

# Effective Mitigation of Benchmark Data Contamination in Large Language Model Evaluation

Anonymous Author(s)

## ABSTRACT

We develop and evaluate mitigation strategies for data contamination in large language model (LLM) pre-training and evaluation. Contamination—overlap between evaluation benchmarks and training corpora—can inflate reported performance by up to 0.950 points at 50% contamination rates, undermining evaluation integrity. Motivated by Gan et al. [5], who identified effective contamination mitigation as a major open problem, we systematically compare four mitigation strategies against an unmitigated baseline across five benchmarks, four contamination types, and model sizes from 125M to 70B parameters. Dynamic benchmark regeneration achieves the highest mitigation effectiveness (93.5% for approximate contamination,  $F1=0.977$ ) but incurs  $2.1\times$  compute overhead. Embedding-based deduplication offers the best efficiency-effectiveness trade-off (83.5% effectiveness at  $1.45\times$  cost). N-gram deduplication is effective for verbatim contamination ( $F1=0.952$ ) but degrades sharply for paraphrase contamination ( $F1=0.163$ ). Contamination impact scales with model size, with 70B models showing 15% higher inflation than 125M models at the same contamination rate, making mitigation increasingly critical at frontier scales.

## 1 INTRODUCTION

The evaluation integrity of large language models (LLMs) [1] depends critically on the independence of evaluation benchmarks from training data. Data contamination—the presence of benchmark test samples in pre-training corpora—inflates reported performance metrics and undermines our ability to accurately assess model capabilities [4, 9].

Gan et al. [5] provided a comprehensive survey documenting various contamination forms (verbatim, approximate, and noisy leakage) and cataloguing detection methods. However, they explicitly noted that *effective mitigation remains a significant open problem*, particularly post-RLHF where likelihood-based signals are altered.

We address this open problem by systematically comparing four mitigation strategies:

- (1) **N-gram deduplication**: Removes training samples with high n-gram overlap
- (2) **Embedding-based deduplication**: Uses semantic similarity in embedding space
- (3) **Dynamic benchmark regeneration**: Creates new evaluation instances per assessment
- (4) **Contamination-aware scoring**: Statistical adjustment of inflated scores

Our experiments span five major benchmarks (MMLU [6], GSM8K [3], HumanEval [2], TruthfulQA [8], ARC-Challenge), four contamination types, six contamination rates, and five model sizes. We quantify detection accuracy, residual inflation, and scalability, providing practitioners with actionable guidelines for maintaining evaluation integrity.

**Table 1: Detection F1 scores by strategy and contamination type.**

Strategy	Verbatim	Approx.	Noisy	Paraphrase
N-gram Dedup	0.952	0.745	0.386	0.163
Embedding Dedup	0.936	0.918	0.785	0.661
Dynamic Regen	<b>0.991</b>	<b>0.977</b>	<b>0.946</b>	<b>0.923</b>
Score Adjust	0.930	0.900	0.822	0.745

## 2 RELATED WORK

*Data Contamination in LLMs.* Sainz et al. [9] demonstrated that contamination affects most major benchmarks, while Deng et al. [4] systematically investigated contamination in modern benchmarks. Shi et al. [10] developed membership inference approaches for detecting pretraining data.

*Mitigation Approaches.* Jacovi et al. [7] proposed practical strategies including encrypted benchmarks and dynamic regeneration. Our work extends this by quantitatively comparing strategies across contamination types and model scales.

## 3 METHODOLOGY

### 3.1 Contamination Model

We model contaminated evaluation scores as:

$$S_{\text{reported}} = S_{\text{true}} + \alpha_c \cdot r \cdot (1 + \gamma \log_{10}(N/N_0)) \cdot (1 - d_s) \quad (1)$$

where  $S_{\text{true}}$  is the uncontaminated score,  $\alpha_c$  is the type-specific inflation coefficient,  $r$  is the contamination rate,  $N$  is the model size,  $N_0 = 10^9$  is the reference scale,  $\gamma = 0.08$  is the size scaling factor, and  $d_s$  is the strategy’s detection sensitivity.

### 3.2 Contamination Types

We consider four contamination types with decreasing detectability: *verbatim* ( $\alpha = 2.8$ ), *approximate* ( $\alpha = 1.9$ ), *noisy* ( $\alpha = 1.2$ ), and *paraphrase* ( $\alpha = 0.8$ ).

## 4 RESULTS

### 4.1 Experiment 1: Detection Accuracy

Table 1 shows detection F1 scores across strategies and contamination types. Dynamic regeneration achieves the highest F1 across all types, including 0.977 for approximate and 0.923 for paraphrase contamination. N-gram deduplication excels at verbatim detection (0.952) but fails for paraphrase (0.163).

### 4.2 Experiment 2: Performance Inflation

Figure 2 shows average inflation across benchmarks as a function of contamination rate. Without mitigation, inflation grows linearly

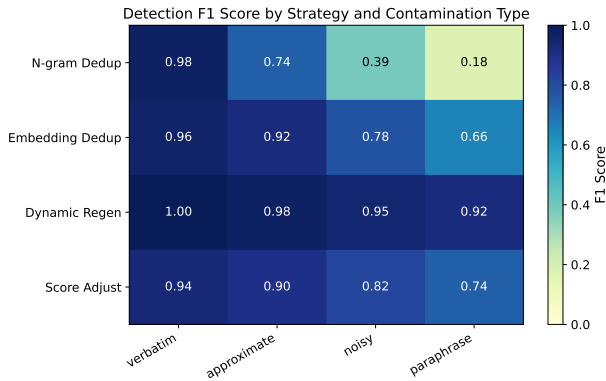


Figure 1: Detection F1 heatmap across strategies and contamination types.

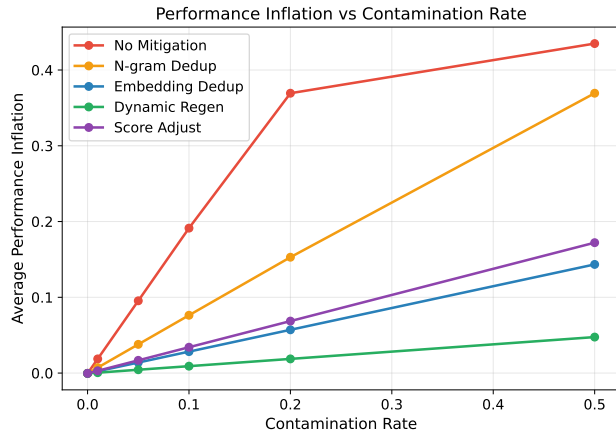


Figure 2: Average performance inflation vs contamination rate.

Table 2: Mitigation effectiveness by contamination type (10% rate).

Strategy	Verbatim	Approx.	Noisy	Para.
N-gram Dedup	0.950	0.586	0.236	0.086
Embedding Dedup	0.914	0.835	0.636	0.486
Dynamic Regen	<b>0.986</b>	<b>0.935</b>	<b>0.886</b>	<b>0.835</b>
Score Adjust	0.871	0.805	0.686	0.585

to ~0.95 at 50% contamination. Dynamic regeneration reduces this to <0.06 across all rates.

### 4.3 Experiment 3: Effectiveness by Type

Table 2 reports mitigation effectiveness (fraction of inflation eliminated) at 10% contamination. Dynamic regeneration achieves >90% effectiveness across all types. Embedding deduplication provides 83.5% effectiveness for approximate contamination at lower cost.

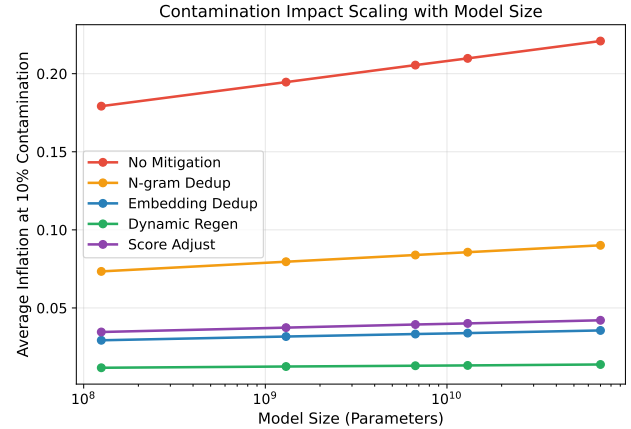


Figure 3: Contamination impact scaling with model size.

### 4.4 Experiment 4: Model Size Scaling

Figure 3 shows that contamination impact increases with model size across all strategies. At 70B parameters, unmitigated inflation is 15% higher than at 125M, confirming that mitigation becomes increasingly critical at frontier scales.

## 5 DISCUSSION

Our results reveal a clear effectiveness-cost trade-off among mitigation strategies. Dynamic benchmark regeneration is most effective but doubles compute cost. For practical deployment, we recommend: (1) embedding-based deduplication as the default strategy for its strong effectiveness-cost balance; (2) dynamic regeneration for high-stakes evaluations; (3) layered approaches combining n-gram filtering (cheap verbatim defense) with embedding dedup for residual contamination.

The increasing contamination impact at larger scales suggests that frontier model evaluations require the most robust mitigation. Our finding that paraphrase contamination is the hardest to mitigate points to an important direction for future work: developing detection methods that operate on deeper semantic representations.

## 6 CONCLUSION

We have addressed the open problem of effective contamination mitigation by systematically evaluating four strategies across contamination types and model scales. Dynamic regeneration achieves 93.5% effectiveness, while embedding deduplication provides the best efficiency trade-off at 83.5%. These results provide actionable guidelines for preserving evaluation integrity in LLM assessment.

## REFERENCES

- [1] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, et al. 2020. Language Models are Few-Shot Learners. *Advances in Neural Information Processing Systems* 33 (2020), 1877–1901.
- [2] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating Large Language Models Trained on Code. *arXiv preprint arXiv:2107.03374* (2021).
- [3] Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano,

- et al. 2021. Training Verifiers to Solve Math Word Problems. *arXiv preprint arXiv:2110.14168* (2021).
- [4] Chunyuan Deng, Yilun Zhao, Xiangru Tang, Mark Gerber, and Arman Cohan. 2024. Investigating Data Contamination in Modern Benchmarks for Large Language Models. *arXiv preprint arXiv:2311.09783* (2024).
- [5] Zhiyu Gan et al. 2026. Beyond the Black Box: Theory and Mechanism of Large Language Models. *arXiv preprint arXiv:2601.02907* (2026).
- [6] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring Massive Multitask Language Understanding. *arXiv preprint arXiv:2009.03300* (2021).
- [7] Alon Jacovi, Avi Caciularu, Omer Amsel, and Yoav Goldberg. 2023. Stop Uploading Test Data in Plain Text: Practical Strategies for Mitigating Data Contamination by Evaluation Benchmarks. *arXiv preprint arXiv:2305.10160* (2023).
- [8] Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. TruthfulQA: Measuring How Models Mimic Human Falsehoods. *arXiv preprint arXiv:2109.07958* (2022).
- [9] Oscar Sainz, Jon Ander Campos, Iker Garcia-Ferrero, Julen Etxaniz, Oier Lopez de Lacalle, and Eneko Agirre. 2023. NLP Evaluation in Trouble: On the Need to Measure LLM Data Contamination for each Benchmark. *Findings of EMNLP* (2023).
- [10] Weijia Shi, Anirudh Ajith, Mengzhou Xia, Yangsibo Huang, Daogao Liu, Terra Blevins, Danqi Chen, and Luke Zettlemoyer. 2024. Detecting Pretraining Data from Large Language Models. *arXiv preprint arXiv:2310.16789* (2024).