230 hyperparameters, different random seeds exhibit varying degrees of breakthroughness and linearity.

231 Another element of tension is with the notion of a breakthrough parameter size. While the scaling
232 laws literature tends to treat depth and width equivalently—as they may be in-distribution—we
233 show that depth and width are, in fact, materially different in how they change the probability of
234 a breakthrough. We discuss this in depth in Section 3.2, where upon inspecting the distribution
235 across scale, we see that increasing depth increases the likelihood of sampling from a “successful”
236 distribution, whereas width improves the average performance of “successful” models.

237 **D.3 How do you expect your submission to affect future work?**

238 When targeting emergent properties, our initial results suggest that depth may be prioritized over
239 width. Furthermore, future work on emergence may involve training multiple seeds—though of
240 course, at very large scales this may be resource-intensive. In general, we hope our results provide
241 further evidence to the susceptibility of research on emergence in LLMs to statistical artifacts like
242 the Texas sharpshooter fallacy. Although it becomes computationally prohibitive to run multiple
243 seeds with the same set of hyperparameters to statistically validate empirical observations, we believe
244 the effect of random variation on out-of-distribution or downstream performance should not be
245 disregarded especially if the variation is not Gaussian and can be presented as multimodal. There
246 are several directions for future work, such as investigating the generalizability of our findings and
247 particularly the presence of bimodal distributions across different tasks, investigating influences
248 of various hyperparameters to the random variation distribution, and studying seed performance
249 correlation in the multi-task setting.