Assessing Variations in Open Datasets for Training Large Language Models: Biases and Benchmarking

Vincent Koc, vincentkoc@ieee.org

Hyperthink, Australia

Abstract:

Open datasets are critical to the development and training of large language models (LLMs). However, variations in dataset composition often introduce biases that can impact model performance and reliability. This Article investigates the nature and extent of these variations, categorizes biases inherent in datasets, and examines their implications on LLM training. We also evaluate benchmarking standards currently employed to measure LLM performance and propose enhancements for a fairer and more inclusive evaluation framework. Through extensive experiments and analyses, we reveal the consequences of dataset heterogeneity and demonstrate practical strategies for mitigating biases. Our findings emphasize the importance of transparent dataset curation and robust benchmarking practices to ensure the ethical development of LLMs.

Keywords: Open Datasets, Large Language Models, Biases, Benchmarking, Dataset Variations, NLP, Dataset Evaluation

I. Introduction

The rapid evolution of large language models (LLMs) has revolutionized natural language processing (NLP) tasks, enabling significant advancements in machine translation, question answering, and content generation. These models rely heavily on large-scale datasets for training, where the quality and composition of the dataset directly influence model performance. Open datasets, celebrated for their accessibility and cost-effectiveness, have become the backbone of LLM development [1]. However, variations in dataset properties such as size, diversity, and representativeness can lead to biases that undermine the fairness, inclusivity, and reliability of LLMs. For instance, skewed representation in training data can result in models that perpetuate stereotypes or exclude minority groups. Additionally, the lack of standardized benchmarks complicates the evaluation of LLMs across diverse use cases [2].

This Article aims to address these challenges by systematically analyzing variations in open datasets, categorizing biases, and evaluating the efficacy of existing benchmarking practices. We conduct experiments to measure the impact of dataset properties on model performance, providing actionable insights for improving dataset curation and evaluation frameworks. Our findings underscore the need for a collaborative effort among researchers, developers, and policymakers to ensure the ethical and effective deployment of LLMs.

II. Background and Related Work

The success of LLMs, such as GPT and BERT, hinges on the availability of extensive and diverse datasets. Researchers have highlighted the critical role of dataset quality in influencing model capabilities and generalization. Open datasets like Common Crawl, Wikipedia, and OpenWebText have been widely adopted, offering diverse sources of information. However, their construction and use in language models are not always transparent [3, 4], leading to issues such as duplication, noise, and bias. Bias in datasets has been extensively studied in the context of machine learning. Representation bias, where certain groups are overrepresented or underrepresented, and selection bias, stemming from non-random sampling of data, are particularly problematic. Annotator bias, introduced during the labeling process, further complicates the issue. These biases not only affect model accuracy but also raise ethical concerns, particularly when LLMs are deployed in sensitive domains such as healthcare or legal services.

Benchmarking standards are equally important but often fail to address dataset-induced biases. Popular benchmarks, including GLUE and SuperGLUE, focus on task-specific performance metrics without considering fairness or inclusivity. Recent efforts, such as the introduction of ethical AI guidelines, aim to bridge these gaps, but a comprehensive solution remains elusive.

III. Assessing Variations in Open Datasets

The composition and diversity of datasets are crucial factors in determining their efficacy for training LLMs. Datasets vary widely in their sources, languages, and topical coverage, which directly influence the generalization capabilities of models trained on them.

Datasets can be broadly categorized by their sources, such as curated encyclopedic content, web crawls, or domain-specific text. For instance, Wikipedia provides high-quality curated text but lacks the linguistic and topical diversity of web-based datasets like Common Crawl. These differences result in significant variation in vocabulary, linguistic constructs, and topical coverage, affecting model training. Temporal diversity is another critical aspect. Datasets constructed over specific time periods may fail to capture evolving language usage, slang, or emerging topics. The inclusion of temporally diverse data ensures that models remain relevant and adaptable to real-world applications.

Quantifying diversity in datasets involves measuring aspects such as language representation, domain coverage, and demographic inclusion. Language diversity, for example, can be assessed through metrics like the proportion of non-English text. Similarly, domain diversity metrics analyze the spread of topics within the dataset. These metrics are essential for evaluating how representative a dataset is of the target population. A dataset lacking in linguistic or demographic diversity may lead to biased models that fail to generalize effectively across various user groups [5].

IV. Methods for Measuring Variations:

Measuring variations in datasets is a multifaceted process that involves both quantitative and qualitative methods. These methods help identify the degree of bias, noise, and heterogeneity present in the datasets [6].

Quantitative methods include statistical metrics such as vocabulary size, word frequency distribution, and sentence length variance. These metrics provide insights into the linguistic richness and structural properties of the dataset. Additionally, demographic analysis can identify overrepresented or underrepresented groups within the data. Another key metric is the bias index, which quantifies disparities in representation or content. For instance, measuring the gender pronoun ratio in text can reveal gender biases, while analyzing the prevalence of specific cultural references may highlight ethnocentric tendencies.

Qualitative methods involve manual inspection of dataset samples to identify patterns or anomalies. This can include reviewing text for offensive or outdated language, assessing topic relevance, and identifying instances of misinformation. Manual reviews are particularly useful for detecting subtle biases that may not be evident through statistical analysis [7]. For example, a qualitative review might uncover narratives that reinforce stereotypes or marginalize minority groups, even if these biases are not statistically significant. Below Figure 1 shows the Examples of the open Datasets.

Dataset	Size (GB)	Languages	Source Types
Common Crawl	100+	Multiple	Web Content
Wikipedia	20	Multiple	Encyclopedia Entries
OpenWebText	40	English	Web Articles

Figure 1: Shows the Examples of open Datasets.

V. Challenges in Standardization

Standardizing open datasets presents significant challenges due to the diverse nature of data sources, languages, and cultural contexts [8]. These challenges must be addressed to ensure fair and equitable use of datasets in LLM training. One of the primary challenges is the absence of universally accepted guidelines for dataset curation. Different organizations and researchers often adopt varying practices for data collection, cleaning, and annotation [9]. This lack of standardization results in inconsistent dataset quality and reliability. Efforts to create shared repositories with standardized metadata and documentation have gained traction. However, achieving global consensus on these standards remains elusive, particularly when accounting for regional and cultural differences in data practices.

Datasets often reflect the cultural and regional biases of their creators, leading to an overrepresentation of dominant languages and perspectives. For instance, English-language datasets are abundant, while low-resource languages are significantly underrepresented.

Addressing this imbalance requires collaborative efforts to collect and document data from diverse linguistic and cultural backgrounds. Standardization must also account for varying data privacy regulations and ethical considerations across regions. For example, data anonymization practices may differ between countries, impacting the availability and usability of datasets for global research.

VI. Biases in Open Datasets

Types of Biases:

- 1. **Representation Bias:** Disproportionate representation of specific demographics, leading to exclusion or overemphasis of certain groups.
- 2. **Selection Bias:** Non-random inclusion of data sources that skew the dataset toward particular domains or perspectives.
- 3. **Annotation Bias:** Subjectivity in labeling data that reflects annotators' cultural or personal biases, which can affect the interpretability of the dataset.

Case Studies of Bias: Analysis of datasets like Common Crawl reveals overrepresentation of English-language content, marginalizing low-resource languages. Similarly, gender biases are prevalent in datasets sourced from online forums or social media platforms. Studies have also identified instances of political bias in datasets sourced from news websites, where certain viewpoints dominate.

Consequences of Bias: Biases manifest as disparities in model predictions, such as gendered pronoun misclassification or culturally inappropriate responses. These issues undermine user trust and limit the applicability of LLMs across diverse contexts. In sensitive applications like legal or medical NLP tasks, biases can lead to significant ethical and operational risks.

Mitigation Strategies: Approaches to mitigate bias include balanced sampling, debiasing algorithms, and post-hoc adjustments to model outputs. Balanced sampling involves proportionate inclusion of diverse demographic and linguistic groups. Debiasing algorithms, such as adversarial training or reweighting, aim to neutralize bias during model training. Post-hoc techniques adjust model outputs to align with fairness criteria. Collaborative efforts to

develop inclusive datasets are also essential, involving stakeholders from varied cultural and linguistic backgrounds.

VII. Benchmarking Standards for LLMs

Benchmarking provides a standardized framework for evaluating model performance. It facilitates comparisons across models and highlights areas for improvement [10]. Benchmarks also serve as a guide for researchers to identify strengths and weaknesses in model capabilities [11]. Blow is the Existing Benchmark, task covered and their limitations.

Benchmark	Tasks Covered	Limitations
GLUE	Text Classification	Lacks fairness metrics
SuperGLUE	NLP Tasks	Focuses on English-only datasets
BigBench	Multitask Learning	Limited diversity in benchmarks

Figure 2: shows the Existing benchmarks.

Incorporating fairness, inclusivity, and robustness into benchmarks can address existing limitations. Metrics that evaluate bias, inclusivity, and adaptability are particularly valuable. For instance, multilingual benchmarks can assess models' proficiency across diverse languages. Benchmarks should also incorporate real-world scenarios, such as cross-cultural communication tasks, to evaluate practical applicability. Collaboration with domain experts can ensure benchmarks reflect ethical and technical considerations [12].

VIII. Experiment Design and Methodology

We selected three widely used datasets: Common Crawl, Wikipedia, and OpenWebText. These datasets vary in size, source, and linguistic diversity. Common Crawl represents an extensive and diverse collection of web-based content [13]. Wikipedia is a curated and factual dataset. OpenWebText bridges the gap with a mix of web content and curated quality. Metrics Used are:

- i. Vocabulary size and overlap.
- ii. Demographic representation.
- iii. Bias metrics, such as word association tests.
- iv. Task-specific performance metrics, including accuracy and F1 scores.

Workflow:

- 1. Preprocessing: Tokenization, deduplication, and noise removal.
 - o Tokenization involves breaking down text into manageable units for model training.
 - o Deduplication ensures that repeated text does not distort model learning.
 - Noise removal eliminates irrelevant or low-quality data, such as corrupted files or nonlinguistic content.
- **2.** Analysis: Measuring variations and biases.
 - Statistical tools and scripts identify key characteristics like word frequency and distribution.
 - Bias detection focuses on analyzing demographic representation and language usage patterns.
- **3.** Model Training: Fine-tuning LLMs on each dataset.
 - o Models are trained using frameworks like PyTorch or TensorFlow to ensure consistent performance comparisons.
 - o Training involves iterative cycles to optimize loss functions and improve output quality.
- **4.** Evaluation: Comparing model outputs using benchmarks.
 - Benchmarks assess task-specific performance and include metrics for diversity and fairness.

Experimental Setup:

The experiments were conducted using a standardized training pipeline to ensure reproducibility. Hyperparameters were aligned across datasets to maintain consistency. Model performance was evaluated on tasks spanning sentiment analysis, question answering, and text generation. Key aspects include:

i. **Hardware Resources**: Experiments utilized NVIDIA A100 GPUs with 40 GB memory per GPU, allowing efficient processing of large datasets. Clusters with multiple GPUs enabled parallel training for reduced computational time.

- ii. **Hyperparameter Settings**: Learning rates were set to 3e-5, with batch sizes of 32 and a sequence length of 512 tokens. Optimization algorithms such as AdamW were used to ensure convergence during training.
- iii. **Dataset Partitioning**: Each dataset was divided into training (80%), validation (10%), and test (10%) sets to maintain consistent evaluation protocols.
- iv. **Performance Metrics**: Models were evaluated on precision, recall, and F1 scores across benchmark tasks, ensuring comprehensive analysis.

IX. Results and Discussion

Dataset	Vocabulary Size	Bias Index	Model Accuracy (%)	F1 Score (%)
Common Crawl	1M+	0.8	85	78
Wikipedia	500K	0.5	88	83
OpenWebText	700K	0.7	86	81

Figure 3: shows the Experimental Result.

The results reveal that Wikipedia datasets, with curated content, lead to higher accuracy and lower bias indices [14]. Common Crawl, while extensive, suffers from higher biases due to unfiltered web sources. OpenWebText strikes a balance but lacks diversity in linguistic representation. The F1 scores indicate the robustness of Wikipedia's dataset in task-specific evaluations. These findings highlight the trade-offs between dataset size, diversity, and quality. Models trained on Common Crawl exhibit higher variability in performance due to noise, while Wikipedia offers consistency at the expense of breadth.

X. Recommendations and Future Directions

Creating datasets with proportional representation across demographics and languages is critical. Collaboration between academia and industry can facilitate this process. Leveraging automated tools for dataset analysis can enhance scalability and efficiency. Benchmarks must evolve to include metrics for fairness and inclusivity. Developing multilingual benchmarks is also essential for global applicability. Incorporating real-world tasks into benchmarks can improve the relevance of evaluations.

Open Questions:

- 1. How can we ensure transparency in dataset curation?
- 2. What role can policymakers play in standardizing benchmarks?
- 3. How can low-resource languages be better represented?
- 4. How can ethical considerations be integrated into benchmarking practices?

Conclusion

This Article underscores the critical role of dataset variations and biases in shaping LLM performance. Through a comprehensive analysis, we demonstrate the need for balanced dataset curation and enhanced benchmarking practices. Addressing these challenges is vital for fostering equitable and effective AI systems. Future efforts must prioritize inclusivity, transparency, and collaboration to ensure that LLMs serve diverse global communities responsibly. Furthermore, the benchmarking of LLMs against diverse datasets is essential to ensure that they perform robustly across different contexts, without favoring certain perspectives or data distributions. The evaluation of multiple metrics—such as fairness, inclusivity, and performance—against these datasets is paramount to identifying potential areas for improvement. Without rigorous benchmarking, LLMs may inadvertently produce biased or unreliable outputs, undermining their usefulness in real-world applications. Creating more diverse and representative datasets, along with comprehensive and transparent evaluation frameworks will be key to reducing bias in AI systems. Future research should focus on developing methods to mitigate dataset biases during the collection, curation, and preprocessing phases, while also exploring techniques to fine-tune models in ways that promote fairness without sacrificing performance.

REFERENCES:

- [1] R. Dutta, "Benchmarking stereotype bias and toxicity in large language models," University of Illinois at Urbana-Champaign, 2024.
- [2] I. O. Gallegos *et al.*, "Bias and fairness in large language models: A survey," *Computational Linguistics*, pp. 1-79, 2024.
- [3] L. Gao *et al.*, "The pile: An 800gb dataset of diverse text for language modeling," *arXiv* preprint arXiv:2101.00027, 2020.

- [4] T. Wu, M. Terry, and C. J. Cai, "Ai chains: Transparent and controllable human-ai interaction by chaining large language model prompts," in *Proceedings of the 2022 CHI conference on human factors in computing systems*, 2022, pp. 1-22.
- [5] C. Xu, S. Guan, D. Greene, and M. Kechadi, "Benchmark Data Contamination of Large Language Models: A Survey," *arXiv preprint arXiv:2406.04244*, 2024.
- [6] A. R. Ives, P. E. Midford, and T. Garland Jr, "Within-species variation and measurement error in phylogenetic comparative methods," *Systematic biology*, vol. 56, no. 2, pp. 252-270, 2007.
- [7] J. Umbrich, S. Neumaier, and A. Polleres, "Quality assessment and evolution of open data portals," in *2015 3rd international conference on future internet of things and cloud*, 2015: IEEE, pp. 404-411.
- [8] H. W. Vesper, G. L. Myers, and W. G. Miller, "Current practices and challenges in the standardization and harmonization of clinical laboratory tests," *The American journal of clinical nutrition*, vol. 104, pp. 907S-912S, 2016.
- [9] A. Aapaoja and H. Haapasalo, "The challenges of standardization of products and processes in construction," in *Proceedings of the 22nd Annual Conference of the International Group for Lean*, 2014: Citeseer, pp. 983-993.
- [10] M.-Y. Chan and S.-M. Wong, "A comparative analysis to evaluate bias and fairness across large language models with benchmarks," 2024.
- [11] X. Li *et al.*, "Benchmarking Bias in Large Language Models during Role-Playing," *arXiv* preprint arXiv:2411.00585, 2024.
- [12] L. Yuan *et al.*, "Revisiting out-of-distribution robustness in nlp: Benchmarks, analysis, and LLMs evaluations," *Advances in Neural Information Processing Systems*, vol. 36, pp. 58478-58507, 2023.
- [13] E. G. Martin, J. Law, W. Ran, N. Helbig, and G. S. Birkhead, "Evaluating the quality and usability of open data for public health research: a systematic review of data offerings on 3 open data platforms," *Journal of Public Health Management and Practice*, vol. 23, no. 4, pp. e5-e13, 2017.
- [14] J. Dhamala *et al.*, "Bold: Dataset and metrics for measuring biases in open-ended language generation," in *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*, 2021, pp. 862-872.