
Static Benchmarks Are Broken: The Case for Dynamic Evaluation of LLMs

Anonymous Authors¹

Abstract

Static, deterministic benchmarks have become the primary tool for measuring large language model (LLM) progress, yet growing evidence suggests they measure memorization rather than genuine capability. Performance on canonical benchmarks such as MMLU and GSM8k degrades sharply under semantics-preserving perturbations, including answer reordering, surface rephrasing, and distractor addition, revealing brittle pattern matching rather than robust understanding. We argue this fragility is not an implementation flaw but a structural consequence of fixed evaluation sets in the era of web-scale training. We advocate for dynamic, synthetically generated benchmarks constructed fresh at evaluation time, making contamination impossible by construction and enabling principled, reproducible evaluation of genuine model capability.

1. Introduction

Benchmarks have become the de facto measure of progress in large language model (LLM) research (Hendrycks et al., 2021; Cobbe et al., 2021). Leaderboard rankings drive research directions, inform product decisions, and serve as proxies for real-world capability (Fourrier et al., 2023; 2024). This reliance rests on a critical assumption: that high performance on a static benchmark reflects the underlying skill it is designed to measure.

We argue this assumption is broken. As LLMs are trained on increasingly vast web-scale datasets, the boundary between training and evaluation data erodes (Mirzadeh et al., 2025; Wang et al., 2024a). Static benchmarks, which consist of fixed question sets with predetermined answers, become inadvertent training targets. Even without deliberate data contamination, models tuned to maximize performance on well-known benchmarks such as MMLU (Hendrycks et al.,

2021) and GSM8k (Cobbe et al., 2021) learn to pattern-match against memorized answer forms rather than develop genuine understanding (Zheng et al., 2024; Pezeshkpour & Hruschka, 2023). This paper presents evidence that this failure mode is already observable, argues that structural features of static benchmarks make it unavoidable, and advocates for a shift toward dynamic, synthetically generated evaluation as the only principled remedy.

2. The Fragility of Static Benchmarks

Multiple-Choice Benchmarks MMLU (Hendrycks et al., 2021) has become the canonical benchmark for measuring LLM knowledge and reasoning across 57 tasks, yet its scores are highly sensitive to evaluation artifacts unrelated to the knowledge being tested. Permuting the order of answer choices causes substantial accuracy drops, with some models losing over 10 percentage points (Pezeshkpour & Hruschka, 2023; Zheng et al., 2024), an effect confirmed at scale (Gupta et al., 2024). Rephrasing questions or introducing linguistic variation causes further degradation even when factual content is unchanged (Wang et al., 2024a). The multiple-choice format itself conflates test-taking heuristics with the target capability, making it a poor proxy for genuine understanding (Li et al., 2024b). A model with genuine knowledge should be invariant to distractor ordering; such sensitivity is the empirical signature of positional bias and memorization, not understanding.

Mathematical Reasoning Benchmarks GSM8k (Cobbe et al., 2021) is the standard benchmark for grade-school mathematical reasoning, widely used to evaluate chain-of-thought arithmetic capabilities (Wei et al., 2022). GSM-Symbolic (Mirzadeh et al., 2025) replaces names and numerical values in GSM8k problems while preserving the underlying mathematical structure. Performance drops significantly across all tested models, with variance increasing as problem complexity grows. A broader robustness evaluation through GSM-Plus confirms systematic degradation across diverse perturbation types, including numerical substitution, question rephrasing, and the addition of irrelevant clauses (Li et al., 2024a). The mathematical structure is identical across original and perturbed variants; only the surface form changes. A model reasoning from first principles should solve both equally well. The performance

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

gap reveals that models are recovering memorized solution templates keyed to surface features, not solving the abstract problem (Mirzadeh et al., 2025; Li et al., 2024a). Both benchmark types exhibit the same failure mode: high performance on the canonical evaluation form, brittleness to semantics-preserving variation, invariance to meaning and sensitivity to form. That this pattern appears independently across benchmark types, research groups, and model families (Pezeshkpour & Hruschka, 2023; Zheng et al., 2024; Gupta et al., 2024; Wang et al., 2024a; Mirzadeh et al., 2025; Li et al., 2024a) suggests it reflects a systematic property of static evaluation rather than isolated benchmark flaws.

3. Static Benchmarks Cannot Fix Themselves

The structural argument applies to any fixed evaluation set regardless of domain or difficulty; the same fragility has been documented in reading comprehension (Jia & Liang, 2017) and leaderboard saturation across benchmark types confirms it is pervasive (Fourrier et al., 2023; 2024) (see Appendix A). The natural response to these findings is to construct harder static benchmarks with stricter contamination controls. MMLU-Pro (Wang et al., 2024b) attempts exactly this by expanding the choice set to ten options and curating harder problems. But this approach treats the symptom rather than the cause.

The fundamental problem is structural: a static benchmark is a closed dataset. Once its contents enter training pipelines, even indirectly through web crawls, distillation, or synthetic data generation, no amount of curation restores its validity as an unbiased measure (Jacovi et al., 2023; Mirzadeh et al., 2025; Wang et al., 2024a). Benchmark updates simply reset the clock; the replacement eventually becomes contaminated in turn, a cycle already evident in the rapid saturation of successive MMLU variants (Fourrier et al., 2023; 2024).

Furthermore, fixed answer sets directly incentivize memorization during fine-tuning. Greedy decoding against a pre-determined correct answer is, mechanistically, next-token prediction on memorized text (Mirzadeh et al., 2025; Zheng et al., 2024). This is not a solvable problem within the static benchmark paradigm.

4. Dynamic Benchmarks as the Path Forward

We advocate for *dynamic benchmarks*: evaluation instances generated fresh at evaluation time from a parameterized generative process, rather than drawn from a fixed held-out set. Dynamic generation is contamination-resistant by construction, since it is not possible to memorize an instance that does not exist until the moment of evaluation (Jacovi et al., 2023; Castillo-Bolado et al., 2024). Two existing benchmarks demonstrate that this is feasible at scale. RULER (Hsieh et al., 2024) evaluates long-context capabil-

ities through synthetically generated needle-in-a-haystack and retrieval tasks, with ground truth derived algorithmically from the generated context. LongReason (Ling et al., 2025) constructs synthetic multi-step reasoning problems through context expansion. Both benchmarks produce meaningful, fine-grained capability signals, providing evidence that synthetic generation does not sacrifice evaluation quality (Hsieh et al., 2024; Ling et al., 2025).

Beyond contamination resistance, dynamic benchmarks offer structural properties that static benchmarks typically cannot provide:

- **Tunable difficulty.** Generative parameters directly control problem complexity, enabling controlled evaluation across the capability spectrum (Hsieh et al., 2024; Ling et al., 2025).
- **Verifiable ground truth.** Domains with algorithmic answers, symbolic reasoning, code execution, and formal logic, admit objective correctness without human annotation (Hsieh et al., 2024; Maheshwari et al., 2024).
- **Reproducibility via seeding.** Dynamic does not mean irreproducible. Versioned random seeds produce identical instance distributions for all evaluators, preserving comparability without exposing instances in advance (Hsieh et al., 2024; Castillo-Bolado et al., 2024).

The reproducibility objection to dynamic benchmarks is therefore an implementation concern, not a principled one: seeded generation gives every research group the same distribution while denying any group advance knowledge of the specific instances. Complementary evidence shows that dynamic and synthetic evaluation frameworks are both practically viable and empirically informative across diverse task settings (Castillo-Bolado et al., 2024; Aluffi et al., 2025; Maheshwari et al., 2024).

5. Implications and Call to Action

The community’s reliance on static benchmarks has produced a problematic feedback loop: models are optimized against known evaluation sets, reported scores rise, and the field interprets this as capability progress (Fourrier et al., 2023; 2024; Wang et al., 2024b). Breaking this loop requires changing what we measure, not merely how carefully we measure it. We call for three concrete shifts: benchmark creators should publish *generative procedures* alongside fixed datasets; evaluation frameworks should adopt versioned seeding as a reproducibility standard; and the community should prioritize domains with algorithmic ground truth for initial dynamic evaluation suites. Real-world deployment does not present models with questions drawn from a fixed academic dataset, and human preference rankings diverge substantially from static benchmark scores (Chiang et al., 2024). Evaluation should reflect that reality.

References

Aluffi, P. A., Zietkiewicz, P., Bazzi, M., Arderne, M., and Murevics, V. Dynamic benchmarking framework for llm-based conversational data capture, 2025. URL <https://arxiv.org/abs/2502.04349>.

Castillo-Bolado, D., Davidson, J., Gray, F., and Rosa, M. Beyond prompts: Dynamic conversational benchmarking of large language models, 2024. URL <https://arxiv.org/abs/2409.20222>.

Chiang, W.-L., Zheng, L., Sheng, Y., Angelopoulos, A. N., Li, T., Li, D., Zhang, H., Zhu, B., Jordan, M., Gonzalez, J. E., and Stoica, I. Chatbot arena: An open platform for evaluating LLMs by human preference, 2024. URL <https://arxiv.org/abs/2403.04132>.

Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., Hesse, C., and Schulman, J. Training verifiers to solve math word problems, 2021. URL <https://arxiv.org/abs/2110.14168>.

Fourrier, C., Habib, N., Lozovskaya, A., Szafer, K., and Wolf, T. OpenLLM Leaderboard V1. https://huggingface.co/docs/leaderboards/en/open_llm_leaderboard/archive, 2023.

Fourrier, C., Habib, N., Lozovskaya, A., Szafer, K., and Wolf, T. OpenLLM Leaderboard V2. https://huggingface.co/spaces/open-llm-leaderboard/open_llm_leaderboard, 2024.

Gupta, V., Pantoja, D., Ross, C., Williams, A., and Ung, M. Changing answer order can decrease mmlu accuracy, 2024. URL <https://arxiv.org/abs/2406.19470>.

Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., and Steinhardt, J. Measuring massive multitask language understanding. In *International Conference on Learning Representations*, 2021. URL <https://arxiv.org/abs/2009.03300>.

Hsieh, C.-P., Sun, S., Krizan, S., Acharya, S., Rekish, D., Jia, F., Zhang, Y., and Ginsburg, B. Ruler: What’s the real context size of your long-context language models?, 2024. URL <https://arxiv.org/abs/2404.06654>.

Jacovi, A., Caciularu, A., Goldman, O., and Goldberg, Y. Stop uploading test data in plain text: Practical strategies for mitigating data contamination by evaluation benchmarks, 2023. URL <https://arxiv.org/abs/2305.10160>.

Jia, R. and Liang, P. Adversarial examples for evaluating reading comprehension systems. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, 2017. URL <https://arxiv.org/abs/1707.07328>.

Li, Q., Cui, L., Zhao, X., Kong, L., and Bi, W. Gsm-plus: A comprehensive benchmark for evaluating the robustness of llms as mathematical problem solvers, 2024a. URL <https://arxiv.org/abs/2402.19255>.

Li, W., Li, L., Xiang, T., Liu, X., Deng, W., and Garcia, N. Can multiple-choice questions really be useful in detecting the abilities of llms?, 2024b. URL <https://arxiv.org/abs/2403.17752>.

Liang, P., Bommasani, R., Lee, T., Tsipras, D., Soylu, D., Yasunaga, M., Zhang, Y., Narayanan, D., Wu, Y., Kumar, A., Newman, B., Yuan, B., Yan, B., Zhang, C., Manning, C. D., Ré, C., Acosta-Navas, D., Hudson, D. A., Zelikman, E., Durmus, E., Ladhak, F., Rong, F., Ren, H., Yao, H., Wang, J., Santhanam, K., Orr, L., Zheng, L., Yuksekogonul, M., Suzgun, M., Kim, N., Guha, N., Chatterji, N., Khattab, O., Henderson, P., Huang, Q., Chi, R., Xie, S. M., Hashimoto, T., Wu, T., Zhang, T., Chaudhary, V., Wang, W., Li, X., Mai, Y., Zhang, Y., and Koreeda, Y. Holistic evaluation of language models, 2022. URL <https://arxiv.org/abs/2211.09110>.

Ling, Z., Liu, K., Yan, K., Yang, Y., Lin, W., Fan, T.-H., Shen, L., Du, Z., and Chen, J. Longreason: A synthetic long-context reasoning benchmark via context expansion, 2025. URL <https://arxiv.org/abs/2501.15089>.

Maheshwari, G., Ivanov, D., and Haddad, K. E. Efficacy of synthetic data as a benchmark, 2024. URL <https://arxiv.org/abs/2409.11968>.

Mirzadeh, I., Alizadeh, K., Shahrokhi, H., Tuzel, O., Bengio, S., and Farajtabar, M. Gsm-symbolic: Understanding the limitations of mathematical reasoning in large language models, 2025. URL <https://arxiv.org/abs/2410.05229>.

Pezeshkpour, P. and Hruschka, E. Large language models sensitivity to the order of options in multiple-choice questions, 2023. URL <https://arxiv.org/abs/2308.11483>.

Wang, W., Jain, S., Kantor, P., Feldman, J., Gallos, L., and Wang, H. Mmlu-sr: A benchmark for stress-testing reasoning capability of large language models, 2024a. URL <https://arxiv.org/abs/2406.15468>.

Wang, Y., Ma, X., Zhang, G., Ni, Y., Chandra, A., Guo, S., Ren, W., Arulraj, A., He, X., Jiang, Z., Li, T., Ku,

M., Wang, K., Zhuang, A., Fan, R., Yue, X., and Chen, W. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark, 2024b. URL <https://arxiv.org/abs/2406.01574>.

Wei, J., Wang, X., Schuurmans, D., Bosma, M., Ichter, B., Xia, F., Chi, E., Le, Q. V., and Zhou, D. Chain-of-thought prompting elicits reasoning in large language models. In *Advances in Neural Information Processing Systems*, 2022. URL <https://arxiv.org/abs/2201.11903>.

Zheng, C., Zhou, H., Meng, F., Zhou, J., and Huang, M. Large language models are not robust multiple choice selectors, 2024. URL <https://arxiv.org/abs/2309.03882>.

Zheng, L., Chiang, W.-L., Sheng, Y., Zhuang, S., Wu, Z., Zhuang, Y., Lin, Z., Li, Z., Li, D., Xing, E. P., Zhang, H., Gonzalez, J. E., and Stoica, I. Judging LLM-as-a-judge with MT-bench and chatbot arena. In *Advances in Neural Information Processing Systems*, 2023. URL <https://arxiv.org/abs/2306.05685>.

A. Summary of Benchmark Fragility Evidence

We focus on MMLU (Hendrycks et al., 2021) and GSM8k (Cobbe et al., 2021) because their ubiquity makes them the highest-value targets for contamination and the most extensively studied cases; together they also cover two distinct evaluation modalities, multiple-choice knowledge and free-form mathematical reasoning, ensuring the fragility pattern is not an artifact of a single format. Table 1 consolidates the perturbation findings from Section 2. The perturbations span three qualitatively distinct categories: *positional* (answer reordering), *linguistic* (question rephrasing), and *symbolic* (variable substitution). That substantial degradation occurs across all three categories, both benchmark types, and independent research groups rules out any explanation specific to a single perturbation strategy, benchmark design, or research team. The pattern is systematic.

B. Limitations of Dynamic Benchmarks

Dynamic benchmarks are not a universal solution, and intellectual honesty requires acknowledging their limitations.

Coverage of open-ended tasks. Automatic generation and verification work best for domains with algorithmic ground truth, as demonstrated by RULER (Hsieh et al., 2024) and LongReason (Ling et al., 2025). Tasks requiring subjective judgment, nuanced argumentation, or creative output are harder to generate and harder to evaluate without human involvement; work on LLM-as-judge evaluation shows that

automated scoring of open-ended outputs systematically diverges from human judgments (Zheng et al., 2023), making human-in-the-loop validation necessary in these settings.

Validity of synthetic instances. A generative procedure may produce instances that are syntactically valid but semantically degenerate, including ambiguous, trivially easy, or unrepresentative problems. The efficacy of synthetic data as a reliable benchmark signal is not guaranteed and requires careful empirical validation (Maheshwari et al., 2024; Aluffi et al., 2025); the absence of human curation eliminates contamination but introduces new failure modes around instance quality.

Comparability across generators. Seeded generation ensures reproducibility within a study, but comparing results across labs using different generators or generator versions remains non-trivial. Holistic evaluation frameworks have shown that evaluation decisions, including metric choice, prompt format, and dataset construction, produce large variance in reported scores across otherwise comparable systems (Liang et al., 2022); dynamic benchmarks face the same risk at the generator level.

C. Static vs. Dynamic Benchmarks: A Property Comparison

Table 2 organizes the structural properties discussed in Section 4 alongside the limitations from Appendix B into a single reference. Each row captures a dimension along which the two paradigms differ, making explicit both the advantages of dynamic evaluation and the trade-offs that must be managed in practice.

Reading the table as a whole reveals an important asymmetry. The rows where static benchmarks are structurally disadvantaged (contamination risk, long-term validity, difficulty control, and ground truth reliability) reflect constraints that cannot be resolved through improved dataset curation or harder problem selection; they are consequences of fixing the evaluation set, not of how carefully it was constructed. The rows where dynamic benchmarks are currently weaker (task coverage and evaluation cost) are engineering and scope challenges that are addressable over time as generative methods and automated verification mature (Castillo-Bolado et al., 2024; Maheshwari et al., 2024).

This asymmetry is the core of the position: the limitations of dynamic benchmarks are tractable, while those of static benchmarks are structural. The goal is not to abandon static benchmarks immediately, but to prioritize domains where dynamic evaluation is already viable and expand from there.

Table 1. Benchmark fragility under semantics-preserving perturbations. In each case the correct answer is unchanged; performance drops reveal sensitivity to surface form rather than underlying capability.

Benchmark	Perturbation	Finding	Source
MMLU	Answer choice reordering	>10 percentage point accuracy drop; some models show up to 13 pp variation across orderings	(Pezeshkpour & Hruschka, 2023; Zheng et al., 2024)
MMLU	Answer choice reordering (scale)	Effect confirmed across a broad range of model families and sizes	(Gupta et al., 2024)
MMLU	Linguistic rephrasing	Significant accuracy degradation despite identical factual content	(Wang et al., 2024a)
GSM8k	Symbolic substitution (names, values)	Consistent performance drop across all tested models; variance increases with problem complexity	(Mirzadeh et al., 2025)
GSM8k	Numerical, rephrasing, irrelevant clauses	Systematic degradation confirmed across all perturbation types	(Li et al., 2024a)

Table 2. Structural comparison of static and dynamic benchmark paradigms across evaluation-relevant properties.

Property	Static Benchmarks	Dynamic Benchmarks
Contamination risk	High; fixed instances can enter training corpora through web crawls, distillation, or synthetic augmentation (Jacovi et al., 2023)	None by construction; instances do not exist prior to evaluation time
Reproducibility	Exact; all evaluators use identical instances	Exact with versioned random seeds; results require seed and generator version to be published (Hsieh et al., 2024)
Difficulty control	Fixed at dataset creation; no post-hoc adjustment	Tunable via generative parameters; enables controlled evaluation across capability levels (Hsieh et al., 2024; Ling et al., 2025)
Ground truth	Human-curated; subject to annotation noise and subjectivity	Algorithmic; derived directly from generation parameters for suitable domains (Maheshwari et al., 2024)
Long-term validity	Degrades as models saturate the fixed instance set (Fourrier et al., 2023; 2024)	Maintained; fresh instances generated at each evaluation
Task coverage	Broad, including open-ended and subjective tasks	Currently limited to domains with verifiable ground truth; open-ended tasks require human-in-the-loop (Zheng et al., 2023)
Evaluation cost	Low; fixed corpus evaluated once and reused	Higher; generation overhead at each evaluation run
Cross-study comparability	High when using the same dataset version	Requires standardized generators and versioning; variance from generator choice is a known risk (Liang et al., 2022)