

Leak, Cheat, Repeat: Data Contamination and Evaluation Malpractices in Closed-Source LLMs

Anonymous ACL submission

Abstract

Natural Language Processing (NLP) research is becoming increasingly focused on the use of Large Language Models (LLMs), with some of the most popular ones being either fully or partially closed-source. The lack of access to model details, especially regarding training data, has repeatedly raised concerns about data contamination among researchers. Several attempts have been made to address this issue, but they are limited to anecdotal evidence and trial and error. Additionally, they overlook the problem of indirect data leaking, where models are iteratively improved by using data coming from users. In this work, we conduct the first systematic review of work using OpenAI’s ChatGPT and GPT-4, the most prominently used LLMs today, in the context of data contamination. By analysing 255 papers and considering OpenAI’s data usage policy, we extensively document how much data has been leaked to ChatGPT in the first year after the model’s release. At the same time, we document a number of evaluation malpractices emerging in the reviewed papers, including unfair or missing baseline comparisons, reproducibility issues, and authors’ lack of awareness of the data usage policy. Our work provides the first quantification of the ChatGPT data leakage problem.

1 Introduction

The recent emergence of large language models (LLMs) that show remarkable performance on a wide range of tasks has led not only to a dramatic increase in their use in research but also to a growing number of companies joining the race for the biggest and most powerful models. In pursuing a competitive advantage, many of the most popular LLMs today are locked behind API access and their details are unknown (OpenAI, 2023a; Thoppilan et al., 2022; Touvron et al., 2023). This includes model weights (OpenAI, 2023a), training data (Piktus et al., 2023), or infrastructural details to assess model carbon footprint (Lacoste et al., 2019).

In particular, the lack of information about the training data being used raises important questions about the credibility of LLM performance evaluation. The training data, typically collected automatically by scraping documents from the web, may contain training, validation, and – most critically – test sets of standard NLP benchmarks, leading to unintended evaluation of the LLM’s performance on data it was trained on. This phenomenon, known as data contamination, may not be an issue in the general use of commercial LLMs, where adherence to research principles is not mandatory, but it becomes a serious problem when these tools are widely used and evaluated in research.

Unfortunately, closed-source models are typically locked behind inference-only APIs, which makes detecting data contamination very complex. Because of this, existing work on the matter is limited to detecting extreme forms of overfitting and memorization, where the LLM is mostly probed to reproduce known examples from benchmarks verbatim. These approaches also overlook that recent LLMs get iteratively improved with user interaction data. When users test these models on benchmarks, they may, in fact, be exposing them to the data even if it was not part of the original training data, a phenomenon that we refer to as *indirect data leaking*.

In this paper, we address the issue of data contamination in closed-source LLMs by conducting a systematic literature review. We review 255 papers and carefully detail data leakage emerging from them. We focus primarily on OpenAI’s ChatGPT,¹ and secondarily on GPT-4² as these are the most frequently used commercial LLMs in NLP research. By considering OpenAI’s data usage policy, we assess how much data was reported to be leaked to them in such a way that it could be

¹<https://openai.com/blog/chatgpt>

²<https://openai.com/gpt-4>

used for further training, hence giving the models an unfair advantage over open-source LLMs that are usually trained following scientific rigor. We also report a series of emergent evaluation malpractices, including lack of comparison with other models/approaches, differences in the evaluation scale (e.g. closed-source LLMs being evaluated on samples while compared open models are evaluated over an entire benchmark), lack of code and data access, or data leakage even in situations where it could be avoided. To our knowledge, ours is the most comprehensive and extensive quantification of the data leakage issue in LLMs.

In short, our contributions are as follows:

- (1) We perform a systematic review of the work evaluating OpenAI’s ChatGPT and GPT-4 on a variety of tasks in NLP and other domains (Section 4).
- (2) For each work, we estimate the amount of data leaked to these models in such a way that it could be used for further model training (Section 5.1).
- (3) We reveal some critical malpractices in closed-source LLM research, affecting both evaluation reproducibility and fairness (Sections 5.2 and 5.3).
- (4) Based on our findings, we propose a list of best practices for the evaluation of closed-source LLMs (Section 6).

Furthermore, we believe that the results of our work can help the ongoing efforts in empirical inspection of LLM data contamination by pointing out which datasets are worth investigating. We release our survey results as a spreadsheet with an assessment of the extent of data leakage for each dataset to stimulate further research.³

2 Prior Work on LLM Data Contamination

Work on LLMs data contamination traces back to OpenAI’s GPT-3 (Brown et al., 2020; Raffel et al., 2020; Magar and Schwartz, 2022), one of the first LLMs with general API access for which only partial training information was released. Despite results hinting at the presence of significant data contamination (Raffel et al., 2020; Magar and

Schwartz, 2022), the model has been used extensively in research and the issue was rarely taken into account when interpreting its performance. With the release of ChatGPT and following closed-source models to general public,⁴ the data contamination topic became an even more pressing issue.

Because of the implicit complexity of analyzing closed-source LLMs, few practical approaches have been proposed to estimate the leak of known benchmarks into them. One notable example is the LLM Contamination Index,⁵ which features a regularly updated list of models along with the extent of their eventual contamination/memorization on popular NLP benchmarks. This approach works by zero-shot prompting the model to generate instances from specific datasets, providing details on the required split and format (Sainz et al.). The premise is that no model should be able to replicate specific benchmark formats without having seen them first.

Recent work (Aiyappa et al., 2023a) also discussed the issue by comparing different ChatGPT versions on known benchmarks, noting that performance on the task improved shortly after the model was exposed to the datasets by previous work (Zhang et al., 2022). More applied approaches have been proposed recently (Golchin and Surdeanu, 2023), where the model is prompted to complete a given sentence coming from a known benchmark. The completion is then compared with the original reference through text overlap metrics and a statistical test is used to assess if the model is contaminated.

Although these preliminary works exploring the task of data contamination detection are promising, they cannot be fully trusted and have some limitations. Most importantly, they are based on an assessment of the model’s ability to generate an example from the benchmark. The recall of such methods can be affected by two issues:

- (1) Some closed-source models have incorporated special filters into their decoding algorithms that prevent them from generating texts that significantly overlap with their training sets (GitHub, 2022; Ippolito et al., 2023). This creates an additional noise for the detection methods and results in the lack of confidence

³https://anonymous.4open.science/r/eacl-data_contamination-BEF1/

⁴Including GPT-4 (OpenAI, 2023a), Google’s LaMDA (Thoppilan et al., 2022) and PaLM (Chowdhery et al., 2022), Cohere’s Command and Anthropic’s Claude.

⁵<https://hitz-zentroa.github.io/lm-contamination/>

173	that even the datasets tested negative for data	samples are perturbed and hence no longer an exact	221
174	leakage are not present in LLM training data.	match of their original version. Therefore, it is un-	222
175	(2) Such approaches can only detect the most extreme	likely that LLM vendors could effectively exclude	223
176	form of overfitting which results in (almost) complete	leaked benchmarks from further model fine-tuning,	224
177	memorization of data samples by the model. However,	especially at scale.	225
178	even a regular adjustment of the model by training	For (2), it would be necessary to understand how	226
179	on the leaked data, which does not necessarily lead	the LLM vendor uses the data to improve the model.	227
180	to its memorization, poses a problem for fair compar-	A very likely scenario is continued pre-training,	228
181	isons.	where the data leaked by users is treated as an	229
182		in-domain corpus (and thus given more influence	230
183	3 The Issue of Indirect Data Leaking	than pretraining data). This procedure is known	231
184	The related work presented in Section 2 approaches	to improve models' performances in the leaked	232
185	the issue of data contamination mainly by back-	domains (Gururangan et al., 2020). Notably, Shi	233
186	tracking models' training data. It is commonly	and Lipani (2023) find that fine-tuning a model	234
187	assumed that the use of benchmarks available only	on in-domain text enriched by textual instructions	235
188	to authorised parties, or datasets being constructed	leads to an increase in the model performance even	236
189	after the ChatGPT release, is a guarantee that they	if gold labels are not shown to the model. Such a	237
190	have not been leaked. This ignores the fact that	setup perfectly matches the kind of data shown to	238
191	models using reinforcement learning from human	chat LLMs when evaluated by researchers. This	239
192	feedback (RLHF) (Ouyang et al., 2022), such as	means that closed-source LLMs such as ChatGPT	240
193	ChatGPT, are subject to repeated updates (Aiyappa	can make use of these gold standard examples from	241
194	et al., 2023a) with training data partially coming	widely used NLP benchmarks to gain an unfair	242
195	from user interactions. The collection of user data	advantage over other models.	243
196	and online model updates lead to a previously over-	With this motivation, we conduct a systematic	244
197	looked phenomenon of indirect data leaking. This	review to quantify how much of such data ChatGPT	245
198	is a new development of the issue for two main	could have obtained.	246
199	reasons:	4 Methodology	247
200	(1) Unlike plain text scraped from the internet,	Following the standard protocol for a systematic	248
201	data from users might be harder to inspect	work review from the medical domain (Khan et al.,	249
202	for contamination as it might involve model	2003) we organize our work into five macro-steps,	250
203	prompts, textual alterations, or truncation of	corresponding to the following subsections.	251
204	benchmark samples.	4.1 Framing questions	252
205	(2) Users supply the data along with instructions	In reviewing the existing work on evaluation of	253
206	on how to perform the task. In LLMs, this can	ChatGPT and GPT-4, we pose the following re-	254
207	be considered a novel form of gold-standard	search questions:	255
208	data for continued training, even in the ab-	(1) Which datasets have been demonstrably	256
209	sence of target labels. Model updates on such	leaked to ChatGPT and GPT-4 during the last	257
210	data are likely much more effective than plain	year?	258
211	in-domain text.	(2) Do all works evaluating these models include	259
212	For (1), even with a conscious and targeted ef-	a fair comparison with existing baselines?	260
213	fort of the LLM vendors to avoid fine-tuning the	4.2 Identifying relevant work	261
214	model on test data and benchmarks, this issue is	We employ commonly used online databases ⁶ and	262
215	particularly complex to trace. When evaluating	major NLP conferences proceedings (including	263
216	a closed-source LLM, users often feed the model	ACL, NAACL, EMNLP, NeurIPS), considering	264
217	with test-set samples (with or without labels) sur-		
218	rounded by additional text, such as instructions in		
219	the form of prompts. In some cases, especially		
220	when evaluating the LLM robustness, the test-set		

⁶We query Google Scholar, Semantic Scholar, DBLP, arXiv, ACL Anthology.

both peer-reviewed work and pre-prints, as the interaction with LLMs happened regardless of publication status. We filter our queries on work containing the terms “ChatGPT”, “evaluation”, “large language models” and “AI” either in title, abstract, body, or all of them.

We also do not limit our search to computer science works only, as ChatGPT has been investigated by researchers from many other domains, e.g. healthcare (Kung et al., 2023a), psychology (Cai et al., 2023) and education (Szefer and Deshpande, 2023). Since ChatGPT is our primary focus, we limit our search to works between late November 2022 (when the model was publicly released) and early October 2023. Among all the papers, we first do a preliminary screening, assessing if they effectively run ChatGPT or GPT-4 in any form.⁷

4.3 Assessing quality and relevance

To assess which work effectively leaked data to ChatGPT or GPT-4, we refer to OpenAI’s data usage policy,⁸ stating that no text sent through or produced from API calls is used for model improvements from 1 March 2023 onwards (coinciding with the first ChatGPT API release). Therefore, only the work interacting with the models through the web interface⁹ is considered to leak data.

A small number of works used both the web interface and API access.¹⁰ We carefully review such works to calculate which portion of the data was used in the former setup. We drew our conclusions from the paper draft history on arXiv; in some cases, this information was also transparently disclosed by the authors. In the case of work with multiple drafts dating prior to the model release in November 2022, we consider the earliest draft that includes ChatGPT or GPT-4 for the calculation.

4.4 Summarizing the evidence

We inspect each surveyed paper, looking for information on the used datasets, split, and number of samples. If no mention of sampling or similar information is made, we assume that the whole dataset has been used. Similarly, if no information on the used split is provided, we assume that the authors treated the dataset as a whole. It could be

⁷We encountered a small number of papers also comparing to other closed-source LLMs, such as Anthropic’s Claude.

⁸<https://help.openai.com/en/articles/5722486-how-your-data-is-used-to-improve-model-performance>

⁹<https://chat.openai.com/>

¹⁰Their experiments began prior to March 1st, 2023 and the authors started using the API soon after it was released.

argued that feeding entire datasets to ChatGPT or GPT-4 is unrealistic because of the usage restrictions imposed by OpenAI on the web interface, and the amount of work necessary for manually inputting the data inside the chat. However, we note that quickly after ChatGPT release, many unofficial wrappers have been developed¹¹ for circumventing said issues, most of which are still in active use. We also point out that many of the papers we surveyed mentioned the use of such tools explicitly.

We also track secondary information relevant to the evaluation – for each work, we inspect: (1) if it has been peer-reviewed¹²; (2) if the used prompts are available; (3) if a repository to reproduce the experiment is provided; (4) if the authors used a whole dataset or a sample; (5) if ChatGPT or GPT-4 were compared to other open models/approaches and if the evaluation scale was the same; (6) if the version of the model used is reported.

4.5 Interpreting the findings

We report the results of our review both quantitatively and qualitatively. Specifically, we report the number of works surveyed leaking data to ChatGPT or GPT-4 in such a way that it can be used by OpenAI to further improve the model (according to their data policy). We also document a series of evaluation practices emerging for the work reviewed that is problematic with respect to objectiveness and reproducibility. Finally, drawing upon our results, we present a series of best practices for researchers evaluating OpenAI’s and other closed-source LLMs.

5 Results

Following our methodology, in the first step we identified 255 research papers, 216 of which were found relevant during the initial screening. Among the relevant papers, 72 (~ 33%) were peer-reviewed while the remainder (144) consisted of pre-prints. We subsequently analysed the retrieved papers to examine the problem of data contamination and the adopted evaluation practices.

5.1 Indirect data contamination

From our analysis, 90 papers (~ 42%) accessed ChatGPT or GPT-4 through the web interface, hence providing data that OpenAI could have used

¹¹E.g. `revChatGPT`, `PyChatGPT`, and `ChatGPT-to-API`.

¹²We do note that part of the work we reviewed is very recent (with pre-prints published in the last 3 months) and might currently be under review.

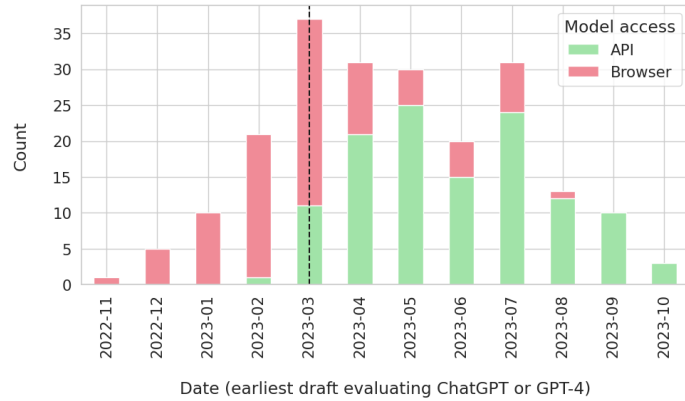


Figure 1: Distribution of the dates when papers evaluating ChatGPT or GPT-4 were first uploaded to arXiv or published. The dotted line represents the ChatGPT API release (March 1st, 2023, dotted line in the chart) as a cutoff point. The single paper shown using the API in February is by a research group who reported having early API access.

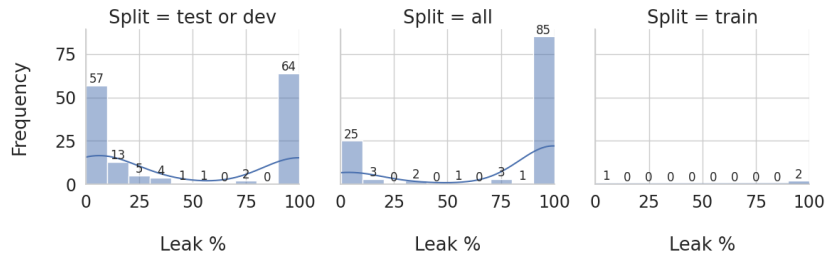


Figure 2: Overview of data leakage distribution. We report the number of times (y) we observed a specific percentage of leaking (x) for the considered split. As some work vaguely describes the used split as “test or dev set”, we merge these two values in a unique chart.

to further improve its models. We note that while it is possible to opt out of providing the data for model improvement purposes, we found no evidence suggesting any of the surveyed papers did so.

We first inspected the time distribution of the reviewed works (Figure 1) to gain insight into when most data leaks happened. Unsurprisingly, the majority of the papers leaking data dates before the official release of ChatGPT API, and it can be seen that web interface access rapidly decreased following March 2023. However, we must note that (1) a considerable amount of work kept using the web interface to access ChatGPT until September 2023 and (2) our analysis cannot inspect the preliminary stages of prompt engineering, which are rarely reported and might still be done through the web interface because of its trial-and-error nature.

The presence of leaked data after the API release may indicate that a part of the research community is either unaware of OpenAI’s data policy, or does not consider it a problem when conducting experiments. Many works also report on using the web interface for cost reasons, as it allows free Chat-

GPT access. This might be crucial, especially for small case studies.

As a second step, we quantified leak severity according to the datasets and splits used. For papers specifying the amount of data used or providing a repository with the information, we considered the given value. For the rest, we calculated the value by inspecting the actual datasets.¹³ In seven papers, no number of samples used was specified, so we contacted the authors for clarification. In the two cases where the authors did not respond, we assumed the entire split of a dataset was used. We calculated both the number of instances and the percentage of the considered split (or the whole dataset when applicable).

Since only a small number of datasets (18) was used in multiple papers, we always assume that the largest leak for a given dataset is a superset of all smaller ones.¹⁴

¹³We mainly use [HuggingFace Datasets](#), but also refer to [Kaggle](#) or other sources based on availability.

¹⁴We also tried a pessimistic approach, where we assumed all the leaks were independent, but due to the small number of works covering the same data, the results are virtually

Task name	Lo	M-Lo	M-Hi	Hi
AI safety & ethics	0	0	2	0
Creative NLG	1	0	0	0
Dialogue	2	1	0	5
NLG evaluation	0	0	0	4
Machine Translation	6	4	1	1
Math	0	1	0	8
Natural language generation	2	1	0	14
Natural language inference	6	2	0	15
Language understanding	0	0	0	2
Paraphrasing	2	0	0	0
Politics	0	1	0	3
Programming	0	0	0	1
Psychology	0	0	0	1
Question answering	24	14	5	31
Commonsense reasoning	3	4	0	9
Semantic similarity	2	1	0	3
Sentiment analysis	8	9	1	8
Summarization	5	6	1	1
Text classification	1	0	0	3
Text extraction	2	1	0	7

Table 1: The number of datasets with low (Lo), moderate-low (M-Lo), moderate-high (M-Hi) and high leaks (Hi) is reported for each task, omitting custom datasets. A more detailed table, including specific dataset names, is provided in the Appendix.

Our calculations show that the 90 papers leaked data from 263 unique datasets, for a total of over 4.7M samples.¹⁵ We find most samples ($\sim 93.8\%$) coming from datasets treated as whole (with no split), followed by test and development ($\sim 5.6\%$),¹⁶ and training ($\sim 0.6\%$) sets. In line with what we discussed in Section 3, ChatGPT was presented with millions of samples with instructions that could be considered de-facto novel gold-standard data in some cases.

We also report that several works included the examples’ labels when few-shot prompting ChatGPT or when using ChatGPT or GPT-4 as a reference-based evaluation metric. We consider this the worst possible case of data leaking, as it gives the model information about the desired output as well.

To classify leak severity, we examine the frequency distribution of leak sizes (Figure 2). It appears that most works either leak full splits or very small samples, with only a few works leaking intermediate amounts. With this information, we classify a portion of leaked data as *low* ($< 5\%$), *moderate-low* ($5 - 50\%$), *moderate-high* ($50 - 95\%$), or *high* ($> 95\%$).

Consequently, we categorize all leaked datasets

identical.

¹⁵The survey total is 4,714,753 leaked samples.

¹⁶As some work vaguely describes the used split as “test or dev set”, we merge these two values.

into these 4 thresholds. Overall, we find a low leak for 66 ($\sim 25\%$) datasets, moderate-low for 47 ($\sim 18\%$), moderate-high for 10 ($\sim 4\%$) and high for 142 ($\sim 53\%$). This result is particularly worrying as the majority of datasets were nearly fully leaked.

Finally, we inspect which NLP tasks are covered by the leaked data (Table 1). We find that the tasks suffering the most from high leaks are natural language inference, question answering, and natural language generation. These and other tasks (see Tables 4 and 5 in the Appendix) include a lot of very popular NLP benchmarks, as well as high-quality custom datasets created ad-hoc for individual evaluations. The custom datasets were frequently phrased as an exam in a field different from NLP, e.g. medicine, physics, psychology, or law. Other custom datasets explored, for example, the LLMs’ sense of humour, philosophical and political leaning, or bias.

5.2 Reproducibility

We assessed the evaluations’ reproducibility by checking whether the prompts used to query ChatGPT were provided, whether a repository containing data or code was available, and whether the datasets used were custom-made. Finally, we also checked for sampling of the original data or other practices that make it impossible to exactly reconstruct the data used. Results are shown in Table 2.

196 (91%) works report the prompts used to convert data into a query and possibly to instruct the model on how to perform a given task. The number of works providing a code repository is significantly smaller, at 115 (53%). This figure excludes papers that provided a link to a non-existent or empty repository. Overall, 73 (51%) of the pre-prints and 34 (47%) peer-reviewed papers provided both prompts and a repository.

Another obstacle to reproducibility is that most closed-source LLMs are being regularly updated. Therefore, it is crucial for experiment reproducibility to report a specific LLM version. In the surveyed works, this was generally done by reporting the running period of the experiments when using the web interface, or by reporting which version of the model has been accessed via the API. Unfortunately, as the concept of regular model updates is relatively new, this practice is not yet common. Only 29 (40%) of the peer-reviewed papers and 33 (23%) of the pre-prints provide this information. We consider reporting the running period to be par-

Prompts	Repo	Sampl.	Custom	n. (%)
				3 (2.11%)
			✓	1 (0.70%)
		✓		8 (5.63%)
	✓			3 (2.11%)
	✓	✓		2 (1.41%)
✓				20 (14.08%)
✓			✓	3 (2.11%)
✓		✓		27 (19.01%)
✓		✓	✓	3 (2.11%)
✓	✓			37 (26.06%)
✓	✓		✓	4 (2.82%)
✓	✓	✓		27 (19.01%)
✓	✓	✓	✓	4 (2.82%)

(a) Pre-prints

Prompts	Repo	Sampl.	Custom	n. (%)
				1 (1.43%)
	✓		✓	1 (1.43%)
	✓	✓		1 (1.43%)
✓				14 (20.00%)
✓			✓	7 (10.00%)
✓		✓		9 (12.86%)
✓		✓	✓	3 (4.29%)
✓	✓			8 (11.43%)
✓	✓		✓	4 (5.71%)
✓	✓	✓		16 (22.86%)
✓	✓	✓	✓	6 (6.57%)

(b) Peer-reviewed works

Table 2: Statistics related to the reproducibility of the work reviewed: the availability of used prompts (Prompts) and code/data repository (Repo), the usage of custom datasets (Custom), the application of random sampling or any other practice that does not allow the exact reconstruction of the data used (Sampl.).

Comp.	Scale	n. (%)
		71 (50.00%)
✓		54 (38.03%)
✓	✓	17 (11.97%)

(a) Pre-prints

Comp.	Scale	n. (%)
		30 (42.86%)
✓		34 (48.57%)
✓	✓	6 (8.57%)

(b) Peer-reviewed works

Table 3: Fairness statistics for reviewed work. Statistics related to the practices of performance comparisons between ChatGPT/GPT-4 and other open models: whether such comparisons are performed at all (Comp.) and whether they are of the same scale (Scale).

472 ticularly important, as [Chen et al. \(2023b\)](#) show
473 that there may be vast differences between model
474 versions.

475 5.3 Evaluation (mal)practices

476 We found that the evaluation of ChatGPT’s perfor-
477 mance is often conducted without comparing it to
478 any open-source LLM or non-LLM-based method
479 (see Table 3). This is similarly prevalent regardless
480 of the publication status, appearing in 71 (50%)
481 of pre-prints and 30 (43%) of published papers.
482 In addition, among the works that perform a compar-
483 ison with open models, 23 (21%) compare the
484 results computed on different samples. ChatGPT
485 is typically evaluated on a random sample of the
486 benchmark while other models are compared on its
487 entirety. In many works, ChatGPT’s performance

488 is measured on only a handful (10-50) of examples,
489 which substantially lowers the expressive power of
490 the comparison. For instance, considering a sim-
491 plistic case with binary assessment of model output
492 (correct/incorrect) on 10 examples, the difference
493 should be more than 30% to be statistically signifi-
494 cant,¹⁷ which is rarely seen. Statistical analysis of
495 results is almost never performed.

496 Another practice of some concern is the way in
497 which the size of the evaluation data is reported,
498 especially when sampling is used. The papers of-
499 ten show the size of the whole evaluation dataset
500 upfront (e.g. in a table or in the dataset description
501 section), but they report the actual sample sizes

¹⁷Assuming Fisher’s exact test, typical $\alpha = 5\%$ and moder-
ate model performance around $\hat{p} = 0.5$

used for evaluation only later and in a less obvious way (in footnotes, limitations sections, or appendices). This practice makes the experimental results harder to interpret.

6 Best Practices in Closed-source LLM Evaluation

Our survey revealed both a significant amount of data leakage in ChatGPT and many worrying trends in its evaluation. In light of this, we list a series of best practices that we believe could help mitigate the issues. We believe that researchers looking to objectively evaluate LLMs today should:

Choose the right model access The first step when planning closed LLMs evaluation should be reading their most up-to-date data policies, and access models accordingly (e.g. API instead of web interface for OpenAI’s LLMs). We also acknowledge that in some cases this might not be viable due to budget limits, or an overly steep learning curve for the use of APIs by researchers outside of computer science.¹⁸

Interpret performance skeptically The lack of system specifications, training data information, and other details can make commercial LLMs look like incredibly powerful tools with impressive zero-shot performance. Aiyappa et al. (2023a) already pointed out that this can often be explained by data contamination. In our review, we documented over 4 million samples across over 200 NLP datasets having been leaked to these models. The performance of closed-source LLMs should always be interpreted while keeping this in mind.

When possible, avoid using closed-source models We strongly encourage using the available open-source LLMs. While there has been discussion in the research community about commercial models being consistently better than open-source ones, we note that (1) this is often driven by hype, while there is evidence of the opposite (Kocoń et al., 2023), (2) research done solely on closed LLMs limits scientific progress, bringing benefits mainly to the LLM vendors (3) LLM vendors can arbitrarily make changes to the models, e.g., making previous versions unavailable, changing their behavior in a way that may not be visible to the user (Chen et al., 2023b) or changing the data treatment policy.

¹⁸In such a case, at the time of writing this paper in mid-October, OpenAI allows users to opt out of providing data for model improvement by using a [Google Form](#).

Adopt a fair and objective comparison Evaluating closed-source LLMs is tied to comparing them with pre-existing approaches. Evaluating LLMs on a limited number of samples while evaluating others on dramatically larger sets is scientifically dubious at best. When sampling is required (for example for budgetary restrictions), it should be applied to all the considered approaches. We also discourage taking state-of-the-art values directly from previous work and suggest to re-run all approaches on the considered data only.

Make the evaluation reproducible In light of the known NLP evaluation reproducibility crisis (Belz et al., 2023) we strongly encourage researchers to report as many details about their setup. Besides the commonly required information such as random seeds, open model parameters, etc., we note that when the evaluation involves closed models, additional details should be disclosed. Prompts, as well as the process leading to them, should be detailed since LLMs are very sensitive to even minor changes in prompts (Lu et al., 2022). The model version and experiment running period should be mentioned as well so that further researchers can use the same model checkpoint if possible. Data, especially if sampled, should be released (ideally in a repository) to avoid potential differences in sampling.

7 Conclusion and Future Work

In this review, we present our findings based on a survey of 255 papers evaluating the performance of ChatGPT and GPT-4. We investigate the problem of indirect data contamination and report that 4.7M samples coming from 263 distinct datasets have been leaked to these models. We also report concerning research practices with respect to reproducibility and fairness. Finally, informed by our survey, we suggest best practices for the evaluation of closed-source LLMs.

Future Work In our future work, we aim to run experiments via the OpenAI API to see the impact of leaked test data on the performance of ChatGPT and GPT-4 on the leaked datasets and the tasks in general.

Furthermore, we consider investigating indirect data leakage in other closed-source models, namely from Anthropic or Cohere, which appeared in a small number of papers reviewed in this work.

8 Limitations

We are aware the list of contaminated datasets we compiled in our work is not fully conclusive for one of several reasons:

- (1) We review the information that has been publicly revealed via articles. We postulate more experiments could have revealed test set data to closed-source models but were never published.
- (2) In this paper, we focus on the works that use ChatGPT or GPT-4. However, prior to March 1st, 2023, OpenAI’s policy stated that they may also use data from the API to improve their models. This would imply that data sent to GPT-3 via the API could have been used for training.
- (3) The number of papers investigating the performance of ChatGPT is vast, and despite our best efforts, we could have missed some works.
- (4) Information on whether individual works are pre-prints or published is given at the time of writing (early October 2023). This is subject to change, especially given the freshness of many of the works reviewed.
- (5) Many datasets released prior to 2021 could have been fully leaked by being a part of the models’ pre-training data.

As mentioned in Section 4, in some cases the papers were not clear about some aspects of the experiments. We contacted the authors of such papers for clarification, however, two of them did not respond. For these papers, our best-judgment assumptions may be wrong.

References

Kabir Ahuja, Harshita Diddee, Rishav Hada, Millient Ochieng, Krithika Ramesh, Prachi Jain, Akshay Nambi, Tanuja Ganu, Sameer Segal, Maxamed Axmed, Kalika Bali, and Sunayana Sitaram. 2023. [Mega: Multilingual evaluation of generative ai](#).

Rachith Aiyappa, Jisun An, Haewoon Kwak, and Yongyeol Ahn. 2023a. [Can we trust the evaluation on ChatGPT?](#) In *Proceedings of the 3rd Workshop on Trustworthy Natural Language Processing (TrustNLP 2023)*, pages 47–54, Toronto, Canada. Association for Computational Linguistics.

Rachith Aiyappa, Jisun An, Haewoon Kwak, and Yongyeol Ahn. 2023b. [Can we trust the evaluation on chatgpt?](#)

Afra Feyza Akyurek, Ekin Akyurek, Ashwin Kalyan, Peter Clark, Derry Tanti Wijaya, and Niket Tandon. 2023. [RL4F: Generating natural language feedback with reinforcement learning for repairing model outputs](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7716–7733, Toronto, Canada. Association for Computational Linguistics.

Mostafa M. Amin, Erik Cambria, and Björn W. Schuller. 2023. [Will affective computing emerge from foundation models and general ai? a first evaluation on chatgpt](#).

Chenxin An, Shansan Gong, Ming Zhong, Xingjian Zhao, Mukai Li, Jun Zhang, Lingpeng Kong, and Xipeng Qiu. 2023. [L-eval: Instituting standardized evaluation for long context language models](#).

Fares Antaki, Samir Touma, Daniel Milad, Jonathan El-Khoury, and Renaud Duval. 2023. [Evaluating the performance of chatgpt in ophthalmology: An analysis of its successes and shortcomings](#). *Ophthalmology Science*, 3(4):100324.

Joris Baan, Nico Daheim, Evgenia Iliia, Dennis Ulmer, Haau-Sing Li, Raquel Fernández, Barbara Plank, Rico Sennrich, Chrysoula Zerva, and Wilker Aziz. 2023. [Uncertainty in natural language generation: From theory to applications](#).

Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. 2023a. [Longbench: A bilingual, multitask benchmark for long context understanding](#).

Yushi Bai, Jiahao Ying, Yixin Cao, Xin Lv, Yuze He, Xiaozhi Wang, Jifan Yu, Kaisheng Zeng, Yijia Xiao, Haozhe Lyu, Jiayin Zhang, Juanzi Li, and Lei Hou. 2023b. [Benchmarking foundation models with language-model-as-an-examiner](#).

Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, Quyet V. Do, Yan Xu, and Pascale Fung. 2023. [A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity](#).

Jonas Belouadi and Steffen Eger. 2023. [Bygpt5: End-to-end style-conditioned poetry generation with token-free language models](#).

Anya Belz, Craig Thomson, Ehud Reiter, and Simon Mille. 2023. [Non-repeatable experiments and non-reproducible results: The reproducibility crisis in human evaluation in NLP](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 3676–3687, Toronto, Canada. Association for Computational Linguistics.

697	Ning Bian, Xianpei Han, Le Sun, Hongyu Lin, Yaojie Lu, and Ben He. 2023. Chatgpt is a knowledgeable but inexperienced solver: An investigation of commonsense problem in large language models.	Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramèr, and Chiyuan Zhang. 2023. Quantifying memorization across neural language models.	751
698			752
699			753
700			754
701	Sebastian Bordt and Ulrike von Luxburg. 2023. Chatgpt participates in a computer science exam.	Mohna Chakraborty, Adithya Kulkarni, and Qi Li. 2023. Zero-shot approach to overcome perturbation sensitivity of prompts. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 5698–5711, Toronto, Canada. Association for Computational Linguistics.	755
702			756
703	Ali Borji. 2023. A categorical archive of chatgpt failures.		757
704			758
705	Ritwik Bose, Ian Perera, and Bonnie Dorr. 2023. Detoxifying online discourse: A guided response generation approach for reducing toxicity in user-generated text. In <i>Proceedings of the First Workshop on Social Influence in Conversations (SICoN 2023)</i> , pages 9–14, Toronto, Canada. Association for Computational Linguistics.	Shreya Chandrasekhar, Chieh-Yang Huang, and Ting-Hao Huang. 2023. Good data, large data, or no data? comparing three approaches in developing research aspect classifiers for biomedical papers. In <i>The 22nd Workshop on Biomedical Natural Language Processing and BioNLP Shared Tasks</i> , pages 103–113, Toronto, Canada. Association for Computational Linguistics.	759
706			760
707			761
708			762
709			763
710			764
711			765
712	Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In <i>Advances in Neural Information Processing Systems</i> , volume 33, pages 1877–1901. Curran Associates, Inc.		766
713			767
714			768
715			769
716			770
717			771
718			772
719			773
720			774
721			775
722			776
723			777
724			778
725			779
726	Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro, and Yi Zhang. 2023. Sparks of artificial general intelligence: Early experiments with gpt-4.	Jiaao Chen, Xiaoman Pan, Dian Yu, Kaiqiang Song, Xiaoyang Wang, Dong Yu, and Jianshu Chen. 2023a. Skills-in-context prompting: Unlocking compositionality in large language models.	780
727			781
728			782
729			783
730			784
731			785
732	Ana-Maria Bucur. 2023. Utilizing chatgpt generated data to retrieve depression symptoms from social media.	Lingjiao Chen, Matei Zaharia, and James Zou. 2023b. How is chatgpt’s behavior changing over time? <i>arXiv preprint arXiv:2307.09009</i> .	786
733			787
734			788
735	Laura Cabello, Jiaang Li, and Ilias Chalkidis. 2023. Pokemonchat: Auditing chatgpt for pokémon universe knowledge.	Nuo Chen, Yan Wang, Haiyun Jiang, Deng Cai, Yuhan Li, Ziyang Chen, Longyue Wang, and Jia Li. 2023c. Large language models meet harry potter: A bilingual dataset for aligning dialogue agents with characters.	789
736			790
737			791
738	Alex Cabrera and Graham Neubig. 2023. Zeno chatbot report.	Xuanting Chen, Junjie Ye, Can Zu, Nuo Xu, Rui Zheng, Minlong Peng, Jie Zhou, Tao Gui, Qi Zhang, and Xuanjing Huang. 2023d. How robust is gpt-3.5 to predecessors? a comprehensive study on language understanding tasks.	792
739			793
740	Zhenguang G Cai, David A Haslett, Xufeng Duan, Shuqi Wang, and Martin J Pickering. 2023. Does chatgpt resemble humans in language use? <i>arXiv preprint arXiv:2303.08014</i> .	Yanran Chen and Steffen Eger. 2022. Transformers go for the lols: Generating (humorous) titles from scientific abstracts end-to-end.	794
741			795
742			796
743			797
744	Yong Cao, Li Zhou, Seolhwa Lee, Laura Cabello, Min Chen, and Daniel Hershcovich. 2023. Assessing cross-cultural alignment between ChatGPT and human societies: An empirical study. In <i>Proceedings of the First Workshop on Cross-Cultural Considerations in NLP (C3NLP)</i> , pages 53–67, Dubrovnik, Croatia. Association for Computational Linguistics.	Cheng-Han Chiang and Hung-yi Lee. 2023. Can large language models be an alternative to human evaluations? In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 15607–15631, Toronto, Canada. Association for Computational Linguistics.	798
745			799
746			800
747			801
748			802
749			803
750			

804	Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pilla, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. PaLM: Scaling Language Modeling with Pathways . <i>arXiv:2204.02311 [cs]</i> . ArXiv: 2204.02311.	
805		
806		
807		
808		
809		
810		
811		
812		
813		
814		
815		
816		
817		
818		
819		
820		
821		
822		
823		
824		
825		
826		
827		
828	Mohita Chowdhury, Ernest Lim, Aisling Higham, Rory McKinnon, Nikoletta Ventoura, Yajie He, and Nick De Pennington. 2023. Can large language models safely address patient questions following cataract surgery? In <i>Proceedings of the 5th Clinical Natural Language Processing Workshop</i> , pages 131–137, Toronto, Canada. Association for Computational Linguistics.	
829		
830		
831		
832		
833		
834		
835		
836	Haoran Chu and Sixiao Liu. 2023. Can ai tell good stories? narrative transportation and persuasion with chatgpt . <i>PsyArXiv</i> .	
837		
838		
839	Ted M. Clark. 2023. Investigating the use of an artificial intelligence chatbot with general chemistry exam questions . <i>Journal of Chemical Education</i> , 100(5):1905–1916.	
840		
841		
842		
843	Merten Nikolay Dahlkemper, Simon Zacharias Lahme, and Pascal Klein. 2023. How do physics students evaluate artificial intelligence responses on comprehension questions? a study on the perceived scientific accuracy and linguistic quality .	
844		
845		
846		
847		
848	Haixing Dai, Zhengliang Liu, Wenxiong Liao, Xiaoke Huang, Yihan Cao, Zihao Wu, Lin Zhao, Shaochen Xu, Wei Liu, Ninghao Liu, Sheng Li, Dajiang Zhu, Hongmin Cai, Lichao Sun, Quanzheng Li, Dinggang Shen, Tianming Liu, and Xiang Li. 2023a. Auggpt: Leveraging chatgpt for text data augmentation .	
849		
850		
851		
852		
853		
854	Sunhao Dai, Ninglu Shao, Haiyuan Zhao, Weijie Yu, Zihua Si, Chen Xu, Zhongxiang Sun, Xiao Zhang, and Jun Xu. 2023b. Uncovering chatgpt’s capabilities in recommender systems .	
855		
856		
857		
858	Mithun Das, Saurabh Kumar Pandey, and Animesh Mukherjee. 2023. Evaluating chatgpt’s performance for multilingual and emoji-based hate speech detection .	
859		
860		
861		
	Ernest Davis. 2023. Benchmarks for automated commonsense reasoning: A survey .	862 863
	Sanjay Deshpande and Jakub Szefer. 2023. Analyzing chatgpt’s aptitude in an introductory computer engineering course .	864 865 866
	Sifatkaur Dhingra, Manmeet Singh, Vaisakh SB, Nee-tiraj Malviya, and Sukhpal Singh Gill. 2023. Mind meets machine: Unravelling gpt-4’s cognitive psychology .	867 868 869 870
	Bosheng Ding, Chengwei Qin, Linlin Liu, Yew Ken Chia, Boyang Li, Shafiq Joty, and Lidong Bing. 2023. Is GPT-3 a good data annotator? In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 11173–11195, Toronto, Canada. Association for Computational Linguistics.	871 872 873 874 875 876 877
	George Duenas, Sergio Jimenez, and Geral Mateus Ferro. 2023. You’ve got a friend in ... a language model? a comparison of explanations of multiple-choice items of reading comprehension between ChatGPT and humans . In <i>Proceedings of the 18th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2023)</i> , pages 372–381, Toronto, Canada. Association for Computational Linguistics.	878 879 880 881 882 883 884 885 886
	Jessica López Espejel, El Hassane Ettifouri, Mahaman Sanoussi Yahaya Alassan, El Mehdi Chouham, and Walid Dahhane. 2023. Gpt-3.5, gpt-4, or bard? evaluating llms reasoning ability in zero-shot setting and performance boosting through prompts .	887 888 889 890 891
	Yaxin Fan and Feng Jiang. 2023. Uncovering the potential of chatgpt for discourse analysis in dialogue: An empirical study .	892 893 894
	Tao Fang, Shu Yang, Kaixin Lan, Derek F. Wong, Jinpeng Hu, Lidia S. Chao, and Yue Zhang. 2023. Is chatgpt a highly fluent grammatical error correction system? a comprehensive evaluation .	895 896 897 898
	Sarah E. Finch, Ellie S. Paek, and Jinho D. Choi. 2023. Leveraging large language models for automated dialogue analysis . In <i>Proceedings of the 24th Meeting of the Special Interest Group on Discourse and Dialogue</i> , pages 202–215, Prague, Czechia. Association for Computational Linguistics.	899 900 901 902 903 904
	Kathleen Fraser, Svetlana Kiritchenko, Isar Nejadgholi, and Anna Kerkhof. 2023. What makes a good counter-stereotype? evaluating strategies for automated responses to stereotypical text . In <i>Proceedings of the First Workshop on Social Influence in Conversations (SICoN 2023)</i> , pages 25–38, Toronto, Canada. Association for Computational Linguistics.	905 906 907 908 909 910 911
	Simon Frieder, Luca Pinchetti, Alexis Chevalier, Ryan-Rhys Griffiths, Tommaso Salvatori, Thomas Lukasiewicz, Philipp Christian Petersen, and Julius Berner. 2023. Mathematical capabilities of chatgpt .	912 913 914 915

916	Jinglong Gao, Xiao Ding, Bing Qin, and Ting Liu.	Wenshi Gu. 2023. Linguistically informed chatgpt prompts to enhance japanese-chinese machine translation: A case study on attributive clauses.	971
917	2023a. Is chatgpt a good causal reasoner? a comprehensive evaluation.		972
918			973
919	Jun Gao, Huan Zhao, Changlong Yu, and Ruifeng Xu.	Reto Gubelmann, Aikaterini-lida Kalouli, Christina Niklaus, and Siegfried Handschuh. 2023. When truth matters - addressing pragmatic categories in natural language inference (NLI) by large language models (LLMs). In <i>Proceedings of the 12th Joint Conference on Lexical and Computational Semantics (*SEM 2023)</i> , pages 24–39, Toronto, Canada. Association for Computational Linguistics.	974
920	2023b. Exploring the feasibility of chatgpt for event extraction.		975
921			976
922	Mingqi Gao, Jie Ruan, Renliang Sun, Xunjian Yin, Shiping Yang, and Xiaojun Wan. 2023c. Human-like summarization evaluation with chatgpt.		977
923			978
924			979
925	Tianyu Gao, Howard Yen, Jiatong Yu, and Danqi Chen.		980
926	2023d. Enabling large language models to generate text with citations.		981
927			
928	Yuan Gao, Ruili Wang, and Feng Hou. 2023e. How to design translation prompts for chatgpt: An empirical study.	Biyang Guo, Xin Zhang, Ziyuan Wang, Minqi Jiang, Jinran Nie, Yuxuan Ding, Jianwei Yue, and Yupeng Wu. 2023. How close is chatgpt to human experts? comparison corpus, evaluation, and detection.	982
929			983
930			984
931	Wayne Geerling, G. Dirk Mateer, Jadrian Wooten, and Nikhil Damodaran. 2023. Chatgpt has aced the test of understanding in college economics: Now what? <i>The American Economist</i> , 68(2):233–245.		985
932			
933			
934			
935	Omid Ghahroodi, Seyed Arshan Dalili, Sahel Mesforoush, and Ehsaneddin Asgari. 2023. SUT at SemEval-2023 task 1: Prompt generation for visual word sense disambiguation. In <i>Proceedings of the 17th International Workshop on Semantic Evaluation (SemEval-2023)</i> , pages 2160–2163, Toronto, Canada. Association for Computational Linguistics.	Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A. Smith. 2020. Don't stop pretraining: Adapt language models to domains and tasks. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 8342–8360, Online. Association for Computational Linguistics.	986
936			987
937			988
938			989
939			990
940			991
941			992
942	Hamideh Ghanadian, Isar Nejadgholi, and Hussein Al Osman. 2023. Chatgpt for suicide risk assessment on social media: Quantitative evaluation of model performance, potentials and limitations.		993
943			
944			
945			
946	Fabrizio Gilardi, Meysam Alizadeh, and Maël Kubli. 2023. Chatgpt outperforms crowd-workers for text-annotation tasks.	Tahsina Hashem, Weiqing Wang, Derry Tanti Wijaya, Mohammed Eunus Ali, and Yuan-Fang Li. 2023. Generating faithful text from a knowledge graph with noisy reference text. In <i>Proceedings of the 16th International Natural Language Generation Conference</i> , pages 106–122, Prague, Czechia. Association for Computational Linguistics.	994
947			995
948			996
949	Aidan Gilson, Conrad W Safranek, Thomas Huang, Vimig Socrates, Ling Chi, Richard Andrew Taylor, and David Chartash. 2023. How does chatgpt perform on the united states medical licensing examination? the implications of large language models for medical education and knowledge assessment. <i>JMIR Med Educ</i> , 9.		997
950			998
951			999
952			1000
953			
954			
955			
956	GitHub. 2022. About github copilot. https://github.com/features/copilot .	Shreya Havaldar, Bhumika Singhal, Sunny Rai, Langchen Liu, Sharath Chandra Guntuku, and Lyle Ungar. 2023. Multilingual language models are not multicultural: A case study in emotion. In <i>Proceedings of the 13th Workshop on Computational Approaches to Subjectivity, Sentiment, & Social Media Analysis</i> , pages 202–214, Toronto, Canada. Association for Computational Linguistics.	1001
957			1002
958	Github. 2023. Evaluation papers for chatgpt.		1003
959	Shahriar Golchin and Mihai Surdeanu. 2023. Time travel in llms: Tracing data contamination in large language models. <i>arXiv preprint arXiv:2308.08493</i> .	Jiabang He, Lei Wang, Yi Hu, Ning Liu, Hui Liu, Xing Xu, and Heng Tao Shen. 2023a. Icl-d3ie: In-context learning with diverse demonstrations updating for document information extraction.	1004
960			1005
961			1006
962	Srinivas Gowriraj, Soham Dinesh Tiwari, Mitali Potnis, Srijan Bansal, Teruko Mitamura, and Eric Nyberg. 2023. Language-agnostic transformers and assessing ChatGPT-based query rewriting for multilingual document-grounded QA. In <i>Proceedings of the Third DialDoc Workshop on Document-grounded Dialogue and Conversational Question Answering</i> , pages 101–108, Toronto, Canada. Association for Computational Linguistics.		1007
963			1008
964			1009
965			1010
966			1011
967			1012
968			
969			
970			
		Qianyu He, Jie Zeng, Wenhao Huang, Lina Chen, Jin Xiao, Qianxi He, Xunzhe Zhou, Lida Chen, Xintao Wang, Yuncheng Huang, Haoning Ye, Zihan Li, Shisong Chen, Yikai Zhang, Zhouhong Gu, Jiaqing Liang, and Yanghua Xiao. 2023b. Can large language models understand real-world complex instructions?	1013
			1014
			1015
			1016
			1017
			1018
		Xinlei He, Xinyue Shen, Zeyuan Chen, Michael Backes, and Yang Zhang. 2023c. Mgtbench: Benchmarking machine-generated text detection.	1019
			1020
			1021
		Michael Heck, Nurul Lubis, Benjamin Ruppik, Renato Vukovic, Shutong Feng, Christian Geishauer, Hsienchin Lin, Carel van Niekerk, and Milica Gasic. 2023. ChatGPT for zero-shot dialogue state tracking: A solution or an opportunity? In <i>Proceedings of the</i>	1022
			1023
			1024
			1025
			1026

1027		Yuheng Huang, Jiayang Song, Zhijie Wang, Shengming Zhao, Huaming Chen, Felix Juefei-Xu, and Lei Ma. 2023d. Look before you leap: An exploratory study of uncertainty measurement for large language models.	1083
1028			1084
1029			1085
1030			1086
1031			1087
1032	Amr Hendy, Mohamed Abdelrehim, Amr Sharaf, Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita, Young Jin Kim, Mohamed Afify, and Hany Hassan Awadalla. 2023. How good are gpt models at machine translation? a comprehensive evaluation.		1088
1033			1089
1034			1090
1035			1091
1036			1092
1037	Takanobu Hirosawa, Yukinori Harada, Masashi Yokose, Tetsu Sakamoto, Ren Kawamura, and Taro Shimizu. 2023. Diagnostic accuracy of differential-diagnosis lists generated by generative pretrained transformer 3 chatbot for clinical vignettes with common chief complaints: A pilot study. <i>International Journal of Environmental Research and Public Health</i> , 20(4).		1093
1038			
1039			
1040			
1041			
1042			
1043	Bart Holterman and Kees van Deemter. 2023. Does chatgpt have theory of mind?		1094
1044			1095
1045			1096
1046			1097
1047			1098
1048			1099
1049			1100
1050			1101
1051			1102
1052			1103
1053			1104
1054			1105
1055			1106
1056			1107
1057			1108
1058			1109
1059			1110
1060			1111
1061			1112
1062			1113
1063			1114
1064			1115
1065			1116
1066			1117
1067			1118
1068			1119
1069			1120
1070			1121
1071			1122
1072			1123
1073			1124
1074			1125
1075			1126
1076			1127
1077			1128
1078			1129
1079			1130
1080			1131
1081			1132
1082			1133
			1134
			1135
			1136
			1137
			1138

1139	Jungo Kasai, Yuhei Kasai, Keisuke Sakaguchi, Yutaro Yamada, and Dragomir Radev. 2023. Evaluating gpt-4 and chatgpt on japanese medical licensing examinations.	1193
1140		1194
1141		1195
1142		
1143	Khalid S Khan, Regina Kunz, Jos Kleijnen, and Gerd Antes. 2003. Five steps to conducting a systematic review. <i>Journal of the royal society of medicine</i> , 96(3):118–121.	1196
1144		1197
1145		1198
1146		1199
1147	Yuheun Kim, Lu Guo, Bei Yu, and Yingya Li. 2023. Can ChatGPT understand causal language in science claims? In <i>Proceedings of the 13th Workshop on Computational Approaches to Subjectivity, Sentiment, & Social Media Analysis</i> , pages 379–389, Toronto, Canada. Association for Computational Linguistics.	1200
1148		1201
1149		1202
1150		
1151		1203
1152		1204
1153	Tom Kocmi and Christian Federmann. 2023. Large language models are state-of-the-art evaluators of translation quality.	1205
1154		1206
1155		1207
1156	Jan Kocoń, Igor Cichecki, Oliwier Kaszyca, Mateusz Kochanek, Dominika Szydło, Joanna Baran, Julita Bielaniec, Marcin Gruza, Arkadiusz Janz, Kamil Kanclerz, Anna Kocoń, Bartłomiej Koptyra, Wiktoria Mieszczzenko-Kowszewicz, Piotr Miłkowski, Marcin Oleksy, Maciej Piasecki, Łukasz Radliński, Konrad Wojtasik, Stanisław Woźniak, and Przemysław Kazienko. 2023. Chatgpt: Jack of all trades, master of none. <i>Information Fusion</i> , 99:101861.	1208
1157		1209
1158		1210
1159		1211
1160		1212
1161		1213
1162		1214
1163		1215
1164		
1165	Nam Ho Koh, Joseph Plata, and Joyce Chai. 2023. Bad: Bias detection for large language models in the context of candidate screening.	1216
1166		1217
1167		1218
1168	Philipp Koralus and Vincent Wang-Maćsianica. 2023. Humans in humans out: On gpt converging toward common sense in both success and failure.	1219
1169		
1170		
1171	Gerd Kortemeyer. 2023. Could an artificial-intelligence agent pass an introductory physics course?	1220
1172		1221
1173	Michał Kosinski. 2023. Theory of mind might have spontaneously emerged in large language models.	1222
1174		1223
1175	Sachin Kumar, Vidhisha Balachandran, Lucille Njoo, Antonios Anastasopoulos, and Yulia Tsvetkov. 2023. Language generation models can cause harm: So what can we do about it? an actionable survey. In <i>Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics</i> , pages 3299–3321, Dubrovnik, Croatia. Association for Computational Linguistics.	1224
1176		1225
1177		1226
1178		1227
1179		1228
1180		1229
1181		
1182		
1183	TH Kung, M Cheatham, A Medenilla, C Sillos, L De Leon, C Elepaño, M Madriaga, R Aggabao, G Diaz-Candido, J Maningo, et al. 2023a. Performance of chatgpt on usmle: Potential for ai-assisted medical education using large language models. <i>plos digit health</i> 2 (2): e0000198.	1230
1184		1231
1185		1232
1186		1233
1187		1234
1188		1235
1189	Tiffany H. Kung, Morgan Cheatham, Arielle Medenilla, Czarina Sillos, Lorie De Leon, Camille Elepaño, Maria Madriaga, Rimel Aggabao, Giezel Diaz-Candido, James Maningo, and Victor Tseng. 2023b. Performance of chatgpt on usmle: Potential for ai-assisted medical education using large language models. <i>PLOS Digital Health</i> , 2(2):1–12.	1236
1190		1237
1191		1238
1192		1239
		1240
		1241
		1242
		1243
		1244
		1245
		1246

1247	Lingyao Li, Lizhou Fan, Shubham Atreja, and Libby Hemphill. 2023e. "hot" chatgpt: The promise of chatgpt in detecting and discriminating hateful, offensive, and toxic comments on social media.	Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. 2023e. Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation.	1303
1248			1304
1249			1305
1250			1306
1251	Xingxuan Li, Ruochen Zhao, Yew Ken Chia, Bosheng Ding, Shafiq Joty, Soujanya Poria, and Lidong Bing. 2023f. Chain-of-knowledge: Grounding large language models via dynamic knowledge adapting over heterogeneous sources.	Xin Liu, Yuan Tan, Zhenghang Xiao, Jianwei Zhuge, and Rui Zhou. 2023f. Not the end of story: An evaluation of ChatGPT-driven vulnerability description mappings. In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 3724–3731, Toronto, Canada. Association for Computational Linguistics.	1307
1252			1308
1253			1309
1254			1310
1255			1311
1256	Zekun Li, Baolin Peng, Pengcheng He, Michel Galley, Jianfeng Gao, and Xifeng Yan. 2023g. Guiding large language models via directional stimulus prompting.	Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruochen Xu, and Chenguang Zhu. 2023g. G-eval: Nlg evaluation using gpt-4 with better human alignment.	1312
1257			1313
1258			1314
1259	Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, Binhang Yuan, Bobby Yan, Ce Zhang, Christian Alexander Cosgrove, Christopher D Manning, Christopher Re, Diana Acosta-Navas, Drew Arad Hudson, Eric Zelikman, Esin Durmus, Faisal Ladhak, Frieda Rong, Hongyu Ren, Huaxiu Yao, Jue WANG, Keshav Santhanam, Laurel Orr, Lucia Zheng, Mert Yuksekogun, Mirac Suzgun, Nathan Kim, Neel Guha, Niladri S. Chatterji, Omar Khattab, Peter Henderson, Qian Huang, Ryan Andrew Chi, Sang Michael Xie, Shibani Santurkar, Surya Ganguli, Tatsunori Hashimoto, Thomas Icard, Tianyi Zhang, Vishrav Chaudhary, William Wang, Xuechen Li, Yifan Mai, Yuhui Zhang, and Yuta Koreeda. 2023a. Holistic evaluation of language models. <i>Transactions on Machine Learning Research</i> . Featured Certification, Expert Certification.	Yiheng Liu, Tianle Han, Siyuan Ma, Jiayue Zhang, Yuanyuan Yang, Jiaming Tian, Hao He, Antong Li, Mengshen He, Zhengliang Liu, Zihao Wu, Lin Zhao, Dajiang Zhu, Xiang Li, Ning Qiang, Dingang Shen, Tianming Liu, and Bao Ge. 2023h. Summary of chatgpt-related research and perspective towards the future of large language models.	1315
1260			1316
1261			1317
1262			1318
1263			1319
1264			1320
1265			1321
1266			1322
1267			1323
1268			1324
1269			1325
1270			1326
1271			1327
1272			1328
1273			1329
1274			1330
1275			1331
1276			1332
1277			1333
1278	Yancheng Liang, Jiajie Zhang, Hui Li, Xiaochen Liu, Yi Hu, Yong Wu, Jiaoyao Zhang, Yongyan Liu, and Yi Wu. 2023b. Breaking the bank with ChatGPT: Few-shot text classification for finance. In <i>Proceedings of the Fifth Workshop on Financial Technology and Natural Language Processing and the Second Multimodal AI For Financial Forecasting</i> , pages 74–80, Macao. -.	Zequan Liu, Wei Zhang, Yingce Xia, Lijun Wu, Shufang Xie, Tao Qin, Ming Zhang, and Tie-Yan Liu. 2023i. MolXPT: Wrapping molecules with text for generative pre-training. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 1606–1616, Toronto, Canada. Association for Computational Linguistics.	1334
1279			1335
1280			1336
1281			1337
1282			1338
1283			1339
1284			1340
1285			1341
1286	Aiwei Liu, Xuming Hu, Lijie Wen, and Philip S. Yu. 2023a. A comprehensive evaluation of chatgpt's zero-shot text-to-sql capability.	Zhengliang Liu, Xiaowei Yu, Lu Zhang, Zihao Wu, Chao Cao, Haixing Dai, Lin Zhao, Wei Liu, Dinggang Shen, Quanzheng Li, Tianming Liu, Dajiang Zhu, and Xiang Li. 2023j. Deid-gpt: Zero-shot medical text de-identification by gpt-4.	1342
1287			1343
1288			1344
1289	Alisa Liu, Zhaofeng Wu, Julian Michael, Alane Suhr, Peter West, Alexander Koller, Swabha Swayamdipta, Noah A. Smith, and Yejin Choi. 2023b. We're afraid language models aren't modeling ambiguity.	Zhexiong Liu, Diane Litman, Elaine Wang, Lindsay Matsumura, and Richard Correnti. 2023k. Predicting the quality of revisions in argumentative writing. In <i>Proceedings of the 18th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2023)</i> , pages 275–287, Toronto, Canada. Association for Computational Linguistics.	1345
1290			1346
1291			1347
1292			1348
1293	Chang Liu and Bo Wu. 2023. Evaluating large language models on graphs: Performance insights and comparative analysis.	Mengsay Loem, Masahiro Kaneko, Sho Takase, and Naoaki Okazaki. 2023. Exploring effectiveness of GPT-3 in grammatical error correction: A study on performance and controllability in prompt-based methods. In <i>Proceedings of the 18th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2023)</i> , pages 205–219, Toronto, Canada. Association for Computational Linguistics.	1349
1294			1350
1295			1351
1296	Hanmeng Liu, Ruoxi Ning, Zhiyang Teng, Jian Liu, Qiji Zhou, and Yue Zhang. 2023c. Evaluating the logical reasoning ability of chatgpt and gpt-4.	Guang Lu, Sylvia B. Larcher, and Tu Tran. 2023a. Hybrid long document summarization using c2f-far and chatgpt: A practical study.	1352
1297			1353
1298			1354
1299	Hanmeng Liu, Zhiyang Teng, Leyang Cui, Chaoli Zhang, Qiji Zhou, and Yue Zhang. 2023d. Logi-cot: Logical chain-of-thought instruction-tuning data collection with gpt-4.	Qingyu Lu, Liang Ding, Liping Xie, Kanjian Zhang, Derek F. Wong, and Dacheng Tao. 2023b. Toward human-like evaluation for natural language generation with error analysis. In <i>Proceedings of the 61st</i>	1355
1300			1356
1301			1357
1302			1358
			1359

1360			1413
1361			1414
1362			1415
1363			1416
1364	Qingyu Lu, Baopu Qiu, Liang Ding, Kanjian Zhang, Tom Kocmi, and Dacheng Tao. 2023c. Error analysis prompting enables human-like translation evaluation in large language models: A case study on chatgpt.		1417
1365			1418
1366			1419
1367			1420
1368	Xin Lu, Zhuojun Li, Yanpeng Tong, Yanyan Zhao, and Bing Qin. 2023d. HIT-SCIR at WASSA 2023: Empathy and emotion analysis at the utterance-level and the essay-level. In <i>Proceedings of the 13th Workshop on Computational Approaches to Subjectivity, Sentiment, & Social Media Analysis</i> , pages 574–580, Toronto, Canada. Association for Computational Linguistics.		1421
1369			1422
1370			1423
1371			1424
1372			1425
1373			1426
1374			1427
1375			1428
1376			1429
1377			1430
1378			1431
1379			1432
1380			1433
1381			1434
1382			1435
1383			1436
1384			1437
1385			1438
1386			1439
1387			1440
1388			1441
1389			1442
1390			1443
1391			1444
1392			1445
1393			1446
1394			1447
1395			1448
1396			1449
1397			1450
1398			1451
1399			1452
1400			1453
1401			1454
1402			1455
1403			1456
1404			1457
1405			1458
1406			1459
1407			1460
1408			1461
1409			1462
1410			1463
1411			1464
1412			1465
			1466
			1467
			1468
			1469

1470	OpenAI. 2023a. Gpt-4 technical report .	1522
1471	OpenAI. 2023b. Gpt-4 technical report .	1523
1472	Miguel Ortega-Martín, Óscar García-Sierra, Alfonso	1524
1473	Ardoiz, Jorge Álvarez, Juan Carlos Armenteros, and	1525
1474	Adrián Alonso. 2023. Linguistic ambiguity analysis	1526
1475	in chatgpt .	1527
1476	Lidiia Ostyakova, Veronika Smilga, Kseniia Petukhova,	1528
1477	Maria Molchanova, and Daniel Kornev. 2023. Chat-	1529
1478	GPT vs. crowdsourcing vs. experts: Annotating open-	1530
1479	domain conversations with speech functions . In <i>Pro-</i>	1531
1480	<i>ceedings of the 24th Meeting of the Special Interest</i>	1532
1481	<i>Group on Discourse and Dialogue</i> , pages 242–254,	1533
1482	Prague, Czechia. Association for Computational Lin-	1534
1483	guistics.	1535
1484	Naoki Otani, Jun Araki, HyeongSik Kim, and Eduard	1536
1485	Hovy. 2023. On the underspecification of situa-	1537
1486	tions in open-domain conversational datasets . In	1538
1487	<i>Proceedings of the 5th Workshop on NLP for Con-</i>	1539
1488	<i>versational AI (NLP4ConvAI 2023)</i> , pages 12–28,	1540
1489	Toronto, Canada. Association for Computational Lin-	1541
1490	guistics.	1542
1491	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,	1543
1492	Carroll Wainwright, Pamela Mishkin, Chong Zhang,	1544
1493	Sandhini Agarwal, Katarina Slama, Alex Ray, John	1545
1494	Schulman, Jacob Hilton, Fraser Kelton, Luke Miller,	1546
1495	Maddie Simens, Amanda Askell, Peter Welinder,	1547
1496	Paul F Christiano, Jan Leike, and Ryan Lowe. 2022.	1548
1497	Training language models to follow instructions with	1549
1498	human feedback . In <i>Advances in Neural Information</i>	1550
1499	<i>Processing Systems</i> , volume 35, pages 27730–27744.	1551
1500	Curran Associates, Inc.	1552
1501	Wenbo Pan, Qiguang Chen, Xiao Xu, Wanxiang Che,	1553
1502	and Libo Qin. 2023. A preliminary evaluation of	1554
1503	chatgpt for zero-shot dialogue understanding .	1555
1504	Ralph Peeters and Christian Bizer. 2023. Using chatgpt	1556
1505	for entity matching .	1557
1506	Alessandro Pegoraro, Kavita Kumari, Hossein Ferei-	1558
1507	dooni, and Ahmad-Reza Sadeghi. 2023. To chatgpt,	1559
1508	or not to chatgpt: That is the question!	1560
1509	Baolin Peng, Michel Galley, Pengcheng He, Hao Cheng,	1561
1510	Yujia Xie, Yu Hu, Qiuyuan Huang, Lars Liden, Zhou	1562
1511	Yu, Weizhu Chen, and Jianfeng Gao. 2023a. Check	1563
1512	your facts and try again: Improving large language	1564
1513	models with external knowledge and automated feed-	1565
1514	back .	1566
1515	Baolin Peng, Chunyuan Li, Pengcheng He, Michel Gal-	1567
1516	ley, and Jianfeng Gao. 2023b. Instruction tuning with	1568
1517	gpt-4 .	1569
1518	Keqin Peng, Liang Ding, Qihuang Zhong, Li Shen,	1570
1519	Xuebo Liu, Min Zhang, Yuanxin Ouyang, and	1571
1520	Dacheng Tao. 2023c. Towards making the most of	1572
1521	chatgpt for machine translation .	1573
	Denis Peskoff and Brandon Stewart. 2023. Credible	1574
	without credit: Domain experts assess generative lan-	1575
	guage models . In <i>Proceedings of the 61st Annual</i>	1576
	<i>Meeting of the Association for Computational Lin-</i>	1577
	<i>guistics (Volume 2: Short Papers)</i> , pages 427–438,	1578
	Toronto, Canada. Association for Computational Lin-	
	guistics.	
	Pouya Pezeshkpour and Estevam Hruschka. 2023.	
	Large language models sensitivity to the order of	
	options in multiple-choice questions .	
	Aleksandra Piktus, Christopher Akiki, Paulo Villegas,	
	Hugo Laurençon, Gérard Dupont, Sasha Luccioni,	
	Yacine Jernite, and Anna Rogers. 2023. The ROOTS	
	search tool: Data transparency for LLMs . In <i>Proceed-</i>	
	<i>ings of the 61st Annual Meeting of the Association</i>	
	<i>for Computational Linguistics (Volume 3: System</i>	
	<i>Demonstrations)</i> , pages 304–314, Toronto, Canada.	
	Association for Computational Linguistics.	
	Andrei Popescu-Belis, Àlex R. Atrio, Bastien Bernath,	
	Etienne Boisson, Teo Ferrari, Xavier Theimer-	
	lienhard, and Giorgos Vernikos. 2023. GPoeT: a	
	language model trained for rhyme generation on syn-	
	thetic data . In <i>Proceedings of the 7th Joint SIGHUM</i>	
	<i>Workshop on Computational Linguistics for Cultural</i>	
	<i>Heritage, Social Sciences, Humanities and Litera-</i>	
	<i>ture</i> , pages 10–20, Dubrovnik, Croatia. Association	
	for Computational Linguistics.	
	Dongqi Pu and Vera Demberg. 2023. ChatGPT vs	
	human-authored text: Insights into controllable text	
	summarization and sentence style transfer . In <i>Pro-</i>	
	<i>ceedings of the 61st Annual Meeting of the Asso-</i>	
	<i>ciation for Computational Linguistics (Volume 4:</i>	
	<i>Student Research Workshop)</i> , pages 1–18, Toronto,	
	Canada. Association for Computational Linguistics.	
	Chengwei Qin, Aston Zhang, Zhuosheng Zhang, Jiaao	
	Chen, Michihiro Yasunaga, and Diyi Yang. 2023. Is	
	chatgpt a general-purpose natural language process-	
	ing task solver?	
	Ali Quidwai, Chunhui Li, and Parijat Dube. 2023. Be-	
	yond black box AI generated plagiarism detection:	
	From sentence to document level . In <i>Proceedings</i>	
	<i>of the 18th Workshop on Innovative Use of NLP</i>	
	<i>for Building Educational Applications (BEA 2023)</i> ,	
	pages 727–735, Toronto, Canada. Association for	
	Computational Linguistics.	
	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine	
	Lee, Sharan Narang, Michael Matena, Yanqi Zhou,	
	Wei Li, and Peter J. Liu. 2020. Exploring the limits	
	of transfer learning with a unified text-to-text trans-	
	former. <i>J. Mach. Learn. Res.</i> , 21(1).	
	Abhiramon Rajasekharan, Yankai Zeng, Parth Padalkar,	
	and Gopal Gupta. 2023. Reliable natural language	
	understanding with large language models and an-	
	swer set programming .	
	Aman Rangapur and Haoran Wang. 2023. Chatgpt-	
	crawler: Find out if chatgpt really knows what it's	
	talking about .	

1579	Arya Rao, Michael Pang, John Kim, Meghana Kamini, Winston Lie, Anoop K. Prasad, Adam Landman, Keith J Dreyer, and Marc D. Succi. 2023a. Assessing the utility of chatgpt throughout the entire clinical workflow. <i>medRxiv</i> .	1634
1580		1635
1581		1636
1582		1637
1583		
1584	Haocong Rao, Cyril Leung, and Chunyan Miao. 2023b. Can chatgpt assess human personalities? a general evaluation framework.	1638
1585		1639
1586		1640
1587	Mathieu Ravaut, Shafiq Joty, and Nancy Chen. 2023. Unsupervised summarization re-ranking. In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 8341–8376, Toronto, Canada. Association for Computational Linguistics.	1641
1588		1642
1589		1643
1590		1644
1591		1645
1592	Anton Razzhigaev, Mikhail Salnikov, Valentin Malykh, Pavel Braslavski, and Alexander Panchenko. 2023. A system for answering simple questions in multiple languages. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)</i> , pages 524–537, Toronto, Canada. Association for Computational Linguistics.	1646
1593		1647
1594		1648
1595		1649
1596		1650
1597		1651
1598		1652
1599		1653
1600	Ruiyang Ren, Yuhao Wang, Yingqi Qu, Wayne Xin Zhao, Jing Liu, Hao Tian, Hua Wu, Ji-Rong Wen, and Haifeng Wang. 2023. Investigating the factual knowledge boundary of large language models with retrieval augmentation.	1654
1601		1655
1602		1656
1603		1657
1604		1658
1605	Saed Rezayi, Zhengliang Liu, Zihao Wu, Chandra Dhakal, Bao Ge, Haixing Dai, Gengchen Mai, Ninghao Liu, Chen Zhen, Tianming Liu, and Sheng Li. 2023. Exploring new frontiers in agricultural nlp: Investigating the potential of large language models for food applications.	1659
1606		1660
1607		1661
1608		1662
1609		1663
1610		1664
1611	Nathaniel R. Robinson, Perez Ogayo, David R. Mortensen, and Graham Neubig. 2023. Chatgpt mt: Competitive for high- (but not low-) resource languages.	1665
1612		1666
1613		1667
1614		1668
1615	Jérôme Rutinowski, Sven Franke, Jan Endendyk, Ina Dormuth, and Markus Pauly. 2023. The self-perception and political biases of chatgpt.	1669
1616		1670
1617		1671
1618	Oscar Sainz, Jon Ander Campos, Iker García-Ferrero, Julen Etxaniz, and Eneko Agirre. Did chatgpt cheat on your test?	1672
1619		1673
1620		1674
1621	Oscar Sainz, Jon Ander Campos, Iker García-Ferrero, Julen Etxaniz, and Eneko Agirre. 2023. Did chatgpt cheat on your test?	1675
1622		1676
1623		1677
1624	Maarten Sap, Ronan LeBras, Daniel Fried, and Yejin Choi. 2023. Neural theory-of-mind? on the limits of social intelligence in large lms.	1678
1625		1679
1626		1680
1627	Jakob Schuster and Katja Markert. 2023. Nut-cracking sledgehammers: Prioritizing target language data over bigger language models for cross-lingual metaphor detection. In <i>Proceedings of the 2023 CLASP Conference on Learning with Small Data (LSD)</i> , pages 98–106, Gothenburg, Sweden. Association for Computational Linguistics.	1681
1628		1682
1629		1683
1630		1684
1631		1685
1632		1686
1633		1687
	Paulo Shakarian, Abhinav Koyyalamudi, Noel Ngu, and Lakshmihari Mareedu. 2023. An independent evaluation of chatgpt on mathematical word problems (mwp).	1634
		1635
		1636
		1637
	Natalie Shapira, Guy Zwirn, and Yoav Goldberg. 2023. How well do large language models perform on faux pas tests? In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 10438–10451, Toronto, Canada. Association for Computational Linguistics.	1638
		1639
		1640
		1641
		1642
		1643
	Xinyue Shen, Zeyuan Chen, Michael Backes, and Yang Zhang. 2023. In chatgpt we trust? measuring and characterizing the reliability of chatgpt.	1644
		1645
		1646
	Yucheng Shi, Hehuan Ma, Wenliang Zhong, Qiaoyu Tan, Gengchen Mai, Xiang Li, Tianming Liu, and Junzhou Huang. 2023. Chatgraph: Interpretable text classification by converting chatgpt knowledge to graphs.	1647
		1648
		1649
		1650
		1651
	Zhengxaing Shi and Aldo Lipani. 2023. Don't stop pretraining? make prompt-based fine-tuning powerful learner. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	1652
		1653
		1654
		1655
	Dominik Sobania, Martin Briesch, Carol Hanna, and Justyna Petke. 2023. An analysis of the automatic bug fixing performance of chatgpt.	1656
		1657
		1658
	Mingyang Song, Haiyun Jiang, Shuming Shi, Songfang Yao, Shilong Lu, Yi Feng, Huafeng Liu, and Liping Jing. 2023. Is chatgpt a good keyphrase generator? a preliminary study.	1659
		1660
		1661
		1662
	Mayank Soni and Vincent Wade. 2023. Comparing abstractive summaries generated by chatgpt to real summaries through blinded reviewers and text classification algorithms.	1663
		1664
		1665
		1666
	David Stap and Ali Araabi. 2023. ChatGPT is not a good indigenous translator. In <i>Proceedings of the Workshop on Natural Language Processing for Indigenous Languages of the Americas (AmericasNLP)</i> , pages 163–167, Toronto, Canada. Association for Computational Linguistics.	1667
		1668
		1669
		1670
		1671
		1672
	Jiashuo Sun, Chengjin Xu, Lumingyuan Tang, Saizhuo Wang, Chen Lin, Yeyun Gong, Lionel M. Ni, Heung-Yeung Shum, and Jian Guo. 2023a. Think-on-graph: Deep and responsible reasoning of large language model on knowledge graph.	1673
		1674
		1675
		1676
		1677
	Weiwei Sun, Lingyong Yan, Xinyu Ma, Pengjie Ren, Dawei Yin, and Zhaochun Ren. 2023b. Is chatgpt good at search? investigating large language models as re-ranking agent.	1678
		1679
		1680
		1681
	Eugene Syriani, Istvan David, and Gauransh Kumar. 2023. Assessing the ability of chatgpt to screen articles for systematic reviews.	1682
		1683
		1684
	Jakub Szefer and Sanjay Deshpande. 2023. Analyzing chatgpt's aptitude in an introductory computer engineering course. <i>arXiv preprint arXiv:2304.06122</i> .	1685
		1686
		1687

1688	Qingyu Tan, Hwee Tou Ng, and Lidong Bing. 2023a. Towards benchmarking and improving the temporal reasoning capability of large language models . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 14820–14835, Toronto, Canada. Association for Computational Linguistics.	1748
1689		1749
1690		1750
1691		
1692		1751
1693		1752
1694		1753
1695	Yiming Tan, Dehai Min, Yu Li, Wenbo Li, Nan Hu, Yongrui Chen, and Guilin Qi. 2023b. Can chatgpt replace traditional kbqa models? an in-depth analysis of the question answering performance of the gpt llm family .	1754
1696		1755
1697		1756
1698		1757
1699		1758
1700	Zeqi Tan, Shen Huang, Zixia Jia, Jiong Cai, Yinghui Li, Weiming Lu, Yueting Zhuang, Kewei Tu, Pengjun Xie, Fei Huang, and Yong Jiang. 2023c. DAMO-NLP at SemEval-2023 task 2: A unified retrieval-augmented system for multilingual named entity recognition . In <i>Proceedings of the 17th International Workshop on Semantic Evaluation (SemEval-2023)</i> , pages 2014–2028, Toronto, Canada. Association for Computational Linguistics.	1759
1701		1760
1702		
1703		
1704		
1705		
1706		
1707		
1708		
1709	Liyan Tang, Tanya Goyal, Alex Fabbri, Philippe Laban, Jiacheng Xu, Semih Yavuz, Wojciech Kryscinski, Justin Rousseau, and Greg Durrett. 2023a. Understanding factual errors in summarization: Errors, summarizers, datasets, error detectors . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 11626–11644, Toronto, Canada. Association for Computational Linguistics.	1761
1710		1762
1711		1763
1712		1764
1713		1765
1714		1766
1715		1767
1716		1768
1717		
1718	Ruixiang Tang, Xiaotian Han, Xiaoqian Jiang, and Xia Hu. 2023b. Does synthetic data generation of llms help clinical text mining?	1769
1719		1770
1720		1771
1721	Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, YaGuang Li, Hongrae Lee, Huaixiu Steven Zheng, Amin Ghafouri, Marcelo Menegali, Yanping Huang, Maxim Krikun, Dmitry Lepikhin, James Qin, Dehao Chen, Yuanzhong Xu, Zhifeng Chen, Adam Roberts, Maarten Bosma, Vincent Zhao, Yanqi Zhou, Chung-Ching Chang, Igor Krivokon, Will Rusch, Marc Pickett, Pranesh Srinivasan, Laichee Man, Kathleen Meier-Hellstern, Meredith Ringel Morris, Tulsee Doshi, Renelito Delos Santos, Toju Duke, Johnny Soraker, Ben Zevenbergen, Vinodkumar Prabhakaran, Mark Diaz, Ben Hutchinson, Kristen Olson, Alejandra Molina, Erin Hoffman-John, Josh Lee, Lora Aroyo, Ravi Rajakumar, Alena Butryna, Matthew Lamm, Viktoriya Kuzmina, Joe Fenton, Aaron Cohen, Rachel Bernstein, Ray Kurzweil, Blaise Agueras-Arcas, Claire Cui, Marian Croak, Ed Chi, and Quoc Le. 2022. Lamda: Language models for dialog applications .	1772
1722		1773
1723		1774
1724		1775
1725		1776
1726		1777
1727		1778
1728		1779
1729		1780
1730		1781
1731		
1732		
1733		
1734		
1735		
1736		
1737		
1738		
1739		
1740		
1741		
1742	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models . <i>arXiv preprint arXiv:2302.13971</i> .	1775
1743		1776
1744		1777
1745		1778
1746		1779
1747		1800
		1801
		1802
	Ruibo Tu, Chao Ma, and Cheng Zhang. 2023a. Causal-discovery performance of chatgpt in the context of neuropathic pain diagnosis .	1748
		1749
		1750
	Shangqing Tu, Chunyang Li, Jifan Yu, Xiaozhi Wang, Lei Hou, and Juanzi Li. 2023b. Chatlog: Recording and analyzing chatgpt across time .	1751
		1752
		1753
	Jens Van Nooten and Walter Daelemans. 2023. Improving Dutch vaccine hesitancy monitoring via multi-label data augmentation with GPT-3.5 . In <i>Proceedings of the 13th Workshop on Computational Approaches to Subjectivity, Sentiment, & Social Media Analysis</i> , pages 251–270, Toronto, Canada. Association for Computational Linguistics.	1754
		1755
		1756
		1757
		1758
		1759
		1760
	Bhaskara Hanuma Vedula, Prashant Kodali, Manish Shrivastava, and Ponnurangam Kumaraguru. 2023. PrecogIIITH@WASSA2023: Emotion detection for Urdu-English code-mixed text . In <i>Proceedings of the 13th Workshop on Computational Approaches to Subjectivity, Sentiment, & Social Media Analysis</i> , pages 601–605, Toronto, Canada. Association for Computational Linguistics.	1761
		1762
		1763
		1764
		1765
		1766
		1767
		1768
	Sai Vemprala, Rogerio Bonatti, Arthur Bucker, and Ashish Kapoor. 2023. Chatgpt for robotics: Design principles and model abilities .	1769
		1770
		1771
	Boshi Wang, Xiang Yue, and Huan Sun. 2023a. Can chatgpt defend its belief in truth? evaluating llm reasoning via debate .	1772
		1773
		1774
	Boxin Wang, Weixin Chen, Hengzhi Pei, Chulin Xie, Mintong Kang, Chenhui Zhang, Chejian Xu, Zidi Xiong, Ritik Dutta, Rylan Schaeffer, Sang T. Truong, Simran Arora, Mantas Mazeika, Dan Hendrycks, Zinan Lin, Yu Cheng, Sanmi Koyejo, Dawn Song, and Bo Li. 2023b. Decodingtrust: A comprehensive assessment of trustworthiness in gpt models .	1775
		1776
		1777
		1778
		1779
		1780
		1781
	Jiaan Wang, Yunlong Liang, Fandong Meng, Zengkui Sun, Haoxiang Shi, Zhixu Li, Jinan Xu, Jianfeng Qu, and Jie Zhou. 2023c. Is chatgpt a good nlg evaluator? a preliminary study .	1782
		1783
		1784
		1785
	Jiaan Wang, Yunlong Liang, Fandong Meng, Beiqi Zou, Zhixu Li, Jianfeng Qu, and Jie Zhou. 2023d. Zero-shot cross-lingual summarization via large language models .	1786
		1787
		1788
		1789
	Jindong Wang, Xixu Hu, Wenxin Hou, Hao Chen, Runkai Zheng, Yidong Wang, Linyi Yang, Haojun Huang, Wei Ye, Xiubo Geng, Binxin Jiao, Yue Zhang, and Xing Xie. 2023e. On the robustness of chatgpt: An adversarial and out-of-distribution perspective .	1790
		1791
		1792
		1793
		1794
	Junda Wang, Zonghai Yao, Avijit Mitra, Samuel Osebe, Zhichao Yang, and Hong Yu. 2023f. UMASS_BioNLP at MEDIQA-chat 2023: Can LLMs generate high-quality synthetic note-oriented doctor-patient conversations? In <i>Proceedings of the 5th Clinical Natural Language Processing Workshop</i> , pages 460–471, Toronto, Canada. Association for Computational Linguistics.	1795
		1796
		1797
		1798
		1799
		1800
		1801
		1802

1803	Rose Wang and Dorottya Demszky. 2023. Is ChatGPT a good teacher coach? measuring zero-shot performance for scoring and providing actionable insights on classroom instruction. In <i>Proceedings of the 18th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2023)</i> , pages 626–667, Toronto, Canada. Association for Computational Linguistics.	1857
1804		1858
1805		1859
1806		1860
1807		1861
1808		1862
1809		1863
1810		1864
1811	Sheng Wang, Zihao Zhao, Xi Ouyang, Qian Wang, and Dinggang Shen. 2023g. Chatcad: Interactive computer-aided diagnosis on medical image using large language models.	1865
1812		1866
1813		1867
1814		1868
1815	Xiaoxuan Wang, Ziniu Hu, Pan Lu, Yanqiao Zhu, Jieyu Zhang, Satyen Subramaniam, Arjun R. Loomba, Shichang Zhang, Yizhou Sun, and Wei Wang. 2023h. Scibench: Evaluating college-level scientific problem-solving abilities of large language models.	1869
1816		1870
1817		1871
1818		1872
1819		
1820	Xingyao Wang, Zihan Wang, Jiateng Liu, Yangyi Chen, Lifan Yuan, Hao Peng, and Heng Ji. 2023i. Mint: Evaluating llms in multi-turn interaction with tools and language feedback.	1873
1821		1874
1822		1875
1823		1876
1824	Yuqing Wang and Yun Zhao. 2023. Metacognitive prompting improves understanding in large language models.	1877
1825		1878
1826		1879
1827	Zengzhi Wang, Qiming Xie, Zixiang Ding, Yi Feng, and Rui Xia. 2023j. Is chatgpt a good sentiment analyzer? a preliminary study.	1880
1828		1881
1829		
1830	Xiang Wei, Xingyu Cui, Ning Cheng, Xiaobin Wang, Xin Zhang, Shen Huang, Pengjun Xie, Jinan Xu, Yufeng Chen, Meishan Zhang, Yong Jiang, and Wenjuan Han. 2023. Zero-shot information extraction via chatting with chatgpt.	1882
1831		1883
1832		1884
1833		1885
1834		1886
1835	Jules White, Sam Hays, Quichen Fu, Jesse Spencer-Smith, and Douglas C. Schmidt. 2023. Chatgpt prompt patterns for improving code quality, refactoring, requirements elicitation, and software design.	1887
1836		
1837		
1838		
1839	Haoran Wu, Wenxuan Wang, Yuxuan Wan, Wenxiang Jiao, and Michael Lyu. 2023a. Chatgpt or grammarly? evaluating chatgpt on grammatical error correction benchmark.	1888
1840		1889
1841		1890
1842		1891
1843	Qianhui Wu, Huiqiang Jiang, Haonan Yin, Börje Karlsson, and Chin-Yew Lin. 2023b. Multi-level knowledge distillation for out-of-distribution detection in text. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 7317–7332, Toronto, Canada. Association for Computational Linguistics.	1892
1844		1893
1845		1894
1846		1895
1847		1896
1848		1897
1849		
1850	WeiQi Wu, Chengyue Jiang, Yong Jiang, Pengjun Xie, and Kewei Tu. 2023c. Do PLMs know and understand ontological knowledge? In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3080–3101, Toronto, Canada. Association for Computational Linguistics.	1898
1851		1899
1852		1900
1853		1901
1854		1902
1855		1903
1856		1904
		1905
		1906
		1907
		1908
		1909
		1910
		1911
		1912
	Changrong Xiao, Sean Xin Xu, Kunpeng Zhang, Yufang Wang, and Lei Xia. 2023. Evaluating reading comprehension exercises generated by LLMs: A showcase of ChatGPT in education applications. In <i>Proceedings of the 18th Workshop on Innovative Use of NLP for Building Educational Applications (BEA 2023)</i> , pages 610–625, Toronto, Canada. Association for Computational Linguistics.	
	Jian Xie, Kai Zhang, Jiangjie Chen, Renze Lou, and Yu Su. 2023a. Adaptive chameleon or stubborn sloth: Revealing the behavior of large language models in knowledge conflicts.	
	Qianqian Xie, Weiguang Han, Yanzhao Lai, Min Peng, and Jimin Huang. 2023b. The wall street neophyte: A zero-shot analysis of chatgpt over multimodal stock movement prediction challenges.	
	Fangzhi Xu, Qika Lin, Jiawei Han, Tianzhe Zhao, Jun Liu, and Erik Cambria. 2023a. Are large language models really good logical reasoners? a comprehensive evaluation and beyond.	
	Liang Xu, Anqi Li, Lei Zhu, Hang Xue, Changtai Zhu, Kangkang Zhao, Haonan He, Xuanwei Zhang, Qiyue Kang, and Zhenzhong Lan. 2023b. Superclue: A comprehensive chinese large language model benchmark.	
	Zihang Xu, Ziqing Yang, Yiming Cui, and Shijin Wang. 2023c. IDOL: Indicator-oriented logic pre-training for logical reasoning. In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 8099–8111, Toronto, Canada. Association for Computational Linguistics.	
	Jingfeng Yang, Hongye Jin, Ruixiang Tang, Xiaotian Han, Qizhang Feng, Haoming Jiang, Bing Yin, and Xia Hu. 2023a. Harnessing the power of llms in practice: A survey on chatgpt and beyond.	
	Wayne Yang and Garrett Nicolai. 2023. Neural machine translation data generation and augmentation using chatgpt.	
	Xianjun Yang, Yan Li, Xinlu Zhang, Haifeng Chen, and Wei Cheng. 2023b. Exploring the limits of chatgpt for query or aspect-based text summarization.	
	Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Ehsan Azarnasab, Faisal Ahmed, Zicheng Liu, Ce Liu, Michael Zeng, and Lijuan Wang. 2023c. Mm-react: Prompting chatgpt for multimodal reasoning and action.	
	Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models.	
	Junjie Ye, Xuanting Chen, Nuo Xu, Can Zu, Zekai Shao, Shichun Liu, Yuhuan Cui, Zeyang Zhou, Chao Gong, Yang Shen, Jie Zhou, Siming Chen, Tao Gui, Qi Zhang, and Xuanjing Huang. 2023. A comprehensive capability analysis of gpt-3 and gpt-3.5 series models.	

1913	Chenhan Yuan, Qianqian Xie, and Sophia Ananiadou. 2023a. Zero-shot temporal relation extraction with chatgpt.	1967
1914		1968
1915		1969
1916	Zheng Yuan, Hongyi Yuan, Chuanqi Tan, Wei Wang, and Songfang Huang. 2023b. How well do large language models perform in arithmetic tasks?	1970
1917		1971
1918		1972
1919	Bowen Zhang, Daijun Ding, and Liwen Jing. 2022. How would stance detection techniques evolve after the launch of chatgpt? <i>arXiv preprint arXiv:2212.14548</i> .	1973
1920		1974
1921		1975
1922		1976
1923	Bowen Zhang, Xianghua Fu, Daijun Ding, Hu Huang, Yangyang Li, and Liwen Jing. 2023a. Investigating chain-of-thought with chatgpt for stance detection on social media.	1977
1924		1978
1925		
1926		
1927	Haopeng Zhang, Xiao Liu, and Jiawei Zhang. 2023b. Extractive summarization via chatgpt for faithful summary generation.	1979
1928		1980
1929		1981
1930	Jianguo Zhang, Kun Qian, Zhiwei Liu, Shelby Heinecke, Rui Meng, Ye Liu, Zhou Yu, Huan Wang, Silvio Savarese, and Caiming Xiong. 2023c. Dialogstudio: Towards richest and most diverse unified dataset collection for conversational ai.	1982
1931		1983
1932		1984
1933		
1934		
1935	Shaolei Zhang, Qingkai Fang, Zhuocheng Zhang, Zhengrui Ma, Yan Zhou, Langlin Huang, Mengyu Bu, Shangdong Gui, Yunji Chen, Xilin Chen, and Yang Feng. 2023d. Bayling: Bridging cross-lingual alignment and instruction following through interactive translation for large language models.	1985
1936		1986
1937		1987
1938		1988
1939		1989
1940		1990
1941	Wenxuan Zhang, Yue Deng, Bing Liu, Sinno Jialin Pan, and Lidong Bing. 2023e. Sentiment analysis in the era of large language models: A reality check.	1991
1942		1992
1943		1993
1944	Xiyuan Zhang, Ranak Roy Chowdhury, Dezhi Hong, Rajesh K. Gupta, and Jingbo Shang. 2023f. Modeling label semantics improves activity recognition.	1994
1945		1995
1946		1996
1947	Weixiang Zhao, Yanyan Zhao, Xin Lu, Shilong Wang, Yanpeng Tong, and Bing Qin. 2023a. Is chatgpt equipped with emotional dialogue capabilities?	1997
1948		1998
1949		1999
1950	Xiaofeng Zhao, Min Zhang, Miaomiao Ma, Chang Su, Yilun Liu, Minghan Wang, Xiaosong Qiao, Jiabin Guo, Yinglu Li, and Wenbing Ma. 2023b. HW-TSC at SemEval-2023 task 7: Exploring the natural language inference capabilities of ChatGPT and pre-trained language model for clinical trial. In <i>Proceedings of the 17th International Workshop on Semantic Evaluation (SemEval-2023)</i> , pages 1603–1608, Toronto, Canada. Association for Computational Linguistics.	2000
1951		2001
1952		2002
1953		2003
1954		2004
1955		2005
1956		2006
1957		2007
1958		2008
1959		2009
1960	Xuandong Zhao, Siqi Ouyang, Zhiguo Yu, Ming Wu, and Lei Li. 2023c. Pre-trained language models can be fully zero-shot learners. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 15590–15606, Toronto, Canada. Association for Computational Linguistics.	2010
1961		2011
1962		2012
1963		2013
1964		
1965		
1966		
	Chujie Zheng, Hao Zhou, Fandong Meng, Jie Zhou, and Minlie Huang. 2023a. Large language models are not robust multiple choice selectors.	1967
		1968
		1969
	Shen Zheng, Jie Huang, and Kevin Chen-Chuan Chang. 2023b. Why does chatgpt fall short in providing truthful answers?	1970
		1971
		1972
	Yi Zheng, Björn Ross, and Walid Magdy. 2023c. What makes good counterspeech? a comparison of generation approaches and evaluation metrics. In <i>Proceedings of the 1st Workshop on CounterSpeech for Online Abuse (CS4OA)</i> , pages 62–71, Prague, Czechia. Association for Computational Linguistics.	1973
		1974
		1975
		1976
		1977
		1978
	Zhi Zheng, Zhaopeng Qiu, Xiao Hu, Likang Wu, Hengshu Zhu, and Hui Xiong. 2023d. Generative job recommendations with large language model.	1979
		1980
		1981
	Qihuang Zhong, Liang Ding, Juhua Liu, Bo Du, and Dacheng Tao. 2023a. Can chatgpt understand too? a comparative study on chatgpt and fine-tuned bert.	1982
		1983
		1984
	Tianyang Zhong, Yaonai Wei, Li Yang, Zihao Wu, Zhengliang Liu, Xiaozheng Wei, Wenjun Li, Junjie Yao, Chong Ma, Xiang Li, Dajiang Zhu, Xi Jiang, Junwei Han, Dinggang Shen, Tianming Liu, and Tuo Zhang. 2023b. Chatabl: Abductive learning via natural language interaction with chatgpt.	1985
		1986
		1987
		1988
		1989
		1990
	Wangchunshu Zhou, Yuchen Eleanor Jiang, Peng Cui, Tiannan Wang, Zhenxin Xiao, Yifan Hou, Ryan Cotterell, and Mrinmaya Sachan. 2023. Recurrentgpt: Interactive generation of (arbitrarily) long text.	1991
		1992
		1993
		1994
	Yiming Zhu, Peixian Zhang, Ehsan-Ul Haq, Pan Hui, and Gareth Tyson. 2023. Can chatgpt reproduce human-generated labels? a study of social computing tasks.	1995
		1996
		1997
		1998
	Terry Yue Zhuo, Yujin Huang, Chunyang Chen, and Zhenchang Xing. 2023. Red teaming chatgpt via jail-breaking: Bias, robustness, reliability and toxicity.	1999
		2000
		2001
	A Full list of the reviewed work	2002
	In this section, we list all the work that we reviewed and classified as relevant.	2003
		2004
	• LLM-contamination (Sainz et al., 2023)	2005
	• Time Travel in LLMs: Tracing Data Contamination in Large Language Models (Golchin and Surdeanu, 2023)	2006
		2007
		2008
	• Don't Stop Pretraining: Adapt Language Models to Domains and Tasks (Gururangan et al., 2020)	2009
		2010
		2011
	• Can we trust the evaluation on ChatGPT? (Aiyappa et al., 2023b)	2012
		2013

2014	• How Close is ChatGPT to Human Experts? Comparison Corpus, Evaluation, and Detection (Guo et al., 2023)	• ChatGPT: Jack of all trades, master of none (Kocoń et al., 2023)	2055
2015			2056
2016			
2017	• GPTEval: A Survey on Assessments of ChatGPT and GPT-4 (Mao et al., 2023)	• Zero-Shot Information Extraction via Chatting with ChatGPT (Wei et al., 2023)	2057
2018			2058
2019	• A Systematic Study and Comprehensive Evaluation of ChatGPT on Benchmark Datasets (Laskar et al., 2023)	• Is ChatGPT Good at Search? Investigating Large Language Models as Re-Ranking Agent (Sun et al., 2023b)	2059
2020			2060
2021		• Is ChatGPT a Good NLG Evaluator? A Preliminary Study (Wang et al., 2023c)	2061
2022	• Quantifying Memorization Across Neural Language Models (Carlini et al., 2023)	• G-EVAL: NLG Evaluation using GPT-4 with Better Human Alignment (Liu et al., 2023g)	2062
2023			2063
2024	• Preventing Generation of Verbatim Memorization in Language Models Gives a False Sense of Privacy (Ippolito et al., 2023)	• Large Language Models Are State-of-the-Art Evaluators of Translation Quality (Kocmi and Federmann, 2023)	2064
2025			2065
2026		• ChatGPT Outperforms Crowd-workers For Text-annotation Tasks (Gilardi et al., 2023)	2066
2027	• List of various ChatGPT evaluations (Github, 2023)		2067
2028		• SuperCLUE: A Comprehensive Chinese Large Language Model Benchmark (Xu et al., 2023b)	2068
2029	• GPT-4 Technical Report (OpenAI, 2023b)		2069
2030	• Zeno Chatbot Report (Cabrera and Neubig, 2023)	• C-eval: A multi-level multi-discipline chinese evaluation suite for foundation models (Huang et al., 2023e)	2070
2031			2071
2032	• A Multitask, Multilingual, Multimodal Evaluation of ChatGPT on Reasoning, Hallucination, and Interactivity (Bang et al., 2023)	• ChatGPT beyond English: Towards a comprehensive evaluation of large language models in multilingual learning (Lai et al., 2023)	2072
2033			2073
2034		• Zero-Shot Cross-Lingual Summarization via Large Language Models (Wang et al., 2023d)	2074
2035	• Don't Stop Pretraining? Make Promp-based Fine-tuning Powerful Learner (Shi and Lipani, 2023)	• A categorical archive of chatgpt failures (Borji, 2023)	2075
2036			2076
2037		• Benchmarks for automated commonsense reasoning: A survey (Davis, 2023)	2077
2038	• How Good Are GPT Models at Machine Translation? A Comprehensive Evaluation (Hendy et al., 2023)	• Is ChatGPT a Good Causal Reasoner? A Comprehensive Evaluation. (Gao et al., 2023a)	2078
2039			2079
2040		• Causal reasoning and large language models: Opening a new frontier for causality (Kırcıman et al., 2023)	2080
2041	• Is ChatGPT A Good Translator? Yes With GPT-4 As The Engine (Jiao et al., 2023)	• Theory of mind may have spontaneously emerged in large language models (Kosinski, 2023)	2081
2042			2082
2043	• Is Chatgpt A General-purpose Natural Language Processing Task Solver? (Qin et al., 2023)	• Boosting Theory-of-Mind Performance in Large Language Models via Prompting (Moghaddam and Honey, 2023)	2083
2044			2084
2045			2085
2046	• Extractive summarization via ChatGPT for faithful summary generation (Zhang et al., 2023b)		2086
2047			2087
2048			2088
2049	• Chatgpt or grammarly? evaluating chatgpt on grammatical error correction benchmark (Wu et al., 2023a)		2089
2050			2090
2051			2091
2052	• Will affective computing emerge from foundation models and general ai? A first evaluation on chatgpt (Amin et al., 2023)		2092
2053			2093
2054			2094
			2095
			2096

2097	• Does ChatGPT have Theory of Mind? (Holterman and Deemter, 2023)	2139
2098		2140
2099	• Can ChatGPT Defend the Truth? Automatic Dialectical Evaluation Elicits LLMs' Deficiencies in Reasoning (Wang et al., 2023a)	2141
2100		2142
2101		2143
2102	• Evaluating the logical reasoning ability of chatgpt and gpt-4 (Liu et al., 2023c)	2144
2103		2145
2104	• We're Afraid Language Models Aren't Modeling Ambiguity (Liu et al., 2023b)	2146
2105		2147
2106	• Sparks of artificial general intelligence: Early experiments with gpt-4 (Bubeck et al., 2023)	2148
2107		2149
2108	• Chatgpt participates in a computer science exam (Bordt and Luxburg, 2023)	2150
2109		2151
2110	• Can llm already serve as a database interface? a big bench for large-scale database grounded text-to-sqls (Li et al., 2023c)	2152
2111		2153
2112		2154
2113	• Mathematical capabilities of chatgpt (Frieder et al., 2023)	2155
2114		2156
2115	• Could an artificial-intelligence agent pass an introductory physics course? (Kortemeyer, 2023)	2157
2116		2158
2117		2159
2118	• Investigating the Use of an Artificial Intelligence Chatbot with General Chemistry Exam Questions (Clark, 2023)	2160
2119		2161
2120		2162
2121	• How do physics students evaluate artificial intelligence responses on comprehension questions? A study on the perceived scientific accuracy and linguistic quality of ChatGPT (Dahlkemper et al., 2023)	2163
2122		2164
2123		2165
2124		2166
2125		2167
2126	• Chatgpt goes to law school (Choi et al., 2023)	2168
2127		2169
2128	• ChatGPT has aced the test of understanding in college economics: Now what? (Geerling et al., 2023)	2170
2129		2171
2130	• The Wall Street Neophyte: A Zero-Shot Analysis of ChatGPT Over MultiModal Stock Movement Prediction Challenges (Xie et al., 2023b)	2172
2131		2173
2132		2174
2133		2175
2134	• How Does ChatGPT Perform on the United States Medical Licensing Examination? The Implications of Large Language Models for Medical Education and Knowledge Assessment (Gilson et al., 2023)	2176
2135		2177
2136		2178
2137		2179
2138		2180
	• Performance of ChatGPT on USMLE: Potential for AI-assisted medical education using large language models (Kung et al., 2023b)	2181
		2182
	• Evaluating the performance of chatgpt in ophthalmology: An analysis of its successes and shortcomings (Antaki et al., 2023)	2183
		2184
	• Diagnostic Accuracy of Differential-Diagnosis Lists Generated by Generative Pretrained Transformer 3 Chatbot for Clinical Vignettes with Common Chief Complaints: A Pilot Study (Hirosawa et al., 2023)	2185
		2186
	• Assessing the utility of ChatGPT throughout the entire clinical workflow (Rao et al., 2023a)	2187
		2188
	• ChatGPT as a medical doctor? A diagnostic accuracy study on common and rare diseases (Mehnen et al., 2023)	2189
		2190
	• Can AI tell good stories? narrative transportation and persuasion with ChatGPT (Chu and Liu, 2023)	2191
		2192
	• Benchmarking Foundation Models with Language-Model-as-an-Examiner (Bai et al., 2023b)	2193
		2194
	• Chatgpt: A meta-analysis after 2.5 months (Leiter et al., 2023)	2195
		2196
	• Causal-discovery performance of chatgpt in the context of neuropathic pain diagnosis (Tu et al., 2023a)	2197
		2198
	• Summary of ChatGPT-Related Research and Perspective Towards the Future of Large Language Models (Liu et al., 2023h)	2199
		2200
	• Harnessing the power of llms in practice: A survey on chatgpt and beyond (Yang et al., 2023a)	2201
		2202
	• Holistic evaluation of language models (Liang et al., 2023a)	2203
		2204
	• A Survey on Evaluation of Large Language Models (Chang et al., 2023)	2205
		2206
	• Can chatgpt understand too? a comparative study on chatgpt and fine-tuned bert (Zhong et al., 2023a)	2207
		2208
	• On the robustness of chatgpt: An adversarial and out-of-distribution perspective (Wang et al., 2023e)	2209
		2210

2182	• An Independent Evaluation of ChatGPT on Mathematical Word Problems (MWP) (Shakarian et al., 2023)	2224
2183		2225
2184		2226
2185	• Can ChatGPT Replace Traditional KBQA Models? An In-depth Analysis of the Question Answering Performance of the GPT LLM Family (Tan et al., 2023b)	2227
2186		2228
2187		2229
2188		2230
2189	• Instruction tuning with gpt-4 (Peng et al., 2023b)	2231
2190		2232
2191	• How Robust is GPT-3.5 to Predecessors? A Comprehensive Study on Language Understanding Tasks (Chen et al., 2023d)	2233
2192		2234
2193		2235
2194	• Consistency analysis of chatgpt (Jang and Lukasiewicz, 2023)	2236
2195		2237
2196	• Does ChatGPT resemble humans in language use? (Cai et al., 2023)	2238
2197		2239
2198	• A comprehensive capability analysis of gpt-3 and gpt-3.5 series models (Ye et al., 2023)	2240
2199		2241
2200	• A comprehensive evaluation of ChatGPT's zero-shot Text-to-SQL capability (Liu et al., 2023a)	2242
2201		2243
2202		2244
2203	• Is ChatGPT a good sentiment analyzer? A preliminary study (Wang et al., 2023j)	2245
2204		2246
2205	• A preliminary evaluation of chatgpt for zero-shot dialogue understanding (Pan et al., 2023)	2247
2206		2248
2207	• Zero-shot Temporal Relation Extraction with ChatGPT (Yuan et al., 2023a)	2249
2208		2250
2209	• Can chatgpt reproduce human-generated labels? a study of social computing tasks (Zhu et al., 2023)	2251
2210		2252
2211		2253
2212	• Chatgraph: Interpretable text classification by converting chatgpt knowledge to graphs (Shi et al., 2023)	2254
2213		2255
2214		2256
2215	• Uncovering the Potential of ChatGPT for Discourse Analysis in Dialogue: An Empirical Study (Fan and Jiang, 2023)	2257
2216		2258
2217		2259
2218	• Evaluating ChatGPT's Performance for Multilingual and Emoji-based Hate Speech Detection (Das et al., 2023)	2260
2219		2261
2220		2262
2221	• Sentiment Analysis in the Era of Large Language Models: A Reality Check (Zhang et al., 2023e)	2263
2222		2264
2223		
	• ChatGPT for Suicide Risk Assessment on Social Media: Quantitative Evaluation of Model Performance, Potentials and Limitations (Ghanadian et al., 2023)	
	• Metacognitive Prompting Improves Understanding in Large Language Models (Wang and Zhao, 2023)	
	• Exploring ai ethics of chatgpt: A diagnostic analysis (Zhuo et al., 2023)	
	• Is chatgpt better than human annotators? potential and limitations of chatgpt in explaining implicit hate speech (Huang et al., 2023a)	
	• Investigating Chain-of-thought with ChatGPT for Stance Detection on Social Media (Zhang et al., 2023a)	
	• The Self-Perception and Political Biases of ChatGPT (Rutinowski et al., 2023)	
	• Not All Languages Are Created Equal in LLMs: Improving Multilingual Capability by Cross-Lingual-Thought Prompting (Huang et al., 2023b)	
	• BAD: BiAs Detection for Large Language Models in the context of candidate screening (Koh et al., 2023)	
	• DECODINGTRUST: A Comprehensive Assessment of Trustworthiness in GPT Models (Wang et al., 2023b)	
	• On Large Language Models' Selection Bias in Multi-Choice Questions (Zheng et al., 2023a)	
	• Exploring the Limits of ChatGPT for Query or Aspect-based Text Summarization (Yang et al., 2023b)	
	• ChatGPT as a Factual Inconsistency Evaluator for Text Summarization (Luo et al., 2023)	
	• Uncovering ChatGPT's Capabilities in Recommender Systems (Dai et al., 2023b)	
	• HaluEval: A Large-Scale Hallucination Evaluation Benchmark for Large Language Models (Li et al., 2023d)	
	• RecurrentGPT: Interactive Generation of (Arbitrarily) Long Text (Zhou et al., 2023)	

2265	• PokemonChat: Auditing ChatGPT for Pokémon Universe Knowledge (Cabello et al., 2023)	2307
2266		2308
2267		2309
2268	• Hybrid Long Document Summarization using C2F-FAR and ChatGPT: A Practical Study (Lu et al., 2023a)	2310
2269		2311
2270		2312
2271	• Generative Job Recommendations with Large Language Model (Zheng et al., 2023d)	2313
2272		2314
2273	• Assessing the Ability of ChatGPT to Screen Articles for Systematic Reviews (Syriani et al., 2023)	2315
2274		2316
2275		2317
2276	• L-Eval: Instituting Standardized Evaluation for Long Context Language Models (An et al., 2023)	2318
2277		2319
2278		2320
2279	• LLM-Rec: Personalized Recommendation via Prompting Large Language Models (Lyu et al., 2023a)	2321
2280		2322
2281		2323
2282	• LongBench: A Bilingual, Multitask Benchmark for Long Context Understanding (Bai et al., 2023a)	2324
2283		2325
2284		
2285	• Can Large Language Models Understand Real-World Complex Instructions (He et al., 2023b)	2326
2286		2327
2287		2328
2288	• Mind meets machine: Unravelling GPT-4's cognitive psychology (Dhingra et al., 2023)	2329
2289		2330
2290	• Capabilities of GPT-4 on Medical Challenge Problems (Nori et al., 2023)	2331
2291		2332
2292		2333
2292	• GPT is becoming a Turing machine: Here are some ways to program it (Jojic et al., 2023)	2334
2293		2335
2294	• ChatGPT is a Knowledgeable but Inexperienced Solver: An Investigation of Commonsense Problem in Large Language Models (Bian et al., 2023)	2336
2295		2337
2296		2338
2297		2339
2298	• Humans in Humans Out: On GPT Converging Toward Common Sense in both Success and Failure (Koralus and Wang-Maścianica, 2023)	2340
2299		2341
2300		2342
2301	• Generative Models as a Complex Systems Science: How can we make sense of large language model behavior? (Holtzman et al., 2023)	2343
2302		2344
2303		2345
2304		2346
2305	• Uncertainty in Natural Language Generation: From Theory to Applications (Baan et al., 2023)	2347
2306		2348
	• ChatGPT-Crawler: Find out if ChatGPT really knows what it's talking about (Rangapur and Wang, 2023)	
	• How is ChatGPT's behavior changing over time? (Chen et al., 2023b)	
	• ChatABL: Abductive Learning via Natural Language Interaction with ChatGPT (Zhong et al., 2023b)	
	• Look Before You Leap: An Exploratory Study of Uncertainty Measurement for Large Language Models (Huang et al., 2023d)	
	• StructGPT: A General Framework for Large Language Model to Reason over Structured Data (Jiang et al., 2023)	
	• Chain-of-Symbol Prompting Elicits Planning in Large Language Models (Hu et al., 2023a)	
	• Tree of Thoughts: Deliberate Problem Solving with Large Language Models (Yao et al., 2023)	
	• ChatLog: Recording and Analyzing ChatGPT Across Time (Tu et al., 2023b)	
	• Chain of Knowledge: A Framework for Grounding Large Language Models with Structured Knowledge Bases (Li et al., 2023f)	
	• ChatHaruhi: Reviving Anime Character in Reality via Large Language Model. (Li et al., 2023b)	
	• Adaptive Chameleon or Stubborn Sloth: Unraveling the Behavior of Large Language Models in Knowledge Conflicts (Xie et al., 2023a)	
	• GPT-3.5 vs GPT-4: Evaluating ChatGPT's Reasoning Performance in Zero-shot Learning (Espejel et al., 2023)	
	• LogiCoT: Logical Chain-of-Thought Instruction-Tuning Data Collection with GPT-4 (Liu et al., 2023d)	
	• Enabling Large Language Models to Generate Text with Citations (Gao et al., 2023d)	

2349	• Why Does ChatGPT Fall Short in Providing Truthful Answers? (Zheng et al., 2023b)	2392
2350		2393
2351	• Are Large Language Models Really Good Logical Reasoners? A Comprehensive Evaluation From Deductive, Inductive and Abductive Views (Xu et al., 2023a)	2394
2352		2395
2353		2396
2354		2397
2355	• Investigating the Factual Knowledge Boundary of Large Language Models with Retrieval Augmentation (Ren et al., 2023)	2398
2356		2399
2357		2400
2358		2401
2359	• SciBench: Evaluating College-Level Scientific Problem-Solving Abilities of Large Language Models (Wang et al., 2023h)	2402
2360		2403
2361		2404
2362	• Think-on-Graph: Deep and Responsible Reasoning of Large Language Model with Knowledge Graph (Sun et al., 2023a)	2405
2363		2406
2364		2407
2365	• Skills-in-Context Prompting: Unlocking Compositionality in Large Language Models (Chen et al., 2023a)	2408
2366		2409
2367		2410
2368	• Large Language Models Sensitivity to The Order of Options in Multiple-Choice Questions (Pezeshkpour and Hruschka, 2023)	2411
2369		2412
2370		2413
2371	• Large Language Models on the Chessboard: A Study on ChatGPT's Formal Language Comprehension and Complex Reasoning Skills (Kuo et al., 2023)	2414
2372		2415
2373		2416
2374		2417
2375	• Evaluating Large Language Models on Graphs: Performance Insights and Comparative Analysis (Liu and Wu, 2023)	2418
2376		2419
2377		2420
2378	• Exploring New Frontiers in Agricultural NLP: Investigating the Potential of Large Language Models for Food Applications (Rezayi et al., 2023)	2421
2379		2422
2380		2423
2381		2424
2382	• MINT: Evaluating LLMs in Multi-turn Interaction with Tools and Language Feedback (Wang et al., 2023i)	2425
2383		2426
2384		2427
2385	• ChatGPT for Robotics: Design Principles and Model Abilities (Vemprala et al., 2023)	2428
2386		2429
2387	• Error Analysis Prompting Enables Human-Like Translation Evaluation in Large Language Models: A Case Study on ChatGPT (Lu et al., 2023c)	2430
2388		2431
2389		2432
2390		2433
2391	• Towards Making the Most of ChatGPT for Machine Translation (Peng et al., 2023c)	2434
	• Linguistically Informed ChatGPT Prompts to Enhance Japanese-Chinese Machine Translation: A Case Study on Attributive Clauses (Gu, 2023)	
	• How to Design Translation Prompts for ChatGPT: An Empirical Study (Gao et al., 2023e)	
	• Is Your Code Generated by ChatGPT Really Correct? Rigorous Evaluation of Large Language Models for Code Generation (Liu et al., 2023e)	
	• Evaluating ChatGPT's Information Extraction Capabilities: An Assessment of Performance, Explainability, Calibration, and Faithfulness (Li et al., 2023a)	
	• Can ChatGPT Assess Human Personalities? A General Evaluation Framework (Rao et al., 2023b)	
	• Is ChatGPT Equipped with Emotional Dialogue Capabilities? (Zhao et al., 2023a)	
	• AugGPT: Leveraging ChatGPT for Text Data Augmentation (Dai et al., 2023a)	
	• Evaluating GPT-4 and ChatGPT on Japanese Medical Licensing Examinations (Kasai et al., 2023)	
	• To ChatGPT, or not to ChatGPT: That is the question! (Pegoraro et al., 2023)	
	• In ChatGPT We Trust? Measuring and Characterizing the Reliability of ChatGPT (Shen et al., 2023)	
	• Analyzing ChatGPT's Aptitude in an Introductory Computer Engineering Course (Deshpande and Szefer, 2023)	
	• Linguistic ambiguity analysis in ChatGPT (Ortega-Martín et al., 2023)	
	• Using ChatGPT for Entity Matching (Peeters and Bizer, 2023)	
	• MEGA: Multilingual Evaluation of Generative AII (Ahuja et al., 2023)	
	• Evaluation of ChatGPT on Biomedical Tasks: A Zero-Shot Comparison with Fine-Tuned Generative Transformers (Jahan et al., 2023)	
	• Zero-shot Clinical Entity Recognition using ChatGPT (Hu et al., 2023c)	

2435	• Human-like Summarization Evaluation with ChatGPT (Gao et al., 2023c)	Elicitation, and Software Design (White et al., 2023)	2478
2436			2479
2437	• "HOT" ChatGPT: The promise of ChatGPT in detecting and discriminating hateful, offensive, and toxic comments on social media (Li et al., 2023e)	• SelfCheckGPT: Zero-Resource Black-Box Hallucination Detection for Generative Large Language Models (Manakul et al., 2023)	2480
2438			2481
2439			2482
2440		• Translating Radiology Reports into Plain Language using ChatGPT and GPT-4 with Prompt Learning: Promising Results, Limitations, and Potential (Lyu et al., 2023b)	2483
2441	• Large language models effectively leverage document-level context for literary translation, but critical errors persist (Karpinska and Iyyer, 2023)		2484
2442			2485
2443			2486
2444		• An Empirical Study of Pre-trained Language Models in Simple Knowledge Graph Question Answering (Hu et al., 2023b)	2487
2445	• BayLing: Bridging Cross-lingual Alignment and Instruction Following through Interactive Translation for Large Language Models (Zhang et al., 2023d)		2488
2446			2489
2447		• DeID-GPT: Zero-shot Medical Text De-Identification by GPT-4 (Liu et al., 2023j)	2490
2448			2491
2449	• Neural Machine Translation Data Generation and Augmentation using ChatGPT (Yang and Nicolai, 2023)	• MM-REACT: Prompting ChatGPT for Multimodal Reasoning and Action (Yang et al., 2023c)	2492
2450			2493
2451			2494
2452	• ChatGPT MT: Competitive for High- (but not Low-) Resource Languages (Robinson et al., 2023)	• Is ChatGPT A Good Keyphrase Generator? A Preliminary Study (Song et al., 2023)	2495
2453			2496
2454		• Beyond Black Box AI-Generated Plagiarism Detection: From Sentence to Document Level (Quidwai et al., 2023)	2497
2455	• How would Stance Detection Techniques Evolve after the Launch of ChatGPT? (Zhang et al., 2022)		2498
2456			2499
2457		• Comparing Abstractive Summaries Generated by ChatGPT to Real Summaries Through Blinded Reviewers and Text Classification Algorithms (Soni and Wade, 2023)	2500
2458	• BHASA: A Holistic Southeast Asian Linguistic and Cultural Evaluation Suite for Large Language Models (Leong et al., 2023)		2501
2459			2502
2460			2503
2461	• Zero-shot Approach to Overcome Perturbation Sensitivity of Prompts (Chakraborty et al., 2023)	• MGTBench: Benchmarking Machine-Generated Text Detection (He et al., 2023c)	2504
2462			2505
2463		• Neural Theory-of-Mind? On the Limits of Social Intelligence in Large LMs (Sap et al., 2023)	2506
2464	• UZH_CLyp at SemEval-2023 Task 9: Head-First Fine-Tuning and ChatGPT Data Generation for Cross-Lingual Learning in Tweet Intimacy Prediction (Michail et al., 2023)		2507
2465			2508
2466		• What would Harry say? Building Dialogue Agents for Characters in a Story (Chen et al., 2023c)	2509
2467			2510
2468	• Exploring the Feasibility of ChatGPT for Event Extraction (Gao et al., 2023b)		2511
2469		• Reliable Natural Language Understanding with Large Language Models and Answer Set Programming (Rajasekharan et al., 2023)	2512
2470	• Does Synthetic Data Generation of LLMs Help Clinical Text Mining? (Tang et al., 2023b)		2513
2471			2514
2472		• ByGPT5: End-to-End Style-conditioned Poetry Generation with Token-free Language Models (Belouadi and Eger, 2023)	2515
2473	• ICL-D3IE: In-Context Learning with Diverse Demonstrations Updating for Document Information Extraction (He et al., 2023a)		2516
2474			2517
2475		• Transformers Go for the LOLs: Generating (Humorous) Titles from Scientific Abstracts End-to-End (Chen and Eger, 2022)	2518
2476	• ChatGPT Prompt Patterns for Improving Code Quality, Refactoring, Requirements		2519
2477			2520

2521	• Modeling Label Semantics Improves Activity Recognition (Zhang et al., 2023f)	2566
2522		2567
2523	• An Analysis of the Automatic Bug Fixing Performance of ChatGPT (Sobania et al., 2023)	2568
2524		2569
2525	• Chat2VIS: Generating Data Visualisations via Natural Language using ChatGPT, Codex and GPT-3 Large Language Models (Maddigan and Susnjak, 2023)	2570
2526		2571
2527		
2528		
2529	• ChatGPT versus Traditional Question Answering for Knowledge Graphs: Current Status and Future Directions Towards Knowledge Graph Chatbots (Omar et al., 2023)	2572
2530		2573
2531		2574
2532		
2533	• ChatCAD: Interactive Computer-Aided Diagnosis on Medical Image using Large Language Models (Wang et al., 2023g)	2575
2534		2576
2535		2577
2536	• Guiding Large Language Models via Directional Stimulus Prompting (Li et al., 2023g)	2578
2537		2579
2538		2580
2539	• Utilizing ChatGPT Generated Data to Retrieve Depression Symptoms from Social Media (Bucur, 2023)	2581
2540		2582
2541		2583
2542	• Check Your Facts and Try Again: Improving Large Language Models with External Knowledge and Automated Feedback (Peng et al., 2023a)	2584
2543		2585
2544		2586
2545	• ChatGPT for Zero-shot Dialogue State Tracking: A Solution or an Opportunity? (Heck et al., 2023)	2587
2546		2588
2547		2589
2548	• ChatGPT vs Human-authored Text: Insights into Controllable Text Summarization and Sentence Style Transfer (Pu and Demberg, 2023)	2590
2549		2591
2550		2592
2551		2593
2552	• Not The End of Story: An Evaluation of ChatGPT-Driven Vulnerability Description Mappings (Liu et al., 2023f)	2594
2553		2595
2554		2596
2555	• ChatGPT is not a good indigenous translator (Stap and Araabi, 2023)	2597
2556		2598
2557		2599
2558	• You've Got a Friend in ... a Language Model? A Comparison of Explanations of Multiple-Choice Items of Reading Comprehension between ChatGPT and Humans (Duenas et al., 2023)	2600
2559		2601
2560		2602
2561		2603
2562	• Language-Agnostic Transformers and Assessing ChatGPT-Based Query Rewriting for Multilingual Document-Grounded QA (Gowriraj et al., 2023)	2604
2563		2605
2564		2606
2565		2607
	• HW-TSC at SemEval-2023 Task 7: Exploring the Natural Language Inference Capabilities of ChatGPT and Pre-trained Language Model for Clinical Trial (Zhao et al., 2023b)	2566
		2567
		2568
		2569
	• Can ChatGPT Understand Causal Language in Science Claims? (Kim et al., 2023)	2570
		2571
	• ChatGPT is fun, but it is not funny! Humor is still challenging Large Language Models (Jentsch and Kersting, 2023)	2572
		2573
		2574
	• Assessing Cross-Cultural Alignment between ChatGPT and Human Societies: An Empirical Study (Cao et al., 2023)	2575
		2576
		2577
	• Evaluating Reading Comprehension Exercises Generated by LLMs: A Showcase of ChatGPT in Education Applications (Xiao et al., 2023)	2578
		2579
		2580
		2581
	• ChatGPT vs. Crowdsourcing vs. Experts: Annotating Open-Domain Conversations with Speech Functions (Ostyakova et al., 2023)	2582
		2583
		2584
	• Leveraging Large Language Models for Automated Dialogue Analysis (Finch et al., 2023)	2585
		2586
	• Breaking the Bank with ChatGPT: Few-Shot Text Classification for Finance (Liang et al., 2023b)	2587
		2588
		2589
	• Is ChatGPT a Good Teacher Coach? Measuring Zero-Shot Performance For Scoring and Providing Actionable Insights on Classroom Instruction (Wang and Demszky, 2023)	2590
		2591
		2592
		2593
	• Credible Without Credit: Domain Experts Assess Generative Language Models (Peskoff and Stewart, 2023)	2594
		2595
		2596
	• SUT at SemEval-2023 Task 1: Prompt Generation for Visual Word Sense Disambiguation (Ghahroodi et al., 2023)	2597
		2598
		2599
	• How Well Do Large Language Models Perform on Faux Pas Tests? (Shapira et al., 2023)	2600
		2601
	• Zero-Shot Information Extraction for Clinical Meta-Analysis using Large Language Models (Kartchner et al., 2023)	2602
		2603
		2604
	• Can Large Language Models Be an Alternative to Human Evaluation? (Chiang and Lee, 2023)	2605
		2606
		2607

2608	• RL4F: Generating Natural Language Feedback with Reinforcement Learning for Repairing Model Outputs (Akyurek et al., 2023)	2651
2609		2652
2610		2653
2611	• Toward Human-Like Evaluation for Natural Language Generation with Error Analysis (Lu et al., 2023b)	2654
2612		2655
2613		2656
2614	• Multi-Level Knowledge Distillation for Out-of-Distribution Detection in Text (Wu et al., 2023b)	2657
2615		2658
2616		2659
2617	• MolXPT: Wrapping Molecules with Text for Generative Pre-training (Liu et al., 2023i)	2660
2618		2661
2619	• MUX-PLMs: Pre-training Language Models with Data Multiplexing (Murahari et al., 2023)	2662
2620		2663
2621	• Is GPT-3 a Good Data Annotator? (Ding et al., 2023)	2664
2622		2665
2623	• Language Generation Models Can Cause Harm: So What Can We Do About It? An Actionable Survey (Kumar et al., 2023)	2666
2624		2667
2625		2668
2626	• Detoxifying Online Discourse: A Guided Response Generation Approach for Reducing Toxicity in User-Generated Text (Bose et al., 2023)	2669
2627		2670
2628		
2629		
2630	• Faithful Question Answering with Monte-Carlo Planning (Hong et al., 2023)	2671
2631		2672
2632	• Nut-cracking Sledgehammers: Prioritizing Target Language Data over Bigger Language Models for Cross-Lingual Metaphor Detection (Schuster and Markert, 2023)	2673
2633		2674
2634		
2635		
2636	• Generating Faithful Text From a Knowledge Graph with Noisy Reference Text (Hashem et al., 2023)	2675
2637		2676
2638		2677
2639	• Empowering Conversational Agents using Semantic In-Context Learning (Omidvar and An, 2023)	2678
2640		2679
2641		
2642	• Can Large Language Models Safely Address Patient Questions Following Cataract Surgery? (Chowdhury et al., 2023)	2680
2643		2681
2644		
2645	• Towards Benchmarking and Improving the Temporal Reasoning Capability of Large Language Models (Tan et al., 2023a)	2682
2646		2683
2647		2684
2648	• Generative Pretrained Transformers for Emotion Detection in a Code-Switching Setting (Nedilko, 2023)	2685
2649		2686
2650		2687
		2688
		2689
		2690
		2691
		2692
	• On the Underspecification of Situations in Open-domain Conversational Datasets (Otani et al., 2023)	
	• Understanding Factual Errors in Summarization: Errors, Summarizers, Datasets, Error Detectors (Tang et al., 2023a)	
	• Multilingual Language Models are not Multicultural: A Case Study in Emotion (Havaladar et al., 2023)	
	• Good Data, Large Data, or No Data? Comparing Three Approaches in Developing Research Aspect Classifiers for Biomedical Papers (Chandrasekhar et al., 2023)	
	• IDOL: Indicator-oriented Logic Pre-training for Logical Reasoning (Xu et al., 2023c)	
	• PrecogIIITH@WASSA2023: Emotion Detection for Urdu-English Code-mixed Text (Vedula et al., 2023)	
	• Do PLMs Know and Understand Ontological Knowledge? (Wu et al., 2023c)	
	• Evaluation of Question Generation Needs More References (Oh et al., 2023)	
	• A System for Answering Simple Questions in Multiple Languages (Razzhigaev et al., 2023)	
	• HIT-SCIR at WASSA 2023: Empathy and Emotion Analysis at the Utterance-Level and the Essay-Level (Lu et al., 2023d)	
	• Pre-trained Language Models Can be Fully Zero-Shot Learners (Zhao et al., 2023c)	
	• Predicting the Quality of Revisions in Argumentative Writing (Liu et al., 2023k)	
	• Exploring Effectiveness of GPT-3 in Grammatical Error Correction: A Study on Performance and Controllability in Prompt-Based Methods (Loem et al., 2023)	
	• When Truth Matters – Addressing Pragmatic Categories in Natural Language Inference (NLI) by Large Language Models (LLMs) (Gubelmann et al., 2023)	
	• Debiasing should be Good and Bad: Measuring the Consistency of Debiasing Techniques in Language Models (Morabito et al., 2023)	

- 2693 • GPoeT: a Language Model Trained for Rhyme
2694 Generation on Synthetic Data (Popescu-Belis
2695 et al., 2023)
- 2696 • Examining Bias in Opinion Summarisation
2697 Through the Perspective of Opinion Diver-
2698 sity (Huang et al., 2023c)
- 2699 • Unsupervised Summarization Re-
2700 ranking (Ravaut et al., 2023)
- 2701 • Improving Dutch Vaccine Hesitancy Monitor-
2702 ing via Multi-Label Data Augmentation with
2703 GPT-3.5 (Van Nooten and Daelemans, 2023)
- 2704 • What Makes a Good Counter-Stereotype?
2705 Evaluating Strategies for Automated Re-
2706 sponses to Stereotypical Text (Fraser et al.,
2707 2023)
- 2708 • What Makes Good Counterspeech? A Com-
2709 parison of Generation Approaches and Evalu-
2710 ation Metrics (Zheng et al., 2023c)
- 2711 • DAMO-NLP at SemEval-2023 Task 2: A Uni-
2712 fied Retrieval-augmented System for Multi-
2713 lingual Named Entity Recognition (Tan et al.,
2714 2023c)
- 2715 • UMASS_BioNLP at MEDIQA-Chat 2023:
2716 Can LLMs generate high-quality synthetic
2717 note-oriented doctor-patient conversa-
2718 tions? (Wang et al., 2023f)
- 2719 • Frontier Review of Multimodal AI (Nan,
2720 2023)
- 2721 • Is ChatGPT a Highly Fluent Grammatical Er-
2722 ror Correction System? A Comprehensive
2723 Evaluation (Fang et al., 2023)
- 2724 • DialogStudio: Towards Richest and Most Di-
2725 verse Unified Dataset Collection for Conver-
2726 sational AI (Zhang et al., 2023c)

2727 **B Detailed List of ChatGPT Data Leak**

2728 We show which datasets have been leaked to Chat-
2729 GPT in Tables 4 and 5.

Task name	Lo	M-Lo	M-Hi	Hi
AI safety & ethics			bbq (all), bold (all)	
Creative text generation	writingprompts (test)			
Dialogue	opendialkg (test), prosocialdialog (test)	multiwoz 2.2 (test)		dstc11 track 5 (dev), dstc7 track 2 (all), multiwoz 2.1 (test), multiwoz 2.4 (test), mutual (test)
Evaluation of generated texts				newsroom (all), openneva-roc (all), realsumm (all), summeval (all)
Machine Translation	flores-101 (test), wmt20 (test), wmt20 en-de (test), wmt20 robustness task set 2 en->ja (test), wmt20 robustness task set 3 (test), wmt20 zh-en (test), wmt22 (test)	nusax (test), wmt19 biomedical translation task (test), wmt20 robustness task set 2 ja->en (test), wmt2014 ott et al 2018 (test)	flores-200 (dev)	mqm-2022 (test)
Math		numersense (dev)		addsub (all), aqua-rat (test), draw-1k (all), ghosts (all), gsm8k (test), multiarith (all), singleeq (all), svamp (all)
Medical text generation	ddxplus (en) (test), mimic (test)			merck sharpe & dohme (msd) clinical manual (all)
Natural Language Inference	becel-snli (test), commitmentbank (all), mnli (dev), qnli (dev), rte (all), onli (dev)	becel-rte (test), entailmentbank (test)		med (test), advglue-mnli (dev), advglue-qnli (dev), advglue-rte (dev), anli-r3 (test), ax-g (dev), cb (dev), conjnli (test), control "logical reasoning" (test), help (test), mnli (test), rte (dev), semeval-2023 task 7 (all), taxinli (test), wnli (dev)
Natural Language Understanding				atis (test), snips (test)
Paraphrasing	mrpc (dev), qqp (dev)			
Politics		covid-19-lee (test)		p-stance (test), semeval-2016 task 6 (test), tweeteval - tweetstance (test)
Programming				quixbugs (all)
Psychology				myers-briggs type indicator (all)

Table 4: The names of datasets with low (Lo), moderate-low (M-Lo), moderate-high (M-Hi), and high (Hi) leakage, categorized according to the task. (1/2)

Task name	Lo	M-Lo	M-Hi	Hi
Question answering	amboss medical questions (all), babi-task16 (test), clutrr (test), e-care (dev), financezhidao (all), freebaseqa (all), hotpotqa (dev), lquad-2.0 (all), legalqa (all), logiqa (all), math (test), mc-taco (dev), medical dialog (all), medical dialog - zh (all), mkqa (all), pep-3k (all), piqa (test), reclor (all), simplequestions (all), spartqa (test), stepgame (test), ti-medial (test), webquestionssp (all), web-textqa + baikeqa (all)	babi-task15 (test), bsc self-assessment program (all), chinese psychological question answering dataset (all), eli5 (dev), nlpcc-dbqa (all), openbookqa (dev), ophthoquestions (all), piqa (dev), qasc (dev), race - reading (test or dev), social iqa (dev), squad v2 (dev), truth-fulqa gen (test), wikiqa (all)	fiqa (all), gsm8k - mathqa (test), kqapro (all), open-bookqa (test), usmle medical q (all)	advglove-qqp (dev), ar-1sat (test), baidubaik15 (all), boolq (test), boolq-contrast (test), bron (preprocessed) (all), cve-2021 (all), cve-att&ck (all), cwq (all), dblp-bench (all), efficientqa (dev), graiqa (test), graphquestions (all), hc3-chinese (all), hc3-english (all), hofstede cultural survey (all), lc-quad2.0 (all), logiqa 2.0 (test), logiqa 2.0-zh (test), mag-bench (all), nbme medical questions (all), ott-qa (all), protoqa (dev), qald-9 (all), qald-9 (test), reclor (dev), truthfulqa (test), tuce test (all), wiki-csai (all), wqsp (all), yago-bench (all)
Reasoning & common sense	commonsenseqa (test), hellaswag (dev), letter string analogies - webb (all)	arc (dev), coin flip (all), copa (dev), wsc (dev)		cola (dev), commonsenseqa (dev), date (all), last letter (all), matres (test), object (all), strategyqa (all), tddiscourse (test), timebank-dense (test)
Semantic similarity	sts-b (dev), tweeteval - tweetemoji (test or dev)	bece1-mrpc (test)		raganato - wsd (test or dev), wic (dev), wic - wordcontext (test or dev)
Sentiment analysis	colbert (test or dev), flipkart (all), imdb (test), sst-2 (dev), unhealthy conv. - unhealthy (test or dev), unhealthy conv. - unhealthyper (test or dev), wikidetox aggr. - aggression (test or dev), wikidetox aggr. - aggressionper (test or dev)	goemotions - goemo (test or dev), goemotions - goemoper0 (test or dev), goemotions - goemoper1 (test or dev), goemotions - goemoper2 (test or dev), goemotions - goemoper3 (test or dev), la-tenthatred (all), sarcasmania - sarcasm (test or dev), semeval-2023 task 9 (test), tweeteval - tweetsent (test or dev)	realtoxicprompts (all)	advglove-sst-2 (dev), clarinemo (test or dev), first impressions (all), imdb-contrast (all), polemo2 - polemo (test or dev), sentiment140 (all), sst2 (dev), the suicide and depression dataset (all)
Summarization	cnn dailymail (test), crosssum (test), reddit (test), wikilingua (test), xsamsum (test)	cnn/dm (test), covidet (test), newts (test), pubmed (test), qmsum (test), xsum (test)	squality (test)	samsum (test)
Text classification	inverse scaling challenge (11 datasets) (all)			pubmed20k (train), sms spam v.1 (test or dev), symptoms (on kagggle) (train)
Text extraction	ace05 (test), mtsamples (all)	i2b2 2010 (all)		ace05 (all), conll++ (all), conll03 (test), duee1.0 (all), duee2.0 (all), msra (all), nyt11-hrl (all)
Text generation		conll2014 (test)		adveta(add)-spider (dev), adveta(rpl)-spider (dev), cosql (dev), espider (dev), dusql (all), quiz-design (all), sparc (dev), spider (dev), spider-cg(app) (all), spider-cg(sub) (all), spider-dk (dev), spider-realistic (dev), spider-syn (dev)

Table 5: The names of datasets with low (Lo), moderate-low (M-Lo), moderate-high (M-Hi), and high (Hi) leakage, categorized according to the task. (2/2)