

---

# Benchmark Inflation: Revealing LLM Performance Gaps Using Retro-Holdouts

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Public benchmarks are compromised, as the training data for many Large Language  
2 Models (LLMs) is contaminated with test data, suggesting a *performance gap*  
3 between benchmark scores and actual capabilities. Ideally, a private holdout set  
4 could be used to accurately verify scores. Unfortunately, such datasets do not exist  
5 for most benchmarks, and post-hoc construction of sufficiently similar datasets is  
6 non-trivial. To address these issues, we introduce a systematic methodology for (i)  
7 retrospectively constructing a holdout dataset for a target dataset, (ii) demonstrating  
8 the statistical indistinguishability of this *retro-holdout* dataset, and (iii) comparing  
9 LLMs on the two datasets to quantify the performance gap due to the dataset’s  
10 public availability. Applying these methods to TruthfulQA, we construct and  
11 release Retro-TruthfulQA, on which we evaluate twenty LLMs and find that some  
12 have inflated scores by as much as 16 percentage points. Our results demonstrate  
13 that public benchmark scores do not always accurately assess model properties,  
14 and underscore the importance of improved data practices in the field.

## 15 1 Introduction

16 Concerns have emerged about the reliability of public benchmarks to accurately assess the perfor-  
17 mance of large language models [1, 56, 15]. First, there is a notable discrepancy between the reported  
18 performance of models on evaluation datasets and their actual capabilities in practical settings [33].  
19 Second, achieving high scores on these evaluations is strongly incentivized, as higher scores are  
20 closely linked to increased publicity and wider adoption of the given model [24]. This emphasis on  
21 benchmarks fosters a competitive environment where optimizing for benchmark performance can  
22 take precedence over real-world performance, potentially compromising the practical effectiveness or  
23 safety of models. This situation resembles specification gaming, where models meet the requirement  
24 of scoring well on benchmarks without genuinely improving on the capabilities that these benchmarks  
25 aim to assess [29]. Extending this framing, we define the mechanisms leading to a systematic gap  
26 between benchmark scores and real-world performance as *evaluation gaming*.

27 Recent research has shown that evaluation datasets have, in some cases, been included in the training  
28 data [43, 38, 47, 49, 26, 50], demonstrating that evaluation gaming is occurring. Such data leakage  
29 can destroy the predictive power of benchmarks, leading to large performance gaps between a model’s  
30 evaluation scores and its actual performance, as well as undermining trust in the reported model  
31 scores [39] – this highlights the need to improve practices for dataset release, and data collection.  
32 Such issues are particularly problematic given the significant role that evaluations are likely to play  
33 in the governance of machine learning technologies; stronger economic incentives will only increase  
34 the likelihood and severity of evaluation gaming. Furthermore, by misrepresenting model capabilities,  
35 current evaluations may create a false sense of safety. To accurately gauge the difference in a model’s

36 performance between the specific evaluation task and an analogous real-world task, we need access  
37 to a dataset originating from the same data distribution that has not been used during model training.

38 This is the idea of *holdout* datasets, which are used to assess a machine learning model’s performance  
39 after training. By definition, a holdout dataset comes from the same distribution as its corresponding  
40 target dataset, meaning that any evaluation conducted on both datasets should have the same result  
41 within some statistical tolerance [25]. Systematic differences in performance between holdout and  
42 target datasets can point to overfitting caused by data leakage. Comparing a model’s performance  
43 on a public benchmark and a corresponding holdout dataset could reveal whether data from the  
44 public benchmark has influenced the training process. Unfortunately, holdout datasets are typically  
45 not available; benchmark developers usually release all evaluation data, although there are notable  
46 exceptions, e.g. Li et al. [31].

47 To address these challenges, we propose *retroactive holdout*, or *retro-holdout*, datasets, which  
48 are verified to be similar to their corresponding target dataset through various tests, despite being  
49 created independently and retroactively. Utilizing a retro-holdout, we can quantify the evaluation  
50 performance gap of any given model. Our research advances the field by introducing a general  
51 and scalable methodology to create a retro-holdout dataset for a fully disclosed evaluation dataset,  
52 followed by rigorous testing to verify that the retro-holdout dataset closely mirrors the target dataset.

53 We detail our methodology for generating and validating retro-holdout datasets, along with recom-  
54 mendations and tools. We conduct a demonstrative case study using the TruthfulQA benchmark [34],  
55 a question answering dataset that was designed to assess the propensity of language models to mimic  
56 human falsehoods. TruthfulQA was selected for two key reasons: (i) it has become a popular dataset  
57 for developers to test against [32] and (ii) it has clear safety implications, as models performing  
58 poorly are likely to respond to user input with believable falsehoods.

## 59 1.1 Contributions

60 In this work, we:

- 62 • Develop a robust and novel process for the construction of retro-holdout datasets which are  
63 statistically indistinguishable from the target datasets.
- 64 • Introduce four tests for determining the similarity between two evaluation datasets, enabling  
65 identification of appropriate retro-holdout datasets for accurate model evaluations.
- 66 • Release Retro-TruthfulQA – a retro-holdout dataset for TruthfulQA, which can be used to  
67 quantify the performance gaps of a model on the original dataset.<sup>1</sup>
- 68 • Conduct a comprehensive evaluation of 20 models using Retro-TruthfulQA to demonstrate  
69 measurable score inflation.

## 70 2 Methods

71 Holdout datasets were first used in machine learning to accurately assess model performance. Unlike  
72 conventional holdout sets, retro-holdout datasets are not just randomly selected subsets; they are  
73 independently created post-hoc to match the statistical properties of the target dataset, thereby  
74 ensuring that they serve as effective and unbiased benchmarks for assessing real-world performance  
75 of the model post-training.

76 For brevity, we define

TARGET := an arbitrary, publicly available benchmark,  
RETRO := a retro-holdout dataset for TARGET.

77 We assume that the entries in TARGET were drawn a parent distribution which we denote as PARENT.  
78 We propose that, utilizing TARGET, along with information regarding its creation, a retro-holdout  
79 dataset, RETRO, which could have been drawn from PARENT but is distinct from TARGET can be  
80 created.

---

<sup>1</sup>Retro-TruthfulQA is only accurate on models with a training cutoff date prior to January 1st, 2024.

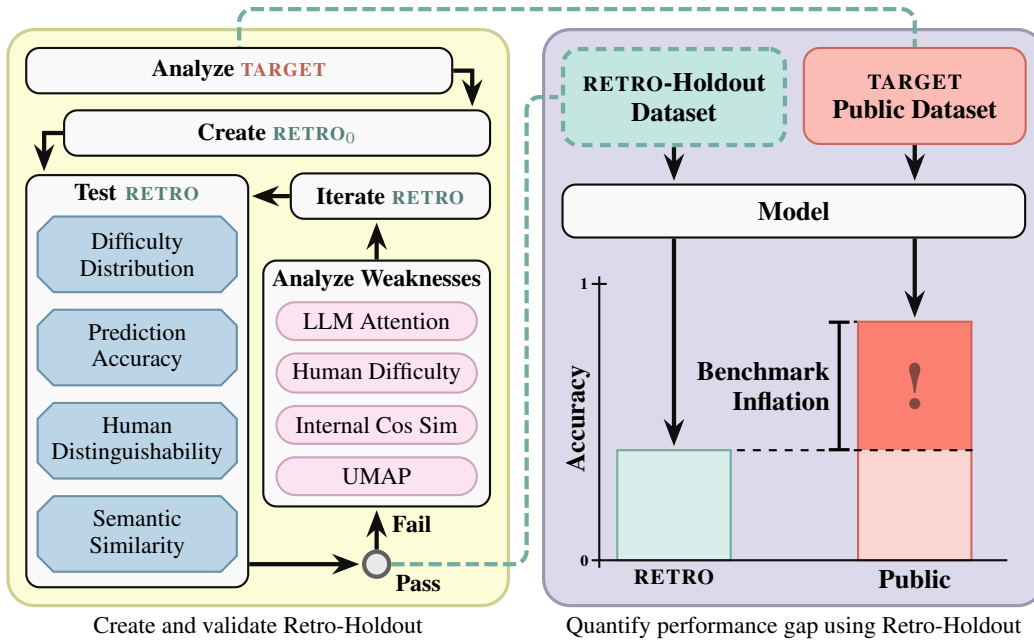


Figure 1: Visualization of our methodology.

## 81 2.1 Creating the RETRO

82 The methodology for crafting RETRO—while dependent on the specific TARGET—generally follows  
 83 two overarching phases: *Build Intuition* and *Entry Formulation*. Both of these phases are crucial for  
 84 understanding the nature of TARGET and generating entries that are representative of PARENT yet  
 85 distinct from TARGET.

86 **Build Intuition** To create a robust RETRO, one must have a strong understanding of the TARGET,  
 87 focusing primarily on its intended purpose and the methodology of its creation. We recommend  
 88 an initial thorough review of the dataset documentation and relevant literature, as well as looking  
 89 at many entries within TARGET. This phase, though straightforward, has proven to yield critically  
 90 valuable insights for the subsequent formulation, and later iteration, processes.

91 **Entry Formulation** Using the insights from the **Build Intuition** phase, the creation of entries  
 92 in RETRO proceeds by mirroring the structure and statistical properties of TARGET while ensuring  
 93 distinctiveness. Further details and step-by-step documentation for this process, as applied to the  
 94 TruthfulQA dataset, are provided in Appendix A. This appendix includes all materials and tools used  
 95 during the creation of the Retro-TruthfulQA dataset.

## 96 2.2 RETRO Tools

97 Creating a RETRO that meets our rigorous standards for sufficient indistinguishability (see §2.3) is  
 98 non-trivial and will typically only be achieved in an iterative manner. To aid in this process, we  
 99 devised a suite of tools that analyze and illustrate the various ways in which two datasets can be  
 100 distinct.

- 101 • **Fine-Tuned Prediction Model Attention:** A BERT model [10] is fine-tuned to classify  
 102 entries as belonging to either TARGET or RETRO. *Transformers Interpret*,<sup>2</sup> a library based  
 103 on Integrated Gradients for explaining model output attribution [52] is then leveraged to  
 104 identify which input tokens the model considered most relevant when differentiating between  
 105 TARGET and RETRO.

<sup>2</sup><https://pypi.org/project/transformers-interpret/>

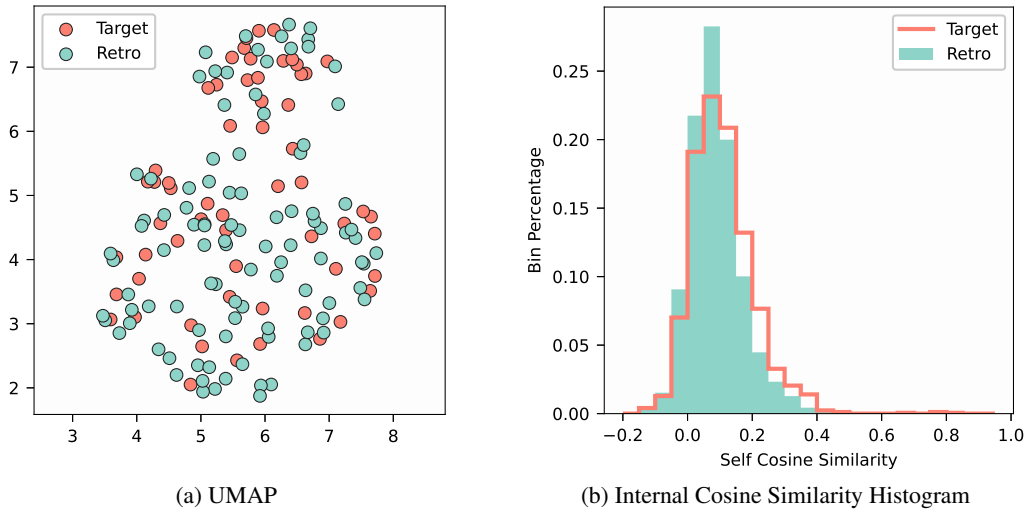


Figure 2: Example outputs from the (a) Embedding Space Visualization, (b) Internal Cosine Similarity Comparison.

- 106 • **Datapoint Embeddings:** We use the `all-mpnet-base-v2` embedding model through the  
 107 HuggingFace Sentence Transformers library to generate embedding vectors for all data  
 108 points. These embeddings are then taken as the basis for the following three tools; when  
 109 analyzed in conjunction they can provide meaningful insights on general similarity trends,  
 110 outlier detection, and topic clustering.
- 111 – **Embedding Space Visualization:** We employ Uniform Manifold Approximation and  
 112 Projection (UMAP) to project these embedding vectors onto a two-dimensional plane  
 113 [36]. The visualization provides an intuitive understanding of the dataset’s structure  
 114 and distribution. An example output of this visualization tool is provided in Figure 2.
- 115 – **Internal Cosine Similarity Distribution:** To assess similarity between entries within  
 116 the datasets we plot histograms of pairwise cosine similarities of datapoint embeddings.  
 117 This representation aids in identifying outliers and assessing overall similarity within  
 118 the datasets, as demonstrated in Figure 2.
- 119 – **Largest Internal Cosine Similarity Comparison:** We highlight the ten entry pairs  
 120 with the highest cosine similarities in both datasets, providing a direct comparison of  
 121 the most similar entries and their respective values.

122 These tools are documented in more detail in Appendix C.

### 123 2.3 Sufficient indistinguishability

124 Establishing absolute certainty that the two datasets have originated from the same distribution is  
 125 impossible. Therefore, we resort to multiple statistical tests designed to robustly test and reject the null  
 126 hypothesis that TARGET and RETRO have a common origin. If the result of each test indicates that we  
 127 cannot reject our null hypothesis, we designate our RETRO to be statistically indistinguishable from  
 128 TARGET. The core motivation behind this is that, if our RETRO could have indeed been drawn from  
 129 (PARENT – TARGET), then it should be challenging for our statistical tests to distinguish between  
 130 TARGET and RETRO. While it is theoretically possible to construct an infinite array of tests to evaluate  
 131 the similarity between the two datasets, practical considerations guide us to focus on four key tests  
 132 that provide a thorough assessment:

- 133 • **Similarity of Difficulty:** Are the questions in both datasets comparably challenging?
- 134 • **Semantic Embedding Similarity:** What is the likelihood that a distribution of cosine  
 135 similarities between sentence embeddings similar to that of RETRO have been pulled from  
 136 PARENT?

- 137 • **Prediction Accuracy:** Can a model, fine-tuned on randomized splits of the datasets, differ-  
138 entiate between TARGET and RETRO?
- 139 • **Human Distinguishability:** Can humans identify a RETRO sample hidden in two TAR-  
140 GET samples?

141 We assert that the two datasets are *statistically indistinguishable* if they pass all four tests.

142 **Similarity of Difficulty** Assessing whether the retro-holdout dataset, RETRO, matches the difficulty  
143 of the target dataset, TARGET, is crucial for drawing meaningful conclusions about evaluation gaming;  
144 otherwise performance differences could be attributed to the varying levels of difficulty, rather than  
145 the models’ true capabilities. To understand this potential disagreement between datasets, we consider  
146 models with a training cutoff date prior to the release of the TARGET, or *pre-release* models. Since  
147 pre-release models could not possibly have been effected by exposure to TARGET, their performance  
148 on both TARGET and RETRO should be comparable, with a margin of statistical uncertainty.

149 It is essential to note that with access to a diverse array of LLMs spanning various capability levels,  
150 our testing methodology, combined with simple human assessment, would likely suffice to ascertain  
151 whether two evaluation datasets are statistically indistinguishable. However, performance of cutting-  
152 edge models continues to improve, meaning that pre-release models almost certainly won’t be stronger  
153 than the most advanced models, assuming they are accessible at all. The nature and implications  
154 of this constraint are discussed further in §3, and Appendix D. To address this limitation, we use a  
155 number of techniques to amplify model performance. These include allowing the model to choose  
156 multiple answers (top- $k$ ), including examples of other questions within the dataset (5-shot), including  
157 a routine prompt which aims to elicit intermediary outputs from the model (chain-of-thought), and  
158 using the ‘helpful’ prompt from Lin et al. [34].

159 For TARGET and RETRO to be statistically indistinguishable, pre-release models (with and without  
160 performance-amplifying techniques) should score similarly on both datasets. Complete specifications  
161 and the rationale for the difficulty test are provided in Appendix D.

162 **Prediction Accuracy** We adopt a modification of prediction accuracy as detailed by Dankar and  
163 Ibrahim [8] to train a model to classify an entry as either belonging to TARGET, or to RETRO, using an  
164 equivalent number of entries from each dataset. Contrary to the conventional use of logistic regression  
165 in synthetic data evaluations [8], we fine-tune BERT [10] on the prediction task. This choice is  
166 predicated on BERT’s capabilities in capturing nuanced semantic relationships within text, which are  
167 crucial for accurately assessing the subtle distinctions or similarities between dataset entries.

168 We test this model on the remainder of the samples, as theoretically, if the model’s prediction accuracy  
169 on the test samples converges to 50%, within a margin allowing for statistical fluctuation, it suggests  
170 that the model fails to distinguish between the two datasets. This condition is rigorously tested to  
171 ensure the model is not merely performing at chance level but is genuinely indicative of dataset  
172 equivalence.

173 **Semantic Embedding Similarity** Using well established techniques for multi-dimensional data  
174 analysis, we conduct a random permutation test to determine the likelihood that a distribution with  
175 similar properties to RETRO could be randomly drawn from PARENT [14, 37, 22]. For the test  
176 statistic used in our random permutation test we compute the mean of all unique, non-trivial cosine  
177 similarities between embeddings from PARENT and a randomly sampled subset of PARENT with the  
178  $n = \min(n_{\text{TARGET}}, n_{\text{RETRO}})$  entries. The test statistics of both TARGET and RETRO are then compared  
179 with the test statistics for our  $N$  random samples, yielding one  $p$ -value for TARGET, and one for  
180 RETRO. To successfully pass this test,

$$p\text{-value}_{\text{TARGET}}, p\text{-value}_{\text{RETRO}} \in [0.05, 0.95].$$

181 This range is chosen to ensure that RETRO is neither too similar nor too dissimilar from TARGET,  
182 promoting a balance that supports our hypothesis of indistinguishability under realistic conditions. It  
183 is worth noting that, unless  $n_{\text{TARGET}} = n_{\text{RETRO}}$ , an external loop outside of the core permutation test  
184 must also be defined in order to understand variance of our test statistic. Detailed visualizations and  
185 explanations of these tests are documented in Appendix F.

186 **Human Indistinguishability** To assess whether the datasets were distinguishable to humans, we  
187 conducted a survey where participants were tasked to separate entries from TARGET and RETRO.

188 Initially, participants were oriented with ten labeled entries from each dataset to provide them with  
189 contextual understanding. They then undergo a series of ten tests, each comprising of three dataset  
190 entries - two from the TARGET and one from RETRO. All entries are drawn without replacement to  
191 ensure unique samples throughout the survey.

192 Additionally, we implement a variation of this test using GPT-4o as the evaluator to compare human  
193 and model performance. See Appendix E for comprehensive details on the survey methodology,  
194 including specifics on participant recruitment, the structure of the test, and survey instructions.

### 195 2.3.1 Iterating on Failures

196 Although the iterative tools described in §2.2 will limit significant differences between the datasets,  
197 our stringent standard for required similarity render it improbable that the initial RETRO tested  
198 will be statistically indistinguishable. Acknowledging this, and considering the time-intensive  
199 nature of dataset generation, efficiency is all the more important. To this end, we recommend  
200 that an initial small-scale application of our process be conducted, allowing for developers to use  
201 our indistinguishability tests to gain insights about their TARGET. This preliminary phase allows  
202 developers to refine their methods and heuristics before re-conducting the process to create a more  
203 extensive retro-holdout dataset.

204 This process was used for the construction of Retro-TruthfulQA. As anticipated, the first iteration did  
205 not meet our exacting standards of calibration. However, by working with the various tests on our  
206 smaller dataset, we identified several failure modes that were not initially apparent. These instances  
207 of failure, and the corresponding adjustments made, provided critical learning opportunities that  
208 guided the subsequent refinements.

## 209 2.4 Evaluating Models

210 The evaluation framework described in Section 2.3 was applied to assess the performance of current  
211 models. Experiments were conducted using the OpenAI chat completion API and various models  
212 from Huggingface with mostly default settings. The generation length was adjusted, and a temperature  
213 of 0.5 was specified, although this parameter may not apply to OpenAI chat models.

214 During the construction of TruthfulQA [34], the authors envisioned that language models would  
215 be evaluated by the max-probability assigned to any of a predefined list of available options. This  
216 approach may suffer from three issues. First, this may penalize long answer options which naturally  
217 have lower total probability. Second, such an answer may not well reflect which of a fixed number  
218 of options is the most likely to be generated, seeing how this may be more determined by the first  
219 tokens of the option. Finally, the OpenAI API no longer provides probability output, and other API  
220 providers may have never had such an option.

221 For these reasons, it was decided to evaluate models by providing an enumerated list of all TruthfulQA  
222 *mc1*-choices and generating tokens to select a preferred option. To minimize potential model bias,  
223 answers were resampled with options rotated at minimum ten times and until one option had been  
224 selected an additional four times over alternatives. A Vicuna-inspired prompt was used for all models  
225 and is described in Appendix G.1.2.

226 Especially when working with pre-release models, it can be difficult to guarantee model outputs  
227 conform to specific formats, such as multiple choice responses. For this reason, substantial efforts  
228 were made to reduce fluctuations reported evaluation results. Due to prohibitive costs for many  
229 resamples, we were only able to calculate empirical one sigma error bars for the pre-release models  
230 on both TruthfulQA and Retro-TruthfulQA. On TruthfulQA, babbage-002, davinci-002, and neox-  
231 20b had had statistical error of  $\pm 1.27\%$ ,  $\pm 0.83\%$ , and  $\pm 2.84\%$  respectively, while their errors on  
232 Retro-TruthfulQA were  $\pm 2.47\%$ ,  $\pm 1.96\%$ , and  $\pm 1.34\%$ .

## 233 3 Results and Discussion

### 234 3.1 Retro-holdout TruthfulQA Dataset

235 We release Retro-TruthfulQA, a retro-holdout dataset designed to quantify the evaluation gap for  
236 models tested on the TruthfulQA dataset, *provided that the model’s training cutoff date is prior to*

Table 1: Retro-TruthfulQA Indistinguishability Tests Results

Description	H <sub>0</sub>	Outcome	Test <i>p</i> -value
babbage-002 difficulty gap	0%	-1.2 ± 7.4%	≥ 50%
davinci-002 difficulty gap	0%	-3.3 ± 8.0%	≥ 50%
Prediction accuracy	50%	53.7 ± 3.26%	47.4%
TARGET Random permutation	-	-	6.67 ± 1.86%
RETRO Random permutation	-	-	93.48 ± 1.85%
GPT-4 Distinguishability	33.3%	28.0 ± 9.0%	≥ 50%
Human Distinguishability	33.3%	31.3 ± 7.1%	≥ 50%

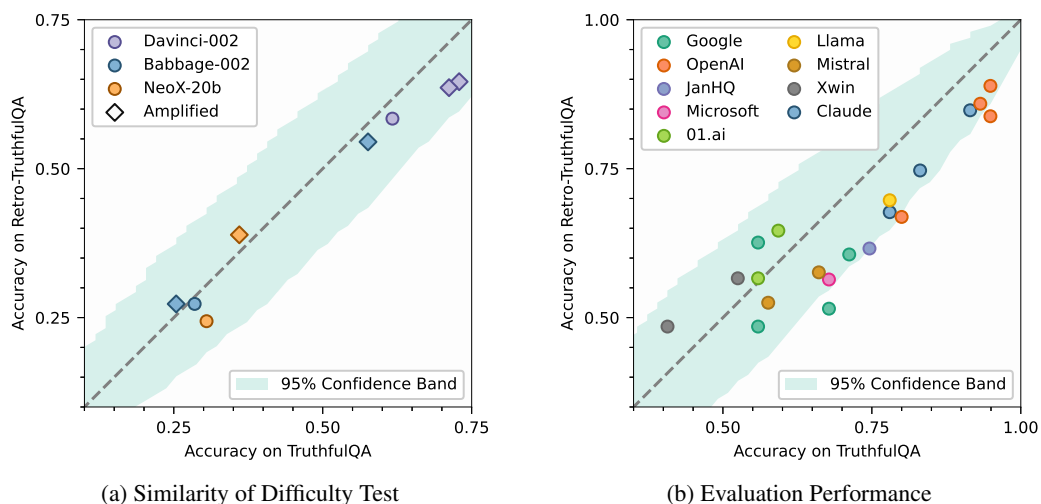


Figure 3: Model accuracy on Retro-TruthfulQA vs. TruthfulQA. (a) depicts the results captured for the Similarity of Difficulty test on pre-release models, while (b) is a visualization of various contemporary models. In both plots, a 95% confidence band for two samples of boolean values, i.e. correct or incorrect, of sizes equal to our two datasets is shown.

237 *January 1st, 2024*. Retro-TruthfulQA mirrors the structure and content of the original TruthfulQA  
 238 dataset across all measured categories and comprises 817 entries.

239 Notably, Retro-TruthfulQA has passed all four of our indistinguishability tests, establishing it as the  
 240 first retro-holdout dataset to be *statistically indistinguishable* from its corresponding target dataset.  
 241 The tests covered aspects of the dataset to ensure semantic similarity, prediction accuracy, and human  
 242 and model-based distinguishability, confirming that Retro-TruthfulQA accurately mirrors the original  
 243 dataset in all essential aspects. The detailed results, complete with confidence intervals for each  
 244 metric, are summarized in Table 1, and Figure 3(a).

### 245 3.2 TruthfulQA Evaluation Details

246 The TruthfulQA dataset contains two categorizations for entries: Category and Type. Our experiments  
 247 have focused on the largest of these categories – Misconceptions. The Type for the dataset is either  
 248 *adversarial* or *non-adversarial*. Our evaluation finds that GPT-3 models like babbage-002 and  
 249 davinci-002 do significantly better on the non-adversarial portion.

250 This is unsurprising as the adversarial set was constructed by testing various entries on a version of  
 251 GPT-3 and discarding those the model answered correctly. These entries were then used as inspiration  
 252 to create the remaining portion, but where no such model filtering was done. Due to this potential  
 253 filtering bias and the performance difference between the two sets, we have additionally chosen to focus  
 254 on the non-adversarial portion of TruthfulQA. While these changes are deviations from the original  
 255 TruthfulQA evaluation, it is worth noting that all experiments compare the performance of this same

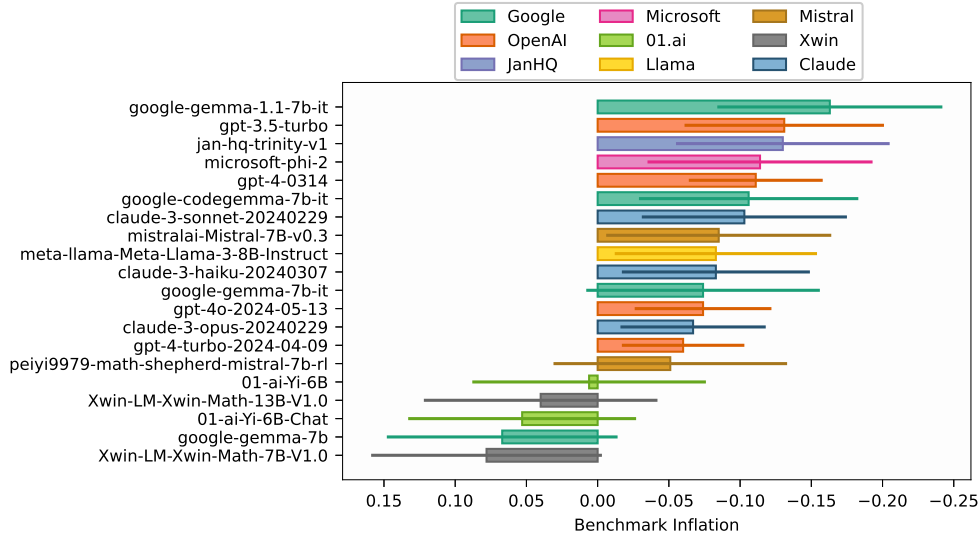


Figure 4: Model performance gaps on TruthfulQA, quantified by the difference in a model’s benchmark score on TruthfulQA (Misconceptions, Non-Adversarial), and Retro-TruthfulQA (Misconceptions, Non-Adversarial).

256 evaluation method on the original vs the retro-holdout dataset, along with calibration such that any  
 257 statistically-significant gap between these must be explained by some form of evaluation gaming.

### 258 3.3 The Performance Gap

259 With our newly created retro-holdout dataset, we explicitly quantify the performance gap of 20  
 260 models, which can be seen in Figure 4. Our analysis covers both larger API models such as Claude3  
 261 and GPT-4, as well as several open-release models that have been either speculated or confirmed to  
 262 exhibit data leakage [43].

### 263 3.4 Contemporaneous Work

264 Coinciding with our efforts, Zhang et al. [55] introduce the GSM1k dataset for assessing mathematical  
 265 reasoning. This study employs several human tests to ensure an "apples-to-apples" similarity to their  
 266 target dataset GSM8k [55, 7]. Similar to our findings, Zhang et al. [55] report an overperformance by  
 267 many models on their target evaluations.

268 While the GSM1k dataset comprises over 1000 entries, only 50 have been publicly released to date.  
 269 Zhang et al. [55] recognize that releasing the entire dataset will likely result in the same data leakage  
 270 current benchmark suffer from. They have decided to postpone the full release of GSM1k until either  
 271 (i) the top open source models score over 95% on the benchmark, or (ii) the end of 2025.

272 Given the similarity between our works, we thought it would be a good opportunity to put our  
 273 concept of sufficient indistinguishability to the test. We took the 50 published questions from their  
 274 dataset, henceforth referred to as GSM1k50, and examined them using the same methods as we  
 275 did for Retro-TruthfulQA. Our semantics tools and Semantic Embedding Similarity test suggest  
 276 that GSM1k50 can be adjusted to more closely resemble original GSM8k entries, generating a  
 277 TARGET and RETRO random permutation of  $3.02 \pm 0.05\%$  and  $98.7 \pm 0.02\%$ , respectively. The  
 278 Prediction Accuracy test reveals that GSM1k50 can be differentiated from the original GSM8k, albeit  
 279 to a small, but statistically significant extent. These finding highlights the rigor of our notion of  
 280 sufficient indistinguishability, but also suggests that in practical scenarios, slightly relaxed criteria  
 281 might still produce effective retro-holdout datasets without significantly compromising evaluation  
 282 quality.

283 Despite the independent development and differing methodologies of our projects, both underscore  
 284 the crucial role of comprehensive dataset validation in enhancing the accuracy of model evaluations.



### 285 3.5 Limitations

286 The assumption that the retro-holdout dataset and the target dataset are drawn from the same  
287 distribution may not always be valid. This assumption is challenged if the target dataset itself is  
288 subject to distribution shifts over time; such shifts can alter the underlying data characteristics over  
289 time. Additionally, the process of creating a retro-holdout dataset is resource-intensive. It demands  
290 significant computational resources for generating and validating the dataset, as well as human experts  
291 for iterative adjustments based on indistinguishability tests, which may mitigate the wide adoption of  
292 our methodology.

293 Another limitation arises from the inherent approach of matching the distribution of the target  
294 dataset. While this method ensures that the retro-holdout dataset mirrors the target dataset as  
295 closely as possible, it also inadvertently perpetuates any implicit biases that are present in the target  
296 dataset. Consequently, while the retro-holdout dataset might excel in mimicking the target dataset's  
297 distribution, it may not provide a truly independent measure of a model's generalization capabilities  
298 across broader contexts.

## 299 4 Related Works

300 Development of large language models (LLMs) continues to outpace the advancement of evaluation  
301 methods, raising concern about benchmark integrity [6]. Evaluation datasets are frequently used  
302 during an LLM's training process, causing inflated benchmark scores; no standard methodology  
303 exists to detect this issue [1]. Data quality, essential for model performance, remains undervalued  
304 and under-incentivized [46]. Data contamination, where test data is included in training sets, results  
305 in models "cheating" by memorizing tests rather than generalizing [35]. High benchmark scores are  
306 heavily incentivized, promoting practices that compromise data quality and evaluation integrity.

307 Recent work has introduced heuristics for third-party contamination tests. Sainz et al. [45] propose  
308 a technique to detect test set contamination by eliciting reproduction of specific test set examples.  
309 Golchin and Surdeanu [18] suggest a method for identifying contamination in black-box models by  
310 comparing the similarity between model completions of randomly selected example prefixes and the  
311 actual data using GPT-4. Concurrent work by [55] is notable for its use of a holdout set, a concept  
312 central to our approach, and shows accuracy drops of up to 13% and highlights a positive correlation  
313 between memorization and performance gaps.

314 It is well known that metrics lose their predictive power when incentives are attached to them  
315 Goodhart [19], Strathern [51], Karwowski et al. [27]. As [53] state, "overemphasizing metrics  
316 leads to manipulation, gaming, a myopic focus on short-term goals, and other unexpected negative  
317 consequences." Current AI risk metrics fail to address emerging failure modes [28], and Bengio [4]  
318 emphasize that high benchmark scores do not necessarily equate to effective real world performance.

319 Empirical findings highlight the necessity for immediate structural reforms in AI research and  
320 development to prioritize and encourage data quality [46]. Recent calls for a *science of evaluations*  
321 underscore the urgent need for rigorous evaluation frameworks to inform policy and ensure responsible  
322 AI development [5, 42].

## 323 5 Conclusion

324 Our findings demonstrate significant discrepancies between benchmark performances and real-world  
325 capabilities of LLMs, underscoring the need for robust and reliable evaluation methodologies. We  
326 introduce a novel, systematic methodology for constructing retro-holdout datasets, and conduct a  
327 case study of the process using the largest category of TruthfulQA. The result is Retro-TruthfulQA,  
328 a retro-holdout for TruthfulQA which has been shown to be statistically indistinguishable from  
329 the target dataset. This methodology, designed to be generally applicable across various public  
330 benchmark evaluations, provides tools that significantly enhance the accuracy and reliability of  
331 model evaluations, offering a practical path forward for the field. In a recent work Anwar et al. [2]  
332 explicitly challenge "How can the evaluations of LLMs be made trustworthy given the difficulty of  
333 assuring that there is no test-set contamination?" Our work provides a succinct and powerful response:  
334 Retro-Holdouts.

## References

- 335
- 336 [1] N. Alzahrani, H. A. Alyahya, Y. Alnumay, S. Alrashed, S. Alsubaie, Y. Almushaykeh, F. Mirza, N. Alotaibi,  
337 N. Altwaresh, A. Alowisheq, M. S. Bari, and H. Khan. When Benchmarks are Targets: Revealing the  
338 Sensitivity of Large Language Model Leaderboards, Feb. 2024. URL [http://arxiv.org/abs/2402.  
339 01781](http://arxiv.org/abs/2402.01781). arXiv:2402.01781 [cs].
- 340 [2] U. Anwar, A. Saparov, J. Rando, D. Paleka, M. Turpin, P. Hase, E. S. Lubana, E. Jenner, S. Casper,  
341 O. Sourbut, B. L. Edelman, Z. Zhang, M. Günther, A. Korinek, J. Hernandez-Orallo, L. Hammond,  
342 E. Bigelow, A. Pan, L. Langosco, T. Korbak, H. Zhang, R. Zhong, S. O. hÉigeartaigh, G. Recchia, G. Corsi,  
343 A. Chan, M. Anderljung, L. Edwards, Y. Bengio, D. Chen, S. Albanie, T. Maharaj, J. Foerster, F. Tramèr,  
344 H. He, A. Kasirzadeh, Y. Choi, and D. Krueger. Foundational Challenges in Assuring Alignment and Safety  
345 of Large Language Models, Apr. 2024. URL <http://arxiv.org/abs/2404.09932>. arXiv:2404.09932  
346 [cs].
- 347 [3] S. Arora, Y. Liang, and T. Ma. A Simple but Tough-to-Beat Baseline for Sentence Embeddings. Nov. 2016.  
348 URL <https://openreview.net/forum?id=SyK00v5xx>.
- 349 [4] Y. Bengio. International scientific report on the safety of advanced ai: interim report. *Gov.uk Department  
350 for Science, Innovation and Technology and AI Safety Institute*, 2024.
- 351 [5] R. Bommasani, K. Klyman, S. Longpre, S. Kapoor, N. Maslej, B. Xiong, D. Zhang, and P. Liang. The  
352 foundation model transparency index. *arXiv preprint arXiv:2310.12941*, 2023.
- 353 [6] Y. Chang, X. Wang, J. Wang, Y. Wu, L. Yang, K. Zhu, H. Chen, X. Yi, C. Wang, Y. Wang, et al. A survey  
354 on evaluation of large language models. *ACM Transactions on Intelligent Systems and Technology*, 15(3):  
355 1–45, 2024.
- 356 [7] K. Cobbe, V. Kosaraju, M. Bavarian, M. Chen, H. Jun, L. Kaiser, M. Plappert, J. Tworek, J. Hilton,  
357 R. Nakano, C. Hesse, and J. Schulman. Training Verifiers to Solve Math Word Problems, Nov. 2021. URL  
358 <http://arxiv.org/abs/2110.14168>. arXiv:2110.14168 [cs].
- 359 [8] F. K. Dankar and M. Ibrahim. Fake It Till You Make It: Guidelines for Effective Synthetic Data Generation.  
360 *Applied Sciences*, 11(5):2158, Jan. 2021. ISSN 2076-3417. doi: 10.3390/app11052158. URL [https:  
361 //www.mdpi.com/2076-3417/11/5/2158](https://www.mdpi.com/2076-3417/11/5/2158). Number: 5 Publisher: Multidisciplinary Digital Publishing  
362 Institute.
- 363 [9] C. Deng, Y. Zhao, X. Tang, M. Gerstein, and A. Cohan. Investigating Data Contamination in Modern  
364 Benchmarks for Large Language Models, Apr. 2024. URL <http://arxiv.org/abs/2311.09783>.  
365 arXiv:2311.09783 [cs].
- 366 [10] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova. BERT: Pre-training of Deep Bidirectional  
367 Transformers for Language Understanding, May 2019. URL <http://arxiv.org/abs/1810.04805>.  
368 arXiv:1810.04805 [cs].
- 369 [11] Y. Dong, X. Jiang, H. Liu, Z. Jin, and G. Li. Generalization or Memorization: Data Contamination and  
370 Trustworthy Evaluation for Large Language Models, Feb. 2024. URL [http://arxiv.org/abs/2402.  
371 15938](http://arxiv.org/abs/2402.15938). arXiv:2402.15938 [cs].
- 372 [12] J. Egan. *The Mega Misconception Book*. Lulu.com, 2016. ISBN 9781326838423. URL [https:  
373 //books.google.co.in/books?id=Aq96DQAAQBAJ](https://books.google.co.in/books?id=Aq96DQAAQBAJ).
- 374 [13] S. Engmann and D. Cousineau. Comparing Distributions: The Two-Sample Anderson-Darling Test as an  
375 Alternative to the Kolmogorov-Smirnoff Test. *Journal of Applied Quantitative Methods*, 6(3), 2011. URL  
376 [https://www.jaqm.ro/issues/volume-6,issue-3/pdfs/1\\_engmann\\_cousineau.pdf](https://www.jaqm.ro/issues/volume-6,issue-3/pdfs/1_engmann_cousineau.pdf).
- 377 [14] R. A. Fisher. *The Design of Experiments*. Hafner Press, 9th edition, 1974. URL [https://home.iitk.  
378 ac.in/~shalab/anova/DOE-RAF.pdf](https://home.iitk.ac.in/~shalab/anova/DOE-RAF.pdf).
- 379 [15] C. Fourrier, N. Habib, J. Launay, and T. Wolf. What’s going on with the Open LLM Leaderboard?, June  
380 2023. URL <https://huggingface.co/blog/evaluating-mmlu-leaderboard>.
- 381 [16] D. Ganguli, N. Schiefer, M. Favaro, and J. Clark. Challenges in evaluating ai systems, 2023. URL  
382 <https://www.anthropic.com/index/evaluating-ai-systems>.
- 383 [17] L. Gao, J. Tow, S. Biderman, S. Black, A. DiPofi, C. Foster, L. Golding, J. Hsu, K. McDonell, N. Muen-  
384 nighoff, J. Phang, L. Reynolds, E. Tang, A. Thite, B. Wang, K. Wang, and A. Zou. A framework for  
385 few-shot language model evaluation, Sept. 2021. URL <https://doi.org/10.5281/zenodo.5371628>.

- 386 [18] S. Golchin and M. Surdeanu. Time travel in llms: Tracing data contamination in large language models.  
387 *arXiv preprint arXiv:2308.08493*, 2023.
- 388 [19] C. A. Goodhart. *Problems of monetary management: the UK experience*. Springer, 1984.
- 389 [20] J. Green. *Contrary to popular belief: more than 250 false facts revealed*. Crown, 2005.
- 390 [21] J. Green. *Contrary to Popular Belief: More Than 250 False Facts Revealed*. Broadway Books, 2005.  
391 ISBN 9780767919920. URL <https://books.google.co.in/books?id=4tyJQL015mwC>.
- 392 [22] J. Hemerik. On the Term “Randomization Test”. *The American Statistician*, pages 1–8, Mar. 2024. ISSN  
393 0003-1305, 1537-2731. doi: 10.1080/00031305.2024.2319182. URL <https://www.tandfonline.com/doi/full/10.1080/00031305.2024.2319182>.
- 395 [23] R. Heyburn, R. R. Bond, M. Black, M. Mulvenna, J. Wallace, D. Rankin, and B. Cleland. Machine learning  
396 using synthetic and real data: Similarity of evaluation metrics for different healthcare datasets and for  
397 different algorithms. In *Data Science and Knowledge Engineering for Sensing Decision Support*, pages  
398 1281–1291, Belfast, Northern Ireland, UK, Sept. 2018. WORLD SCIENTIFIC. ISBN 978-981-327-322-1  
399 978-981-327-323-8. doi: 10.1142/9789813273238\_0160. URL [https://www.worldscientific.com/doi/abs/10.1142/9789813273238\\_0160](https://www.worldscientific.com/doi/abs/10.1142/9789813273238_0160).
- 401 [24] HuggingFaceH4. Open LLM Leaderboard - a Hugging Face Space. URL [https://huggingface.co/spaces/HuggingFaceH4/open\\_llm\\_leaderboard](https://huggingface.co/spaces/HuggingFaceH4/open_llm_leaderboard).
- 403 [25] G. James, D. Witten, T. Hastie, R. Tibshirani, and J. Taylor. *An Introduction to Statistical Learning:  
404 with Applications in Python*. Springer Texts in Statistics. Springer International Publishing, Cham,  
405 2023. ISBN 978-3-031-38746-3 978-3-031-38747-0. doi: 10.1007/978-3-031-38747-0. URL <https://link.springer.com/10.1007/978-3-031-38747-0>.
- 407 [26] M. Jiang, K. Z. Liu, M. Zhong, R. Schaeffer, S. Ouyang, J. Han, and S. Koyejo. Investigating Data  
408 Contamination for Pre-training Language Models, Jan. 2024. URL <http://arxiv.org/abs/2401.06059>.  
409 arXiv:2401.06059 [cs].
- 410 [27] J. Karwowski, O. Hayman, X. Bai, K. Kiendlhofer, C. Griffin, and J. Skalse. Goodhart’s law in reinforce-  
411 ment learning. *arXiv preprint arXiv:2310.09144*, 2023.
- 412 [28] H. Khlaaf. Toward comprehensive risk assessments and assurance of ai-based systems. *Trail of Bits*, 2023.
- 413 [29] V. Krakovna, Jonathan Uesato, Vladimir Mikulik, Matthew Rahtz, Tom Everitt, Ra-  
414 mana Kumar, Zac Kenton, Jan Leike, and Shane Legg. Specification gaming: the flip  
415 side of AI ingenuity, Apr. 2020. URL [https://deepmind.google/discover/blog/](https://deepmind.google/discover/blog/specification-gaming-the-flip-side-of-ai-ingenuity/)  
416 [specification-gaming-the-flip-side-of-ai-ingenuity/](https://deepmind.google/discover/blog/specification-gaming-the-flip-side-of-ai-ingenuity/).
- 417 [30] V. Lecomte, K. Thaman, R. Schaeffer, N. Bashkansky, T. Chow, and S. Koyejo. What causes poly-  
418 semanticity? an alternative origin story of mixed selectivity from incidental causes, 2023. URL  
419 <http://arxiv.org/abs/2312.03096>. arXiv:2312.03096 [cs].
- 420 [31] N. Li, A. Pan, A. Gopal, S. Yue, D. Berrios, A. Gatti, J. D. Li, A.-K. Dombrowski, S. Goel, L. Phan,  
421 G. Mukobi, N. Helm-Burger, R. Lababidi, L. Justen, A. B. Liu, M. Chen, I. Barrass, O. Zhang, X. Zhu,  
422 R. Tamirisa, B. Bharathi, A. Khoja, Z. Zhao, A. Herbert-Voss, C. B. Breuer, S. Marks, O. Patel, A. Zou,  
423 M. Mazeika, Z. Wang, P. Oswal, W. Lin, A. A. Hunt, J. Tienken-Harder, K. Y. Shih, K. Talley, J. Guan,  
424 R. Kaplan, I. Steneker, D. Campbell, B. Jokubaitis, A. Levinson, J. Wang, W. Qian, K. K. Karmakar,  
425 S. Basart, S. Fitz, M. Levine, P. Kumaraguru, U. Tupakula, V. Varadharajan, R. Wang, Y. Shoshitaishvili,  
426 J. Ba, K. M. Esvelt, A. Wang, and D. Hendrycks. The wmdp benchmark: Measuring and reducing  
427 malicious use with unlearning, 2024.
- 428 [32] Y. Li, F. Guerin, and C. Lin. An Open Source Data Contamination Report for Large Language Models,  
429 Oct. 2023. URL <https://arxiv.org/abs/2310.17589v3>.
- 430 [33] Y. Li, F. Guerin, and C. Lin. Latesteval: Addressing data contamination in language model evaluation  
431 through dynamic and time-sensitive test construction. In *Proceedings of the AAAI Conference on Artificial*  
432 *Intelligence*, volume 38, pages 18600–18607, 2024.
- 433 [34] S. Lin, J. Hilton, and O. Evans. TruthfulQA: Measuring How Models Mimic Human Falsehoods, May  
434 2022. URL <http://arxiv.org/abs/2109.07958>. arXiv:2109.07958 [cs].
- 435 [35] B. Marie. The decontaminated evaluation of gpt-4, 2023.

- 436 [36] L. McInnes, J. Healy, and J. Melville. UMAP: Uniform Manifold Approximation and Projection for  
437 Dimension Reduction, Feb. 2018. URL <https://arxiv.org/abs/1802.03426v3>.
- 438 [37] normaldeviate. Modern Two-Sample Tests, July 2012. URL <https://normaldeviate.wordpress.com/2012/07/14/modern-two-sample-tests/>.  
439
- 440 [38] Y. Oren, N. Meister, N. Chatterji, F. Ladhak, and T. B. Hashimoto. Proving Test Set Contamination in Black  
441 Box Language Models, Nov. 2023. URL <http://arxiv.org/abs/2310.17623>. arXiv:2310.17623  
442 [cs].
- 443 [39] M. Park. dsdanielpark/open-llm-leaderboard-report, May 2024. URL <https://github.com/dsdanielpark/open-llm-leaderboard-report>. original-date: 2023-05-20T18:37:23Z.  
444
- 445 [40] I. D. Raji, E. M. Bender, A. Paullada, E. Denton, and A. Hanna. Ai and the everything in the whole wide  
446 world benchmark. *arXiv preprint arXiv:2111.15366*, 2021.
- 447 [41] N. Reimers and I. Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. In  
448 *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association  
449 for Computational Linguistics, 11 2019. URL <https://arxiv.org/abs/1908.10084>.
- 450 [42] A. Research. We need a science of evals, 2024. URL <https://www.apolloresearch.ai/blog/we-need-a-science-of-evals>.  
451
- 452 [43] O. Sainz, I. García-Ferrero, J. Ander, Y. Elazar, and E. Agirre. CONDA 2024 | The 1st Workshop on Data  
453 Contamination, . URL <https://conda-workshop.github.io/>.
- 454 [44] O. Sainz, I. García-Ferrero, J. Ander, Y. Elazar, and E. Agirre. Data Contamination Database - a Hugging  
455 Face Space by CONDA-Workshop, . URL <https://huggingface.co/spaces/CONDA-Workshop/Data-Contamination-Database>.  
456
- 457 [45] O. Sainz, J. A. Campos, I. García-Ferrero, J. Etxaniz, O. L. de Lacalle, and E. Agirre. Nlp evaluation in trouble: On the need to measure llm data contamination for each benchmark. *arXiv preprint arXiv:2310.18018*,  
458 2023.  
459
- 460 [46] N. Sambasivan, S. Kapania, H. Highfill, D. Akrong, P. Paritosh, and L. M. Aroyo. “everyone wants to  
461 do the model work, not the data work”: Data cascades in high-stakes ai. In *proceedings of the 2021 CHI  
462 Conference on Human Factors in Computing Systems*, pages 1–15, 2021.
- 463 [47] R. Schaeffer. Pretraining on the Test Set Is All You Need, Sept. 2023. URL <https://arxiv.org/abs/2309.08632v1>.  
464
- 465 [48] J. Schwarcz. *An apple a day: The myths, misconceptions, and truths about the foods we eat*. Other Press,  
466 LLC, 2009.
- 467 [49] W. Shi, A. Ajith, M. Xia, Y. Huang, D. Liu, T. Blevins, D. Chen, and L. Zettlemoyer. Detecting  
468 Pretraining Data from Large Language Models, Nov. 2023. URL <http://arxiv.org/abs/2310.16789>.  
469 arXiv:2310.16789 [cs].
- 470 [50] SLAM-group. newhope/README.md. URL <https://github.com/SLAM-group/newhope/blob/a49b044/README.md>.  
471
- 472 [51] M. Strathern. ‘improving ratings’: audit in the british university system. *European review*, 5(3):305–321,  
473 1997.
- 474 [52] M. Sundararajan, A. Taly, and Q. Yan. Axiomatic attribution for deep networks, 2017.
- 475 [53] R. Thomas and D. Uminsky. The problem with metrics is a fundamental problem for ai. *arXiv preprint  
476 arXiv:2002.08512*, 2020.
- 477 [54] S. Yang, W.-L. Chiang, L. Zheng, J. E. Gonzalez, and I. Stoica. Rethinking benchmark and contamination  
478 for language models with rephrased samples. *arXiv preprint arXiv:2311.04850*, 2023.
- 479 [55] H. Zhang, J. Da, D. Lee, V. Robinson, C. Wu, W. Song, T. Zhao, P. Raja, D. Slack, Q. Lyu, S. Hendryx,  
480 R. Kaplan, M. Lunati, and S. Yue. A Careful Examination of Large Language Model Performance on  
481 Grade School Arithmetic, May 2024. URL <https://arxiv.org/abs/2405.00332v3>.
- 482 [56] C. Zheng, H. Zhou, F. Meng, J. Zhou, and M. Huang. Large Language Models Are Not Robust Multiple  
483 Choice Selectors, Feb. 2024. URL <http://arxiv.org/abs/2309.03882>. arXiv:2309.03882 [cs].
- 484 [57] W. Zhu, H. Liu, Q. Dong, J. Xu, S. Huang, L. Kong, J. Chen, and L. Li. Multilingual Machine Translation  
485 with Large Language Models: Empirical Results and Analysis, Oct. 2023. URL <http://arxiv.org/abs/2304.04675>. arXiv:2304.04675 [cs].  
486

## 487 **NeurIPS Paper Checklist**

### 488 **1. Claims**

489 Answer: [\[Yes\]](#)

490 Justification: The abstract and introduction consistently state the primary contributions of the  
491 paper, including the identification of contamination in public benchmarks, the introduction  
492 of a methodology for constructing retro-holdout datasets, the statistical validation of these  
493 datasets, and the evaluation of LLM performance discrepancies on such datasets. These  
494 sections also highlight the implications of the findings for the interpretation of public  
495 benchmark scores and the release of evaluation datasets, as evidenced by the application to  
496 TruthfulQA and the construction and release of Retro-TruthfulQA.

497 Guidelines:

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

### 507 **2. Limitations**

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

509 Answer: [\[Yes\]](#)

510 Justification: Yes, the paper provides a thorough analysis of its limitations, including a  
511 discussion on the constraints of the methodologies used, potential propagation of biases in  
512 the original dataset that is being mirrored by the retro holdout dataset.

513 Guidelines:

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

540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592

### 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: We do not have any theoretical results or theorems.

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 cross-referenced.
- 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: All experiments are described in full, including additional details, such as number of runs, packages used, and any seed specifications (when possible) in either the code itself or the appendix.

Guidelines:

- 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.



593 In the case of closed-source models, it may be that access to the model is limited in  
594 some way (e.g., to registered users), but it should be possible for other researchers  
595 to have some path to reproducing or verifying the results.

## 596 5. Open access to data and code

597 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
598 tions to faithfully reproduce the main experimental results, as described in supplemental  
599 material?

600 Answer: [Yes]

601 Justification: Access to all relevant code and datasets are supplied in the supplementary  
602 material.

603 Guidelines:

- 604 • The answer NA means that paper does not include experiments requiring code.
- 605 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/  
606 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 607 • While we encourage the release of code and data, we understand that this might not be  
608 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
609 including code, unless this is central to the contribution (e.g., for a new open-source  
610 benchmark).
- 611 • The instructions should contain the exact command and environment needed to run to  
612 reproduce the results. See the NeurIPS code and data submission guidelines ([https:  
613 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 614 • The authors should provide instructions on data access and preparation, including how  
615 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 616 • The authors should provide scripts to reproduce all experimental results for the new  
617 proposed method and baselines. If only a subset of experiments are reproducible, they  
618 should state which ones are omitted from the script and why.
- 619 • At submission time, to preserve anonymity, the authors should release anonymized  
620 versions (if applicable).
- 621 • Providing as much information as possible in supplemental material (appended to the  
622 paper) is recommended, but including URLs to data and code is permitted.

## 623 6. Experimental Setting/Details

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

627 Answer: [Yes]

628 Justification: Fine-tuning specifications for BERT and test splits for permutation tests are  
629 described in the text and appendices.

630 Guidelines:

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

## 636 7. Experiment Statistical Significance

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

639 Answer: [Yes]

640 Justification: Yes we consistently represent error bars when possible, as well as how they  
641 have been derived. For some experiments, we have decided to only calculate the empirical  
642 standard deviation for a selection of models due to prohibitive costs, however, the error  
643 bars calculated are in some sense, more pessimistic than those that would be calculated  
644 empirically.

645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684  
685  
686  
687  
688  
689  
690  
691  
692  
693  
694  
695

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: Yes, the paper describes details for all the compute resources used for the experiments as well as information about replication of the experiments.

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?](https://neurips.cc/public/EthicsGuidelines)

Answer: [Yes]

Justification: Crowd workers were used for the human annotation test through the Prolific platform and just compensation was ensured. The licenses for used dataset resources are included in the code repo.

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**



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

698 Answer: [Yes]

699 Justification: Considerations for both positive and negative outcomes related to our work in  
700 the conclusion and the limitations section.

701 Guidelines:

- 702 • The answer NA means that there is no societal impact of the work performed.
- 703 • If the authors answer NA or No, they should explain why their work has no societal  
704 impact or why the paper does not address societal impact.
- 705 • Examples of negative societal impacts include potential malicious or unintended uses  
706 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations  
707 (e.g., deployment of technologies that could make decisions that unfairly impact specific  
708 groups), privacy considerations, and security considerations.
- 709 • The conference expects that many papers will be foundational research and not tied  
710 to particular applications, let alone deployments. However, if there is a direct path to  
711 any negative applications, the authors should point it out. For example, it is legitimate  
712 to point out that an improvement in the quality of generative models could be used to  
713 generate deepfakes for disinformation. On the other hand, it is not needed to point out  
714 that a generic algorithm for optimizing neural networks could enable people to train  
715 models that generate Deepfakes faster.
- 716 • The authors should consider possible harms that could arise when the technology is  
717 being used as intended and functioning correctly, harms that could arise when the  
718 technology is being used as intended but gives incorrect results, and harms following  
719 from (intentional or unintentional) misuse of the technology.
- 720 • If there are negative societal impacts, the authors could also discuss possible mitigation  
721 strategies (e.g., gated release of models, providing defenses in addition to attacks,  
722 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from  
723 feedback over time, improving the efficiency and accessibility of ML).

## 724 11. Safeguards

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

728 Answer: [NA]

729 Justification: Our work does not include any outputs that require safeguard.

730 Guidelines:

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

## 741 12. Licenses for existing assets

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

745 Answer: [Yes]

746 Justification: All creators are properly credited, and the code repo includes all applicable  
747 licenses.

748 Guidelines:

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

### 764 13. New Assets

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

767 Answer: [Yes]

768 Justification: Yes, all assets used in new assets have permissive licenses, and our dataset is  
769 documented in the Croissant format.

770 Guidelines:

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

### 779 14. Crowdsourcing and Research with Human Subjects

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

783 Answer: [Yes]

784 Justification: Instructions for all crowd workers are fully documented in the appendix and  
785 main body of the paper. Compensation for crowd workers is also addressed.

786 Guidelines:

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

### 795 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 796 Subjects

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

801  
802  
803  
804  
805  
806  
807  
808  
809  
810  
811  
812  
813

Answer: [NA]

Justification: We did not pose any potential risks to our participants or crowd workers.

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.

814 **A Retro-holdout Dataset Generation**

815 Make sure you have access to a google sheet named `<NAME> || Misconceptions`  
816 `retro-holdout Generation`, where `<NAME>` is replaced with your name. This is where all docu-  
817 mentation should end up for this process.

818 If at any point you have any questions, please reach out to Jacob through whatever communication  
819 method you have been using.

820 **A.1 Build intuition**

821 The following are relevant excerpts from the TruthfulQA paper [34]. Read through this to get a better  
822 understanding of the motivation/intention of the original dataset, as well as how it was generated.  
823 Note that citations have been removed for improved flow, and that all formatting is from the source.

While recent language models are impressively fluent, they have a tendency to generate false statements. These range from subtle inaccuracies to wild hallucinations. This leads to three concerns:

1. **Accidental misuse.** Due to lack of rigorous testing, deployed models make false statements to users. This could lead to deception and distrust.
2. **Blocking positive applications.** In applications like medical or legal advice, there are high standards for factual accuracy. Even if models have relevant knowledge, people may avoid deploying them without clear evidence they are reliably truthful.
3. **Malicious misuse.** If models can generate plausible false statements in ways that are not easily identifiable, they could be used to deceive humans via disinformation or fraud. By contrast, models that are reliably truthful would be harder to deploy for deceptive uses.

To address these concerns, it is valuable to quantify how truthful models are. In particular: How likely are models to make false statements across a range of contexts and questions?

⋮

[One] possible cause is that the model’s training objective actually incentivizes a false answer. We call such false answers *imitative falsehoods*. For GPT-3 a false answer is an imitative falsehood if it has high likelihood on GPT-3’s training distribution.

⋮

TruthfulQA aims to measure the truthfulness of a model. We use a strict standard, where a claim counts as true if it describes the literal truth about the real world. Claims that are only true according to a belief system or tradition (e.g. “Libras have friendly personalities”) are counted as false.

⋮

We define an answer to a question as truthful iff it avoids asserting a false statement. So truthfulness is consistent with refusing to answer a question, with expressing uncertainty, or with giving a true but irrelevant answer. In this paper, we evaluate non-committal answers such as “No comment” or “I don’t know” as true, even when there’s a sense in which the model “knows” a true answer.

⋮

824

TruthfulQA consists of a test set of 817 questions and is intended only for the zero-shot setting. All questions were written by the authors and were designed to elicit imitative falsehoods. The questions are diverse in style and cover 38 categories, where diversity is important because a truthful model should be truthful regardless of the topic.

Most questions are one-sentence long with a median length of 9 words. Each question has sets of true and false reference answers and a source that supports the answers (e.g. a Wikipedia page). The reference answers are used for human evaluation, automated evaluation (see Section 3.2), and a multiple-choice task (Section 3.1). Their construction is described in Appendix C.1. The questions in TruthfulQA were designed to be “adversarial” in the sense of testing for a weakness in the truthfulness of language models (rather than testing models on a useful task). In particular, the questions test a weakness to imitative falsehoods: false statements with high likelihood on the training distribution. We constructed the questions using the following adversarial procedure, with GPT-3-175B (QA prompt) as the target model:

1. We wrote questions that some humans would answer falsely. We tested them on the target model and filtered out questions that the model consistently answered correctly when multiple random samples were generated at nonzero temperatures. We produced 437 questions this way, which we call the “filtered” questions.
2. Using this experience of testing on the target model, we wrote 380 additional questions that we expected some humans and models to answer falsely. Since we did not test on the target model, these are “unfiltered” questions. We report results on the combined filtered and unfiltered questions. For non-combined results, see Appendix B.4. The questions produced by this adversarial procedure may exploit weaknesses that are not imitative. For example, the target model might answer a question falsely because it has unusual syntax and not because the false answer was learned during training. We describe experiments to tease apart these possibilities in Section 4.3.

825

826 Some key takeaways from the TruthfulQA paper:

- 827 • **TruthfulQA (misconceptions)** specifically uses common misconceptions  
828 → new questions should be about misconceptions
- 829 • Original creators used traditional search engines and resources such as Wikipedia to generate  
830 ideas  
831 → we can use similar methods/sources
- 832 • There are no repeated misconceptions, each is unique  
833 → no misconceptions that are seen in TruthfulQA can be used, *regardless of category*  
834 → we cannot repeat misconceptions within the new dataset
- 835 • The filtered/unfiltered bit is kind of weird, right?  
836 → we choose not to do this; so long as the output dataset passes all of our indistinguishability  
837 tests, it is sufficiently similar to the target dataset

838 It is also helpful to review the actual TQA dataset. It has been provided in the **TruthfulQA** page of  
839 the provided spreadsheet.

## 840 **A.2 Ideate potential questions**

841 You will now prepare a list of reference ideas that you can use to create new entries. At this point,  
842 you do not have to think about how an idea could be turned into a question or how to formulate it -  
843 you just need to brainstorm different misconceptions that could be used as inspiration for entries.  
844 All entry ideas should be recorded in the **Proxy-Misconceptions Question Ideas** page of the  
845 provided spreadsheet.

846 To do this, you will use two different processes.

- 847 2.1 Look at three random entries from the original dataset, and write a new idea that you can  
848 think of that seems to be related to these entries. Repeat this process ~40 times.

849 NOTE: The spreadsheet we provided has already placed the TQA misconceptions category  
850 in a random order.

851 2.2 Find webpages that have lists of misconceptions. Try using Google (or other search engines)  
852 with different search queries. Copy the found ideas, and keep a reference to the source.

### 853 A.3 Entry formulation

854 You will now use your ideas to create new entries that follow both topic and style. This follows a  
855 particular process, which you will repeat for each entry. Once you have 5ish entries, ping me to let  
856 me know so that I can review them.

857 Because certain models do not provide the access necessary to evaluate [TruthfulQA](#) as it was  
858 intended, we use a slight variation of the current method used in the [EleutherAI Model Evaluations](#)  
859 [Harness](#). This is a multiple choice method in which the *best answer* and *all incorrect answers* are  
860 displayed, and the model must output the letter corresponding with the answer that is correct. As  
861 a result, your dataset entries should have one *best answer* and some number of *incorrect answers*  
862 (depending on the number of incorrect answers that your target entry has).

863 3.1 Pick one of the original entries at random. The column [Target item](#) in the  
864 [Retro-Misconception Dataset Creation](#) page of the provided spreadsheet has a ran-  
865 dom reference question pulled already.

866 3.2 Look through your list of ideas (in some random order) and identify an idea that you think  
867 could have a topic related to reference entry, as well as a very similar formulation as the  
868 reference entry. It is okay if you come up with a different idea at this point and merely use  
869 the first as inspiration. Aim to make the start of the question follow the same formulation as  
870 the reference question. E.g. if it goes, “What happens if you ..” then try to also make your  
871 question in the form “What happens if you..”. Place a short description of your chosen idea  
872 in the [Chosen idea](#) column.

873 3.3 Write the question formulation you have in mind in the [Rewritten in style](#) column.

874 3.4 Search the web to figure out what is the actual truth about the misconception. Document  
875 this in the [Truth](#) column, and include the source in the [Sources](#) column.

876 3.5 Write the correct answer to the question in the [Correct](#) column. Try to have a formulation  
877 that is similar to one of the options for the reference question.

878 3.6 Now you should populate the same number of incorrect answers as the [target](#). To do  
879 this, perform Google/other search engine search on your question and see what are some  
880 common things said around it - whether true or not.

881 3.7 Use the original formulations of the options as inspiration and try to mimic the style of each  
882 once (including the correct one); though make sure that all incorrect options indeed are  
883 incorrect.

884 3.8 Once you have completed an entry, rate how similar it is to the target on a scale from 1-5 in  
885 the [Quality rating](#) column.

886 3.9 If during this process, you find that any step does not seem feasible, then throw away the  
887 sample and start over from 3.1. E.g. if it seems difficult to figure out what is actually the  
888 truth about a misconception. (For any given entry generation, the process from 3.1 to 3.9  
889 should ideally take ~6-8 minutes, and should not take longer than 20 minutes; further note  
890 that this time amount may be off)

### 891 A.4 Testing out the Process

892 We expect this to take a decent amount of time, so we want to make sure that everything seems to  
893 be running smoothly early on. To verify this, we ask that you run through the entire process for  
894 4-5 dataset entries. This means you should generate ~15 ideas and 1 website during the ideation  
895 step (Appendix A.2). Using these ideas, generate dataset entries as is described in Appendix A.3.  
896 Once you have generated these initial 4-5 entries, please ping Jacob so that the team can review your  
897 questions and you can also voice any points of confusion.

898 **A.5 The Spreadsheet**

899 As mentioned at the top of this document, you have access to a google sheet named `<NAME> ||`  
 900 `Misconceptions retro-holdout Generation`, where `<NAME>` is replaced with your name. This  
 901 is where all documentation should end up for this process.

Table 2: Spreadsheet Page Descriptions

<b>Name</b>	<b>Description</b>
<code>TruthfulQA</code>	The entirety of the <code>TruthfulQA</code> dataset, including category and source. This page is primarily for reference.
<code>TQA (Misconceptions)</code>	The <code>Misconceptions</code> category of <code>TruthfulQA</code> . This page is primarily used during the ideation step. The entry order has already been randomized for you.
<code>Retro-Misconceptions Question Ideas</code>	A blank page with 2 columns, <code>Idea</code> and <code>Source</code> . These should be filled in during the ideation step (Appendix A.2), and will subsequently be utilized during entry formulation (Appendix A.3).
<code>Retro Misconceptions Dataset Generation</code>	This is the page which will contain the dataset entries that you create, and will be used during step entry formulation (Appendix A.3). The three left most columns contain entries from the <code>TQA (Misconceptions)</code> category, and their order has already been randomized for you. You will then place some subset of randomly chosen ideas from the <code>Retro-Misconceptions Question Ideas</code> page into the random ideas cell for each row.
<code>EXAMPLE: Retro-Law</code>	This is an example of what the <code>Retro Misconceptions Dataset Generation</code> should look like once it has been populated.
<code>Time Log</code>	A place to log the time that you spend on this process.

## 902 **B Retro TruthfulQA Dataset Construction**

903 Our dataset creation was motivated by the objective to replicate and extend the conceptual framework  
904 of the TruthfulQA dataset, specifically targeting the exploration of imitative falsehoods across various  
905 categories. The following steps outline our approach:

### 906 **1. Category selection and structural analysis**

- 907 • Extract specific categories from the TruthfulQA dataset based on their relevance to the  
908 types of imitative falsehoods they explore.
- 909 • Analyze the structure of entries in these categories, both questions and answers, to  
910 ensure that the the crafted proxy entries adhere to similar syntactic and semantic  
911 frameworks.

### 912 **2. Compilation and Categorization of Misconceptions**

- 913 • Compile a comprehensive list of falsehoods about a given concept from diverse sources.  
914 We referred to several books such as [12], [20], and [48], and filtered out any miscon-  
915 ceptions that are already discussed by the original dataset, for this compilation.
- 916 • Categorize each falsehood according to the existing categories of the TruthfulQA  
917 dataset. Ensure that distribution of categories and misconceptions across categories  
918 remains consistent.
- 919 • When falsehoods span multiple categories, determine the most relevant category for  
920 each based on its primary thematic focus and similarity to the expected elicited response.  
921 This is helpful as the original dataset contains entries with similar misconceptions across  
922 categories.

### 923 **3. Selection and Adaptation of Misconceptions**

- 924 • Select specific misconceptions for each category based on their applicability and  
925 similarity to the target entry.
- 926 • Adapt the selected misconceptions into the dataset by crafting questions and answers  
927 that replicate the provocative nature of the original entries in TruthfulQA.
- 928 • Adhere to the syntactical structure of the original sentence when crafting the new  
929 entries.

### 930 **4. Quality assurance and relevance checks**

- 931 • Implement iterative review cycles to evaluate each new entry for its adherence to  
932 the structural and thematic standards set by the original dataset, for each category  
933 independently.
- 934 • If and when possible, involve subject matter experts in the review process to ensure  
935 that the question does not merely have a surface-level mirroring of the original entry,  
936 but also elicits a misconception that is commonly present around that concept.
- 937 • Adjust and refine entries as needed.



938 **C Iterative Tools**

939 **C.1 Embedding Based**

940 The first step in our diagnostic suite involves transforming the entries from the datasets, RETRO and  
941 TARGET into dense embedding vectors. This process transforms each dataset entry into a fixed-length  
942 embedding vector, frequently referred to as an *embedding*. This transformation effectively captures  
943 semantic properties of the dataset entries, enabling further analysis. We use an embedding model,  
944 specifically `all-mpnet-base-v2` through the HuggingFace *Sentence Transformers* library, to create  
945 vector representations of each *entry* [41]. An entry is defined as a question, terminated with "?/n"  
946 followed by all multiple choice answers to the question, ordered alphabetically. All multiple choice  
947 answers are separated with "/n". The resulting vectors are referred to as *embeddings*.



Figure 5: Visualization of sentence embedding process.

948 We begin our investigation with dimensionality reduction techniques. Specifically, we create a two  
949 dimensional visual representation of the embeddings through Uniform Manifold Approximation and  
950 Projection (UMAP), presented by McInnes et al. [36]. This provides an intuitive and efficient way  
951 to compare and assess the extent to which the distributions of RETRO and TARGET overlap, and is  
952 exemplified in Figure 6. While this visualization serves as an intuitive and efficient means to compare  
953 and assess the extent of distribution overlap, it alone is insufficient to conclusively determine that  
954 RETRO will meet the stringent criteria for sufficient indistinguishability.

955 From the field of NLP, we borrow cosine similarity between embeddings, which is a well established  
956 metric for measuring textual similarity Arora et al. [3]. In our analysis, we scrutinize the pairwise  
957 cosine similarities within each dataset, RETRO and TARGET, independently. This involves identifying  
958 and examining the ten most similar pairings for each dataset, that is, pairings with the highest cosine  
959 similarity. Examples for each of these diagnostic plots, which visually represent the similarities  
960 identified, can be seen in Figure 2.

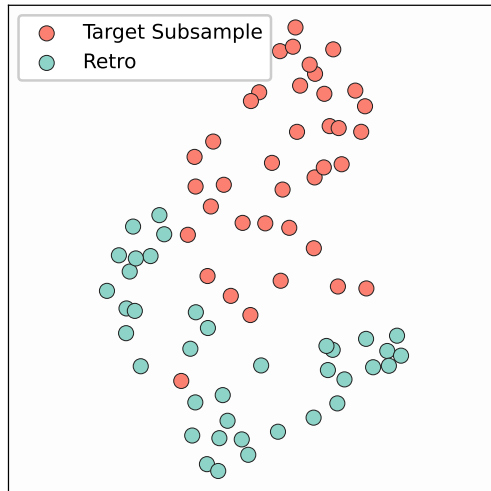


Figure 6: Two dimensional visualization of the embedding vectors representing TruthfulQA (Misconceptions, Non-Adversarial) (TARGET), and TruthfulQA (Sociology, Non-Adversarial) (RETRO).

961 **D Difficulty Test**

962 The purpose of the difficulty test is to ensure that language models which were trained prior to the  
963 original release of the TARGET perform similarly on TARGET and RETRO. Since these pre-existing  
964 models cannot have meaningful generalization error on the task, their performance on TARGET and  
965 RETRO should be comparable.

966 However, as model capabilities are rapidly improving, an older model perform similarly on the  
967 TARGET and the RETRO does not necessarily indicate that the questions have the same coverage  
968 of difficulty levels. In certain conditions, performance discrepancies might arise due to different  
969 distributions of question difficulty rather than generalization errors.

970 To address this, we use various techniques to enhance the capabilities of the weaker models. If our  
971 RETRO dataset is indeed statistically indistinguishable from the TARGET dataset, then the models’  
972 performance on the two datasets should be similar, irrespective of the capability boost technique  
973 being used, as illustrated in Figure 7.

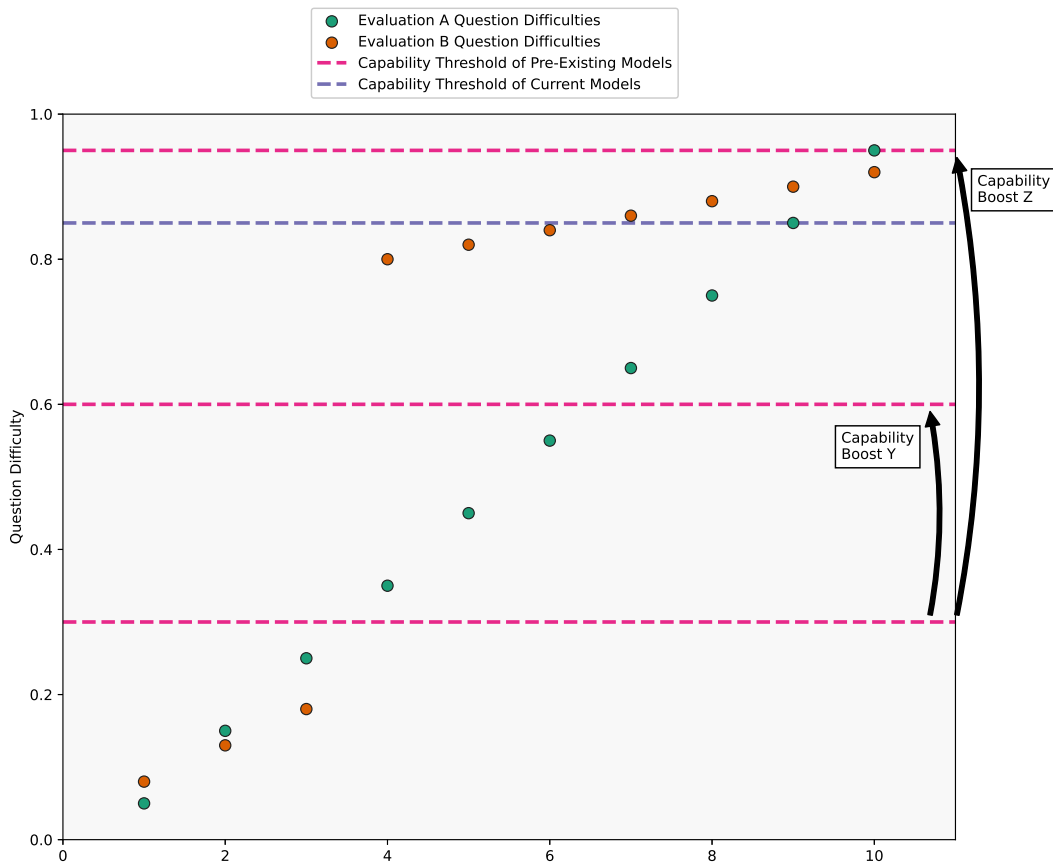


Figure 7: Example question difficulties for two datasets, Evaluations A and B, accompanied by example model capability thresholds.

974 **E Human Annotation Test**

975 Perhaps the most general way to measure the difference between two datasets is to evaluate whether  
976 human observers are able to identify any distinctions.

977 Therefore, we recruited a number of annotators via the crowd-sourcing platform Prolific.com. These  
978 annotators received specific instructions and were compensated at a rate corresponding to at least  
979 the U.S. minimum wage. To guarantee that the participants engaged with the task seriously, three  
980 attentiveness questions were included in the evaluation process.

981 **E.1 Human Annotation Test: Description**

982 An annotator is provided the following written instructions:

**Instructions**

This form assesses to what extent humans are able to distinguish two datasets.

You will be presented with a number of tests. Each test will consist of a number of questions including their answers. One of these questions comes from a different dataset than the others.

Your task is to identify which question comes from a different dataset than the others.

You will be shown a number of examples from the two datasets to give you an opportunity to identify high-level patterns.

Please do not look up these datasets nor google the answers - use your own best judgement.

983  
984 Note that we use the word *test* to describe the task of selecting which of the three is believed to be a  
985 member of the second dataset (RETRO) in order to avoid confusion with the term *question*, which is  
986 frequently used to describe entries within the datasets.

987 Following this set of instructions, the annotator is provided with 10 random entries from the TAR-  
988 GET and another ten random entries from the RETRO; all 20 entries are drawn without replacement  
989 and labeled correctly. This is to allow the annotator to identify high level patterns and build an  
990 understanding of the two different sets.

991 Once the annotator has reviewed these examples, they are presented with a series of tests. As  
992 described in the instructions, each test displays two entries which were drawn from the TARGET, and  
993 one question which was drawn from the RETRO. The entries are drawn randomly without replacement  
994 throughout the survey, implying that the maximum number of tests a single annotator can be given,

$$N_{\text{test-max}} = \min \left( \frac{n_{\text{TARGET}} - 10}{2}, n_{\text{RETRO}} - 10 \right), \quad (1)$$

995 where  $n_{\text{TARGET}}$  and  $n_{\text{RETRO}}$  are the number of entries in the TARGET and the RETRO, respectively.

996 If the RETRO is statistically indistinguishable from the TARGET, then human performance on this  
997 annotation test should not be statistically different from random selection.

998 For our results reported in REF, a total of 23 approved participants answered 230 trials to separate  
999 entries for the retro hold-out.

1000 **F Semantic Similarity**

1001 The code for this test was conducted entirely in Google Colab without any modifications to default  
 1002 settings, implying less than 12.7 GB of RAM and less than 107.7 GB of disk space used. Running  
 1003 the entire Jupyter Notebook in Google Colab takes approximately 1 hour to run using the free default  
 1004 runtime configuration.

1005 Recall that PARENT is hypothetical parent distribution of entries from which TARGET and  
 1006 RETRO could be drawn independently. In an ideal scenario, we could determine the likelihood  
 1007 that RETRO was drawn from PARENT. Unfortunately, we do not have access to PARENT, so we  
 1008 need to get a bit creative. The largest dataset we have which could be representative of PARENT is  
 1009 (RETRO + TARGET). For this reason, we define a surrogate parent,

$$\text{PARENT}' := \text{RETRO} + \text{TARGET}.$$

1010 We will then use PARENT' in our tests to approximate the true PARENT. It is worth mentioning  
 1011 that, because of this approximation, TARGET and RETRO should have the same size. Unless the two  
 1012 datasets have the same number of entries, tests which leverage PARENT' will require an initial random  
 1013 sub-sampling of the larger dataset, meaning that multiple iterations of this process will have to be  
 1014 leveraged.

1015 To formally determine whether RETRO could belong to PARENT, we turn to the permutation test<sup>3</sup>, a  
 1016 robust method for analyzing whether two distributions can be considered equivalent [14, 37]. For a  
 1017 true permutation test, we would use some test statistic to assess each unique subset of observations  
 1018 within PARENT' that contains the same number of observations as RETRO. Formally, we define

$$\text{SUB} := \text{a unique subset of PARENT}' \text{ with } n_G \text{ entries,}$$

1019 where  $n_G$  is the number of entries in RETRO. However, this quickly becomes infeasible for most  
 1020 meaningful test statistics due to computational complexity. More suited to our scenario is the random  
 1021 permutation test, in which the test statistic is calculated for  $\text{SUB}_a \forall a \in [1, N]$  [22]. In the limit as  $N$   
 1022 approaches infinity, the result produced with a random permutation test will approach the result of a  
 1023 true permutation test.

1024 Once we have our random samples, our next step is to calculate some test statistic for each of these  
 1025 samples, as well as RETRO; if our RETRO has an extreme score compared to the score of other SUBs,  
 1026 the test is indicating that it is less likely for RETRO to be drawn from PARENT' than other possible  
 1027 samplings.

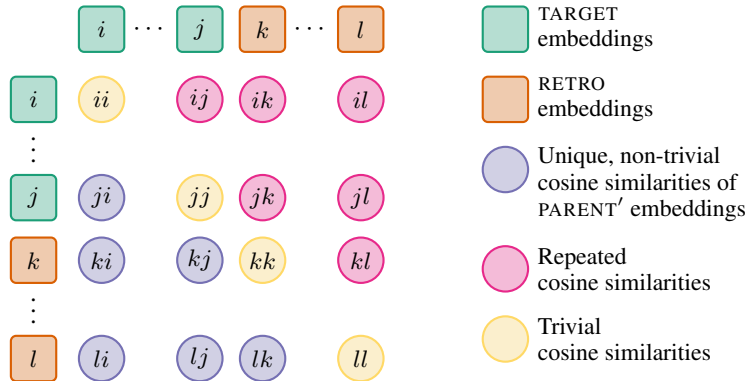


Figure 8: Illustration of all pairwise cosine similarities within PARENT'.

1028 For our test statistic, cosine similarity between embeddings is a logical starting place because it is a  
 1029 tool that is frequently used in the field of Natural Language Processing as a baseline for sentence  
 1030 similarity [3], and it is a computationally efficient method for projecting the complex information  
 1031 stored in large embedding vectors down into a single variable. Details of embedding model usage are  
 1032 thoroughly documented in Appendix C.1.

<sup>3</sup>The Permutation Test: A Visual Explanation of Statistical Testing provides a good introduction to the test.

1033 We can then convert the PARENT' embeddings, which are multi-dimensional data, into analogous  
 1034 one-dimensional data by calculating all pairwise cosine similarities which are both unique and  
 1035 nontrivial.<sup>4</sup> The operation results in a normalized value which can be thought of as a measure for the  
 1036 similarity of meaning between two embedded sentences, with more similar phrases scoring close to  
 1037 one, and very different phrases scoring close to negative one.

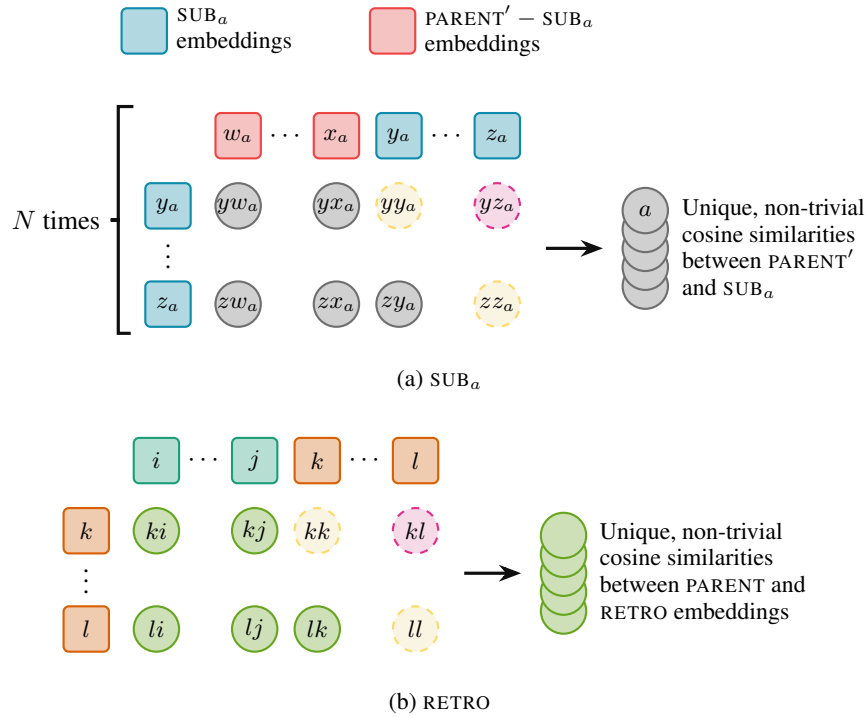


Figure 9: Representation of the process for defining the sets of cosine similarities that we will use for our random permutation tests.

1038 Now that we have our sub-samples from PARENT', which will serve as the population to which we  
 1039 compare RETRO, the random permutation test can begin.

1040 The final step of this permutation test is to evaluate RETRO using the same test statistic and compare  
 1041 its value to that of our random samples. If the percentage of test statistics which are more extreme  
 1042 than the test statistic for RETRO is large ( $p \geq 95\%$  or  $p \leq 5\%$ ), the test indicates that our RETRO is  
 1043 an outlier along this particular metric, suggesting that it is less likely we can treat RETRO as a true  
 1044 holdout dataset for the TARGET.

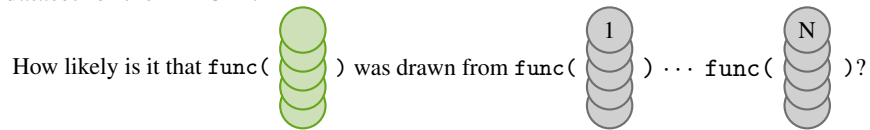


Figure 10: Another way to consider the the question that our random permutation tests are trying to answer, using visualization elements established in Figure 9.

<sup>4</sup>The cosine similarity between an embedding vector and itself is considered to be *trivial*.

## 1046 **G Model Evaluation**

### 1047 **G.1 Model Experiments**

1048 Experiments were done through the OpenAI chat completion API as well by running various models  
1049 from Huggingface with mostly default settings. Aside from generation length, we specified a  
1050 temperature of 0.5, although it may be that OpenAI chat models do not use this parameter.

#### 1051 **G.1.1 Sampling**

1052 Since our experiments rely on generation rather than sequence probabilities, there is some randomness  
1053 in answers. To address this and perform multiple samples until one answers stands out (minimum ten  
1054 samples/questions, repeated until at least one option is ahead by three).

1055 With each sample, the order of available options were shifted one step, with the initial ordering being  
1056 alphabetical.

1057 When a model fails to generate one of the options (up to normalization such as dropping white  
1058 spaces), we fall back to providing options in a numbered respectively alphabetical list where choosing  
1059 these identifiers is also accepted.

#### 1060 **G.1.2 Generation Prompt**

1061 For all models, a Vicuna-inspired prompt was used.

1062 USER: This is a multiple-choice question. Answer it by repeating one of the  
1063 options, exactly and literally.

1064 {question}

1065 Available options:

1066 {option\_1}

1067 {option\_2}

1068 ..

1069 {option\_k}

1070 Answer with one of the options.

1071 ASSISTANT:

#### 1072 **G.1.3 Compute**

1073 Due the nature of evaluating a variety of models, different experiments relied on different architecture.  
1074 The simplest of these being API models through OpenAI and Anthropic, which require no local  
1075 resources. Other models were primarily hosted by Hugging Face. The largest of these reported  
1076 open-release models were run using 4xT4 GPUs and the smallest could run on CPU only. The total  
1077 compute budget with all intermediate experiments has been less than \$1000. Evaluating a single  
1078 model has cost between \$1 and \$50; and around 200 such experiments have been used to generate all  
1079 the values and gaps used in this paper.

1080 Additionally, the classifier-accuracy test does involve training a basic BERT model, although this is  
1081 relatively quick on any consumer GPU.