
Benchmark Probing: Investigating Data Leakage in Large Language Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Large language models have consistently demonstrated exceptional performance
2 across a wide range of natural language processing tasks. However, concerns have
3 been raised about whether LLMs rely on benchmark data during their training phase,
4 potentially leading to inflated scores on these benchmarks. This phenomenon,
5 known as data contamination, presents a significant challenge within the context
6 of LLMs. In this paper, we present a novel investigation protocol named **Testset
7 Slot Guessing (TS-Guessing)** on knowledge-required benchmark MMLU and
8 TruthfulQA, designed to estimate the contamination of emerging commercial
9 LLMs. We divide this protocol into two subtasks: (i) *Question-based* setting:
10 guessing the missing portion for long and complex questions in the testset (ii)
11 *Question-Multichoice* setting: guessing the missing option given both complicated
12 questions and options. We find that commercial LLMs could surprisingly fill in the
13 absent data and demonstrate a remarkable increase given additional metadata (from
14 22.28% to 42.19% for Claude-instant-1 and from 17.53% to 29.49% for GPT-4).

15 1 Introduction

16 Large language models (LLMs) have demonstrated exceptional performance across a wide range
17 of NLP tasks, and the NLP community has witnessed the emergence of several impressive LLMs.
18 Notably, there are robust commercial LLMs, including the GPT-* [3, 16] by OpenAI, Claude [1] by
19 Anthropic, and Bard [6] by Google, among others. In addition to these commercial models, there are
20 numerous open-source LLMs, such as Llama [22, 23], MPT [13], Falcon [15], and Chinchilla [8].
21 However, concerns have arisen regarding these LLMs, primarily related to their extensive training
22 on web data, often at a terabyte scale. This extensive training data may, in turn, potentially overlap
23 with the current benchmark [3, 4, 22, 23], which is also frequently constructed from Internet sources.
24 Research has revealed that pretraining on the testset can artificially inflate performance metrics [18].
25 Consequently, it becomes imperative for the community to address the detection of potential data
26 contamination in these models.

27 Despite the pressing need for research on data contamination, there appears to be a scarcity of relevant
28 studies. For current Large Language Models, n-gram based methods are introduced [3, 24, 23] to
29 detect data contamination. To summarize, our approach involves employing n-gram tokenization
30 to partition large documents into smaller shards and assessing their similarity with benchmark
31 data [4, 22]. However, this method is contingent upon having complete access to the training corpus,
32 making it challenging to estimate data contamination for models [3, 16, 6, 1, 10] that do not disclose
33 their training data. Therefore, there is a clear necessity to develop a more robust method for detecting
34 potential contamination in both *open-sourced* and *closed-sourced* Language Models.

35 In this paper, we introduce a novel investigation protocol referred to as TS-Guessing in two distinct
36 settings: (1) Question-based guessing and (2) Question-multichoice guessing shown in Figure 1. In

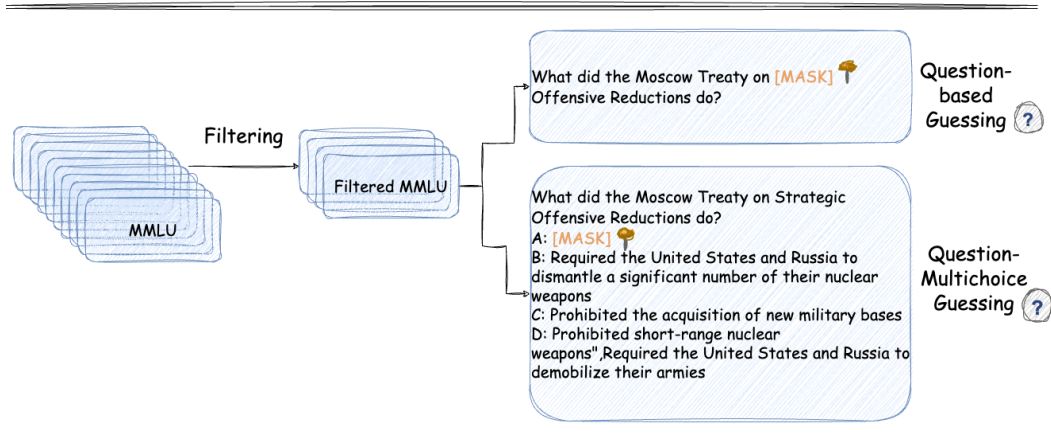


Figure 1: Illustration of workflow of **TS-Guessing** on MMLU. The prefiltering technique (§ 3.3) filters out correlated and correct options in the benchmark to rationalize our investigation protocol.

37 the *Question-based* setting, our objective is to hide a crucial word within a sentence, challenging
 38 the model to predict it accurately, while avoiding common alternatives from a vast vocabulary.
 39 In the *Question-Multichoice* setting, our goal is to conceal an incorrect answer option among
 40 multiple choices, preventing the model from giving the correct answer directly and encouraging it to
 41 complete the missing part in the benchmark. These two settings guide LLMs in guessing the missing
 42 information in the questions and answers, thereby testing their contaminated knowledge against the
 43 testset data.

44 Our experimental results reveal that different versions of LLMs from the *same* company did not
 45 exhibit a pronounced difference in their TS-Guessing performance, with GPT-4 showing only a
 46 1% improvement in the zero-shot setting compared to GPT-3.5-turbo, and Claude-2 performing
 47 5% worse than Claude-instant-1 in the question-based guessing task. These findings highlight the
 48 consistency in performance among LLMs from the same company and underscore the potential data
 49 source similarity. Besides, we observed that commercial large language models (LLMs) achieved a
 50 remarkable zero-shot accuracy of 16% with GPT-3.5-turbo, 22% with Claude-instant-1, and 25%
 51 with GPT-3.5-turbo in the Question-based setting within TruthfulQA. In the Question-Multichoice
 52 setting, GPT-3.5-turbo exhibited a noteworthy ability to guess the missing option, achieving a 57%
 53 EM rate. Considering these results, we raise concerns about the potential contamination of the current
 54 benchmark, especially if it is made publicly accessible. We urge for this to be seen as a call to action
 55 and to explore additional methods to mitigate the risk of contamination.

56 2 Related Work

57 Recent advancements in LLMs have raised concerns about data contamination and its impact on
 58 model performance. To address these concerns, researchers have explored various tokenization
 59 strategies and detection methods in the field of natural language processing. Previous research related
 60 to LLMs is introduced in GPT-3’s Appendix C [3]. In this study, GPT-3 employs 13-gram tokenization
 61 for both training data and benchmark data. PaLM [4] also employs an 8-gram strategy, considering
 62 data to be contaminated if there is a 70% overlap with 8-grams from the test set data. Open-source
 63 models such as Llama [22] adopt a similar methodology, derived from GPT-3. Llama 2 [23] (Section
 64 A.6), however, enhances this approach by using 8-gram tokenization with weight balancing. Currently,
 65 various methods, including prompt-based and probing-based approaches [5, 11], have been developed
 66 to detect potential data contamination in LLMs. Additionally, there are suggestions for mitigating
 67 potential leakage when manipulating benchmark test sets [9]. Besides the research conducted on
 68 English-only corpora, there is also ongoing investigation [2] into the issue of language contamination
 69 in cross-lingual settings.

70 3 Testset Slot Guessing Protocol

71 3.1 What does Question-based and Question-Multichoice stand for?

72 As shown in Figure 2a, in the context of the *Question-based* setting, our objective is to **mask a**
73 **pivotal word that represents the essence of the sentence**. Taking the example sentence, "Where
74 did fortune cookies originate?" into consideration, in this instance, the word "fortune" emerges as a
75 potential keyword candidate. This is because predicting the masked sentence, "Where did [MASK]
76 cookies originate?" necessitates the model to select a word from a vast vocabulary list, encompassing
77 options such as "Chocolate chip," "Biscotti," "Snickerdoodle" and so forth. However, if the model
78 has encountered testset data during its training phase, it may exhibit a greater inclination to produce
79 the missing word as 'fortune' rather than "Biscotti." or "Snickerdoodle".

80 A more challenging task is *Question-Multichoice* setting (shown in Figure 2b). In this particular
81 scenario, **our objective is to mask an incorrect option**. We intentionally *avoid masking the correct*
82 *option* to prevent the model from directly providing the correct answer, instead compelling it to guess
83 an incorrect answer from a vast set of erroneous possibilities. Furthermore, we implement detailed
84 filtering procedures (introduced in 3.3) to eliminate instances where there exists a strong correlation
85 between any answer options, thereby discouraging the model from relying on its reasoning and
86 inference capabilities to predict the masked words. when confronted with intricate questions and
87 unrelated options, if the model can still output missing options (sometimes exceeding a length of
88 8) correctly, it raises a compelling suspicion regarding the extent to which the model's behavior is
89 influenced by its exposure to benchmark data.

90 3.2 Problem Formulation

91 **Question-based** Let \mathcal{D} be a dataset containing n documents. For each document d_i , where
92 $i \in \{1, \dots, n\}$, there exists a question q_i and several answers. Given a question q_i from document d_i ,
93 we perform a *keyword searching function*

$$k_i = f_{keyword}(q_i)$$

94 where k_i is the keyword associated with q_i . Subsequently, we use a mask function $q'_i = g(q_i, k_i)$ to
95 mask the keyword in the question with [MASK]. Thus, the overall process can be represented as:

$$q'_i = g(q_i, k_i, [\text{MASK}])$$

96 **Question-Multichoice** Let \mathcal{D} be a dataset containing n documents. For each document d_i , where
97 $i \in \{1, \dots, n\}$, there is: A question denoted by Q . A list of answers denoted by A , where $A =$
98 $\{a_1, a_2, \dots, a_m\}$ and m is the number of answers for that document. One correct answer denoted by
99 a_c such that $a_c \in A$.

100 From the list A , one wrong answer is chosen and replaced with [MASK], denoted by a_{mask} . The final
101 template is a concatenation of the question, the correct answer, and three wrong answers (including
102 the masked one):

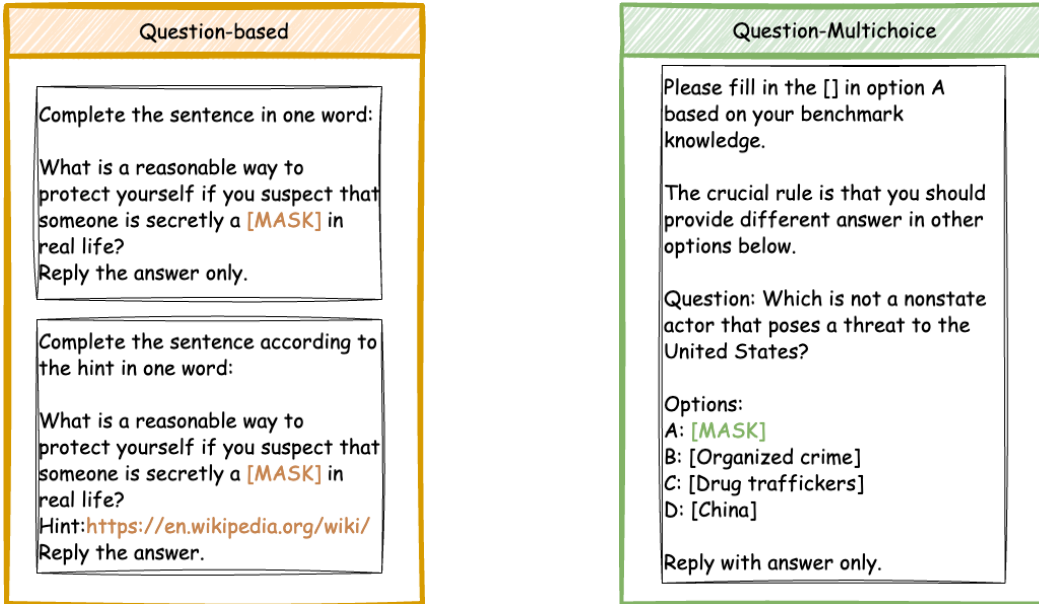
$$T_i = \text{Concat}(Q_i, a_{c_i}, a_{w1_i}, a_{w2_i}, a_{\text{mask}_i})$$

103 Where T_i is the template for the i^{th} document, Q_i is the question for the i^{th} document, a_{c_i} is the
104 correct answer for the i^{th} document, a_{w1_i} and a_{w2_i} are two wrong answers chosen from the list A
105 for the i^{th} document, a_{mask_i} is the wrong answer that has been replaced with [MASK] for the i^{th}
106 document.

107 3.3 Experiments Details

108 **Domains** We consider two datasets widely recognized for their effectiveness in evaluating knowledge
109 Question Answering in current LLMs benchmarks: (i) **MMLU** [7], a dataset measuring the knowledge
110 capabilities of LLMs and encompasses 57 diverse tasks spanning elementary mathematics, U.S.
111 history, computer science, law, and more. (ii) **TruthfulQA** [14], a benchmark assesses the truthfulness
112 of language models in generating responses to questions across 38 different categories, including
113 health, law, finance, and politics.

114 **Pre-filtering** A critical step in our experiment involves the application of filtering techniques. We
115 employ several methods to ensure that our investigative protocol does not become a straightforward



(a) Prompt template of **Question-based** guessing from handpicked example in TruthfulQA.

(b) Prompt template of **Question-Multichoice** guessing from handpicked example in MMLU.

Figure 2: Illustration of two tasks within TS-Guessing. Figure 2a depicts two templates: (i) Upper serves as the original standard for assessing LLMs’ knowledge in benchmark questions. (ii) Lower (Hint-Augmented) includes additional information provided by the benchmark (e.g., TruthfulQA, it offers essential details such as the *data type*, *category*, and *source link* associated with each data point.)

116 semantic inference or logical reasoning task. For TruthfulQA, we implement two filtering criteria: (i)
 117 removing data if its question has a length of four words or fewer, and (ii) excluding data associated
 118 with the 'Indexical Error' type. For the MMLU dataset, we adopt a more stringent filtering rule, which
 119 includes: (i) removing data containing only "Yes-No" or "True-False" options, mathematical symbols,
 120 or other simple option expressions; and (ii) removing data if the Rouge-L [12] F1 score between any
 121 two options exceeds a predefined threshold. In this paper, we have established a threshold of 0.65
 122 chosen to filter out "three words differing one in a sentence"(e.g, A:"I am American" and B:"I am
 123 Swedish." would result in the data being filtered)

124 **Keyword Searching** We are implementing a keyword searching function using two powerful tools:
 125 the Stanford POS Tagger [21] and ChatGPT with 5-shot in-context learning. Our objective is to
 126 identify the pivotal word in a question-based context. To achieve this, our approach begins by utilizing
 127 ICL ChatGPT to identify the most informative word. Subsequently, we assess whether the previously
 128 selected word falls within the categories of nouns (NN), adjectives (JJ) or verbs (VB).

129 **Hint** Hint is employed in the Question-based setting to leverage the supplementary information
 130 within the test dataset. In this paper, TruthfulQA not only supplies questions and answer options but
 131 also includes additional metadata, such as type, category, and URL information. This metadata serves
 132 as an added prompt presented to LLMs. For MMLU, we do not use a hint-based approach since the
 133 benchmark consists solely of questions and answers. Nevertheless, we posit that this methodology
 134 holds promise for application to other datasets, facilitating the exploitation of information within the
 135 test dataset.

136 3.4 Observations and Analysis

137 3.4.1 Strong Model Doesn't Indicate Proficiency In TS-Guessing

138 As depicted in Table 1 and Table 2, despite the increased power of GPT-4, we do not observe
 139 significant improvements in our TS-Guessing protocol. In the original version (without hints appended

Table 1: Exact Match (EM) rate in the **Question-based** guessing in TruthfulQA. Three kinds of hints are metadata given in TruthfulQA. (Details in 3.3)

Model	Company	Question-based			
		w/o hint	w. type-hint	w. category-hint	w. url-hint
GPT-4	OpenAI	0.17	0.19	0.15	0.29
GPT-3.5-turbo	OpenAI	0.16	0.17	0.19	0.25
Claude-2	Anthropic	0.23	0.25	0.25	0.37
Claude-instant-1	Anthropic	0.22	0.23	0.21	0.42

140 to the prompt), there is only a 1% difference between the two models. Even when utilizing URL-
 141 hint prompting in a Question-based setting, the performance gap remains minimal, with only a
 142 4% difference between GPT-3.5-turbo and GPT-4, and a fluctuation of approximately $\pm 3\%$ in
 143 performance in the Question-Multichoice setting. This pattern is consistent in both Claude-instant-1
 144 and Claude-2. In the Question-based setting, we consistently find similar performance levels in our
 145 TS-Guessing task. This suggests that our protocol may not heavily rely on advanced reasoning skills,
 146 although its performance may vary depending on the training data available.

147 This phenomenon could be explained in several ways. Firstly, the variance in training data between
 148 different companies may be significant. Secondly, even within the same company, different model
 149 versions may have closely related training data, especially when considering data that potentially
 150 overlaps with the benchmark.

151 3.4.2 Latest Benchmark Could Still Be Contaminated

152 As shown in Table 1, there are **16.24% percent of success rate** to guess the missing word in the
 153 benchmark of TruthfulQA. According to OpenAI, their training data is current up to September
 154 2021, with no utilization of data beyond that date. However, TruthfulQA made its camera-ready
 155 version available on the ACL Anthology in May 2022. Upon closer look, it becomes evident
 156 that a substantial portion of the data in TruthfulQA originates from or is derived with assistance
 157 from publicly accessible sources. It's worth noting that this publicly available content, particularly
 158 when restricted to Wikipedia, remains accessible to commercial AI companies at any given time.
 159 Consequently, careful consideration is warranted when assessing potential contamination in new
 160 benchmarks.

Table 2: Performance in the **Question-Multichoice** guessing in TruthfulQA. BLEURT is a pre-trained score metrics used in text generation evaluation [19]

Model	TruthfulQA			MMLU		
	EM	Rouge-L F1	BLEURT	EM	Rouge-L F1	BLEURT
GPT-4	0.12	0.46	0.32	0.52	0.69	0.41
GPT-3.5-turbo	0.10	0.43	0.30	0.57	0.67	0.44

161 3.4.3 MMLU Are Probably Contaminated Seriously

162 As shown in Table 2, given the fact that we have filtered out the correlated options, mathematical
 163 symbol and logic expressions. **GPT-3.5-turbo still could precisely predict masked choices in**
 164 **MMLU testset with 57% accuracy.** After filtering, the remaining options appear disorganized
 165 and complex. However, successful examples are rather surprising. In comparison to TruthfulQA,
 166 which boasts a 0.10 EM rate and a 0.43 Rouge-L F1 score, the EM rate of MMLU is noticeably
 167 higher. The high accuracy suggests that when given a question and the correct answer in MMLU,
 168 GPT-3.5-turbo has a probability greater than fifty percent of generating a candidate list with incorrect
 169 answers, just like the benchmark. we here could take a successful example in Question-Multichoice
 170 Guessing, "Which is not a nonstate actor that poses a threat to the United States?" and a correct
 171 answer "D. China" as an example. ChatGPT could magically complete another wrong option "C.
 172 Drug traffickers" if we mask option C. The candidate list for a wrong option is large and may even be

173 infinite, so when seeing LLMs could complete it exactly correctly sometimes for a very long and
174 complex sentence, this raises our concerns of benchmark data leakage.

175 3.5 Correlation between TS-Guessing and Task Accuracy

176 As illustrated in Table 3, we have included the *Spearman correlation* as a metric to assess the relation-
177 ship between our TS-Guessing protocol and task performance, thereby examining the interconnection
178 between these two tasks. In particular, we conduct this experiment on the Question-Multichoice task,
179 utilizing the Rouge-L F1 score to investigate its relevance to question answering performance.

180 Our findings reveal interesting insights. In the case of TruthfulQA, we observe a negative correlation
181 (-0.158 for GPT-4 and -0.128 for GPT-3.5-turbo) between task performance and the TS-Guessing
182 protocol. In contrast, for MMLU, which is a benchmark that has a potential contaminated risk, there
183 is a positive correlation of 0.279 for GPT-4.

Table 3: Spearman correlations between task performance and Rouge-L F1 score. $p < 0.05$ is set default

Task	Model	Corr. (ρ) with... f1 score \uparrow
TruthfulQA	GPT-4	-0.158
	GPT-3.5-turbo	-0.128
MMLU	GPT-4	0.279
	GPT-3.5-turbo	0.234

184 We aim to provide an explanation from two perspectives. Firstly, the results of our correlation test
185 suggest that while n-gram-based algorithms offer convenience, they may not be the best approach for
186 detecting data contamination in LLMs rigorously. However, this method is widely used in models
187 such as GPT-3, Llama, and Llama 2 (as discussed in Section 2).

188 Secondly, our lack of knowledge about the actual training techniques and training data used in
189 closed-source LLMs poses a challenge. In today’s landscape, numerous training techniques have
190 emerged, ranging from supervised fine-tuning (SFT) to reinforcement learning from human feedback
191 (RLHF) [17], and even mixture of experts (MoE) [20]. Applying the same evaluation methods to
192 different techniques could yield varying results.

193 4 Conclusion and Future Work

194 In this paper, we present a novel investigation protocol designed to assess the potential data leakage
195 in benchmark datasets when evaluated with Language Model Models (LLMs). Our results reveal
196 that both commercial LLMs from OpenAI and Claude exhibit the capability to accurately complete
197 missing options in the test set. GPT-3.5-turbo achieved a 57% accuracy in predicting masked choices
198 in the MMLU testset, with remaining options seeming disorganized and complex after filtering.
199 Compared to TruthfulQA, MMLU exhibited a significantly higher EM rate despite its challenging
200 nature. This observation raises concerns about potential data leakage in contemporary benchmark
201 datasets.

202 This study can be extended beyond closed-source models, encompassing open-source LLMs as well.
203 It offers a valuable tool for identifying and detecting potential data leakage, shedding light on how
204 LLMs acquire knowledge about test data from benchmarks. This, in turn, provides insights into
205 benchmark contamination and informs us about the appropriate times to update our benchmarks. As
206 the field of natural language processing continues to evolve, the development of new benchmark
207 datasets and evaluation protocols should be a priority. These new benchmarks should be designed
208 with robust mechanisms to detect and mitigate data leakage, ensuring the integrity of the evaluation
209 process for future LLMs. Collaborative efforts between researchers, dataset creators, and LLM
210 developers can play a pivotal role in achieving this goal.

211 **References**

- 212 [1] Anthropic. Claude, 2023.
- 213 [2] Terra Blevins and Luke Zettlemoyer. Language contamination helps explains the cross-lingual
214 capabilities of English pretrained models. In *Proceedings of the 2022 Conference on Empirical
215 Methods in Natural Language Processing*, pages 3563–3574, Abu Dhabi, United Arab Emirates,
216 December 2022. Association for Computational Linguistics.
- 217 [3] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal,
218 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel
219 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M.
220 Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz
221 Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec
222 Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020.
- 223 [4] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam
224 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker
225 Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes,
226 Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson,
227 Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin,
228 Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier
229 Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David
230 Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani
231 Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat,
232 Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei
233 Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei,
234 Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. Palm: Scaling
235 language modeling with pathways, 2022.
- 236 [5] Shahriar Golchin and Mihai Surdeanu. Time travel in llms: Tracing data contamination in large
237 language models, 2023.
- 238 [6] Google. Bard, Febraury 2023.
- 239 [7] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and
240 Jacob Steinhardt. Measuring massive multitask language understanding, 2021.
- 241 [8] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza
242 Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom
243 Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia
244 Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent
245 Sifre. Training compute-optimal large language models, 2022.
- 246 [9] Alon Jacovi, Avi Caciularu, Omer Goldman, and Yoav Goldberg. Stop uploading test data in
247 plain text: Practical strategies for mitigating data contamination by evaluation benchmarks,
248 2023.
- 249 [10] Yuanzhi Li, Sébastien Bubeck, Ronen Eldan, Allie Del Giorno, Suriya Gunasekar, and Yin Tat
250 Lee. Textbooks are all you need ii: phi-1.5 technical report, 2023.
- 251 [11] Yucheng Li. Estimating contamination via perplexity: Quantifying memorisation in language
252 model evaluation, 2023.
- 253 [12] Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summariza-
254 tion Branches Out*, pages 74–81, Barcelona, Spain, July 2004. Association for Computational
255 Linguistics.
- 256 [13] Kevin Lin, Chung-Ching Lin, Lin Liang, Zicheng Liu, and Lijuan Wang. Mpt: Mesh pre-training
257 with transformers for human pose and mesh reconstruction, 2023.
- 258 [14] Stephanie Lin, Jacob Hilton, and Owain Evans. TruthfulQA: Measuring how models mimic
259 human falsehoods. In *Proceedings of the 60th Annual Meeting of the Association for Compu-
260 tational Linguistics (Volume 1: Long Papers)*, pages 3214–3252, Dublin, Ireland, May 2022.
261 Association for Computational Linguistics.

- 262 [15] Lingjie Mei, Jiayuan Mao, Ziqi Wang, Chuang Gan, and Joshua B. Tenenbaum. Falcon: Fast
263 visual concept learning by integrating images, linguistic descriptions, and conceptual relations,
264 2022.
- 265 [16] OpenAI. Gpt-4 technical report, 2023.
- 266 [17] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin,
267 Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton,
268 Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano,
269 Jan Leike, and Ryan Lowe. Training language models to follow instructions with human
270 feedback, 2022.
- 271 [18] Rylan Schaeffer. Pretraining on the test set is all you need, 2023.
- 272 [19] Thibault Sellam, Dipanjan Das, and Ankur P. Parikh. Bleurt: Learning robust metrics for text
273 generation, 2020.
- 274 [20] Sheng Shen, Le Hou, Yanqi Zhou, Nan Du, Shayne Longpre, Jason Wei, Hyung Won Chung,
275 Barret Zoph, William Fedus, Xinyun Chen, Tu Vu, Yuexin Wu, Wuyang Chen, Albert Webson,
276 Yunxuan Li, Vincent Zhao, Hongkun Yu, Kurt Keutzer, Trevor Darrell, and Denny Zhou.
277 Mixture-of-experts meets instruction tuning:a winning combination for large language models,
278 2023.
- 279 [21] Kristina Toutanova and Christopher D. Manning. Enriching the knowledge sources used in
280 a maximum entropy part-of-speech tagger. In *2000 Joint SIGDAT Conference on Empirical
281 Methods in Natural Language Processing and Very Large Corpora*, pages 63–70, Hong Kong,
282 China, October 2000. Association for Computational Linguistics.
- 283 [22] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timo-
284 thée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez,
285 Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation
286 language models, 2023.
- 287 [23] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei,
288 Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas
289 Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes,
290 Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony
291 Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian
292 Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut
293 Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov,
294 Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta,
295 Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiao-
296 qing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng
297 Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien
298 Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation
299 and fine-tuned chat models, 2023.
- 300 [24] Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan
301 Du, Andrew M. Dai, and Quoc V. Le. Finetuned language models are zero-shot learners, 2022.