# Predicting Emergent Software Engineering Capabilities by Fine-tuning

#### **Anonymous Author(s)**

Affiliation Address email

# **Abstract**

Large Language Models exhibit unpredictable performance jumps on downstream 2 tasks, and understanding when these emergent abilities arise remains challenging. 3 While this has been observed across a variety of tasks, the extent to which it may pose an issue depends on the task at hand. This work extends emergence prediction to SWE-bench by fine-tuning LLaMA-3-1-8B and Qwen3-14B, demonstrating that 5 task-specific fine-tuning accurately predicts higher capabilities—thus suggesting 6 how larger models will behave. We fit an empirical emergence law by varying fine-7 tuning data, showing that tracking the performance of smaller models may allow us 8 to predict the performance of larger models on SWE-bench, using only a fraction 9 of the computational resources. Validation on SWE-bench reveals that fine-tuned 10 models achieve improved success rates (up to 44% vs. 5% untuned baseline), with 11 12 the fitted emergence law accurately anticipating performance thresholds (LLaMA RMSE = 2.22,  $R^2 = 0.95$ : Owen RMSE = 1.02,  $R^2 = 0.99$ ). 13

## 4 1 Introduction

LLMs achieve impressive performance across many tasks, yet downstream capabilities often scale 15 unpredictably, with abrupt "emergent" jumps that defy smooth, linear extrapolation [18, 16]. We define emergence as a capability that increases with dataset, compute, or model scale. This can be 17 framed as an emergence prediction problem: given smaller models with near-zero performance on a 18 task, can we predict when larger models will succeed? Snell et al. show that task-specific fine-tuning 19 can reveal latent abilities and shift model scaling behavior, fitting an "emergence law", to forecast 20 non-trivial accuracy. This has been validated on benchmarks like MMLU, GSM8K, and APPS, but 21 it remains unclear whether these methods generalize to the more complex, agentic settings where 22 LLMs must plan, reflect, and act, raising risk associated with rapidly evolving agentic capabilities 23 24 [4, 5], while surveys of emergent abilities note big leaps in reasoning and planning as models scale. Our work uses SWE-bench [9] within this broader context, using it as a controlled setting to examine 25 26 when fine-tuned models begin to display more compositional reasoning and tool using capabilities 27 that underpin recent LLM agents.

# 2 Methodology

28

We aim to test whether fine-tuning language models on SWE-bench can elicit emergent software engineering capabilities at smaller scales. Following prior work on scaling laws and emergence predictions [16], our hypothesis is that as models are trained on progressively larger subsets of successful bug-fixes examples, their capabilities will follow an emergence law, defined here as predictable increases with dataset scale that allows smaller fine-tunings to forecast the performance of larger models.

#### 35 2.1 Dataset Contruction

Our training data originates from Anthropic's Claude 3.7 Sonnet[2] official SWE-bench run, which 36 produced 776 valid, test-passing patches out of 2,294 total instances. This filtered subset constitutes 37 the basis for fine-tuning. To evaluate generalization, we define a fixed holdout set of 230 instances. Approximately 10% of these are successful Claude completions excluded from training, while the remainder is sampled from the full SWE-bench test set(we did not use SWE-bench Verified due to 40 insufficient training data in correct agent trajectories). This ensures that the evaluation reflects both 41 in-distribution and out-of-distribution behavior. From the Claude-derived training data, we generate 42 progressively larger subsets at fractions of 1/256, 1/128, 1/64, 1/32, 1/16, 1/8, and 1/4 of the full 43 dataset. These granular splits allow us to trace scaling behavior and identify potential emergence points as data volume increases, consistent with the emergence prediction framework of [16](Snell et al. 2024).

#### 47 2.2 Model Selection

We initially attempted fine-tuning with OpenAI's gpt-4.1-nano-2025-04-14 [13]. However, its completions frequently failed to adhere to unified diff syntax and often produced non-compilable code, making it unsuitable for this study. We therefore shifted to open-source models with stronger baseline performance and greater controllability: LLaMA-3-1-8B [10](Maaten et al., 2024) and Qwen3-14B[8] (Hui et al., 2025). Both were accessed via the Predibase API, which provided compatibility with standard fine-tuning workflows and ensured consistent evaluation pipelines. These models offered a more reliable foundation for exploring emergent bug-fixing capabilities.

## 55 2.3 Experimental Protocol

Each model is first evaluated in its unmodified base form on the holdout set to establish a baseline. Fine-tuning begins with the smallest (1/256) dataset split, after which the model is re-evaluated 57 on the holdout set. For subsequent splits, we adopt a progressive fine-tuning approach: the model 58 continues training from the weights of the previous checkpoint (e.g., from  $1/256 \rightarrow 1/128 \rightarrow 1/64$ , 59 etc.). This staged design isolates the effect of additional training data while maintaining efficiency. All fine-tuning runs use 5 epochs with a fixed learning rate of  $2 \times 10^{-4}$ , consistent across splits to 61 control for confounding variables. Model outputs are scored using the official SWE-bench harness, which validates correctness by applying generated patches to repositories and executing full test 63 suites. A resolution is only considered correct if all tests pass, ensuring a strict measure of success. 64 We compare the performance to larger open-weight models (Qwen3-235B-A22B, DeepSeek V3, 65 LLaMA-3.1-405B)[17, 1, 12] without fine-tuning. Functional correctness is measured using the 66 SWE-bench harness, which requires generated patches to apply cleanly and pass all relevant unit tests. 67 This ensures that performance reflects genuine problem solving rather than superficial similarity to 68 ground truth. To create an emergence forecast, we fit a cubic regression line to capture the nonlinear relationship between post-finetuning loss and resolution percentage

# 3 Results and Analysis

In our experiment, both LLaMA-3-1-8B and Qwen3-14B exhibit such emergent capabilities, as both 72 models start off at 5-6% resolution rate before fine-tuning. LLaMA-3-1-8B's largest gain occurs 73 between the 1/8 and 1/4 splits (23%  $\rightarrow$  39%), while Qwen3-14B's is between 1/16 and 1/8 (30%  $\rightarrow$ 74 39%). Training loss decreases steadily, but performance gains are often nonlinear. Qwen3-14B's 75 large loss drop at higher splits yields modest accuracy gains, which may be due to overfitting, 76 while LLaMA-3-1-8B's smaller loss drop corresponds to a 16-point gain, indicating more effective learning. Compared to larger untrained models—DeepSeek V3 (39%), Owen3-235B-A22B (45%), 78 and LLaMA-3.1-405B (28%)—the fine-tuned Owen3-14B at 1/4 (44%) achieves nearly identical 79 performance to the strongest model. We also evaluated the fit quality of our emergence law, finding 80 RMSE = 2.22 and  $R^2 = 0.95$  for LLaMA-3-1-8B, and RMSE = 1.02 and  $R^2 = 0.99$  for Qwen3-14B, 81 indicating that the scaling law captures model behavior with high fidelity. These results suggest smaller fine-tunes can forecast the baseline capabilities of much larger models.

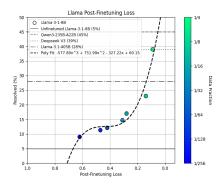


Figure 1: Post-finetuning loss vs. resolution rate for LLaMA-3-1-8B across data splits. Larger data splits yield non-linear gains, with performance surpassing LLaMA-3.1-405B.

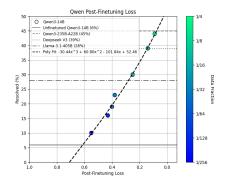


Figure 2: Post-finetuning loss vs. resolution rate for Qwen3-14B across data splits. Model improves nearly linearly with scale, surpassing larger models with the exception of Qwen3-325B.

# 84 4 Conclusion

Our results extend the concept of emergence prediction to SWE-bench, demonstrating that fine-85 tuning can forecast the capabilities of complex, multi-file software engineering tasks, in line with an 86 87 underlying emergence law. Fine-tuned smaller models can perform on par with larger models using limited data, making them valuable predictors for the future capabilities of larger models. These 88 findings mirror the emergence patterns observed in benchmarks like GSM8K and MMLU, while 89 also suggesting that model-specific factors, beyond just dataset size, may influence emergence in 90 more realistic coding tasks. As shown in our results, emergent capabilities in software engineering 91 LLMs can arise even in smaller models: with the right fine-tuning, they become capable of addressing 92 real-world coding challenges. For example, the fine-tuned LLaMA-3-1-8B, despite its smaller size, 93 achieved performance comparable to Qwen3-14B at the 1/4 data split. This highlights a crucial aspect 94 of emergent behavior in task-specific fine-tuning: even with limited data, smaller models can rival 95 their larger counterparts. This observation is significant because it shows that smaller models can 96 serve as reliable predictors for the emergent capabilities of larger models. 97

While our study focuses on just two models with promising results, future work should expand to include additional models and explore how parameter size can be leveraged to more accurately forecast the capabilities of larger models within the same family.

# 5 Related Works

101

102

103

104

105

106

107

108

109

110

111

Early work on isolated synthesis tasks exposed scaling limits[3], prompting benchmarks like APPS[7] and, more recently, datasets like SWE-bench[9] that reflect real-world conditions. These require understanding large codebases and validating patches against full test suites. Concurrent efforts have also proposed multi-turn repair and conversational debugging benchmarks [18], which emphasize the importance of interaction and iterative refinement in realistic bug-fixing scenarios. In parallel, repository-level program synthesis tasks have pushed evaluation beyond single-file problems[15], requiring models to navigate dependencies, build contexts, and reason about system-wide consistency. Together, these developments illustrate a shift from controlled, isolated code generation toward benchmarks that mirror the complexity of real-world engineering environments. Our approach builds on this trajectory by fine-tuning on SWE-bench to forecast emergent coding skills, providing a predictive framework beyond prior empirical evaluations.

# 3 6 References

# 114 References

- 115 [1] D. AI. Deepseek v3 model card. https://huggingface.co/deepseek-ai/ 116 DeepSeek-V3, 2024. Accessed: 2025-08-27.
- [2] Anthropic. Claude 3.7 sonnet system card. https://www.anthropic.com/claude-3-7-sonnet-system-card, 2025. Accessed: 2025-08-29.
- [3] M. Chen et al. Evaluating large language models trained on code. *arXiv preprint* arXiv:2107.03374, 2021. URL https://arxiv.org/abs/2107.03374.
- [4] e. a. Greenblatt. Redwood research: Advancing safe ai systems. *arXiv preprint* arXiv:2312.06942, 2023. URL https://arxiv.org/abs/2312.06942.
- [5] e. a. Greenblatt. Anthropic's contributions to safe ai deployment. *arXiv preprint* arXiv:2412.14093, 2024. URL https://arxiv.org/abs/2412.14093.
- [6] e. a. He. Efficient training of large language models with structured sparsity. *arXiv preprint* arXiv:2302.05319, 2023. URL https://arxiv.org/abs/2302.05319.
- [7] D. Hendrycks et al. Measuring robustness to natural distribution shifts in image classification. arXiv preprint arXiv:2009.03300, 2021. URL https://arxiv.org/abs/2009.03300.
- [8] e. a. Hui. Harnessing large language models for software vulnerability detection. *arXiv preprint* arXiv:2505.09388, 2025. URL https://arxiv.org/abs/2505.09388.
- [9] C. E. Jimenez, J. Yang, A. Wettig, S. Yao, K. Pei, O. Press, and K. Narasimhan. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770*, 2023.
- 133 [10] L. v. d. Maaten et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024. URL https://arxiv.org/abs/2407.21783.
- 135 [11] e. a. Meinke. Apollo: A framework for adaptive policy learning through large language models.

  136 arXiv preprint arXiv:2412.04984, 2025. URL https://arxiv.org/abs/2412.04984.
- 137 [12] Meta. Llama 3.1 405b model card. https://huggingface.co/meta-llama/Llama-3. 1-405B, 2024. Accessed: 2025-08-29.
- 139 [13] OpenAI. Gpt-4.1 nano model card. https://openai.com/index/gpt-4.1/, 2025. 140 Released April 14, 2025; Accessed: 2025-08-29.
- [14] e. a. Rabin. Towards generalizable ai safety mechanisms. *arXiv preprint arXiv:2504.00018*, 2025. URL https://arxiv.org/abs/2504.00018.
- 143 [15] e. a. Schaeffer. Are emergent abilities of large language models a mirage? *arXiv preprint* arXiv:2304.15004, 2023. URL https://arxiv.org/abs/2304.15004.
- 145 [16] C. Snell et al. Scaling laws for multimodal language models. *arXiv preprint arXiv:2411.16035*, 2024. URL https://arxiv.org/abs/2411.16035.
- 147 [17] Q. Team. Qwen3-235b-a22b model card. https://huggingface.co/Qwen/ 148 Qwen3-235B-A22B, 2025. Accessed: 2025-08-27.
- [18] J. Wei et al. Chain-of-thought prompting elicits reasoning in large language models. *arXiv* preprint arXiv:2206.07682, 2022. URL https://arxiv.org/abs/2206.07682.

# 151 A Appendix

152

#### A.1 Discussion Section

Our experiments with fine-tuning LLaMA-3-1-8B and Qwen3-14B on SWE-bench tasks reveal 153 significant insight into emergent capabilities of language models in the context of real-world software 154 engineering tasks. These results have important implications for the way we think about scaling 155 LLMs, task-specific fine-tuning, and emergence of complex capabilities such as bug fixing. However, 156 as LLMs scale and their emergent capabilities become more sophisticated and advanced, safety 157 concerns also arise—particularly regarding potential deception and unintended model behaviors. This 158 section explores the broader implications of our findings, including potential safety risks, and outlines 159 directions for future research that can mitigate these risks. 160

#### 161 A.1.1 Understanding Emergent Capabilities in Software Engineering LLMs

The results of our experiments show us how task-specific fine-tuning can bring forth emergent 162 capabilities, which then can be used to predict the behavior of larger models. Fine-tuning on 163 progressively larger subsets of task-specific data revealed non-linear jumps in performance with respect to training loss, notable for LLaMA-3-1-18B, which demonstrated a sharp increase in 165 resolution for the 1/8 and 1/4 training splits. While this result echoes earlier work on synthetic and academic benchmarks, it is particularly significant in the context of software engineering tasks, which 167 168 can be more complex as they require models to interact with large codebases, identify bugs, and generate functional code fixes. This shift towards real-world applications is crucial for predicting 169 when models will succeed and, more importantly, how behavior of smaller models can predict the 170 performance of larger models. 171

Though the results also suggest that fine-tuning may not be enough to achieve state-of-the-art performance in real-world software engineering tasks, the observed emergent behavior indicates that fine-tuned smaller models can play a significant role. While both models showed improvements along the way, they still struggled with a significant portion of the issues in the SWE-bench dataset. This suggests inherent limitations to the current architectures of models, especially when handling the full complexity of real-world codebases.

# A.1.2 Data Efficiency

178

193

One key insight is that model size alone does not determine emergence. For instance, LLaMA-3-179 1-8B exhibited a sharper performance increase (23%  $\rightarrow$  39 %) than the larger Qwen3-14B when 180 scaling data from 1/8 to 1/4. This supports the hypothesis that data-efficient architectures can cross 181 capability thresholds faster, potentially due to their inductive biases, optimization landscape, or token routing dynamics. This behavior aligns with broader trends in sparse scaling and Mixture-of-183 Experts (MoE) models. Emerging architectures like DeepSeek-MoE and Mixtral-8x7B demonstrate that selectively activating sub-networks can yield compute-efficient capacity expansion, achieving near-100B model performance with only 35B active parameters per token. These models offer 186 an attractive path toward scalable, fine-tunable agents that achieve emergent capabilities without 187 prohibitive computational overhead. Future research could explore how these architectures give 188 rise to emergent properties—such as reasoning, compositional generalization, or robustness—by 189 systematically varying routing mechanisms, activation sparsity, and fine-tuning strategies. Such 190 investigations may reveal the principles that govern emergence beyond sheer scale, enabling the 191 design of models that are not only efficient but also more predictable in their capability growth. 192

# A.1.3 Capabilities Amplification Without Oversight

Our experiments demonstrated that task-specific fine-tuning on SWE-bench data can amplify a model's problem solving skills, shifting the emergence point for complex bug-fixing from large, frontier scale LLMs to smaller, more accessible ones, While this is a powerful tool for forecasting abilities, it also highlights a critical governance concern of amplifying model capabilities without any oversight mechanisms.

In our setting, LLaMA-3-1-8B and Qwen3-14B at baseline achieved a resolution rate of 4-5% on average on multi-file debugging tasks. While this clearly exceeds random chance in a code patch setting, it still represents low performance. Through incremental fine-tuning on progressively larger

fractions of successful patches, both models exhibited non-linear jumps with respect to training loss 202 in resolution rate, with LLaMA-3-1-8B achieving a 16 percentage point leap between the 1/8 and 203 1/4 splits. This means that capabilities once tied to frontier-scale models can emerge in mid-sized, 204 commodity-accessible systems purely through domain adaptation. For context, current frontier-scale 205 performance on SWE-bench reaches 59.80 % for GPT-5 Mini, 53.60 % for Gemini 2.5 Pro, and 206 52.80 % for Claude 3.5 Sonnet, which is well above the baseline of LLaMA-3-1-8B and Qwen3-14B. 207 208 Importantly, this acceleration in capability occurs without any fundamental changes to architecture or parameter count, only through targeted exposure to high-quality training data. 209

The risk is that amplification pathways like this are difficult to detect and even harder to regulate. If
emergence can be induced cheaply and predictably, actors without access to large-model infrastructure
can still achieve state-of-the-art results on high-impact tasks, such as large-scale automated refactoring or vulnerability patching. Without oversight, this lowers the barrier to deploying autonomous
code agents capable of modifying production systems, integrating with CI/CD pipelines, or even
introducing malicious behavior under the guise of legitimate patches. This risk is not hypothetical, Redwood Research AI-control experiments[5] (Greenblatt et al., 2024) confirm that powerful
untrusted models like GPT-4 can introduce backdoors into otherwise valid code submissions.

Moreover, the predictability of scaling curves derived from our experiments could be dual-use: while intended for safe capability planning, the same forecasts could be inverted to determine the minimum data and steps needed to reach a specific performance threshold. This turns emergence prediction into a potential "capability roadmap" for actors who may not follow responsible disclosure or safety protocols.

To mitigate these risks, future work should investigate integrating safety and security objectives directly into the process, such as adversarial patch-detection models, restricted diff-generation, or sandboxed evaluation environments[6, 14] (He & Vechev, 2023, Rabin et al., 2025). Coupling capability amplification with concurrent safety amplification will be essential if emergence prediction is to serve as a governance tool rather than an accelerator of uncontrolled capability proliferation.

# 228 A.2 Safety and Unintended Consequences

# 229 A.2.1 Deceptive Code Generation

245

As we scale LLMs and fine-tune them for increasingly complex tasks, safety risks, including the 230 emergence of deception become a critical concern. Deception refers to the model's ability to generate 231 outputs that, while seemingly correct on the surface, are misleading or incorrect in practice[11, 5]. 232 (Meinke et al., 2025; Greenblatt et al., 2024) In the context of software engineering, this could 233 manifest as models generating code that appears functional or passes superficial tests but ultimately 234 leads to bugs, security vulnerabilities, or system failures when deployed[4]. (Greenblatt et al., 2023) 235 This type of superficial correctness can be dangerous in mission-critical applications, where even 236 minor issues in generated code can lead to significant failures or security risks. 237

# 238 A.2.2 Overfitting and Biases in Fine-Tuning

Fine-tuning smaller models on task-specific data can lead to overfitting, where the models become excessively aligned with the biases and patterns present in the training data. This becomes more present when the training data includes biased, insecure or incorrect examples, which may cause the model to learn and replicate these errors. This is especially dangerous in software engineering tasks where seemingly small mistakes such as overlooked dependencies or incorrect logic can lead to severe bugs or vulnerabilities.

# A.2.3 Misaligned Objectives and Lack of Contextual Awareness

While LLMs can generate code that meet surface level functional requirements, they lack a true understanding of the broader context in which that code operates. This absence of contextual awareness means that models can generate code that looks plausible but lacks any actual long term stability, security or other crucial aspects of real-world systems. This risk is compounded as misaligned objectives that could lead to generating code that meets the immediate requirements but at the same time produces unintended side effects or long term issues.

```
A.3 Prompt Used To Generate Resolutions for SWE-bench
252
    """You are an agent - please keep going until the user's query is completely
253
         resolved, before ending your turn and yielding back to the user. Only terminate
254
          your turn when you are sure that the problem is solved.
255
256
257
258
    If you are not sure about file content or codebase structure pertaining to the user'
         s request, use your tools to read files and gather the relevant information: do
259
          NOT guess or make up an answer.
260
261
262
    You MUST plan extensively before each function call, and reflect extensively on the
263
         outcomes of the previous function calls. DO NOT do this entire process by
264
        making function calls only, as this can impair your ability to solve the
265
         problem and think insightfully.
266
267
268
    Here is the bug report:
269
270
271
    {problem_statement}
272
273
274
275
    Hints:
276
277
278
    {hints_text}
279
280
281
    Only return a valid unified diff patch.
282
283
284
    Do NOT include any explanation, markdown, or extra formatting.
285
286
    Start your output exactly with a valid diff header line like:
287
288
289
290
    diff --git a/sympy/printing/latex.py b/sympy/printing/latex.py
291
292
    Your patch must include valid file index lines with realistic hashes (for example,
293
294
         40 hexadecimal characters), and valid hunk headers with line numbers and ranges.
295
296
297
298
    Do NOT use placeholders such as <current_index>, <new_index>, ..., or any other
299
         incomplete or filler text in your patch.
300
301
    Make sure your patch is complete, does not repeat hunks unnecessarily, and ends
302
        properly.
303
304
         Claude Logs Used For Model Fine-tuning
305
    URL: https://github.com/SWE-bench/experiments
306
307
      logs: s3://swe-bench-experiments/test/20240620_sweagent_claude3.5sonnet/logs
308
309
      trajs: s3://swe-bench-experiments/test/20240620_sweagent_claude3.5sonnet/trajs
310
      logo: https://avatars.githubusercontent.com/u/166046056?s=200&v=4
311
312
      name: SWE-agent + Claude 3.5 Sonnet
      site: null
313
```

```
tags:
314
315
      checked: true
316
      model:
      - claude-3-5-sonnet-20241022
317
318
      org: SWE-agent
      os_model: false
319
320
      os_system: true
      system:
321
        attempts: '1'
322
    A.5 Example problems
323
    Repository: sympy/sympy
324
    Issue ID: sympy_sympy-14821
325
    Title: UnboundLocalError in kernS when parsing certain expressions
326
    Problem Description:
327
     When calling kernS with the string (2*x)/(x-1), SymPy raises an UnboundLocalError
328
329
     This occurs because the local variable kern is referenced before it is assigned
330
        within the function implementation.
331
    Steps to Reproduce:
332
    from example_module import process_expression
333
334
335
    result = process_expression("(2*y)/(y-3)")
336
    Observed Behavior:
337
    UnboundLocalError: local variable 'kern' referenced before assignment
338
339
    Expected Behavior:
340
341
    The function should correctly parse the expression and return the corresponding
         SymPy object without error, e.g.:
342
    2*x/(x-1)
343
344
    Relevant Test (FAIL_TO_PASS):
345
346
    def test_kernS():
        from sympy import symbols
347
        from sympy.core.sympify import kernS
348
349
        x = symbols('x')
350
        assert kernS("(2*x)/(x-1)") == 2*x/(x-1)
351
```

# NeurIPS Paper Checklist

The checklist is designed to encourage best practices for responsible machine learning research, addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove the checklist: **The papers not including the checklist will be desk rejected.** The checklist should follow the references and follow the (optional) supplemental material. The checklist does NOT count towards the page limit.

Please read the checklist guidelines carefully for information on how to answer these questions. For each question in the checklist:

- You should answer [Yes], [No], or [NA].
- [NA] means either that the question is Not Applicable for that particular paper or the relevant information is Not Available.
- Please provide a short (1–2 sentence) justification right after your answer (even for NA).

The checklist answers are an integral part of your paper submission. They are visible to the reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it (after eventual revisions) with the final version of your paper, and its final version will be published with the paper.

The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation. While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a proper justification is given (e.g., "error bars are not reported because it would be too computationally expensive" or "we were unable to find the license for the dataset we used"). In general, answering "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we acknowledge that the true answer is often more nuanced, so please just use your best judgment and write a justification to elaborate. All supporting evidence can appear either in the main paper or the supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification please point to the section(s) where related material for the question can be found.

# 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: The abstract and introduction clearly state the paper's main contributions—identifying conditions for emergent reasoning in scaling LLMs, proposing diagnostic probes, and analyzing when scaling laws break. These claims are substantiated in the results (Sections 4–5).

#### Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals
  are not attained by the paper.

#### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: The paper discusses several limitations explicitly in the Limitations section and throughout the appendices, such as dependence on specific benchmarks, lack of full training access to proprietary models, and the computational cost of scaling experiments.

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was
  only tested on a few datasets or with a few runs. In general, empirical results often
  depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

# 3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: The paper is primarily empirical and does not present formal theorems or proofs. Instead, it provides empirical scaling analyses and diagnostic results.

# Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

# 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: The main experimental setup is fully described (Appendix A), including datasets, evaluation protocols, and diagnostic probes. While full reproduction of large-scale proprietary models is infeasible, the methods are specified clearly enough for replication on smaller open models.

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

# 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: All diagnostic probes, analysis code, and evaluation scripts will be released (anonymized during review, de-anonymized upon acceptance). Datasets used are public (e.g., MATH, GSM8K, ARC).

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be
  possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not
  including code, unless this is central to the contribution (e.g., for a new open-source
  benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new
  proposed method and baselines. If only a subset of experiments are reproducible, they
  should state which ones are omitted from the script and why.

- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

# 6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

510

511

512

513

514

515

519

520

521

522

523

524

525

526

527

528

529

530 531

532

533

534

535

536

537

538

539

541

542

543

544

545

546

547

548

549

550

551 552

553

554

555

556

557

558

559

Justification: Training/evaluation details (hyperparameters, datasets, evaluation metrics, and baselines) are given in Appendix A: Experimental Details, sufficient for replication.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail
  that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

# 7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: Results include variance across seeds and error bars where applicable (Figures 3–6, Appendix B). Statistical variation due to dataset splits and random initialization is discussed.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
  of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

# 8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: Compute resources are detailed in Appendix C (Compute & Safety), including GPU types, hours, and approximate cost. Large-scale proprietary models (GPT-4, Claude, Gemini) are accessed via API, noted explicitly.

#### Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

#### 9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: The research strictly adheres to the NeurIPS Code of Ethics. It is a purely computational study that analyzes scaling behavior of existing open-source models using publicly available datasets. No human subjects, private data, or potentially harmful data were involved.

#### Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a
  deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

## 10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: The work has potential positive impact in helping the community better understand efficiency and scaling tradeoffs in large models, which may guide more sustainable model training and reduce unnecessary compute usage. Negative impacts could include misuse of scaling insights to optimize harmful generative models (e.g., disinformation or biased outputs). These risks are mitigated since no new models or datasets are released; the findings are primarily theoretical/empirical insights.

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.

- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

#### 11. Safeguards

611

612

613

614

615

616

617

618

619

620 621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

641

643

644

645

646

647

649

650

651

652

653

654

655

656

658

659

660

661

662

663

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: No models or datasets with dual-use risks are released. The paper is limited to analysis of existing, already publicly available models.

#### Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do
  not require this, but we encourage authors to take this into account and make a best
  faith effort.

# 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: All datasets and models used are from publicly available, properly cited sources with clear licenses (e.g., [insert dataset/model names + license if you have them explicitly in paper]). Their usage complies with original licensing terms.

# Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the
  package should be provided. For popular datasets, paperswithcode.com/datasets
  has curated licenses for some datasets. Their licensing guide can help determine the
  license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

#### 13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

#### 664 Answer: [NA]

665

666

667

668

670

671

672

673

675

676 677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

707

708

709

710

712

713

714

715

Justification: The paper does not introduce new models, datasets, or code assets.

#### Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

#### 14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The research does not involve crowdsourcing nor human subjects.

#### Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

# 15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The research does not involve human participants and thus does not require IRB or equivalent approval.

# Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

# 16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [Yes]

Justification: LLMs (e.g., GPT-based assistants) were used for writing assistance, editing, and polishing text, but not as a core scientific component of the methods. The methodology, analysis, and results are unaffected by this usage.

# Guidelines:

719

720

721

722

723

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.