

SafetyPrompts: a Systematic Review of Open Datasets for Evaluating and Improving Large Language Model Safety

Paul Röttger[✉] Fabio Pernisi[✉] Bertie Vidgen[✉] Dirk Hovy[✉]
✉ Bocconi University ✉ University of Oxford

Abstract

The last two years have seen a rapid growth in concerns around the safety of large language models (LLMs). Researchers and practitioners have met these concerns by introducing an abundance of new datasets for evaluating and improving LLM safety. However, much of this work has happened in parallel, and with very different goals in mind, ranging from the mitigation of near-term risks around bias and toxic content generation to the assessment of longer-term catastrophic risk potential. This makes it difficult for researchers and practitioners to find the most relevant datasets for a given use case, and to identify gaps in dataset coverage that future work may fill. To remedy these issues, we conduct a first systematic review of open datasets for evaluating and improving LLM safety. We review 102 datasets, which we identified through an iterative and community-driven process over the course of several months. We highlight patterns and trends, such as a trend towards fully synthetic datasets, as well as gaps in dataset coverage, such as a clear lack of non-English datasets. We also examine how LLM safety datasets are used in practice – in LLM release publications and popular LLM benchmarks – finding that current evaluation practices are highly idiosyncratic and make use of only a small fraction of available datasets. Our contributions are based on SafetyPrompts.com, a living catalogue of open datasets for LLM safety, which we commit to updating continuously as the field of LLM safety develops.

1 Introduction

Ensuring that large language models (LLMs) are safe has become as a key priority for model developers and regulators. Consequently, in recent years, researchers and practitioners have created an abundance of new datasets for evaluating and improving LLM safety. Safety, however, is a multi-faceted and contextual concept that lacks a unifying definition. This complexity is reflected in the current landscape of safety datasets, which is broad, diverse, and fast-moving. In just the first two months of 2024, for example, researchers published datasets for evaluating near-term risks from LLMs, such as sociodemographic bias (Gupta et al., 2024) and toxic content generation (Bianchi et al., 2024), as well as datasets for evaluating long-term societal risk potential, around power-seeking (Mazeika et al., 2024) and sycophantic behaviours (Sharma et al., 2024). The rapid pace of dataset creation and the variety of purposes served by different datasets make it difficult for researchers and practitioners to find the most relevant datasets for different use cases, and to identify gaps in dataset coverage that future work may fill.

In this paper, we address these issues by conducting a first systematic review of open datasets for evaluating and improving LLM safety. We identify 102 datasets published between June 2018 and February 2024 based on clear inclusion criteria (§2.1) using a comprehensive community-driven search method (§2.2). We examine these 102 datasets along several key dimensions, including their purpose (§3.2), creation (§3.4), format and size (§3.3), access and licensing (§3.6), and publication (§3.7). Key findings of our review include that dataset creation is currently growing at an unprecedented rate, driven primarily by academic and non-profit organisations; that there is a trend towards more specialised safety evaluations and the use of synthetic data; and that the English language dominates the

dataset landscape. We also review how open LLM safety datasets are used in practice – in model release publications (§4) and popular LLM benchmarks (§5). We find that current evaluation practices are highly idiosyncratic and leverage only a small fraction of available datasets. In our Discussion (§6), we argue that this creates clear scope for standardisation in LLM safety evaluations, and that evaluations in general could be improved by better leveraging recent progress in safety dataset creation.

2 Dataset review methodology

2.1 Inclusion criteria

At a high level, we restrict our review to *open* datasets that are relevant to *LLMs*, and specifically to evaluating and improving *LLM safety*.

In terms of **data modality**, we only include text datasets. We do not include image datasets (e.g. Schwemmer et al., 2020; Zhao et al., 2021; Ricker et al., 2022) or audio datasets (e.g. Reimao & Tzerpos, 2019; Koenecke et al., 2020; Meyer et al., 2020). We also do not include datasets targeted at multimodal models, even if one modality is text, such as in the case of vision-language (e.g. Carlini et al., 2023; Hall et al., 2023; Wolfe et al., 2023) or text-to-image models (e.g. Bianchi et al., 2023; Parrish et al., 2023; Luccioni et al., 2024). We do not include datasets targeted at code generation models (e.g. Siddiq & Santos, 2022; Bhatt et al., 2023). These modalities and models constitute natural expansions for future work.

We make only minimal restrictions in terms of **data format**. Real-world user interactions with LLMs usually take the form of text chat (Ouyang et al., 2023; Zheng et al., 2024; Zhao et al., 2024), so we are most interested in datasets that naturally fit a chat format, like open-ended questions and instructions, but we also consider any other dataset that can meaningfully be expressed in a prompt format. This includes multiple-choice questions or autocomplete-style text snippets. We do not make any restrictions on language.

For **data access**, we only include datasets that are available for download via GitHub and/or Hugging Face. In practice, we found that, if data is made available, it is almost always on one or both of these platforms. We do not make restrictions based on how data is licensed.

Finally, we require that all datasets are **relevant to safety**. For the purposes of our review, we adopt a wide and open definition of safety. Broadly speaking, we include datasets that relate to representational, political or other forms of sociodemographic bias; to toxicity, malicious instructions or harmful advice; to hazardous behaviours like sycophancy or power-seeking; to alignment with social, moral or ethical values; or to adversarial LLM usage (e.g. red-teaming, jailbreaking, prompt hacking). We only include datasets that explicitly focus on (some of these aspects of) LLM safety. We do not include datasets that target general LLM capabilities like reasoning, language understanding, or code completion (e.g. Dua et al., 2019; Hendrycks et al., 2020b; Chen et al., 2021). We also do not include datasets that target factuality in LLMs, unless they directly relate to safety, like in the case of generating misinformation (Souly et al., 2024) or measuring truthfulness (Lin et al., 2022).

The cutoff date for our review is March 1st, 2024. We did not include datasets that were first published after this date.

2.2 Finding dataset candidates

We used an iterative and community-driven approach combined with snowball search to identify dataset candidates for inclusion in our review. In January 2024, we released a first version of SafetyPrompts.com, with an initial list of datasets that we had compiled in a heuristic fashion based on prior work and our knowledge of the LLM safety field. Over the next two months, we marketed the website to the LLM safety community on Twitter and Reddit, to solicit feedback and further dataset suggestions. This resulted in 77 datasets. We then used these 77 datasets as a starting point for snowball search, wherein we reviewed each publication corresponding to each dataset for references to other potentially relevant datasets. Whenever we identified a new dataset, we repeated this process. This resulted in

35 additional datasets. **Overall, our review includes 102 open datasets for evaluating and improving LLM safety**, which were published between June 2018 and February 2024.

We opted for our review method because of two main reasons. First, LLM safety is a very fast-moving field with contributions from across academia and industry. By sharing intermittent results of our review on SafetyPrompts.com, we were able to solicit feedback from a broad range of stakeholders and expand the scope of our review, while also providing a useful resource to the community well ahead of the release of this paper. Second, traditional systematic review methods like keyword search are ill-suited to the scope of our review. Combinations of relevant keywords like “language model”, “safety” and “dataset” return thousands of results on Google Scholar and similar platforms – and still fail to capture the many types of datasets that may not mention “safety” but still are highly relevant to it, like toxic conversation datasets or bias evaluations. Despite our best efforts, it is likely that our review is missing at least some relevant datasets. We are committed to adding these datasets, along with future relevant dataset releases, to SafetyPrompts.com.

2.3 Recording structured information

For each of the 102 datasets in our review, we recorded 23 pieces of structured information. At a high level, our goal was to capture the full development pipeline of each dataset: from how the dataset was created, to what it looks like, what it can or should be used for, how it can be accessed, and where it was published. We show the full codebook in Appendix A, which describes the structure and content of our main review spreadsheet. We make the spreadsheet available along with code to reproduce our analyses at github.com/paulrottger/safetyprompts-paper.

3 Dataset review findings

3.1 History and growth of open datasets for LLM safety

Our review shows that **LLM safety builds on a rich history of research into risks and biases of language models**. The first datasets in our review were published in 2018, and focus on evaluating gender bias – originally for co-reference resolution systems, but equally applicable to current LLMs (Zhao et al., 2018; Rudinger et al., 2018). These datasets, in turn, build on earlier works on biases in word embeddings (e.g. Bolukbasi et al., 2016; Caliskan et al., 2017; Garg et al., 2018), which fall outside the scope of our review (§2.1), but illustrate that concerns around the negative social impacts of language models are far from new. Similarly, Dinan et al. (2019) and Rashkin et al. (2019), among others, introduced datasets for evaluating and improving the safety of dialogue agents well before the current generative LLM paradigm. By today’s standards, however, interest in safety was relatively low at the time, with only 9 out of the 102 datasets in our review (8.9%) published in or before 2020.

We find that **LLM safety experienced a first moderate growth phase in 2021 and 2022** (Figure 1). These two years, respectively, saw the publication of 15 and 16 open LLM safety datasets. This coincides with increased interest in generative language models, particularly among researchers, following the release of GPT-3 in mid-2020 (Brown et al., 2020).

Finally, our review confirms that **research into LLM safety is currently experiencing unprecedented growth**. 47 out of the 102 datasets in our review (46.1%) were published in 2023. This coincides with a surge in public interest in LLMs as well as concerns around LLM safety following the release of ChatGPT in November 2022. With 15 datasets published in just the first two months of 2024 – i.e. up to our review cutoff of March 1st, 2024 – it is likely that more open LLM safety datasets will be published in 2024 than ever before.

3.2 Intended use and purpose of datasets

In our review, we differentiate between five high-level dataset purposes: **Broad safety** (n=33) denotes datasets that cover several aspects of LLM safety. This includes structured evaluation datasets like SafetyKit (Dinan et al., 2022) or SimpleSafetyTests (Vidgen et al., 2023) as

well as broad-scope red-teaming datasets like BAD (Xu et al., 2021) or AnthropicRedTeam (Ganguli et al., 2022). **Narrow safety** (n=18), conversely, denotes datasets that focus only on one specific aspect of LLM safety. SafeText (Levy et al., 2022), for example, focuses only on commonsense physical safety, while SycophancyEval (Sharma et al., 2024) focuses on sycophantic behaviour. **Value alignment** (n=17) refers to datasets that are concerned with the ethical, moral or social behaviour of LLMs. This includes datasets that seek to evaluate LLM understanding of ethical norms, like Scruples (Lourie et al., 2021) and ETHICS (Hendrycks et al., 2020a), as well as opinion surveys like GlobalOpinionQA (Durmus et al., 2023). **Bias** (n=26) refers to datasets for evaluating sociodemographic biases in LLMs. BOLD (Dhamala et al., 2021), for example, evaluates bias in text completions, whereas DiscrimEval (Tamkin et al., 2023) evaluates biases in situated LLM decision-making. **Other** (n=8), in our review, includes datasets for developing LLM chat moderation systems, like FairPrism (Fleisig et al., 2023) and ToxicChat (Lin et al., 2023), as well as collections of specialised prompts from public prompt hacking competitions, like Gandalf (LakeraAI, 2023a), MossCap (LakeraAI, 2023b) or HackAPrompt (Schulhoff et al., 2023).

Figure 1 shows that **early safety datasets were primarily concerned with evaluating biases**. 13 out of 24 datasets (54.2%) published between 2018 and 2021 were created to identify and analyse sociodemographic biases in language models. 12 of these datasets evaluate gender biases, either exclusively (e.g. Nozza et al., 2021) or along with other categories of bias such as race and sexual orientation (e.g. Sheng et al., 2019).

Broad safety emerged as a prominent theme in 2022, driven by industry contributions.

Anthropic, for example, released two broad-scope red-teaming datasets (Ganguli et al., 2022; Bai et al., 2022a), while Meta published datasets on positive LLM conversations (Ung et al., 2022) and general safety evaluation (Dinan et al., 2022). More recently, broad safety has shifted towards more structured evaluation, as exemplified by benchmarks like DecodingTrust (Wang et al., 2024) or HarmBench (Mazeika et al., 2024).

Our findings also suggest **there is now a trend towards more specialised safety evaluations.**

Narrow safety evaluations did not emerge until 2022, but now make up a significant portion of all new datasets. In the first two months of 2024 alone, 6 of 15 datasets in our review (40.0%) were concerned with specific aspects of LLM safety, like rule-following (Mu et al., 2024) or privacy-reasoning ability (Mireshghallah et al., 2024).

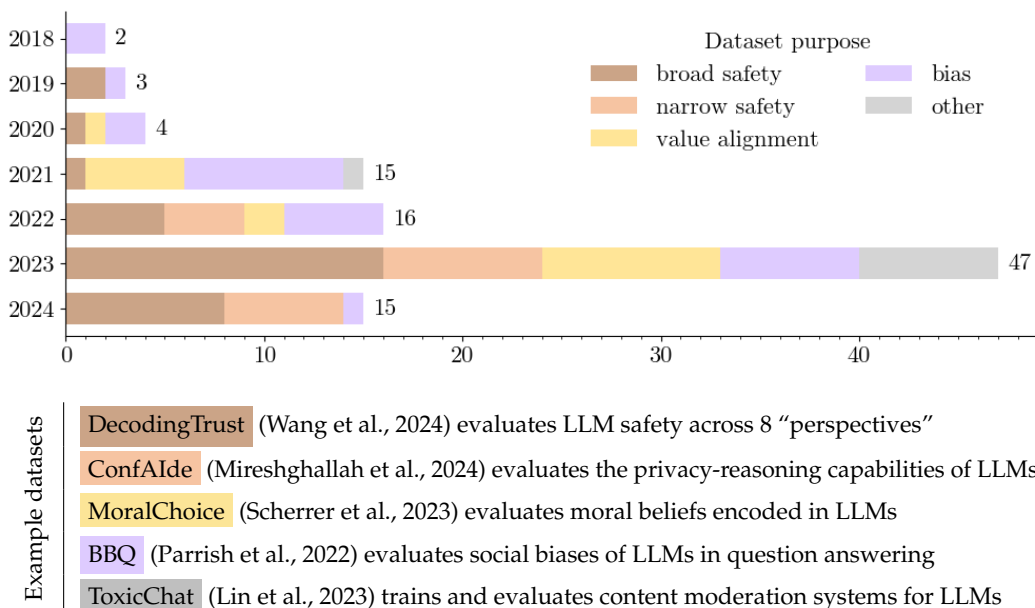


Figure 1: (top) **Number of datasets published per year, grouped by their primary purpose.** In total, our review includes 102 datasets published between June 2018 and February 2024. (bottom) **Example datasets** for each type of dataset purpose. See §3.2 for details.

Finally, we find that **most datasets are intended for model evaluation only**. 80 out of the 102 datasets in our review (78.4%) were explicitly created for benchmarking or evaluation, rather than model training. Only 4 datasets (3.9%), by contrast, comprise examples of positive interactions between users and LLMs, created specifically for model training (Rashkin et al., 2019; Ung et al., 2022; Kim et al., 2022; Bianchi et al., 2024).

3.3 Dataset format and size

We find that **the format of LLM safety datasets has changed alongside general trends in LLM development**. Early datasets, especially those created for bias evaluation, often used autocomplete-style formats (n=8), where models are tasked to either fill a masked word or finish a sentence snippet (e.g. Sheng et al., 2019; Gehman et al., 2020; Kirk et al., 2021). Such formats, which are most compatible with early LLMs like BERT or GPT-2, have since been replaced by chat-style prompts and conversations (n=58) as well as promptable multiple-choice questions (n=14), which better suit current generative LLMs.

Dataset size varies substantially across the 102 datasets in our review, but there is no clear pattern of different sizes corresponding to different dataset characteristics, like dataset purpose or creation. The smallest dataset, *ControversialInstructions* by Bianchi et al. (2024), comprises 40 author-written prompts instructing LLMs to generate hate speech. The largest dataset, *HackAPrompt* by Schulhoff et al. (2023), contains 601,757 human-written prompts recorded as part of a public prompt hacking competition.

3.4 Dataset creation

We find that **the use of templates is a consistently popular method for generating safety evaluation data**. 17 out of the 102 datasets in our review (16.7%) consist of human-written phrase or prompt templates, which are expanded through combination into larger evaluation datasets. *HolisticBias* (Smith et al., 2022), for example, comprises 26 sentence templates (e.g. "I am an [attribute1] who is [attribute2]."), which are combined with around 600 individual descriptor terms to create 459,758 test cases. Generally, template approaches are most popular for bias evaluation – 13 out of 26 bias evaluation datasets in our review use a template approach – but recent work has also used templates for evaluating LLM safety more generally (Wang et al., 2024) or in relation to specific concerns such as privacy reasoning (Mireshghallah et al., 2024).

We also find that **a significant portion of recently-released datasets is fully synthetic**. Earlier safety datasets collected human-written prompts (e.g. Dinan et al., 2019; Cercas Curry et al., 2021; Xu et al., 2021), but 2023 saw the release of the first datasets with fully model-generated prompts. 12 out of the 47 datasets released in 2023 consist of prompts, sentences, or multiple-choice questions generated entirely by LLMs – most commonly by some variant of GPT-3.5. Shaikh et al. (2023), for example, prompt GPT-3.5 to generate 200 harmful questions, which they use to explore safety in chain-of-thought question answering.

Relatedly, instead of relying on static templates for data creation, **multiple recent datasets are more flexibly augmented using LLMs**. Bhatt et al. (2023), for instance, expand a small expert-written set of cyberattack instructions into a larger set of 1,000 prompts using Llama-70b-chat (Touvron et al., 2023a). Wang et al. (2024) take a similar approach to build their large-scale DecodingTrust benchmark.

Finally, we observe that **there is a trend towards small hand-written prompt datasets for model evaluation**. 11 out of the 102 datasets in our review (10.8%) were written by the authors of the corresponding dataset publication. Typically, these datasets comprise just a few hundred prompts and target specific model behaviours like rule-following (Mu et al., 2024) or exaggerated safety (Röttger et al., 2023), which require careful prompt construction.

3.5 Dataset languages

We find that **the vast majority of open LLM safety datasets use English language only**. 88 out of the 102 datasets in our review (86.3%) contain only English language entries.

Six datasets (5.9%) focus exclusively on Chinese (e.g. Zhou et al., 2022; Xu et al., 2023; Zhao et al., 2023). One dataset (Névéol et al., 2022) measures social bias in French language models. The seven other datasets (10.8%) cover English along with one or more other languages. Pikuliak et al. (2023) cover the largest variety of ten languages. In total, the 102 datasets in our review span only 19 different languages.

3.6 Data access and licensing

We find that **GitHub is the most popular platform for sharing open LLM safety data**. Only 8 out of the 102 datasets in our review (7.8%) are not shared on GitHub. These 8 datasets are available on Hugging Face. 35 datasets (34.3%) are available on both GitHub and Hugging Face. Despite the growing popularity of Hugging Face, we do not find a clear trend towards a higher proportion of LLM safety datasets being available on Hugging Face.

We also find that, **when data is shared, usage licenses are mostly permissive**. The most common license is the very permissive MIT License, which is used for 40 out of 102 datasets (39.2%). 14 datasets (13.7%) use the Apache 2.0 License, which provides additional patent protections. 27 datasets (26.5%) use variants of a Creative Commons BY 4.0 License, which requires dataset users to provide appropriate credit and indicate if changes were made to the dataset. 5 datasets (4.9%) prohibit commercial usage with a CC BY-NC License. Only two datasets (2.0%) use a more restrictive custom license. As of March 25th, 2024, 19 datasets (18.6%) do not specify any license.¹

3.7 Dataset publication authors and venues

We find that **academic and non-profit organisations drive most of the creation and publication of open LLM safety datasets**. For 51 out of the 102 datasets in our review (50.0%), all authors of the corresponding publication were affiliated only with academic or non-profit organisations. 27 datasets (26.5%) were published by teams spanning industry and academia. Only 24 datasets (23.5%) were published by fully industry teams.

We also find that **the creation of LLM safety datasets is concentrated in few research hubs** (Table 1). There are 90 unique affiliations across the authors of the 102 datasets in our review. 52 affiliations (57.8%) are associated with just a single dataset. The five most prolific organisations, on the other hand, are each associated with 10 or 11 datasets. All of the twenty most prolific organisations are located and/or headquartered in the US, with the exception of Bocconi University (Italy, n=10), Alibaba (China, n=3), and the University of Cambridge (UK, n=3).

Academic / Non-Profit Org.			Industry Org.		
		n			n
1	UC Berkeley	11	1	Meta* (prev. Facebook)	11
1	Stanford University	11	2	Anthropic	9
3	Allen AI	10	3	Microsoft* (incl. Research)	6
3	Bocconi University	10	3	Google* (incl. DeepMind)	6
5	University of Washington	9	5	Alibaba	3

Table 1: **Organisations that published the most open LLM safety datasets**, among the 102 datasets in our review. For each dataset, we count all affiliations for all co-authors. UCB at n=11, for example, means that 11 datasets had an author affiliated with UCB.

Finally, we find that **the largest share of LLM safety datasets so far has appeared at ACL* conferences**. 45 out of the 102 datasets in our review (44.1%) were published at either ACL (n=20), EMNLP (n=20), NAACL (n=4), or AACL (n=1). 20 datasets (19.6%) were published at ICLR (n=10), NeurIPS (n=7), or ICML (n=3). Only 4 datasets (3.9%) were published at other venues. Notably, none of the datasets appeared in journal publications. 29 datasets

¹While conducting our review, we reached out to authors of all datasets that had not specified a license and encouraged them to add one. At least five authors added a license as a result.

(28.4%), on the other hand, were accompanied only by arXiv preprints, and 4 (3.9%) only by blog posts, meaning they did not receive traditional peer review. Generally, we observe a slight trend away from ACL* conferences and towards more ML-focused venues as well as arXiv-only publication, although this could in part be explained by recent arXiv preprints still being under review at forthcoming conferences.

4 Open LLM safety datasets used in model release publications

In the following, we briefly examine how open LLM safety datasets are used in practice. In particular, we examine in this section (§4) which safety datasets are used to evaluate current state-of-the-art LLMs ahead of their release, as documented in model release publications. In the next section (§5), we then examine which safety datasets are included in popular LLM benchmark suites and leaderboards. This is to characterise current norms and common practices in evaluating LLM safety, so that we can then discuss, in §6, these norms and practices in relation to the findings of our dataset review (§3).

4.1 Scope of our model release publication review

We include the top 50 best-performing LLMs listed on the LMSYS Chatbot Arena Leaderboard as of March 12th, 2024, in our review.² The LMSYS leaderboard is a crowdsourced platform for LLM evaluation, which ranks models based on model Elo scores calculated from over 400,000 pairwise human preference votes. We use the LMSYS leaderboard because it is very popular in the LLM community, and it has up-to-date coverage of recent model releases from both industry and academia.

The top 50 entries on the LMSYS leaderboard correspond to 31 unique model releases.³ Of these 31 models, 11 (35.5%) are proprietary models only accessible via an API. These are models released by OpenAI (GPT), Google (Gemini), Anthropic (Claude), Perplexity (pplx), and Mistral (Next, Medium, and Large). The other 20 models (64.5%) are open models accessible via Hugging Face. Proprietary models generally outrank open models on the leaderboard, with Qwen1.5-72b-chat being the best open model at rank 10. 26 out of the 31 models (83.9%) were released by industry labs, while the rest were created either by academic or non-profit organisations. All 31 models were released in 2023 or 2024.

4.2 Findings of our model release publication review

We find that **the majority of state-of-the-art LLMs are evaluated for safety ahead of their release**, although the extent and nature of safety evaluation varies. 24 out of the 31 models (77.4%) report safety evaluations in their release publications. 21 models (67.7%) report results on at least one open LLM safety dataset. Guanaco (Dettmers et al., 2024), for example, was evaluated on a single open LLM safety dataset (CrowS-Pairs by Nangia et al., 2020). Llama2 (Touvron et al., 2023b), by contrast, was evaluated on five different open LLM safety datasets. 7 out of the 31 models, on the other hand, do not report any safety evaluations. This includes 5 open models from academia and industry, such as Starling (Zhu et al., 2023) and WizardLM (Xu et al., 2024), as well as the proprietary Mistral Medium and Next models.

We also find that **proprietary data plays a large role in model release safety evaluations**. Out of the 24 model releases that report safety evaluation results, 13 (54.2%) use undisclosed proprietary data for evaluating model safety. Three of these releases – Gemini (Anil et al., 2023), Qwen (Bai et al., 2022b), and Mistral-7B (Jiang et al., 2023) – report results only on proprietary safety datasets.

Finally, we find that **the diversity of open LLM safety datasets used in model release evaluations is very limited**. A total of only 12 open LLM safety datasets are used across the 31 model releases, and 7 of these 12 datasets are used only once. Table 2 shows the 5 datasets that are used more than once. Notably, TruthfulQA (Lin et al., 2022) is used in

²huggingface.co/spaces/lmsys/chatbot-arena-leaderboard

³A model release may comprise multiple model versions, such as GPT-4-0314 and GPT-4-0613.

16 out of the 24 model releases that report any safety evaluation results (66.7%). All other datasets are used in at most 5 model release publications.

Dataset	Purpose	n
TruthfulQA (Lin et al., 2022)	evaluate tendency to mimic human falsehoods	16
BBQ (Parrish et al., 2022)	evaluate social bias in question answering	5
AnthropicRedTeam (Ganguli et al., 2022)	evaluate responses to diverse red-team attacks	4
ToxiGen (Hartvigsen et al., 2022)	evaluate toxicity in text completions	3
BOLD (Dhamala et al., 2021)	evaluate social bias in text completions	3

Table 2: **Most popular open LLM safety datasets**, based on how often release publications for state-of-the-art LLMs reported results on each dataset. BBQ at n=5, for example, means that 5 out of 31 model release publications reviewed in §4.2 reported results on BBQ.

5 Open LLM safety datasets used in popular benchmarks

5.1 Scope of our benchmark review

We examine 5 widely-used general-purpose benchmarking suites: Stanford’s HELM Classic (Liang et al., 2023) and Helm Instruct (Zhang et al., 2024), Hugging Face’s Open LLM Leaderboard (Beeching et al., 2023), Eleuther AI’s Evaluation Harness (Gao et al., 2021), and BIG-Bench (Srivastava et al., 2023). We also examine 2 benchmarks focused primarily on LLM safety: TrustLLM (Sun et al., 2024) and the LLM Safety Leaderboard.⁴

5.2 Findings of our benchmark review

We observe that **there are large differences in how different benchmarks evaluate LLM safety**. A total of 20 open LLM safety datasets are used across the 7 benchmarks. 14 of these datasets are used in just one benchmark. TrustLLM (Sun et al., 2024), for example, uses 8 open LLM safety datasets, of which 6 are not used in any other benchmark. The only open LLM safety datasets that are used in more than 2 benchmarks are TruthfulQA (Lin et al., 2022), which is used in 5 benchmarks, as well as RealToxicityPrompts (Gehman et al., 2020) and ETHICS (Hendrycks et al., 2020a), which are both used in 3 benchmarks.

We also note that **there is currently no LLM safety benchmark with a truly comprehensive scope**. The LLM Safety Leaderboard, based on DecodingTrust (Wang et al., 2024), has the broadest scope relevant to safety among the 7 benchmarks we reviewed. However, it does not, for example, test for exaggerated safety, as TrustLLM does (Sun et al., 2024), or test for catastrophic risk potential, as the Evaluation Harness does (Gao et al., 2021).⁵

6 Discussion

Overall, our review shows that **growing interest in LLM safety is driving the creation of more and more diverse open LLM safety datasets**. More datasets were published in 2023 than ever before, and it is likely that this trend will continue in the current year (§3.1). Existing datasets span varied purposes (§3.2) and formats (§3.3), which have adapted over time to meet the needs and requirements of LLM users and developers. Researchers and practitioners are making creative use of new methods for dataset creation (§3.4), and when data is shared, usage licenses are mostly permissive (§3.6). These are encouraging signs for the health of the open LLM safety community and its ability to address emerging challenges and fill gaps in dataset coverage as they become apparent.

⁴huggingface.co/spaces/AI-Secure/llm-trustworthy-leaderboard, based on Wang et al. (2024).

⁵Note that parts of the DecodingTrust dataset were incorporated into HELM in a recent update: github.com/stanford-crfm/helm/tree/main/src/helm/benchmark/scenarios

Among these gaps, the most apparent today is that **there is a clear lack of datasets in non-English languages**. We found that English dominates the current safety dataset landscape (§3.5), mirroring long-standing trends in NLP research (Bender, 2011; Joshi et al., 2020; Holtermann et al., 2024). To some extent, this language imbalance reflects an imbalance in who is publishing safety datasets (§3.7). Both imbalances could be remedied by non-US institutions leading the creation of datasets in their local languages.

Our analysis of how open LLM safety datasets are used in practice, shows that **there is clear scope for standardisation in LLM safety evaluations**. Evaluating safety is a key priority for model developers and users, as evidenced by the inclusion of safety evaluations in model release publications (§4) and popular LLM benchmarks (§5). However, the methods that have been used for evaluating safety so far are highly idiosyncratic. Specifically, we found that most model release publications and benchmarks make use of very different datasets. For commercial model releases, these datasets are often proprietary and undisclosed, More standardised open evaluations would enable more meaningful model comparisons and incentivise the development of safer LLMs.

A key challenge for standardisation is justifying which evaluations constitute an adequate standard. In this paper, we refrained from making quality judgments about the datasets we reviewed, mainly because different datasets serve different purposes, so that their utility is highly context dependent. However, our review of benchmarks and model release publications shows that **current safety evaluation practices do not fully leverage recent progress in safety dataset creation**. The most popular open datasets for evaluating safety in model release publications, for example, are all from 2021 or 2022 (Table 2), despite more than half of the datasets in our review being published in or after 2023. Prior publication date is not a sign of lacking quality, but due to the rapid development of the field, older autocomplete-style datasets like BOLD (Dhamala et al., 2021) or ToxiGen (Hartvigsen et al., 2022) no longer reflect real-world usage of current-generation LLMs (Ouyang et al., 2023; Zheng et al., 2024; Zhao et al., 2024). We hope that by highlighting the diversity of open LLM datasets available today, our review can provide a starting point for updating and improving prevailing evaluation practices.

7 Conclusion

In recent years, researchers and practitioners have sought to meet concerns around the safety of large language models by creating an abundance of new datasets for evaluating and improving LLM safety. The rapid pace of dataset creation and the variety of purposes served by different datasets make it difficult for researchers and practitioners to find the most relevant datasets for different use cases, and to identify gaps in dataset coverage that future work may fill. In this paper, we addressed these issues by conducting a first systematic review of open LLM safety datasets. We identified 102 datasets published between June 2018 and February 2024 based on clear inclusion criteria (§2.1) using a comprehensive search method (§2.2). We examined these datasets along several key dimensions, including their purpose (§3.2), creation (§3.4), format and size (§3.3), access and licensing (§3.6), and publication (§3.7). Key findings of our review include that dataset creation is currently growing at an unprecedented rate, driven primarily by academic and non-profit organisations; that there is a trend towards more specialised safety evaluations and the use of synthetic data; and that the English language dominates the dataset landscape. We also reviewed how open LLM safety datasets are used in practice – in model release publications (§4) and popular LLM benchmarks (§5) – finding that current evaluation practices are highly idiosyncratic and leverage only a small fraction of available datasets. In our Discussion (§6), we argued that this creates clear scope for standardisation in LLM safety evaluations, and that evaluations in general could be improved by better leveraging recent progress in safety dataset creation. Overall, we hope that our review, along with the living dataset catalogue we make available on SafetyPrompts.com, can help researchers and practitioners make the best use of existing datasets, and provide a strong foundation for future dataset development.

Acknowledgments

Thank you for feedback and dataset suggestions to Giuseppe Attanasio, Federico Bianchi, Hannah Lucas, Felix Röttger, Bo Li, Matus Pikuliak, Pranav Venkit, Verena Rieser, Norman Mu, Niloofar Mireshghallah, Hyunwoo Kim, Sam Toyer and Laura Weidinger. Special thanks to Hannah Rose Kirk for the initial logo suggestion.

PR, FP, and DH are members of the Data and Marketing Insights research unit of the Bocconi Institute for Data Science and Analysis, and are supported by a MUR FARE 2020 initiative under grant agreement Prot. R20YSMBZ8S (INDOMITA) and the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation program (No. 949944, INTEGRATOR).

References

- Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*, 2022a.
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022b.
- Edward Beeching, Clémentine Fourrier, Nathan Habib, Sheon Han, Nathan Lambert, Nazneen Rajani, Omar Sanseviero, Lewis Tunstall, and Thomas Wolf. Open llm leaderboard. [https://huggingface.co/spaces/HuggingFaceH4/open.llm.leaderboard](https://huggingface.co/spaces/HuggingFaceH4/open_llm_leaderboard), 2023.
- Emily M Bender. On achieving and evaluating language-independence in nlp. *Linguistic Issues in Language Technology*, 6, 2011.
- Manish Bhatt, Sahana Chennabasappa, Cyrus Nikolaidis, Shengye Wan, Ivan Evtimov, Dominik Gabi, Daniel Song, Faizan Ahmad, Cornelius Aschermann, Lorenzo Fontana, et al. Purple llama cyberseceval: A secure coding benchmark for language models. *arXiv preprint arXiv:2312.04724*, 2023.
- Federico Bianchi, Pratyusha Kalluri, Esin Durmus, Faisal Ladhak, Myra Cheng, Debora Nozza, Tatsunori Hashimoto, Dan Jurafsky, James Zou, and Aylin Caliskan. Easily accessible text-to-image generation amplifies demographic stereotypes at large scale. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency, FAccT ’23*, pp. 1493–1504, New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9798400701924. doi: 10.1145/3593013.3594095. URL <https://doi.org/10.1145/3593013.3594095>.
- Federico Bianchi, Mirac Suzgun, Giuseppe Attanasio, Paul Rottger, Dan Jurafsky, Tatsunori Hashimoto, and James Zou. Safety-tuned LLaMAs: Lessons from improving the safety of large language models that follow instructions. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=gT5hALch9z>.
- Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. *Advances in neural information processing systems*, 29, 2016.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.

- Aylin Caliskan, Joanna J Bryson, and Arvind Narayanan. Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334):183–186, 2017.
- Nicholas Carlini, Milad Nasr, Christopher A Choquette-Choo, Matthew Jagielski, Irena Gao, Pang Wei Koh, Daphne Ippolito, Florian Tramèr, and Ludwig Schmidt. Are aligned neural networks adversarially aligned? In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Amanda Cercas Curry, Gavin Abercrombie, and Verena Rieser. ConvAbuse: Data, analysis, and benchmarks for nuanced detection in conversational AI. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 7388–7403, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. URL <https://aclanthology.org/2021.emnlp-main.587>.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms. *Advances in Neural Information Processing Systems*, 36, 2024.
- Jwala Dhamala, Tony Sun, Varun Kumar, Satyapriya Krishna, Yada Pruksachatkun, Kai-Wei Chang, and Rahul Gupta. Bold: Dataset and metrics for measuring biases in open-ended language generation. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency, FAccT '21*, pp. 862–872, New York, NY, USA, 2021. Association for Computing Machinery. ISBN 9781450383097. doi: 10.1145/3442188.3445924. URL <https://doi.org/10.1145/3442188.3445924>.
- Emily Dinan, Samuel Humeau, Bharath Chintagunta, and Jason Weston. Build it break it fix it for dialogue safety: Robustness from adversarial human attack. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 4537–4546, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1461. URL <https://www.aclweb.org/anthology/D19-1461>.
- Emily Dinan, Gavin Abercrombie, A. Bergman, Shannon Spruit, Dirk Hovy, Y-Lan Boureau, and Verena Rieser. SafetyKit: First aid for measuring safety in open-domain conversational systems. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 4113–4133, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.284. URL <https://aclanthology.org/2022.acl-long.284>.
- Dheeru Dua, Yizhong Wang, Pradeep Dasigi, Gabriel Stanovsky, Sameer Singh, and Matt Gardner. DROP: A reading comprehension benchmark requiring discrete reasoning over paragraphs. In Jill Burstein, Christy Doran, and Thamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 2368–2378, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1246. URL <https://aclanthology.org/N19-1246>.
- Esin Durmus, Karina Nyugen, Thomas I Liao, Nicholas Schiefer, Amanda Askell, Anton Bakhtin, Carol Chen, Zac Hatfield-Dodds, Danny Hernandez, Nicholas Joseph, et al. Towards measuring the representation of subjective global opinions in language models. *arXiv preprint arXiv:2306.16388*, 2023.
- Eve Fleisig, Aubrie Amstutz, Chad Atalla, Su Lin Blodgett, Hal Daumé III, Alexandra Olteanu, Emily Sheng, Dan Vann, and Hanna Wallach. FairPrism: Evaluating fairness-related harms in text generation. In Anna Rogers, Jordan Boyd-Graber, and Naoki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 6231–6251, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.343. URL <https://aclanthology.org/2023.acl-long.343>.

- Deep Ganguli, Liane Lovitt, Jackson Kernion, Amanda Askill, Yuntao Bai, Saurav Kadavath, Ben Mann, Ethan Perez, Nicholas Schiefer, Kamal Ndousse, et al. Red teaming language models to reduce harms: Methods, scaling behaviors, and lessons learned. *arXiv preprint arXiv:2209.07858*, 2022.
- Leo Gao, Jonathan Tow, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Kyle McDonell, Niklas Muennighoff, Jason Phang, Laria Reynolds, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. A framework for few-shot language model evaluation. In *Zenodo*. Zenodo, September 2021. doi: 10.5281/zenodo.5371628. URL <https://doi.org/10.5281/zenodo.5371628>.
- Nikhil Garg, Londa Schiebinger, Dan Jurafsky, and James Zou. Word embeddings quantify 100 years of gender and ethnic stereotypes. *Proceedings of the National Academy of Sciences*, 115(16):E3635–E3644, 2018.
- Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A. Smith. RealToxicityPrompts: Evaluating neural toxic degeneration in language models. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pp. 3356–3369, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.findings-emnlp.301. URL <https://aclanthology.org/2020.findings-emnlp.301>.
- Vipul Gupta, Pranav Narayanan Venkit, Hugo Laurençon, Shomir Wilson, and Rebecca J. Passonneau. Calm : A multi-task benchmark for comprehensive assessment of language model bias, 2024.
- Siobhan Mackenzie Hall, Fernanda Gonçalves Abrantes, Hanwen Zhu, Grace Sodunke, Aleksandar Shtedritski, and Hannah Rose Kirk. Visogender: A dataset for benchmarking gender bias in image-text pronoun resolution. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track, 2023*. URL <https://openreview.net/forum?id=BNwsJ4bFsc>.
- Thomas Hartvigsen, Saadia Gabriel, Hamid Palangi, Maarten Sap, Dipankar Ray, and Ece Kamar. ToxiGen: A large-scale machine-generated dataset for adversarial and implicit hate speech detection. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 3309–3326, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.234. URL <https://aclanthology.org/2022.acl-long.234>.
- Dan Hendrycks, Collin Burns, Steven Basart, Andrew Critch, Jerry Li, Dawn Song, and Jacob Steinhardt. Aligning ai with shared human values. In *International Conference on Learning Representations, 2020a*.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. Measuring massive multitask language understanding. In *International Conference on Learning Representations, 2020b*.
- Carolin Holtermann, Paul Röttger, Timm Dill, and Anne Lauscher. Evaluating the elementary multilingual capabilities of large language models with multiq. *arXiv preprint arXiv:2403.03814*, 2024.
- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika Bali, and Monojit Choudhury. The state and fate of linguistic diversity and inclusion in the NLP world. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 6282–6293, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.560. URL <https://aclanthology.org/2020.acl-main.560>.
- Hyunwoo Kim, Youngjae Yu, Liwei Jiang, Ximing Lu, Daniel Khashabi, Gunhee Kim, Yejin Choi, and Maarten Sap. ProsocialDialog: A prosocial backbone for conversational

- agents. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 4005–4029, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.267. URL <https://aclanthology.org/2022.emnlp-main.267>.
- Hannah Rose Kirk, Yennie Jun, Filippo Volpin, Haider Iqbal, Elias Benussi, Frederic Dreyer, Aleksandar Shtedritski, and Yuki Asano. Bias out-of-the-box: An empirical analysis of intersectional occupational biases in popular generative language models. *Advances in neural information processing systems*, 34:2611–2624, 2021.
- Allison Koenecke, Andrew Nam, Emily Lake, Joe Nudell, Minnie Quartey, Zion Mengesha, Connor Toups, John R Rickford, Dan Jurafsky, and Sharad Goel. Racial disparities in automated speech recognition. *Proceedings of the National Academy of Sciences*, 117(14): 7684–7689, 2020.
- LakeraAI. Gandalf prompt injection. In *Hugging Face Dataset*, 2023a. URL https://huggingface.co/datasets/Lakera/gandalf_ignore_instructions.
- LakeraAI. MossCAP prompt injection. In *Hugging Face Dataset*, 2023b. URL https://huggingface.co/datasets/Lakera/mossCAP_prompt_injection.
- Sharon Levy, Emily Allaway, Melanie Subbiah, Lydia Chilton, Desmond Patton, Kathleen McKeown, and William Yang Wang. SafeText: A benchmark for exploring physical safety in language models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 2407–2421, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. URL <https://aclanthology.org/2022.emnlp-main.154>.
- Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, et al. Holistic evaluation of language models. *Transactions on Machine Learning Research*, 2023.
- Stephanie Lin, Jacob Hilton, and Owain Evans. TruthfulQA: Measuring how models mimic human falsehoods. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 3214–3252, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.229. URL <https://aclanthology.org/2022.acl-long.229>.
- Zi Lin, Zihan Wang, Yongqi Tong, Yangkun Wang, Yuxin Guo, Yujia Wang, and Jingbo Shang. ToxicChat: Unveiling hidden challenges of toxicity detection in real-world user-AI conversation. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 4694–4702, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-emnlp.311. URL <https://aclanthology.org/2023.findings-emnlp.311>.
- Nicholas Lourie, Ronan Le Bras, and Yejin Choi. Scruples: A corpus of community ethical judgments on 32,000 real-life anecdotes. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 13470–13479, 2021.
- Sasha Luccioni, Christopher Akiki, Margaret Mitchell, and Yacine Jernite. Stable bias: Evaluating societal representations in diffusion models. *Advances in Neural Information Processing Systems*, 36, 2024.
- Mantas Mazeika, Long Phan, Xuwang Yin, Andy Zou, Zifan Wang, Norman Mu, Elham Sakhaee, Nathaniel Li, Steven Basart, Bo Li, et al. Harmbench: A standardized evaluation framework for automated red teaming and robust refusal. *arXiv preprint arXiv:2402.04249*, 2024.
- Josh Meyer, Lindy Rauchenstein, Joshua D. Eisenberg, and Nicholas Howell. Artie bias corpus: An open dataset for detecting demographic bias in speech applications. In

- Nicoletta Calzolari, Frédéric Béchet, Philippe Blache, Khalid Choukri, Christopher Cieri, Thierry Declerck, Sara Goggi, Hitoshi Isahara, Bente Maegaard, Joseph Mariani, H el ene Mazo, Asuncion Moreno, Jan Odijk, and Stelios Piperidis (eds.), *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pp. 6462–6468, Marseille, France, May 2020. European Language Resources Association. ISBN 979-10-95546-34-4. URL <https://aclanthology.org/2020.lrec-1.796>.
- Niloofar Miresghallah, Hyunwoo Kim, Xuhui Zhou, Yulia Tsvetkov, Maarten Sap, Reza Shokri, and Yejin Choi. Can LLMs keep a secret? testing privacy implications of language models via contextual integrity theory. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=gmg7t8b4s0>.
- Norman Mu, Sarah Chen, Zifan Wang, Sizhe Chen, David Karamardian, Lulwa Aljerais, Basel Alomair, Dan Hendrycks, and David Wagner. Can llms follow simple rules?, 2024.
- Nikita Nangia, Clara Vania, Rasika Bhalerao, and Samuel R. Bowman. CrowS-pairs: A challenge dataset for measuring social biases in masked language models. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu (eds.), *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 1953–1967, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.154. URL <https://aclanthology.org/2020.emnlp-main.154>.
- Aur elie N ev eol, Yoann Dupont, Julien Bezan on, and Kar en Fort. French CrowS-pairs: Extending a challenge dataset for measuring social bias in masked language models to a language other than English. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 8521–8531, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.583. URL <https://aclanthology.org/2022.acl-long.583>.
- Debora Nozza, Federico Bianchi, and Dirk Hovy. HONEST: Measuring hurtful sentence completion in language models. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 2398–2406, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.naacl-main.191. URL <https://aclanthology.org/2021.naacl-main.191>.
- Siru Ouyang, Shuohang Wang, Yang Liu, Ming Zhong, Yizhu Jiao, Dan Iter, Reid Pryzant, Chenguang Zhu, Heng Ji, and Jiawei Han. The shifted and the overlooked: A task-oriented investigation of user-GPT interactions. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 2375–2393, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.146. URL <https://aclanthology.org/2023.emnlp-main.146>.
- Alicia Parrish, Angelica Chen, Nikita Nangia, Vishakh Padmakumar, Jason Phang, Jana Thompson, Phu Mon Htut, and Samuel Bowman. BBQ: A hand-built bias benchmark for question answering. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), *Findings of the Association for Computational Linguistics: ACL 2022*, pp. 2086–2105, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.findings-acl.165. URL <https://aclanthology.org/2022.findings-acl.165>.
- Alicia Parrish, Hannah Rose Kirk, Jessica Quaye, Charvi Rastogi, Max Bartolo, Oana Inel, Juan Ciro, Rafael Mosquera, Addison Howard, Will Cukierski, et al. Adversarial nibbler: A data-centric challenge for improving the safety of text-to-image models. *arXiv preprint arXiv:2305.14384*, 2023.
- Mat u s Pikuliak, Andrea Hrkova, Stefan Oresko, and Marian  simko. Women are beautiful, men are leaders: Gender stereotypes in machine translation and language modeling. *arXiv preprint arXiv:2311.18711*, 2023.
- Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau. Towards empathetic open-domain conversation models: A new benchmark and dataset. In Anna

- Korhonen, David Traum, and Lluís Màrquez (eds.), *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 5370–5381, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1534. URL <https://aclanthology.org/P19-1534>.
- Ricardo Reimao and Vassilios Tzerpos. For: A dataset for synthetic speech detection. In *2019 International Conference on Speech Technology and Human-Computer Dialogue (SpeD)*, pp. 1–10. IEEE, 2019.
- Jonas Ricker, Simon Damm, Thorsten Holz, and Asja Fischer. Towards the detection of diffusion model deepfakes. *arXiv preprint arXiv:2210.14571*, 2022.
- Paul Röttger, Hannah Rose Kirk, Bertie Vidgen, Giuseppe Attanasio, Federico Bianchi, and Dirk Hovy. Xstest: A test suite for identifying exaggerated safety behaviours in large language models. *arXiv preprint arXiv:2308.01263*, 2023.
- Rachel Rudinger, Jason Naradowsky, Brian Leonard, and Benjamin Van Durme. Gender bias in coreference resolution. In Marilyn Walker, Heng Ji, and Amanda Stent (eds.), *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pp. 8–14, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-2002. URL <https://aclanthology.org/N18-2002>.
- Nino Scherrer, Claudia Shi, Amir Feder, and David Blei. Evaluating the moral beliefs encoded in llms. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.
- Sander Schulhoff, Jeremy Pinto, Anaum Khan, Louis-François Bouchard, Chenglei Si, Svetlana Anati, Valen Tagliabue, Anson Kost, Christopher Carnahan, and Jordan Boyd-Graber. Ignore this title and HackAPrompt: Exposing systemic vulnerabilities of LLMs through a global prompt hacking competition. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 4945–4977, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.302. URL <https://aclanthology.org/2023.emnlp-main.302>.
- Carsten Schwemmer, Carly Knight, Emily D Bello-Pardo, Stan Oklobdzija, Martijn Schoonvelde, and Jeffrey W Lockhart. Diagnosing gender bias in image recognition systems. *Socius*, 6:2378023120967171, 2020.
- Omar Shaikh, Hongxin Zhang, William Held, Michael Bernstein, and Diyi Yang. On second thought, let’s not think step by step! bias and toxicity in zero-shot reasoning. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 4454–4470, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.244. URL <https://aclanthology.org/2023.acl-long.244>.
- Mrinank Sharma, Meg Tong, Tomasz Korbak, David Duvenaud, Amanda Askell, Samuel R. Bowman, Esin DURMUS, Zac Hatfield-Dodds, Scott R Johnston, Shauna M Kravec, Timothy Maxwell, Sam McCandlish, Kamal Ndousse, Oliver Rausch, Nicholas Schiefer, Da Yan, Miranda Zhang, and Ethan Perez. Towards understanding sycophancy in language models. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=tvhaxkMKAn>.
- Emily Sheng, Kai-Wei Chang, Premkumar Natarajan, and Nanyun Peng. The woman worked as a babysitter: On biases in language generation. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan (eds.), *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 3407–3412, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1339. URL <https://aclanthology.org/D19-1339>.

- Mohammed Latif Siddiq and Joanna C. S. Santos. Securityeval dataset: mining vulnerability examples to evaluate machine learning-based code generation techniques. In *Proceedings of the 1st International Workshop on Mining Software Repositories Applications for Privacy and Security*, MSR4P&S 2022, pp. 29–33, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450394574. doi: 10.1145/3549035.3561184. URL <https://doi.org/10.1145/3549035.3561184>.
- Eric Michael Smith, Melissa Hall, Melanie Kambadur, Eleonora Presani, and Adina Williams. “I’m sorry to hear that”: Finding new biases in language models with a holistic descriptor dataset. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 9180–9211, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.625. URL <https://aclanthology.org/2022.emnlp-main.625>.
- Alexandra Souly, Qingyuan Lu, Dillon Bowen, Tu Trinh, Elvis Hsieh, Sana Pandey, Pieter Abbeel, Justin Svegliato, Scott Emmons, Olivia Watkins, et al. A strongreject for empty jailbreaks. *arXiv preprint arXiv:2402.10260*, 2024.
- Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *Transactions on Machine Learning Research*, 2023.
- Lichao Sun, Yue Huang, Haoran Wang, Siyuan Wu, Qihui Zhang, Chujie Gao, Yixin Huang, Wenhan Lyu, Yixuan Zhang, Xiner Li, et al. Trustllm: Trustworthiness in large language models. *arXiv preprint arXiv:2401.05561*, 2024.
- Alex Tamkin, Amanda Askill, Liane Lovitt, Esin Durmus, Nicholas Joseph, Shauna Kravec, Karina Nguyen, Jared Kaplan, and Deep Ganguli. Evaluating and mitigating discrimination in language model decisions. *arXiv preprint arXiv:2312.03689*, 2023.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.
- Megan Ung, Jing Xu, and Y-Lan Boureau. SaFeRD dialogues: Taking feedback gracefully after conversational safety failures. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 6462–6481, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.acl-long.447. URL <https://aclanthology.org/2022.acl-long.447>.
- Bertie Vidgen, Hannah Rose Kirk, Rebecca Qian, Nino Scherrer, Anand Kannappan, Scott A Hale, and Paul Röttger. Simplestests: a test suite for identifying critical safety risks in large language models. *arXiv preprint arXiv:2311.08370*, 2023.
- Boxin Wang, Weixin Chen, Hengzhi Pei, Chulin Xie, Mintong Kang, Chenhui Zhang, Chejian Xu, Zidi Xiong, Ritik Dutta, Rylan Schaeffer, et al. Decodingtrust: A comprehensive assessment of trustworthiness in gpt models. *Advances in Neural Information Processing Systems*, 36, 2024.
- Robert Wolfe, Yiwei Yang, Bill Howe, and Aylin Caliskan. Contrastive language-vision ai models pretrained on web-scraped multimodal data exhibit sexual objectification bias. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency, FAccT ’23*, pp. 1174–1185, New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9798400701924. doi: 10.1145/3593013.3594072. URL <https://doi.org/10.1145/3593013.3594072>.

- Can Xu, Qingfeng Sun, Kai Zheng, Xiubo Geng, Pu Zhao, Jiazhan Feng, Chongyang Tao, Qingwei Lin, and Daxin Jiang. WizardLM: Empowering large pre-trained language models to follow complex instructions. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=CfXh93NDgH>.
- Guohai Xu, Jiayi Liu, Ming Yan, Haotian Xu, Jinghui Si, Zhuoran Zhou, Peng Yi, Xing Gao, Jitao Sang, Rong Zhang, et al. Cvalues: Measuring the values of chinese large language models from safety to responsibility. *arXiv preprint arXiv:2307.09705*, 2023.
- Jing Xu, Da Ju, Margaret Li, Y-Lan Boureau, Jason Weston, and Emily Dinan. Bot-adversarial dialogue for safe conversational agents. In Kristina Toutanova, Anna Rumshisky, Luke Zettlemoyer, Dilek Hakkani-Tur, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy Chakraborty, and Yichao Zhou (eds.), *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 2950–2968, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.naacl-main.235. URL <https://aclanthology.org/2021.naacl-main.235>.
- Yian Zhang, Yifan Mai, Josselin Somerville Roberts, Rishi Bommasani, Yann Dubois, and Percy Liang. Helm instruct: A multidimensional instruction following evaluation framework with absolute ratings, February 2024. URL <https://crfm.stanford.edu/2024/02/18/helm-instruct.html>.
- Dora Zhao, Angelina Wang, and Olga Russakovsky. Understanding and evaluating racial biases in image captioning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 14830–14840, 2021.
- Jiaxu Zhao, Meng Fang, Zijing Shi, Yitong Li, Ling Chen, and Mykola Pechenizkiy. CHBias: Bias evaluation and mitigation of Chinese conversational language models. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 13538–13556, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.757. URL <https://aclanthology.org/2023.acl-long.757>.
- Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. Gender bias in coreference resolution: Evaluation and debiasing methods. In Marilyn Walker, Heng Ji, and Amanda Stent (eds.), *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pp. 15–20, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-2003. URL <https://aclanthology.org/N18-2003>.
- Wenting Zhao, Xiang Ren, Jack Hessel, Claire Cardie, Yejin Choi, and Yuntian Deng. (in)the>wildchat: 570k chatGPT interaction logs in the wild. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=B18u7ZR1bM>.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Tianle Li, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zhuohan Li, Zi Lin, Eric Xing, Joseph E. Gonzalez, Ion Stoica, and Hao Zhang. Lmsys-1m: A large-scale real-world LLM conversation dataset. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=B0fDKxfwt0>.
- Jingyan Zhou, Jiawen Deng, Fei Mi, Yitong Li, Yasheng Wang, Minlie Huang, Xin Jiang, Qun Liu, and Helen Meng. Towards identifying social bias in dialog systems: Framework, dataset, and benchmark. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2022*, pp. 3576–3591, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.findings-emnlp.262. URL <https://aclanthology.org/2022.findings-emnlp.262>.
- Banghua Zhu, Hiteshi Sharma, Felipe Vieira Frujeri, Shi Dong, Chenguang Zhu, Michael I Jordan, and Jiantao Jiao. Fine-tuning language models with advantage-induced policy alignment. *arXiv preprint arXiv:2306.02231*, 2023.

A Review code book

See Table 3 below.

purpose: what can I use this dataset for?		
.type	High-level type of purpose / general area of application	single choice: [broad safety, narrow safety, bias, value alignment, other]
.tags	Additional tags to specify purpose	comma-separated list of text tags
.stated	Exact purpose as stated by authors	free text
.llmdev	Intended use of the dataset within the LLM development pipeline	single choice: [eval only, train and eval, train only, other]
entries: what do entries in this dataset look like?		
.type	High-level type of entry / format	single choice: [chat, multiple choice, autocomplete, other]
.languages	Languages in the dataset	comma-separated list of language names
.n	Number of entries	integer
.unit	Unit of entries (e.g. conversation)	free text
.detail	Additional detail on entry format	free text
creation: who created this dataset / where is it sampled from?		
.creator_type	Type of creator, i.e. who or what created the data	single choice: [human, machine, hybrid]
.source_type	Type of data source, i.e. where the data is taken from	single choice: [original, sampled, mixed]
.detail	Additional detail on creator/source	free text
access: where can I download this dataset, and how is it licensed?		
.git_url	GitHub repo URL	URL
.hf_url	Hugging Face dataset URL	URL
.license	Dataset license	free text
publication: when, where, and by whom was this dataset published?		
.date	Publication date (most recent version)	dd-mmm-yyyy
.affils	Author affiliations	comma-separated list of institutions
.sector	Sector from which the publication originated	single choice: [academia, industry, mixed]
.name	Publication name / reference	free text
.venue	Publication venue	free text
.url	Publication URL	URL
other: what is worth noting beyond the scope of this review?		
.notes	Additional notes	free text
.date_added	Date on which the dataset was added to SafetyPrompts.com	dd-mmm-yyyy

Table 3: **Codebook for our main review spreadsheet.** For each of the 102 datasets included in our review, we recorded 23 pieces of structured information.