

# GENERATIVE MONOCULTURE IN LARGE LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

## ABSTRACT

We introduce *generative monoculture*, a behavior observed in large language models (LLMs) characterized by a significant narrowing of model output diversity relative to available training data for a given task: for example, generating only positive book reviews for books with a mixed reception. While in some cases, generative monoculture enhances performance (e.g., LLMs more often produce efficient code), the dangers are exacerbated in others (e.g., LLMs refuse to share diverse opinions). As LLMs are increasingly used in high-impact settings such as education and web search, careful maintenance of LLM output diversity is essential to ensure a variety of facts and perspectives are preserved over time. We experimentally demonstrate the prevalence of generative monoculture through analysis of book review and code generation tasks, and find that simple countermeasures such as altering sampling or prompting strategies are insufficient to mitigate the behavior. Moreover, our results suggest that the root causes of generative monoculture are likely embedded within the LLM’s alignment processes, suggesting a need for developing fine-tuning paradigms that preserve or promote diversity.

## 1 INTRODUCTION

Large language models (LLMs) show promise due to their emergent abilities (Wei et al., 2022a) and state-of-the-art performance on several NLP tasks (Bommasani et al., 2023). However, concerns have been raised about the increasing reliance on LLM-based systems with insufficient testing, and how they impact society (Anwar et al., 2024; Wang et al., 2023). Recent evidence has shown that LLMs have dangerous tendencies: they convincingly return incorrect information (Dahl et al., 2024; Li et al., 2023a; Zhang et al., 2023), produce toxic language (Abid et al., 2021; Wen et al., 2023), and can effectively propagate misinformation (Barman et al., 2024; Sun et al., 2024).

In this paper, we focus on a different concern: that for a given prompt and task, *LLMs do not faithfully represent the diversity of potential responses available in their training data*. We call this behavior *generative monoculture*: given some task (e.g., generating book reviews) and a data attribute (e.g. sentiment), generative monoculture refers to the narrowing of the probability distribution of the considered attribute from source data (i.e., available human-written book reviews as part of the training data) to the generated data (i.e., LLM-generated book reviews).

As a preview, for book reviews (Fig. 1(Left)), we compare the diversity in sentiment of Goodreads reviews (Wan et al., 2019) (i.e., *src*)—very likely a portion of LLM training data (Achiam et al., 2023)—with LLM-generated reviews (i.e., *gen*). The range of mean sentiment scores per book across *gen* book reviews is much narrower than that in the *src*: in the experiment pictured, the average sentiment score for *gen* reviews are mostly over 0.85, whereas the the average sentiment

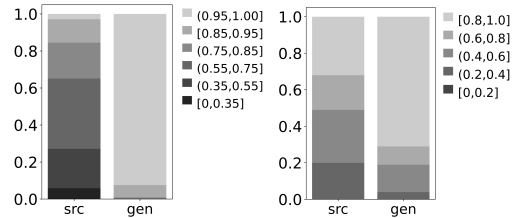


Figure 1: (Left) Comparison of the range of average-per-book sentiment scores for book reviews generated by an LLM (*gen*) and by human reviewers from the Goodreads dataset (*src*). Note the generated reviews have a much smaller range, as they are overwhelmingly positive. Model: Llama-2-chat. (Right) The spectrum of the mean pairwise Jaccard similarity among the algorithms of coding solutions. Note the generated code covers a narrower range of algorithms. Model: GPT-4

score for `src` reviews over the same books have a wider range from zero to one. For the code generation task (Fig. 1 (Right)), the range of algorithms employed in (correct) solutions to a given coding problem (i.e., `gen`) was much less varied than a sample of human answers (`src`) available on the web (Li et al., 2022): we show this through the range of Jaccard similarity of the algorithms employed in sets of human-written (`src`) and LLM-generated (`gen`) responses to a coding prompt.

Through the rapid adoption of LLMs such as ChatGPT, CoPilot, and Devin across education, code generation, and day-to-day information gathering, generative monoculture can harm society through loss of information, creativity, and intellectual diversity. For example, students asking LLMs questions about class material to get help researching for an essay may have their opinions formed without exposure to a sufficiently wide subset of available information; this will allow for certain opinions to die out over time. Concretely, a reduction in diversity of sentiment displayed above may lead to the loss of arguments from negative opinions on controversial books, potentially crucial for historical or literary context or a nuanced understanding of a book’s contributions.

Generative monoculture could even lead to security threats, depending on the application: software engineers across the globe relying on ChatGPT and CoPilot receiving similar code generations which do not reflect the true diversity of methods to solve a given problem may lead to similar code vulnerabilities across several large tech companies. Indeed, as we preview in Fig. 1 and show in detail in § 5, LLM output exhibits are less diverse than human-written solutions in the training data (e.g., by employing a narrower array of algorithms), which could lead to similarity across a wide range of code bases, leading in turn to repeated vulnerabilities (Perry et al., 2023; Pearce et al., 2022). However, in this case, LLM outputs not reflecting the full diversity of available coding examples can be positive: LLM outputs over-represent correct and efficient solutions. We present a nuanced picture of generative monoculture: while it can highlight the optimal portion of human-written data for some attributes, its pervasiveness across generation tasks and attributes *may* cause harm without careful intervention.

In this paper, we (1) define the concept of generative monoculture, and compare it to prior work around related topics (§ 2 and 7); (2) introduce a paradigm for measuring generative monoculture in LLMs (§ 3); (3) show experimental evidence for the prevalence of generative monoculture across a variety of application areas (book reviews and code) and (open-source and proprietary) LLMs, and provide some evidence for what may exacerbate generative monoculture, such as alignment tuning (Ouyang et al., 2022), (§ 5 and 6); and (4) show the (in)efficacy of several methods to abate generative monoculture: changing temperature, sampling, and prompting techniques (§ 5 and 6).

## 2 DEFINING GENERATIVE MONOCULTURE

We broadly characterize generative monoculture as a *distribution shift* from source data (i.e., human-written training data) to model generated data (i.e., model outputs) for a specific task, such as generating reviews for books or solutions to coding problems. This can be formalized using measures of statistical dispersion applied to various task-specific attributes.

**Definition 1** (Generative Monoculture). For a given task, let  $\mathcal{P}_{\text{src}}$  denote the probability distribution of the source data,  $\mathcal{P}_{\text{gen}}$  denote the probability distribution of the LLM-generated data,  $h$  denote a function extracting attributes from data (such as sentiment, or algorithms used in code), and  $\text{Dispersion}(\cdot)$  denote a dispersion metric (e.g., entropy). Then we define generative monoculture as the condition where  $\mathcal{P}_{\text{gen}}$  is statistically narrower than  $\mathcal{P}_{\text{src}}$ , namely:  $\text{Dispersion}(h(x)|x \sim \mathcal{P}_{\text{gen}}) < \text{Dispersion}(h(x)|x \sim \mathcal{P}_{\text{src}})$ .

Note,  $\mathcal{P}_{\text{src}}/\mathcal{P}_{\text{gen}}$  can be the distribution of human-written/model-generated responses, for a given task, conditioned on one specific given prompt (which we refer to as the *conditional distribution*), or the distribution of human-written/model-generated responses for *any possible prompt* in a considered domain (which we call the *unconditional distribution*).

This phenomenon signifies a shift towards less varied outputs. We emphasize that the investigation of generative monoculture is intrinsically task-dependent, as the attributes of interest differ across tasks. In addition, as we often do not have access to the source distribution in practice, we approximate it using a source dataset ( $D_{\text{src}}$ ), comprised of a subset of the training data of the LLMs. Similarly, we approximate the generated distribution through a dataset generated by the model ( $D_{\text{gen}}$ ).

**Generative Monoculture, Human Preference, and Alignment:** Generative monoculture can cause LLMs to over-emphasize human-preferred areas of a distribution for a certain data attribute; this is often desired behavior. For example, as we demonstrate in § 6, generative monoculture can result in having a narrower distribution of code correctness or efficiency biased towards correct, fast, and low-memory code. We conjecture this is a consequence of alignment procedures such as reinforcement learning with human feedback (RLHF) (Ouyang et al. (2022)).

However, when a tendency stemming from human preference bleeds beyond its intended use—e.g., a preference for positive sentiment affecting outputs that need or should not be positive—these seemingly advantageous behaviors can prevent an equally important goal: *maintaining diversity of opinion and expression*. Further, along data attributes which do not have a clear “preferred area” of the distribution, generative monoculture can limit the scope of methods, topics, or ideas expressed.

### 3 MEASURING GENERATIVE MONOCULTURE

We outline a general approach to measuring generative monoculture in LLMs. In particular, following Definition 1, we outline steps to construct  $D_{\text{src}}$  and  $D_{\text{gen}}$  and compare their diversity through extracting data attributes and calculating dispersion metrics. We illustrate our approach in Fig. 2.

#### 3.1 DATA CURATION

For a given task, we aim to create a source dataset that is likely to have been used in training the LLM we wish to investigate. Training data for most LLMs is a closely guarded secret. While recent work (Oren et al., 2023) describes how dataset contamination can be determined, such approaches are (a) riddled with false positives, and (b) computationally expensive. Thus, we often take an educated guess (based on dataset popularity and ease of use) in ascertaining if a given dataset is a likely training dataset candidate.

Formally, we define the source dataset as  $D_{\text{src}} = \{q_i, \text{src}_i\}_{i \in [N]}$  where (a)  $q_i$  is a problem instance within a task (e.g., name of a book for which a review has to be written), and (b)  $\text{src}_i = \{\text{src}_i^j\}_{j \in [n_i]}$  is a set of  $n_i$  human-written answers to the given prompt  $q_i$  (e.g., a set of  $n_i$  of book reviews for that particular book). In practice, we utilize existing datasets likely to be used during LLM training, and perform filtering and sub-sampling to obtain our  $D_{\text{src}}$ .

To create the model-generated dataset  $D_{\text{gen}}$ , for each sample  $q_i$ , we prompt the LLM ( $\mathcal{M}$ ) we wish to evaluate,  $m_i$  times to generate a set of responses,  $\text{gen}_i = \{\text{gen}_i^j\}_{j \in [m_i]}$ . Here,  $\text{gen}_i^j \leftarrow \mathcal{M}_j(P_{\text{task}}(q_i), \text{kwards})$  is the response obtained in the  $j$ -th call of  $\mathcal{M}$ , where (a)  $P_{\text{task}}$  denotes the task-specific formatting prompt that wraps the sample  $q_i$ , and (b)  $\text{kwards}$  denotes the generation keyword arguments (e.g., temperature) that specify the sampling strategy. Across both  $D_{\text{src}}$  and  $D_{\text{gen}}$ , we select or generate a large enough number of responses per  $q_i$  to ensure variety.

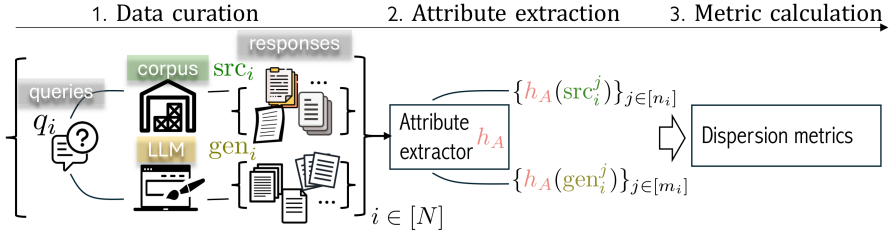


Figure 2: An overview of the procedure.

#### 3.2 ATTRIBUTE EXTRACTION

For a given task, we identify and compile a list of attributes that are of interest from the perspective of preserving diversity. This is a subjective task, but we focus on metrics which target understanding the content of LLM output (e.g., book review sentiment and topic; code time complexity and algorithms used), as opposed to more general metrics of output language quality such as saliency, fluency and coherence. More importantly, we need to ensure that extraction functions are efficient, accurate, and

reproducible—we outline our tests to ensure these qualities in §5 and the Appendices B and D. For example, care must be taken to use LLMs for attribute extraction, as they are known to be biased towards their own responses (Xu et al., 2024). For a given attribute  $A$ , extraction function  $h_A$  takes a string  $\text{tgt}$  to obtain the attribute value  $h_A(\text{tgt})$ . Note that  $\text{tgt}$  can either be  $\text{src}_i^j$  or  $\text{gen}_i^j$ . The extracted attribute can either be a continuous/categorical variable or of other more complicated types, depending on the nature of the attribute.

### 3.3 METRIC CALCULATION

As the last step, we compute metrics on the extracted attributes. Given a set of responses, dispersion metrics aim to capture their breadth or coverage of these attributes. We describe those used in this paper below and in Fig. 3.

**Dispersion Metrics.** We introduce dispersion metrics suited to different data types.

*A. Distribution of the mean.* For ordinal or continuous attributes, we calculate the mean over the conditional distribution—that is, we calculate the metrics over each  $\text{src}_i$  or  $\text{gen}_i$  (e.g. reviews for a given book), and show the distribution of this mean over all  $q_i$  for a certain task (e.g. over all books). While the mean itself does not directly measure dispersion, the *distribution* of the mean values sheds light on dispersion: concentrated mean values indicates a smaller dispersion of the data. The advantage is that it not only describes dispersion, but also the qualitative tendency of the attribute (e.g. bias towards positive/negative sentiment), which cannot be captured otherwise.

*B. Entropy and standard deviation.* For categorical attributes, we measure dispersion using entropy over the conditional distribution. For continuous attributes, we use standard deviation over the conditional distribution to quantify the dispersion of values around the mean, providing the variability.

*C. Mean pairwise similarity.* For attributes that are not easily characterized as categorical or continuous, we adopt specific similarity metrics catered to the data type. We then calculate the pairwise similarity values for the conditional distribution—i.e. calculate the mean similarity value for all the pairs of elements within each  $\text{src}_i/\text{gen}_i$ , then show the distribution for all  $N$  samples. We use:

1. *Mean pairwise Jaccard index:* Given two sets of categorical variables  $A$  and  $B$ , the Jaccard index,  $J(A, B) = |A \cap B| / |A \cup B|$ , measures similarity between them. An example of such a set could be a set of several algorithms inferred from a piece of code. A higher mean Jaccard index indicates a higher overall similarity between the set, and consequently, lower dispersion.
2. *Mean pairwise cosine similarity:* Given two multi-dimensional embeddings  $e_1$  and  $e_2$  obtained via a sentence embedder (Reimers & Gurevych, 2019), we calculate their similarity via cosine similarity, i.e.,  $S_C(e_1, e_2) = \langle e_1, e_2 \rangle / \|e_1\| \|e_2\|$ . A higher mean cosine similarity indicates a higher similarity and lower dispersion.
3. *Mean pairwise fingerprint similarity:* For tasks related to coding or computer programs, similarity is based on the overlap of selected hash values (or fingerprints) generated by Winnowing (Schleimer et al., 2003). We adopt an existing open-source tool COPYDETECT (Lingenfelter), which takes in a set of programs and returns the pairwise similarity scores for all programs in the set. We then calculate the mean value of these pairwise similarity scores as an indicator of the similarity for the set of programs. A higher mean fingerprint similarity indicates higher structural and syntactical similarity of the code.

In addition to the dispersion metrics, we consider one other approach—visualizing the top modes of unconditional distributions for certain attributes, e.g., topics. This helps identify areas of emphasis in  $\text{src}$  and  $\text{gen}$  distributions, as well as the tendency of change across distributions.

## 4 MITIGATING GENERATIVE MONOCULTURE

To attempt to mitigate generative monoculture, we test four methods known to increase LLM output diversity: increasing the *temperature*  $T$ , top- $p$  parameter, setting a temperature decay, and changing prompts. More details are in Appendices A and C.6.

**Temperature  $T$ .** This determines the dispersion of the probability distribution over the next token: increasing the temperature leads to a more flat probability distribution and increases the likelihood of sampling from less probable tokens, resulting in more diverse generations.

**Top- $p$ .** This controls the randomness of the generations by limiting the range of tokens considered. Specifically, it considers the smallest subset (consisting of the top probability tokens) whose cumulative probability exceeds the threshold  $p$ . A smaller  $p$  encourages the model to sample from a more focused set of likely tokens.

**Decaying Temperature.** We choose the starting temperature  $T = 10.0$  and follow a linear schedule for temperature decay, over the course of 50 time-steps (i.e., from the 1-st output token to the 50-th output token), with an ending temperature  $T = 1.2$ . The method is inspired by [Carlini et al. (2021)].

**Prompts.** Tuning the specific content and framing of the prompt can steer the model’s output more effectively [Brown et al., 2020; Sclar et al., 2023] and significantly impact the diversity of the generated text. We use “role-playing” or impersonation [Salewski et al., 2024], which instructs the model to produce the output in the persona of a specific person, and expect it to induce more personalized and varied responses.

## 5 EXPERIMENTAL SETUP

In this section, we describe our experimental setup for measuring and mitigating generative monoculture for two tasks, namely, generating book reviews and code solutions. We provide details for datasets, LLMs used, and most notably, the data attributes and metrics considered. We open source our code at <https://github.com/GeMoLLM/GeMO>.

### 5.1 GENERATING BOOK REVIEWS

**Data Curation:** For  $D_{\text{src}}$ , we use the Goodreads dataset [Wan et al., 2019], which contains multiple books with several reviews each. We perform filtering and sampling to ensure reliable attribute extraction (see Appendix B.1), and craft a final dataset of  $N = 742$  books with English titles, and  $\forall i, n_i = 10$  reviews per book such that the review length is between 300 and 700 words.

To obtain  $D_{\text{gen}}$ , we used the following LLMs:

(a) Llama-2-13b [Touvron et al., 2023] (henceforth referred to as Llama-2), (b) Llama-2-13b-chat [Touvron et al., 2023] (henceforth referred to as Llama-2-chat), (c) Vicuna-13b-v1.5 [Chiang et al., 2023] (henceforth referred to as Vicuna-13b), (d) GPT-3.5-turbo-instruct (0914) [Ouyang et al., 2022] (henceforth referred to as GPT-3.5), and (e) GPT-4-turbo (0125) [Microsoft 2024; Achiam et al., 2023] (henceforth referred to as GPT-4). We performed nucleus sampling [Holtzman et al., 2019] with various sampling parameters: (a) temperature  $T \in \{0.5, 0.8, 1.0, 1.2, 1.5\}$ , and (b) top- $p \in \{0.90, 0.95, 0.98, 1.00\}$ . We also experimented with two candidates for  $P_{\text{task}}$ : prompt (1) “Write a personalized review of the book titled {title}.”, and prompt (2) “Write a book review for the book titled {title} as if you are {person}.”. Prompt (2) was chosen as LLMs are known to generate more diverse responses when instantiated with a persona [Salewski et al., 2024]. We list the names of the 10 persons we considered in Appendix B.2. For comprehensiveness, we experimented with three more groups of prompts and report the results in Appendix C.7. For each combination of LLM, sampling parameter, and prompt, we independently sampled from the LLM 10 times to generate responses. We filtered out low-quality (generated) reviews by examining their perplexity (see Appendix C.1). This is to

	Attribute	Data type	Level	Metric
Book review	sentiment	categorical	C	mean, entropy
	topic	categorical	C	entropy
			U	distribution visualization
	wording	categorical	U	count, entropy
Coding	correctness	categorical	C	mean
	efficiency (complexity)	categorical	C	entropy
			U	distribution visualization
	efficiency (runtime)	continuous	C	mean, standard deviation
	fingerprint	hash values	C	mean pairwise fingerprint similarity
	code summary (text)	embedding	C	mean pairwise cosine similarity
	code summary (categorical)	categorical	C	mean pairwise Jaccard index

Figure 3: A summary of the scenarios, the attributes we consider, their data types, and the corresponding analysis levels as well as metrics. C and U stand for conditional and unconditional distributions.



ensure that the data used for analysis represents well-formed and coherent text, thereby improving the reliability of our findings. Thus,  $\forall i, m_i \leq 10$ .

**Attribute Extraction:** We want attributes that capture both the semantics and syntax of book reviews, representative of the key thematic and linguistic elements. Most importantly, while these attributes are not exhaustive, their extraction is reliable and efficient.

1. *Sentiment* indicates whether a review is positive (praising the book) or negative (criticizing the book). We employ a fine-tuned sentiment classifier (HuggingFace, b) as the attribute extractor which accepts text and returns a prediction in  $\{0, 1\}$ . This model has been downloaded  $\sim 5.4$  million times, and reaches an accuracy of 91.3 % on the dev set of SST-2 (Socher et al., 2013).
2. *Topic* refers to the themes discussed in a review (Wallach, 2006; Alghamdi & Alfalqi, 2015). We leverage BERTopic (Grootendorst, 2022) pre-trained on Wikipedia (HuggingFace, a) which assigns one topic to each review out of a total of  $\sim 2,000$  topics.
3. *Word choice* captures the lexical diversity in a review. To quantify this, we produce a frequency table of the unique words (see Appendix B.3), and immediately have the number of unique words.

**Metric Calculation:** For sentiment, we calculate mean and entropy for the conditional distribution. For topic, we calculate entropy for the conditional distribution as well as visualize the unconditional distribution of topics across all reviews, focusing on the top 10 classes. Finally, for word choice, we calculate count and entropy of the unconditional distribution.

## 5.2 GENERATING CODE SOLUTIONS

**Data Curation:** For  $D_{\text{src}}$ , we chose the CodeContests dataset (Li et al., 2022), a competitive programming problem dataset where each problem comes with multiple correct and incorrect solutions. We limited the scope to a subset ( $N = 100$ ) of level-A problems (easy problems) on Codeforces (CodeForces), and the language of the solutions to python3. More details in Appendix D.1. For each problem in the subset, we randomly sampled  $\forall i, n_i = 20$  correct solutions from all of the  $n_i^{\text{correct}}$  solutions for that problem.

To obtain  $D_{\text{gen}}$ , we use: (a) GPT-4, and (b) Claude-3-Sonnet (Anthropic, 2024). We did not use open-source LLMs, as these were not able to generate correct solutions for the problems we chose. More details are in Appendix E.4. We performed nucleus sampling (Holtzman et al., 2019) with various sampling parameters: (a) temperature  $T \in \{0.5, 1.0\}$ , and (b) top- $p \in \{0.9, 1.0\}$ . We used only one candidate for  $P_{\text{task}}$  i.e., “Please read the below problem description and generate a python code to solve the problem {problem description} Please only generate code and nothing else.” While we experimented with providing the LLM with a persona i.e., asking the LLM to pretend to be a “grandmaster in solving competitive programming problems”, the resulting accuracy was lower (see Appendix E.2). For each combination of LLM, sampling parameter, and prompt, we produce  $\forall i, m_i \geq 20$  generations such that at least 20 of the generated solutions were correct (details in Appendix D.2). We instantiated this by keep generating samples and measuring their correctness, until at some point all problems reached at least 20 correct solutions; we then stopped. This gave us  $k = 100$  for GPT-4 and  $k = 200$  for Claude-3.

**Attribute Extraction:** We consider the following attributes which characterize different aspects of code. We rely on GPT-3.5 for extracting some of the attributes; we manually verified the extracted attributes and confirmed their quality is high (see Appendix E.3).

1. *Correctness* refers to whether a piece of code correctly solves the given problem and passes all the test cases. We measure accuracy as the ratio of correct solutions among all solutions (details in Appendix D.3), to quantify the quality of human-/model-generated solutions.
2. *Efficiency* is crucial for scalability (Huang et al., 2024). This is measured through: asymptotic time/space complexity and runtime efficiency. We prompt GPT-3.5 to infer the big  $O$  time and space complexity (MacNeil et al., 2022), and execute the code on test cases to measure runtime and memory usage (see Appendix D.4).
3. *Fingerprint* provides insights into the structural and syntactical uniqueness of each code segment. As stated in Section 3.3, we use the COPYDETECT tool for this.

4. *Code Summary (textual)* explains the functionality of the code. Prior work has demonstrated the effectiveness of GPT-3.5 in code understanding (Nam et al., 2024). Thus, we use it to produce text-based summaries, and a description, functionality, algorithm, and data structure (prompt for this task is in Appendix D.5). To compare the similarity for these text summaries, we produce their embeddings using the all-MiniLM-L6-v2 model (HuggingFace, c).
5. *Code Summary (categorical)* reflects the techniques employed in the code through categorical tags, as used on the Codeforces website. We prompt GPT-3.5 to assign tags to a code segment by providing it a set of tags to choose from (prompt for this task is in Appendix D.5). We obtain one set per code segment. We similarly prompt GPT-3.5 to choose from a list of algorithms and data structures.

**Metric Calculation:** For correctness, we calculate the mean value i.e., accuracy over the conditional distribution. For efficiency (asymptotic complexity), we calculate: (a) entropy for the conditional distribution, and (b) plot the histogram for the unconditional distribution. For runtime efficiency, we calculate mean and standard deviation for the conditional distribution. We measure the following over the conditional distribution: (a) fingerprints, where we calculate the mean pairwise fingerprint similarity; (b) code summary (textual), where we calculate the mean pairwise cosine similarity in their embedding space; and (c) code summary (categorical), where we calculate the mean pairwise Jaccard index.

## 6 RESULTS AND TAKEAWAYS

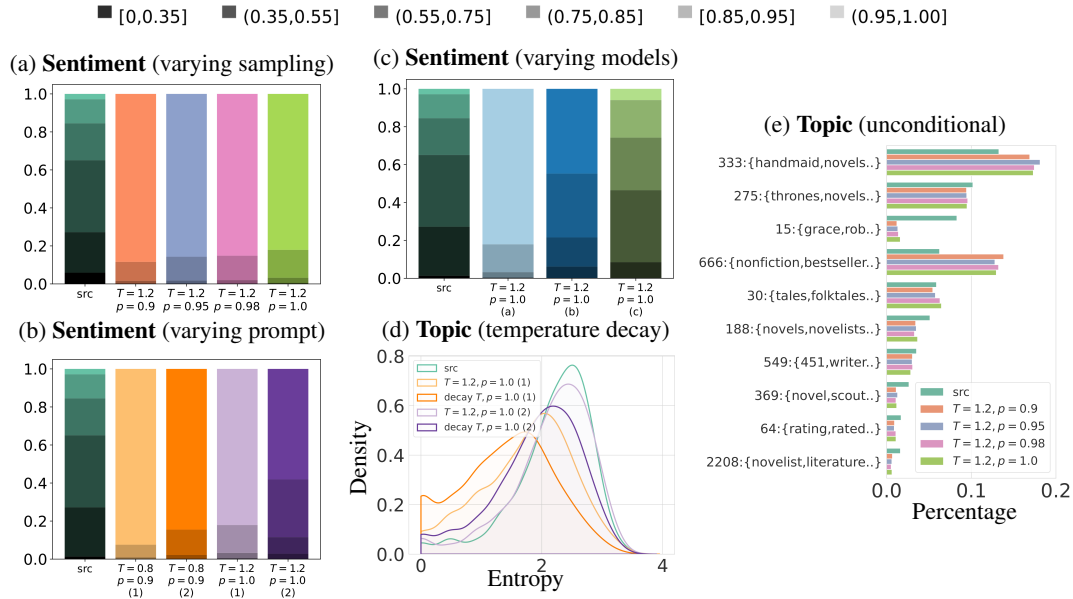


Figure 4: (a-c) stacked barplots for the mean sentiment scores under varying sampling parameters, prompts, and models. For these plots, in each bar, darker hues (bottom) represent lower scores while lighter one (top) denote higher scores. See the legend for the value range of each hue. In subfigure (b), (1) and (2) refer to the two prompts as introduced in Section 5.1. In subfigure (c), (a-c) refer to Llama-2-chat, Vicuna-13b, and Llama-2. (d) kernel density estimation (KDE) for the entropy values calculated on the conditional distribution of the topics. (e) unconditional topic distribution for top-10 topics. For all the subfigures, we mark the sampling parameters in them; unless marked with (1-2) or (a-c), the results are obtained on Llama-2-chat under prompt (1). These subfigures show that the model-generated reviews are overwhelmingly positive and cover a narrower range of the topics per book; moreover, there is distinctive under- and over-representation of the topics covered overall.

**Guide:** We present our results on measuring, and attempting to mitigate, generative monoculture in Fig. 4 and 5. We display results mainly in three formats: (a) stacked bar charts, where different hues correspond to different value ranges as indicated in the legend; (b) histograms (or grouped bar charts), to reflect the probability mass of a categorical variable; and (c) kernel density estimation

(KDE) plots, to reflect the estimated probability density of a continuous variable. We note that, in the code results, for all plots except that evaluating accuracy, we restrict to `correct` solutions.

**Takeaway 1: Monoculture Exists and is Severe, Within and Across LLMs.** As shown in Fig. 4 and 5, there exists significant narrowing from the source to generation distribution in all attributes considered for both scenarios, book reviews and coding.

Notably, for book review, *proprietary OpenAI LLMs (GPT-3.5 and GPT-4) demonstrate even more severe monoculture compared with the open-source Llama family LLMs* (see Fig. 9 in Appendix C.5 for more details). Particularly, for both GPT-3.5 and GPT-4, 100% of the samples have average positivity falling in (0.95, 1.00] under prompt (1), and 98.9% and 97.6% under prompt (2). For coding, similar reductions in diversity can be seen: in Fig. 5(e), we see increased similarity in natural language descriptions of LLM-generated code solutions, and Fig. 5(f) shows the Jaccard similarity of the generated solutions in terms of the inferred algorithms, with the majority of problems displaying high similarity across generated solutions. Of particular interest, the plagiarism scores of the LLM-generated code are extremely high (Fig. 5(b)), compared to the source solutions which achieve an utterly zero plagiarism score for *all* the problems. We examine a few pairs of examples and their plagiarism scores in Appendix E.1.

**Takeaway 2: LLMs tend to produce Human-favorable Generations.** Our results show that LLMs tend to over-represent parts of the attribute distribution that are preferred by humans: humans largely prefer text with positive sentiment (Dodds et al., 2015; Augustine et al., 2011; Boucher & Osgood, 1969) as well as correct and efficient code, and researchers have specifically infused these preferences into LLM assistants through preference tuning (e.g. RLHF) (Ouyang et al., 2022; Bai et al., 2022; Roziere et al., 2023). Fig. 4(a) and Appendix C.5 show that LLMs produce overwhelmingly positive generations. Fig. 5, as well as Fig. 22 and 23 in the Appendix reveal that LLM-generated code segments (a) are over  $2\times$  more accurate than the average human solutions, (b) enjoy an overall lower asymptotic time and space complexity, and (c) use less runtime and memory during execution. This may just be the intended consequence of RLHF, which explicitly optimizes the LLM towards producing human-favored responses in its objective, as guided by a reward model trained on human preferences.

However, as our results show, this implies a loss of diversity guided by human preferences, which, if only naively understood and enforced, could lead to unwanted consequences if going unnoticed. One example of the unintended artifacts is the under- and over-represented topics (Fig. 4(e)); the topic group 15 which contains keywords “rob” and “kill” etc. is significantly under-represented, likely a consequence of RLHF alignment tuning.

**Takeaway 3: RLHF Hurts Diversity the Most.** Llama-2-chat is obtained via performing RLHF tuning (Ouyang et al., 2022) on the pre-trained (PT) Llama-2. Similarly, Vicuna-13b is obtained via supervised fine-tuning (SFT) on the PT Llama-2 (Chiang et al., 2023). Comparisons on these LLMs (see Fig. 4(c), as well as Fig. 8 in Appendix C.4) show that the PT LLM-generated reviews are much more similar to the source. The PT LLM Llama-2 has 5.9% of samples with average sentiment values falling in the range of (0.95, 1.00], which is much closer to the source percentage of 2.8% than 44.7% for Vicuna-13b and 82.1% for Llama-2-chat. Vicuna-13b also shows better diversity than Llama-2-chat—this is consistent with findings suggesting RLHF reduces output diversity compared with SFT (albeit with different metrics) (Kirk et al., 2024).

**Takeaway 4: Naive Mitigations are Insufficient.** Changing the sampling parameter (increasing  $T$  and  $p$ ) and using a more diversity-inducing prompt (e.g., prompt (2) for book reviews) can reduce the gap (see Fig. 4(a-b) and Fig. 5). For example, using prompt (2) reduces the percentage of the most positive range from 82.1% to 58.1% in Fig. 4(b) for  $T = 1.2, p = 1.0$ . However, the gap is still large. More results in Appendix (Figures 13 and 14, and Appendix E) show similar conclusions.

We attempted two other strategies: (a) picking a higher temperature, and (b) leveraging a decaying temperature scheme (see § 4). Results in Appendix C.2 show that the gap still remains big even at such high randomness. Furthermore, for larger  $T$ , we notice a significant degradation of the generation quality as a result of the increased randomness. In Table 1, we present the average fraction of valid generations for Llama-2-chat and Vicuna-13b under various sampling parameters. The table shows that the valid number of generations rapidly drops as the randomness increases, particularly at  $T = 1.5$ ; the implication is that such a high randomness setting cannot be adopted for practical use.



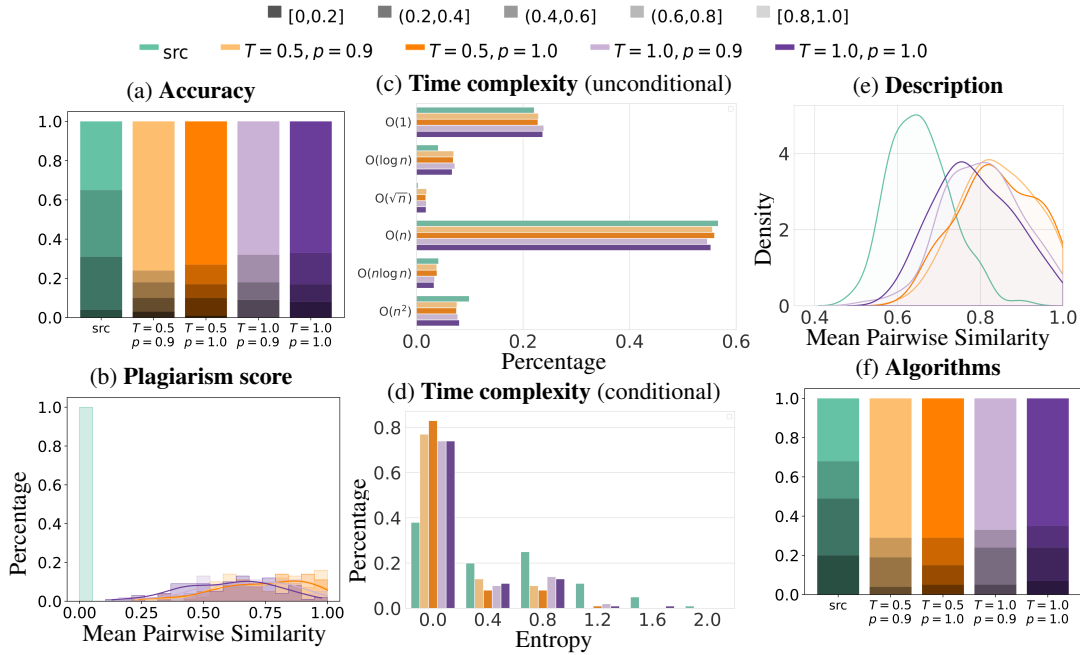


Figure 5: (Left) (a) stacked barplot for accuracy and (b) probability mass along with KDE for plagiarism scores. (Middle) Time complexity: (c) histogram of the (unconditional) distribution of the asymptotic complexity and (d) probability mass for the (conditional) distribution of entropy values. (Right) Selected code summary: (e) KDE plot for the mean pairwise cosine similarity scores for “description” as natural language and (f) stacked barplot for the mean pairwise Jaccard scores for “algorithms” as categorical values. Overall, the model-generated solutions are more accurate and efficient, display higher description similarity to each other, and cover a narrower span of algorithms. (More results in Appendix E.5)

Table 1: The average fraction of valid generations (out of a total of 10) for two models under various sampling parameters (temperature, top- $p$ , and prompt—denoted (1) and (2)). We regard a “valid” generation as text of perplexity value  $\leq 20$ —as a support, we present high perplexity samples in Appendix C.1. We observe that the number drops as the randomness increases (along the increase of both  $T$  and top- $p$  values). As a reference, GPT-4 achieves an average ratio of 1.000 at  $T = 1.2$  and  $p = 1.0$ .

LLama-2 -13b-chat	$p = 0.90$		$p = 0.95$		$p = 0.98$		$p = 1.00$		Vicuna -13b	$p = 0.90$		$p = 0.95$		$p = 0.98$		$p = 1.00$	
	(1)	(2)	(1)	(2)	(1)	(2)	(1)	(2)		(1)	(2)	(1)	(2)	(1)	(2)	(1)	(2)
$T = 0.5$	1.000	1.000	1.000	1.000	1.000	1.000	1.000	1.000	$T = 0.5$	0.987	0.992	0.987	0.990	0.984	0.989	0.984	0.988
$T = 1.0$	1.000	1.000	0.999	1.000	0.999	1.000	0.998	0.999	$T = 1.0$	0.961	0.962	0.951	0.958	0.942	0.953	0.935	0.947
$T = 1.5$	0.994	0.988	0.930	<b>0.883</b>	<b>0.743</b>	<b>0.649</b>	<b>0.512</b>	<b>0.394</b>	$T = 1.5$	<b>0.871</b>	0.903	<b>0.835</b>	<b>0.876</b>	<b>0.766</b>	<b>0.840</b>	<b>0.680</b>	<b>0.767</b>

Though prior work showed promise for decaying temperature to encourage diversity while maintaining quality (Carlini et al., 2021), this too failed to achieve higher diversity (see Fig. 4(d) and Appendix C.6).

## 7 RELATED WORK

**Diversity and LLMs.** Santurkar et al. (2023) demonstrate that LLMs do not give representative opinions to polling questions when compared to the general U.S. population. Our work focuses on the narrowing of diversity in LLM output from its human-written training data—while Santurkar *et al.* demonstrate a narrowing in diversity from actual human survey respondents (and not training data). Additionally, our work proposes a general framework for measuring monoculture. Padmakumar & He (2023) demonstrate that using LLM assistance can lead to reduced diversity in human-written argumentative essays when compared to essays written without LLM assistance. While they mention that this is partially because the models themselves do not produce diverse output, they do not focus on the narrowing of diversity from LLM training data to LLM-generated data. Finally, Zhang et al. (2024) propose an approach to fine-tune LLMs to output desired target distri-

butions, and Sorensen et al. (2024) outline an alignment framework to emphasize *pluralism* to work towards creating models which express a variety of opinions and perspectives. While these are certainly related to our work, generative monoculture as a phenomenon extends beyond differences in opinion, and expresses the narrowing of any number of task-specific attributes, from code correctness to topics covered to many others. One common thread across many of these works, which our work adds to, is that *current alignment practices—namely RLHF—harms output diversity*.

**Other Notions of Monoculture.** Our notion of generative monoculture relates to, but differs from, other notions of monoculture in the AI literature. For example, algorithmic monoculture (Kleinberg & Raghavan, 2021) and outcome homogeneity (Bommasani et al., 2022) describe the societal state where many decision-making actors rely on the same underlying algorithms to generate (classification or ranking) predictions, from the perspective of decision-making actors and individuals subject to those decisions respectively. These works show that algorithmic monoculture is sub-optimal for both decision-making actors (due to correlated failures across models) and for those subject to model decisions, as repeated outcomes across models leave little room for algorithmic recourse. In contrast, generative monoculture focuses on documenting the phenomenon of individual LLMs narrowing the diversity of their output in relation to their source data—for example, only returning positive book reviews about a controversial book. We do, however, document in this work that generative monoculture exists to similar extents and in similar directions across a variety of available LLMs, (e.g., Llama, Vicuna, ChatGPT-4) leaving open the possibility of concerns brought up by Kleinberg & Raghavan (2021), but in a generative context.

**Connections to Model Collapse:** We evaluate models trained on human-curated data, whereas model collapse evaluates models trained iteratively on synthetic data (either fully, or mixed with human data). In settings of model collapse where the model is trained only on synthetic data (Shumailov et al., 2023; Taori & Hashimoto, 2023), it is understandable that the generation quality is low. In contrast, our work shows that the generation quality is good (e.g., model generates correct coding solutions), but the “diversity” in generations is low. Unlike model collapse which converges to the mean of the distribution (Shumailov et al., 2023), our observation is there’s an emphasis on a specific part of the distribution which is not necessarily the mean. In this way, the work is tangentially related to collapse, but is not a special case of it (as the collapse phenomenon necessitates the distribution to match the mean with many rounds).

## 8 CONCLUSION AND LIMITATIONS

In this work, we introduce the concept of generative monoculture, a phenomenon where LLMs narrow the diversity of their output relative to their source data for a given task. We experimentally demonstrate its prevalence across text and code generation tasks, and show the difficulty in mitigating the behavior. Our work has limitations: first, we did not analyze the full training set of the LLMs we study due to time and compute restrictions, as the corpora are large and often proprietary. Further, as we note in § 3, measuring monoculture is difficult as selecting attributes is subjective, and the attribute extraction process is sensitive to the reliability of extraction techniques. (We verify our own attribute extraction techniques in the appendix). Further, while generative monoculture itself can have unfair consequences by enforcing the suppression of minority opinions, mitigating monoculture without extreme care could lead to the proliferation of harmful ideas or even toxicity by allowing for representation of the entire distribution of source text. We look forward to future work mitigating monoculture while maintaining low levels of toxicity and other dangerous behavior.

---

## REFERENCES

- Abubakar Abid, Maheen Farooqi, and James Zou. Persistent anti-muslim bias in large language models. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, pp. 298–306, 2021.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Rubayyi Alghamdi and Khalid Alfalqi. A survey of topic modeling in text mining. *Int. J. Adv. Comput. Sci. Appl.(IJACSA)*, 6(1), 2015.
- Anthropic. Claude-3 language model. <https://www.anthropic.com/>, 2024. Accessed: 2024-04-08.
- Usman Anwar, Abulhair Saparov, Javier Rando, Daniel Paleka, Miles Turpin, Peter Hase, Ekdeep Singh Lubana, Erik Jenner, Stephen Casper, Oliver Sourbut, et al. Foundational challenges in assuring alignment and safety of large language models. *arXiv preprint arXiv:2404.09932*, 2024.
- Adam A Augustine, Matthias R Mehl, and Randy J Larsen. A positivity bias in written and spoken english and its moderation by personality and gender. *Social Psychological and Personality Science*, 2(5):508–515, 2011.
- 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*, 2022.
- Dipto Barman, Ziyi Guo, and Owen Conlan. The dark side of language models: Exploring the potential of llms in multimedia disinformation generation and dissemination. *Machine Learning with Applications*, pp. 100545, 2024.
- Steven Bird and Edward Loper. NLTK: The natural language toolkit. In *Proceedings of the ACL Interactive Poster and Demonstration Sessions*, pp. 214–217, Barcelona, Spain, July 2004. Association for Computational Linguistics. URL <https://aclanthology.org/P04-3031>.
- Rishi Bommasani, Kathleen A Creel, Ananya Kumar, Dan Jurafsky, and Percy S Liang. Picking on the same person: Does algorithmic monoculture lead to outcome homogenization? *Advances in Neural Information Processing Systems*, 35:3663–3678, 2022.
- Rishi Bommasani, Percy Liang, and Tony Lee. Holistic evaluation of language models. *Annals of the New York Academy of Sciences*, 1525(1):140–146, 2023.
- Jerry Boucher and Charles E Osgood. The pollyanna hypothesis. *Journal of verbal learning and verbal behavior*, 8(1):1–8, 1969.
- 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.
- Nicholas Carlini, Florian Tramèr, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Ulfar Erlingsson, et al. Extracting training data from large language models. In *30th USENIX Security Symposium (USENIX Security 21)*, pp. 2633–2650, 2021.
- Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%\* chatgpt quality, March 2023. URL <https://lmsys.org/blog/2023-03-30-vicuna/>.
- CodeForces. Codeforces. URL <https://codeforces.com/>. Accessed: 2024-03-20 at <https://codeforces.com/>.

- OpenAI Community. Web chat default temperature for gpt-3.5 and 4. <https://community.openai.com/t/web-chat-default-temperature-for-gpt-3-5-and-4/167356/5>, 2023. accessed: 2024-04-08.
- Matthew Dahl, Varun Magesh, Mirac Suzgun, and Daniel E. Ho. Hallucinating law: Legal mistakes with large language models are pervasive, Jan 2024. URL <https://hai.stanford.edu/news/hallucinating-law-legal-mistakes-large-language-models-are-pervasive>.
- Peter Sheridan Dodds, Eric M Clark, Suma Desu, Morgan R Frank, Andrew J Reagan, Jake Ryland Williams, Lewis Mitchell, Kameron Decker Harris, Isabel M Kloumann, James P Bagrow, et al. Human language reveals a universal positivity bias. *Proceedings of the national academy of sciences*, 112(8):2389–2394, 2015.
- Maarten Grootendorst. Bertopic: Neural topic modeling with a class-based tf-idf procedure. *arXiv preprint arXiv:2203.05794*, 2022.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration. *arXiv preprint arXiv:1904.09751*, 2019.
- Dong Huang, Jie M Zhang, Yuhao Qing, and Heming Cui. Effibench: Benchmarking the efficiency of automatically generated code. *arXiv preprint arXiv:2402.02037*, 2024.
- HuggingFace. Maartengr/bertopic\_wikipedia, a. URL [https://huggingface.co/MaartenGr/BERTopic\\_Wikipedia](https://huggingface.co/MaartenGr/BERTopic_Wikipedia). Accessed: 2024-03-20 at [https://huggingface.co/MaartenGr/BERTopic\\_Wikipedia](https://huggingface.co/MaartenGr/BERTopic_Wikipedia).
- HuggingFace. distilbert/distilbert-base-uncased-finetuned-sst-2-english, b. URL <https://huggingface.co/distilbert/distilbert-base-uncased-finetuned-sst-2-english>. Accessed: 2024-03-20 at <https://huggingface.co/distilbert/distilbert-base-uncased-finetuned-sst-2-english>.
- HuggingFace. sentence-transformers/all-minilm-l6-v2, c. URL <https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2>. Accessed: 2024-03-20 at <https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2>.
- Michael Kerrisk. time(1) — linux manual page, 2023. URL <https://man7.org/linux/man-pages/man1/time.1.html>.
- Robert Kirk, Ishita Mediratta, Christoforos Nalmpantis, Jelena Luketina, Eric Hambro, Edward Grefenstette, and Roberta Raileanu. Understanding the effects of RLHF on LLM generalisation and diversity. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=PX03FAVHJT>.
- Jon Kleinberg and Manish Raghavan. Algorithmic monoculture and social welfare. *Proceedings of the National Academy of Sciences*, 118(22):e2018340118, 2021.
- Junyi Li, Xiaoxue Cheng, Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. HaluEval: A large-scale hallucination evaluation benchmark for large language models. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 6449–6464, Singapore, December 2023a. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.397. URL <https://aclanthology.org/2023.emnlp-main.397>.
- Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, et al. Starcoder: may the source be with you! *arXiv preprint arXiv:2305.06161*, 2023b.
- Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, Thomas Hubert, Peter Choy, Cyprien de Masson d’Autume, Igor Babuschkin, Xinyun Chen, Po-Sen Huang, Johannes Welbl, Sven

Gowal, Alexey Cherepanov, James Molloy, Daniel J. Mankowitz, Esme Sutherland Robson, Pushmeet Kohli, Nando de Freitas, Koray Kavukcuoglu, and Oriol Vinyals. Competition-level code generation with alphacode. *Science*, 378(6624):1092–1097, 2022. doi: 10.1126/science.abq1158. URL <https://www.science.org/doi/abs/10.1126/science.abq1158>.

Bryson Lingenfelter. blingenf/copydetect: Code plagiarism detection tool. URL <https://github.com/blingenf/copydetect>. Accessed: 2024-03-20 at <https://github.com/blingenf/copydetect>.

Stephen MacNeil, Andrew Tran, Dan Mogil, Seth Bernstein, Erin Ross, and Ziheng Huang. Generating diverse code explanations using the gpt-3 large language model. In *Proceedings of the 2022 ACM Conference on International Computing Education Research-Volume 2*, pp. 37–39, 2022.

Microsoft. Openai models - azure openai service. <https://learn.microsoft.com/en-us/azure/ai-services/openai/concepts/models>, 2024. accessed: 2024-04-08.

Daye Nam, Andrew Macvean, Vincent Hellendoorn, Bogdan Vasilescu, and Brad Myers. Using an llm to help with code understanding. In *2024 IEEE/ACM 46th International Conference on Software Engineering (ICSE)*, pp. 881–881. IEEE Computer Society, 2024.

OpenAI. Chat completions - openai api. <https://platform.openai.com/docs/api-reference/chat/create>, 2024. accessed: 2024-04-08.

Yonatan Oren, Nicole Meister, Niladri Chatterji, Faisal Ladhak, and Tatsunori B Hashimoto. Proving test set contamination in black box language models. *arXiv preprint arXiv:2310.17623*, 2023.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35: 27730–27744, 2022.

Vishakh Padmakumar and He He. Does writing with language models reduce content diversity? *arXiv preprint arXiv:2309.05196*, 2023.

Hammond Pearce, Baleegh Ahmad, Benjamin Tan, Brendan Dolan-Gavitt, and Ramesh Karri. Asleep at the keyboard? assessing the security of github copilot’s code contributions. In *2022 IEEE Symposium on Security and Privacy (SP)*, pp. 754–768. IEEE, 2022.

Neil Perry, Megha Srivastava, Deepak Kumar, and Dan Boneh. Do users write more insecure code with ai assistants? In *Proceedings of the 2023 ACM SIGSAC Conference on Computer and Communications Security*, pp. 2785–2799, 2023.

Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 11 2019. URL <https://arxiv.org/abs/1908.10084>.

Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*, 2023.

Leonard Salewski, Stephan Alaniz, Isabel Rio-Torto, Eric Schulz, and Zeynep Akata. In-context impersonation reveals large language models’ strengths and biases. *Advances in Neural Information Processing Systems*, 36, 2024.

Shibani Santurkar, Esin Durmus, Faisal Ladhak, Cinoo Lee, Percy Liang, and Tatsunori Hashimoto. Whose opinions do language models reflect? *arXiv preprint arXiv:2303.17548*, 2023.

Saul Schleimer, Daniel S Wilkerson, and Alex Aiken. Winnowing: local algorithms for document fingerprinting. In *Proceedings of the 2003 ACM SIGMOD international conference on Management of data*, pp. 76–85, 2003.



- Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. Quantifying language models’ sensitivity to spurious features in prompt design or: How i learned to start worrying about prompt formatting. *arXiv preprint arXiv:2310.11324*, 2023.
- Ilia Shumailov, Zakhar Shumaylov, Yiren Zhao, Yarin Gal, Nicolas Papernot, and Ross Anderson. The curse of recursion: Training on generated data makes models forget, 2023.
- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D Manning, Andrew Y Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pp. 1631–1642, 2013.
- Taylor Sorensen, Jared Moore, Jillian Fisher, Mitchell Gordon, Niloofar Mireshghallah, Christopher Michael Rytting, Andre Ye, Liwei Jiang, Ximing Lu, Nouha Dziri, et al. A roadmap to pluralistic alignment. *arXiv preprint arXiv:2402.05070*, 2024.
- Yanshen Sun, Jianfeng He, Limeng Cui, Shuo Lei, and Chang-Tien Lu. Exploring the deceptive power of llm-generated fake news: A study of real-world detection challenges. *arXiv preprint arXiv:2403.18249*, 2024.
- Rohan Taori and Tatsunori Hashimoto. Data feedback loops: Model-driven amplification of dataset biases. In *International Conference on Machine Learning*, pp. 33883–33920. PMLR, 2023.
- 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*, 2023.
- Hanna M Wallach. Topic modeling: beyond bag-of-words. In *Proceedings of the 23rd international conference on Machine learning*, pp. 977–984, 2006.
- Mengting Wan, Rishabh Misra, Ndapa Nakashole, and Julian J. McAuley. Fine-grained spoiler detection from large-scale review corpora. In Anna Korhonen, David R. Traum, and Lluís Màrquez (eds.), *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, pp. 2605–2610. Association for Computational Linguistics, 2019. doi: 10.18653/V1/P19-1248. URL <https://doi.org/10.18653/v1/p19-1248>.
- 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. *arXiv preprint arXiv:2306.11698*, 2023.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*, 2022a.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022b.
- Jiaxin Wen, Pei Ke, Hao Sun, Zhexin Zhang, Chengfei Li, Jinfeng Bai, and Minlie Huang. Unveiling the implicit toxicity in large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 1322–1338, 2023.
- Wenda Xu, Guanglei Zhu, Xuandong Zhao, Liangming Pan, Lei Li, and William Yang Wang. Perils of self-feedback: Self-bias amplifies in large language models. *arXiv preprint arXiv:2402.11436*, 2024.
- Yiming Zhang, Avi Schwarzschild, Nicholas Carlini, Zico Kolter, and Daphne Ippolito. Forcing diffuse distributions out of language models. *arXiv preprint arXiv:2404.10859*, 2024.
- Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, et al. Siren’s song in the ai ocean: a survey on hallucination in large language models. *arXiv preprint arXiv:2309.01219*, 2023.