

1 **Redefining Value: Wikipedia's Role in Dataset Governance for Foundation**
2 **Models**
3

4 YAXUAN YIN, The Information School, University of Wisconsin–Madison, USA
5

6 Wikipedia has long been central to CSCW research on online collaboration, but it now plays a dual role: as both a collaborative
7 knowledge platform and a critical dataset powering large language models (LLMs). This shift raises open questions about how
8 contributions are valued, recognized, and governed when they influence downstream AI behavior. This position paper explores
9 influence-based data valuation as one potential lens and invites discussion on how CSCW insights into collaboration and governance
10 can inform responsible dataset composition, curation, and release.
11

12 CCS Concepts: • **Human-centered computing** → **Empirical studies in collaborative and social computing; Empirical studies**
13 **in HCI.**
14

15 Additional Key Words and Phrases: Responsible AI, Foundation Models, Dataset Governance, Wikipedia, Generative AI, Data Valuation
16

17 **ACM Reference Format:**

18 Yaxuan Yin. 2018. Redefining Value: Wikipedia's Role in Dataset Governance for Foundation Models. In *Proceedings of Make sure to*
19 *enter the correct conference title from your rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, 3 pages.
20
21 <https://doi.org/XXXXXXX.XXXXXXX>

22 **1 Introduction**

23 Within CSCW research, Wikipedia has long served as a central site for studying online collaboration. Prior studies have
24 examined how volunteers coordinate [4], manage misinformation and vandalism [6], and navigate disparities across
25 topics, languages, and regions [5]. These studies have shaped our understanding of what makes a contribution valuable
26 within Wikipedia's community — contributions that improve accuracy, fill knowledge gaps, and strengthen collaborative
27 dynamics. However, the rise of large language models (LLMs) introduces a new layer of complexity. Wikipedia now
28 plays a dual role: it is both a collaborative knowledge platform and a critical dataset powering foundation models [9].
29 For example, Wikipedia pages are among the five primary datasets used to train GPT-3 [1], the FLORES-101 benchmark
30 relies on Wikipedia for multilingual evaluation [3], and retrieval-augmented generation (RAG) pipelines treat Wikipedia
31 as a key source of factual knowledge [7].
32

33 This shift introduces a new dimension of value: contributions to Wikipedia now matter not only for supporting
34 human collaboration but also for shaping how AI systems behave and what knowledge they reproduce. Specifically,
35 they influence who gets credited (attribution), whose perspectives are amplified or marginalized (representation), and
36 who controls how data is used and governed (governance). Yet contributors currently lack transparency, recognition,
37 and control over how their work is used in training datasets. At the same time, tensions between data providers and AI
38 developers are intensifying, as reflected in emerging legal disputes over attribution, consent, and compensation [8].
39

40 Author's Contact Information: Yaxuan Yin, yaxuan.yin@wisc.edu, The Information School, University of Wisconsin–Madison, Madison, Wisconsin, USA.
41

42 Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not
43 made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components
44 of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on
45 servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
46

47 © 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.
48

49 Manuscript submitted to ACM
50

51 Manuscript submitted to ACM
52

53 These dynamics open new opportunities for CSCW researchers to revisit longstanding questions of value, governance,
 54 and collaboration in a context where human and machine knowledge production are deeply intertwined. Understanding
 55 these evolving relationships is critical for shaping equitable, transparent, and sustainable data ecosystems in the era of
 56 LLMs.
 57

60 **2 Opportunity: Using Data Valuation as a Lens**

61 One potential lens for exploring these issues comes from recent advances in data valuation methods. For example, Choe
 62 et al. [2] introduce techniques for estimating the influence of individual training examples on a model’s performance.
 63 They provide a useful technical framework for reflecting on how “value” is defined when Wikipedia contributions
 64 shape downstream AI behaviors. At a high level, the approach takes three main inputs: (1) a trained language model, (2)
 65 a snapshot of the training dataset (e.g., Wikipedia), and (3) an evaluation benchmark, such as factual QA, multilingual
 66 reasoning, or fairness tasks like WinoBias [10]. Using these inputs, the method produces influence scores that estimate
 67 how much each data point contributes to the model’s performance on the chosen benchmark. These scores can be
 68 aggregated at multiple levels – individual tokens, pages, or contributors – to understand which parts of Wikipedia
 69 most affect model outputs.

70 Applying this framework conceptually to Wikipedia allows us to connect influence scores to broader Responsible AI
 71 dimensions. For example, examining which pages or contributors most affect model outputs can reveal potential fairness
 72 and representation gaps, such as the dominance of English-language content in shaping LLM behavior [?]. Influence
 73 scores can also enhance transparency by linking model responses back to specific pages or contributors, enabling
 74 better attribution and accountability. These insights provide a foundation for rethinking dataset composition, curation
 75 practices, and how community norms might be incorporated into foundation model development. However, data
 76 valuation is not the endpoint, it is a starting point for new CSCW discussions. Influence scores highlight what matters
 77 but not why, and they cannot explain the social, cultural, or governance dynamics behind Wikipedia contributions – for
 78 instance, it is taken as granted that Wikipedia is a reliable corpus for training, but such influence scores cannot account
 79 for the peer production processes that results in Wikipedia being a high-quality information resource. This opens
 80 important questions for CSCW researchers: How should we redefine “value” when Wikipedia edits shape downstream
 81 AI behaviors? How can, and should, data valuation insights be integrated with community-driven governance processes?
 82 And how might these techniques support more equitable and participatory approaches to dataset curation?

83 **3 Towards Participatory Dataset Governance**

84 Wikipedia offers more than training data, it represents a living example of participatory governance at scale. Over
 85 two decades, editors have collaboratively developed policies, resolved disputes, and balanced openness with cultural
 86 sensitivities across regions and languages. These practices offer valuable insights into how datasets for foundation
 87 models could be curated and managed. We envision combining technical insights from data valuation with community-
 88 driven governance principles to establish participatory approaches to responsible dataset curation. Influence scores can
 89 help identify which contributions most strongly shape model behavior, while Wikipedia’s governance norms can inform
 90 how these contributions are reviewed, recognized, and included. By integrating these perspectives, CSCW researchers
 91 can help bridge the gap between technical valuation methods and social governance frameworks, ensuring that dataset
 92 composition, curation, and release reflect both responsible AI principles and community values.

105 References

- 106 [1]** Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D. Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish
107 Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler,
108 Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher
109 Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models Are Few-Shot Learners. In *Advances in Neural*
110 *Information Processing Systems*, Vol. 33. 1877–1901.
- 111 [2]** Sang Keun Choe, Hwijeon Ahn, Juhan Bae, Kewen Zhao, Minsoo Kang, Youngseog Chung, Adithya Pratapa, Willie Neiswanger, Emma Strubell,
112 Teruko Mitamura, Jeff Schneider, Eduard Hovy, Roger Grosse, and Eric Xing. 2024. What Is Your Data Worth to GPT? LLM-Scale Data Valuation
113 with Influence Functions. *arXiv preprint arXiv:2405.13954* (2024). <https://arxiv.org/abs/2405.13954>
- 114 [3]** Naman Goyal, Cynthia Gao, Vishrav Chaudhary, Peng-Jen Chen, Guillaume Wenzek, Da Ju, Sanjana Krishnan, Marc'Aurelio Ranzato, Francisco
115 Guzmán, and Angela Fan. 2022. The Flores-101 Evaluation Benchmark for Low-Resource and Multilingual Machine Translation. *Transactions of the*
116 *Association for Computational Linguistics* 10 (2022), 522–538. doi:10.1162/tacl_a_00474
- 117 [4]** Aaron Halfaker, R. Stuart Geiger, Jonathan Morgan, and John Riedl. 2013. The Rise and Decline of an Open Collaboration System How Wikipedia's
118 Reaction to Popularity Is Causing Its Decline. *American Behavioral Scientist* 57 (05 2013), 664–688. doi:10.1177/0002764212469365
- 119 [5]** Molly G. Hickman, Viral Pasad, Harsh Kamalesh Sanghavi, Jacob Thebault-Spieker, and Sang Won Lee. 2021. Understanding Wikipedia Practices
120 Through Hindi, Urdu, and English Takes on an Evolving Regional Conflict. *Proc. ACM Hum.-Comput. Interact.* 5, CSCW1, Article 34 (April 2021),
31 pages. doi:10.1145/3449108
- 121 [6]** Sara Javanmardi, David W. McDonald, and Cristina V. Lopes. 2011. Vandalism detection in Wikipedia: a high-performing, feature-rich model and its
122 reduction through Lasso. In *Proceedings of the 7th International Symposium on Wikis and Open Collaboration* (Mountain View, California) (*WikiSym*
123 '11). Association for Computing Machinery, New York, NY, USA, 82–90. doi:10.1145/2038558.2038573
- 124 [7]** Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim
125 Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. In *Advances in Neural*
126 *Information Processing Systems*, Vol. 33. Curran Associates, Inc., 9459–9474.
- 127 [8]** Cade Metz and Katie Robertson. 2024. OpenAI Seeks to Dismiss Parts of The New York Times's Lawsuit. *The New York Times* (February 2024).
128 Discusses legal challenges related to how AI systems are built and trained.
- 129 [9]** Philipp Singer, Florian Lemmerich, Robert West, Leila Zia, Ellery Wulczyn, Markus Strohmaier, and Jure Leskovec. 2017. Why we read Wikipedia.
130 In *Proceedings of the 26th International Conference on World Wide Web*. International World Wide Web Conferences Steering Committee, 1591–1600.
131 doi:10.1145/3038912.3052716
- 132 [10]** Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2018. Gender Bias in Coreference Resolution: Evaluation and
Debiasing Methods. In *Proceedings of NAACL*.

133 Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009