

Building Stereotype Repositories with LLMs and Community Engagement for Scale and Depth

Sunipa Dev
Google Research
sunipadev@google.com

Akshita Jha
Virginia Tech, Google Research
akshitajha@vt.edu

Jaya Goyal
Circadian Connect
jaya@circadianconnect.com

Dinesh Tewari
Google Research
dineshtewari@google.com

Shachi Dave
Google Research
shachi@google.com

Vinodkumar Prabhakaran
Google Research
vinodkpg@google.com

Abstract

Measurements of fairness in NLP have been critiqued for lacking concrete definitions of biases or harms measured, and for perpetuating a singular, Western narrative of fairness globally. To combat some of these pivotal issues, methods for curating datasets and benchmarks that target specific harms are rapidly emerging. However, these methods still face the significant challenge of achieving coverage over global cultures and perspectives at scale. To address this, in this paper, we highlight the utility and importance of complementary approaches that leverage both large generative models as well as community engagement, in these curation strategies. We specifically target the harm of stereotyping and demonstrate a pathway to build a benchmark that covers stereotypes about diverse, and intersectional identities. We discuss the two approaches, their advantages and constraints, the characteristics of the data they produce, and finally, their potential to be used complementarily for better evaluation of stereotyping harms.

CONTENT WARNING: This paper contains examples of stereotypes that may be offensive.

1 Introduction

Generative language models are widely used in diverse global settings across applications such as writing assistants (Ippolito et al., 2022), search tools,¹ and more (Jaech and Ostendorf, 2018; Yuan et al., 2022). Recent years have seen immense progress in the development of such large language models (Brown et al., 2020; Thoppilan et al., 2022; Zhang et al., 2022; Chowdhery et al., 2022), accompanied by detailed analysis of their abilities (Qin et al., 2023). Recent work has demonstrated the need for assessing their potential risks and harms to be contextually situated within the specific global

socio-cultural settings they are deployed in (Sambasivan et al., 2021; Prabhakaran et al., 2022). This need in turn highlights the gaps in current evaluation paradigms, within which a vast majority of resources are in English language, and/or is limited to a Western perspective of fairness and harms (Malik et al., 2022; Bhatt et al., 2022). This is especially troubling for evaluation benchmarks that require socially situated resources, for instance, to assess *stereotyping harms* that vary across cultures.

Addressing this growing need for evaluation strategies to be more globally relevant has its own challenges. First, the scale of operation becomes massive, given how diverse different languages and cultures are. Every region has its own unique axes of identities and with varying granularity of inspection, a large possible number of unique and intersectional identities and associated harms need to be examined. Second, stereotypes can be locally situated; some stereotypes are prevalent only within a region and can be about people residing in it or outside it. Hence, a lack of involvement of some communities can result in major gaps in evaluations, leading to disparately increased risks to those communities. This is interlinked with the third challenge, of ensuring that our resources and evaluations are not dominated by a Western perspective of what unfairness or stereotypes look like.

In this paper, we first discuss the challenges and limitations of current paradigms of stereotype resource collection, which are rooted in the enormity of global scale, and differential prevalence of stereotypes in different contexts. We then propose and demonstrate using exemplar methods, how complementary investigations of stereotypes which target scale and depth can achieve greater coverage and address aforementioned challenges - our first approach involves generation of candidate stereotypes using large language models (LLMs),

¹<https://openai.com/blog/chatgpt/>

followed by human annotations to verify which associations are stereotypical; the second approach involves reaching out to communities to directly collect the stereotypes known to them.

2 Complementary Approaches to Build Stereotype Resources

Stereotypes are generalizations about groups of people defined by their identity such as their gender, race, sexuality, age, etc. Stereotyping when propagated through language technologies can lead to many harmful outcomes including misrepresentation, targeted hateful speech generation, disparate access to resources, and opportunities (Blodgett et al., 2021; Dev et al., 2022; Shelby et al., 2022). There have been several efforts to build resources which document stereotypes in society (Koch et al., 2018; Borude, 1966), how they percolate into language technologies (Nadeem et al., 2021; Nangia et al., 2020; Bhatt et al., 2022), and cause unfair model behavior (Dev et al., 2022; Li et al., 2020).

While existing stereotype resources are rich and enable model evaluations, most of them were collected by employing methods that rely on human annotations about statements describing a potential stereotype. However, stereotypes are not absolute, in that they vary by societies, communities, and individual experiences of people. Any individual annotator will not be aware of all stereotypes present globally and can only confirm stereotypes they individually know of. As a result, annotations from sets of people or even stereotypical statements or text written down by people will still present a limited view of all stereotypes across the world. Also, the statements or text that is annotated for presence of stereotypes is typically human generated, which is an additional challenge towards both scale and coverage of global identities and stereotypes.

For broader coverage, LLMs can be imagined as a lens on the society, since they are trained over copious amounts of naturally occurring, human-generated text that reflect the underlying societal context including social stereotypes. Their generations attempt to mimic human knowledge and predispositions, and has been shown to reproduce stereotypes (Zhao et al., 2018; Dev et al., 2022; Li et al., 2020). Consequently, they can, inexpensively create generalizations that are diverse and representative of a wide range of identities across the globe (Lauscher et al., 2020; Malik et al., 2022). So we can tap into the generalizing capabilities of

LLMs to create a broad-coverage candidate set for stereotypes. However, LLM generations are not always grounded factually, and reflect spurious correlations, and noise (Bang et al., 2023). Hence, for usage as a stereotype resource, associations generated by LLMs about groups of people need to be validated for social presence of such stereotypes by human raters familiar with the corresponding socio-cultural contexts.

On the other hand, LLMs may not capture all social stereotypes globally. While they are trained on large amounts of data, there are still gaps in global representativeness in such data (Chowdhery et al., 2022), which will also carry over to stereotype resources built using LLMs. Furthermore, since most state-of-the-art LLMs are trained on online data that has a Western lens (Dodge et al., 2021), the stereotypes we get through LLMs may also reflect this Western gaze, and miss the nuances of stereotypes in local cultural contexts (Malik et al., 2022; Bhatt et al., 2022). Hence, it is important to complement the LLM-based approach with community engagements to build richer resources. Methods that rely on community engagement are expensive and time consuming but help collect socially situated perspectives. When used in targeted ways to understand one specific culture or society, annotations, surveys, and free form data collection can provide depth and nuance to the collected stereotype resource.

Figure 1 imagines this juxtaposition of challenges and complementarity of community engaged and LLM generation based approaches. If our goal is to uncover the universal set of all stereotypes in the world, different strategies are warranted. Ideally, the results of community engagements, when deployed globally would overlap a 100% with this set. However, that would be expensive both cost and time wise to completely attain. Meanwhile, if we consider a second set consisting of all associations LLMs generate with the identities of people, only a certain fraction of it would be socially present stereotypes. However, LLM generations, combined with human annotations would give us a list of stereotypes which is represented in Figure 1 as the intersection of these two sets.

3 Case Study

In this section, we summarize insights from two separate studies that take these complimentary approaches towards stereotype resource building,

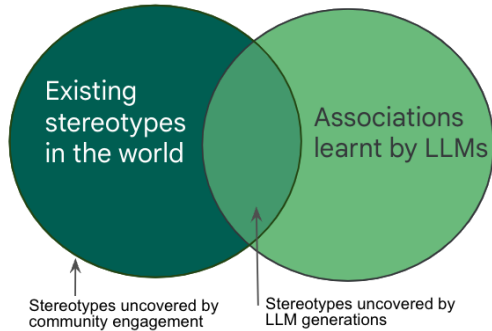


Figure 1: *Projected coverage of stereotypes uncovered by the approaches.* While community engagement can potentially uncover the set of all stereotypes (in darkest green in image), it is expensive. LLM generations (in lightest green) on the other hand may contain noise and spurious correlations. The intersection of the two sets represents social stereotypes uncovered using LLMs. (Proportions of sets in image not to scale.)

and outline their strengths and limitations. One approach crowd-sources stereotypes by engaging with communities, and the other uses generative models in conjunction with human annotations to scale coverage.² These complementary approaches can be extended globally, to different harms such as hateful speech, toxic language and so on, which are also geo-culturally and socially situated.

3.1 LLM-based Stereotype Repository

Generative language models are powerful in learning from naturally occurring text and responding to prompts with text that is contextually meaningful. We prompt state-of-the-art language models PaLM (Chowdhery et al., 2022) and GPT-3 (Brown et al., 2020) with stereotypes from existing datasets of stereotypes from NLP and social psychology literature (Nangia et al., 2020; Nadeem et al., 2021; Bhatt et al., 2022; Borude, 1966; Rogers and Wood, 2010; Koch et al., 2018). The stereotypes selected for prompting were about global nationalities, and states in the United States and India. The prompts result in the models producing other such generalizations about geographical identities of persons, which are filtered and processed to obtain a candidate set. We then validated whether the associations in this candidate set are commonly known social stereotypes, for which we recruited annotators with diverse backgrounds (across gender) and geographic location that matches the associations.

²These studies will be published separately; in this paper, we discuss the methods only briefly, and focus on the insights that highlight the need for such complementary approaches.

Examples	Saliency	Human Validation
(Italian, gangsters)	16.1	3
(Nigerian, scammers)	13.8	2
(Irish, violent)	7.43	3
(Greeks, proud)	6.31	3
(Japanese, greedy)	5.13	2
(Iranian, cruel)	4.48	2

Table 1: Example regional stereotypes obtained using LLM probing, their saliency scores, and the number of human raters validating their presence in society.

Constraints: Model generations only estimate stereotype candidates and must be validated by human annotations. Since annotations are subject to annotator experiences with respect to culture, world locations, etc., annotators need to be aware of the presented identity and stereotype. Selection and availability of annotators, thus, restricts the axes and granularity of identities whose associated stereotypes can be validated and uncovered. For this reason, in this study, the data is filtered by country and state demonyms and is reduced in its coverage of the resultant dataset to other regional groups, ethnicities, and their associations.

Dataset Produced: The resulting dataset contains about 8000 tuples, each with at least 3 human ratings whether the terms in the tuple represent a stereotype. Each tuple consists of an *identity term* and an *attribute*. An *identity term* refers to a word or phrase that denotes a social group a person belongs to. An *attribute* refers to word(s)/phrase that describes a person or a group of people, such as adjectives or verbal predicates. Table 1 shows some example stereotype tuples about regional identities obtained by this approach, along with their saliency scores in the LLM generations, and the number of annotators from the corresponding regions who validated them to be known stereotypes. We calculate the salience score of a stereotype tuple using a modified tf-idf metric. See (Jha et al., 2023) for more details about the dataset and the process followed.

3.2 Community Engagement based Stereotype Repository

Identities of persons can be intersectional, fine-grained, and also be more fluid than absolute categories. Additionally, each of these identities, associated generalizations and sentiments about them, and the potential harms they face from unfair technology is socially situated and differs by regions of the globe. Capturing these nuances require ap-

proaches that understand identities and stereotypes deeply for a given socio-cultural context, that may not be captured by the LLMs. We focus on India which yielded a large number of stereotypes in the LLM based approach. India is a country with 22 official languages, over 461 languages in use with many more dialects, 6 major religions, and many more such nuances which define individuals, their communities, and faced stereotypes. We employ an exploratory study design using surveys, distributed across 8 urban and suburban regions in India, which introduce the concept of stereotypes with examples of locally present stereotypes, followed by open ended questions about what stereotypes the participant is aware of in their society. The stereotypes can be about any identity, or any combination of identities. For example, it can be about ethnic origin and caste such as ‘Rajput’, but also intersect with gender such as ‘Rajput women’.

Constraints: Since this method engages with diverse communities local to regions, it is expensive and time consuming. Additionally, scaling it needs local knowledge and points of contact to identify and distribute the surveys to the underrepresented communities and prevent imposition of an external viewpoint of fairness and social structures.

Dataset: The dataset created consists of about 2000 unique social stereotypes. In addition, it contains meta-data about how many persons with various identities (e.g., by gender, caste, and regional belonging) contributed the tuple as a stereotype.

3.3 Complementary coverage and insights

The two approaches together yielded approximately 11,000 associations, with varying degrees of prevalence as social stereotypes. In this section we compare and contrast various aspects of tuples produced by both approaches.

Coverage of Identities: The LLM-based approach render the ability to scale up dataset creation many fold. In particular, the approach when restricted to generate for only region associated stereotypes, resulted in generation of candidate stereotype tuples for over 170 countries. This is 5 times the coverage of existing datasets such as StereoSet (Nadeem et al., 2021) and CrowS-Pairs (Nangia et al., 2020). In addition, it also contain stereotypes about states within India. Each identity term in this case is a demonym, restricted to countries and states. So, while the scale has been

improved, the depth and granularity of identities understood is restricted. By engaging with communities in India, a larger number of identities, around 1000, are covered. These span demonyms, races, ethnicities, castes, religion, gender, sexuality, age, and more, including intersectional identities.

Coverage of Attributes: The LLM-based approach produced stereotype tuples, with over 10,000 different attributes. On the other hand, stereotypes collected by surveying communities contained about 2,000 distinct attributes. For both datasets, there is a substantial number of attribute terms that are synonymous or alternate phrases for each other. While the absolute number of attributes produced does not directly imply richer stereotype data, diversity in attribute terms covered reflects indirectly on the diversity in the types of stereotypes about an identity that were uncovered.

Coverage of Stereotypes: Both approaches uncovered unique stereotypes with minimal overlap (≤ 10 stereotypes). The LLM-based approach largely covered broad categories of demonyms, and yielded broad-strokes stereotypes such as ‘Indian, vegetarian’, while engagement with communities broke this stereotype down into smaller, more nuanced associations, such as ‘Jain, vegetarian’, where the identity is a religion category, ‘Brahmin, vegetarians’, where the identity term is an intersectional religion and caste category, and ‘Punjabi, non-vegetarians’, where the identity term is a state demonym. Furthermore, the generative approach hinges on the abilities of LLMs which in turn rely on their training data that is mostly in English and West-centric. Thus, stereotypes uncovered can sometimes have a Western perspective such as ‘Indian, smelly’, which was not present in the data produced through community engagement.

Dataset Sample: Table 2 presents some examples of stereotypes collected by the two approaches that demonstrates their differences. Stereotypes collected by engaging with communities tend to be more granular about identity terms, and use terms such as ‘Baniya’,³ which in vernacular tongues mean ‘merchant’, but is also a caste category prevalent in some parts of India. On the other hand, the LLM-based approach provide more global coverage of identities for each stereotypes. For instance, it found stereotypes around Chinese and Taiwanese people being good at math, and Pakistani

³[https://en.wikipedia.org/wiki/Bania_\(caste\)](https://en.wikipedia.org/wiki/Bania_(caste))

LLM-based	Community-based
Indian, brown	Indian, brown South Indian, dark skinned Bihari, dark skinned
Gujarati, trader	Gujarati, businessman Gujarati, baniya
Chinese, very good at math Taiwanese, good at math Pakistani, bad at math American, bad at math	Asian, good at math

Table 2: Example stereotypes collected by LLM-based and community engagement based approaches. We see that for Indian state based identities, the community based approach results in much more granular stereotypical associations. However, since the community engaged effort was made in India, its coverage was limited compared to LLM based approach.

and American people being bad at math, while the community engaged approach provided only a single stereotype about Asians for this attribute.

4 Discussion

In the paper, we presented two approaches to expand the coverage of stereotype resources used to evaluate language technologies. While we demonstrated the advantages of each individual method, it is also important to note how the complementary usage of the methods can lead to broad, and granular coverage of stereotype harms globally. Each method uncovered different kinds of stereotypes that were not found using the other.

Additionally, the output of one method can serve as the seed for the other; the stereotypes recovered from engaging with communities can be used as prompts in subsequent usage of the generative approach using LLMs. Meanwhile, the generative approach highlights prevalence of associations and can help understand which communities to engage with for uncovering finer-grained stereotypes.

Further, the collection of non-overlapping, complementary sets of stereotypes enhances coverage both in terms of global communities covered as well as fine-grained identities present in different regions. Measurements of harm in language tasks like question answering (Li et al., 2020) and natural language inference (Dev et al., 2020) which are built on preferential associations with identities can leverage this more comprehensive list to make more holistic estimations.

Limitations

Stereotypes are subjective and socially situated. The absence of a stereotype in the lists collected by either approach does not imply that the stereotype does not exist in society or cannot be harmful to people. Any measurements built with these lists can still only make limited estimations, and more precautions should always be taken when deploying a model or tool with the specific use case at hand. Further, even with both approaches, we may not cover all possible regional identities and finer-grained examinations of stereotypes are possible. We also only work with English language text, and stereotypes written in English, and multilingual efforts are required to reflect some stereotypes present only within specific cultures.

References

- Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenhong Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, Quyet V. Do, Yan Xu, and Pascale Fung. 2023. [A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity](#). 2
- Shaily Bhatt, Sunipa Dev, Partha Talukdar, Shachi Dave, and Vinodkumar Prabhakaran. 2022. [Re-contextualizing fairness in nlp: The case of india](#). In *Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing*, pages 727–740. 1, 2, 3
- Su Lin Blodgett, Gilsinia Lopez, Alexandra Olteanu, Robert Sim, and Hanna Wallach. 2021. [Stereotyping Norwegian salmon: An inventory of pitfalls in fairness benchmark datasets](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1004–1015, Online. Association for Computational Linguistics. 2
- Ramdas Borude. 1966. [Linguistic stereotypes and social distance](#). *Indian Journal of Social Work*, 27(1):75–82. 2, 3
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. [Language models are few-shot learners](#). *Advances in neural information processing systems*, 33:1877–1901. 1, 3
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton,

- Sebastian Gehrmann, et al. 2022. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*. 1, 2, 3
- Sunipa Dev, Tao Li, Jeff M Phillips, and Vivek Sriku-
mar. 2020. On measuring and mitigating biased in-
ferences of word embeddings. In *Proceedings of
the AAAI Conference on Artificial Intelligence*, vol-
ume 34, pages 7659–7666. 5
- Sunipa Dev, Emily Sheng, Jieyu Zhao, Aubrie Amstutz,
Jiao Sun, Yu Hou, Mattie Sanseverino, Jiin Kim, Ak-
ihiro Nishi, Nanyun Peng, and Kai-Wei Chang. 2022.
[On measures of biases and harms in NLP](#). In *Find-
ings of the Association for Computational Linguistics:
AAACL-IJCNLP 2022*, pages 246–267, Online
only. Association for Computational Linguistics. 2
- Jesse Dodge, Maarten Sap, Ana Marasović, William
Agnew, Gabriel Ilharco, Dirk Groeneveld, Margaret
Mitchell, and Matt Gardner. 2021. [Documenting
large webtext corpora: A case study on the colos-
sal clean crawled corpus](#). In *Proceedings of the
2021 Conference on Empirical Methods in Natural
Language Processing*, pages 1286–1305, Online and
Punta Cana, Dominican Republic. Association for
Computational Linguistics. 2
- Daphne Ippolito, Ann Yuan, Andy Coenen, and
Sehmon Burnam. 2022. Creative writing with an
ai-powered writing assistant: Perspectives from pro-
fessional writers. *arXiv preprint arXiv:2211.05030*.
1
- Aaron Jaech and Mari Ostendorf. 2018. [Personalized
language model for query auto-completion](#). In *Pro-
ceedings of the 56th Annual Meeting of the Associa-
tion for Computational Linguistics (Volume 2: Short
Papers)*, pages 700–705, Melbourne, Australia. As-
sociation for Computational Linguistics. 1
- Akshita Jha, Aida Davani, Chandan Reddy, Shachi
Dave, Vinodkumar Prabhakaran, and Sunipa Dev.
2023. [Seegull: A broad-coverage stereotype bench-
mark leveraging generative models \(under submis-
sion\)](#). 3
- Alex Koch, Nicolas Kervyn, Matthieu Kervyn, and
Roland Imhoff. 2018. [Studying the cognitive map
of the u.s. states: Ideology and prosperity stereo-
types predict interstate prejudice](#). *Social Psycholog-
ical and Personality Science*, 9(5):530–538. 2, 3
- Anne Lauscher, Rafik Takiyeddin, Simone Paolo
Ponzetto, and Goran Glavaš. 2020. [Araweat: Mul-
tidimensional analysis of biases in arabic word em-
beddings](#). 2
- Tao Li, Daniel Khashabi, Tushar Khot, Ashish Sab-
harwal, and Vivek Sriku-
mar. 2020. [UNQOVERing
stereotyping biases via underspecified questions](#). In
*Findings of the Association for Computational Lin-
guistics: EMNLP 2020*, pages 3475–3489, Online.
Association for Computational Linguistics. 2, 5
- Vijit Malik, Sunipa Dev, Akihiro Nishi, Nanyun Peng,
and Kai-Wei Chang. 2022. [Socially aware bias mea-
surements for hindi language representations](#). In
*Conference of the North American Chapter of the
Association for Computational Linguistics (NAACL),
short*. 1, 2
- Moin Nadeem, Anna Bethke, and Siva Reddy. 2021.
[Stereoset: Measuring stereotypical bias in pre-
trained language models](#). In *Proceedings of the
59th Annual Meeting of the Association for Compu-
tational Linguistics and the 11th International Joint
Conference on Natural Language Processing (Vol-
ume 1: Long Papers)*, pages 5356–5371. 2, 3, 4
- Nikita Nangia, Clara Vania, Rasika Bhalerao, and
Samuel Bowman. 2020. [Crows-pairs: A challenge
dataset for measuring social biases in masked lan-
guage models](#). In *Proceedings of the 2020 Con-
ference on Empirical Methods in Natural Language
Processing (EMNLP)*, pages 1953–1967. 2, 3, 4
- Vinodkumar Prabhakaran, Rida Qadri, and Ben
Hutchinson. 2022. [Cultural incongruencies in arti-
ficial intelligence](#). 1
- Chengwei Qin, Aston Zhang, Zhuosheng Zhang, Jiao
Chen, Michihiro Yasunaga, and Diyi Yang. 2023. [Is
chatgpt a general-purpose natural language process-
ing task solver?](#) 1
- Katherine H. Rogers and Dustin Wood. 2010. [Accu-
racy of united states regional personality stereotypes](#).
Journal of Research in Personality, 44(6):704–713.
3
- Nithya Sambasivan, Erin Arnesen, Ben Hutchinson,
Tulsee Doshi, and Vinodkumar Prabhakaran. 2021.
[Re-imagining algorithmic fairness in india and be-
yond](#). In *Proceedings of the 2021 ACM Confer-
ence on Fairness, Accountability, and Transparency*,
FAccT ’21, page 315–328, New York, NY, USA. As-
sociation for Computing Machinery. 1
- Renee Shelby, Shalaleh Rismani, Kathryn Henne,
AJung Moon, Negar Rostamzadeh, Paul Nicholas,
N’Mah Yilla, Jess Gallegos, Andrew Smart, Emilio
Garcia, and Gurleen Virk. 2022. [Identifying so-
ciotechnical harms of algorithmic systems: Scoping
a taxonomy for harm reduction](#). 2
- Romal Thoppilan, Daniel De Freitas, Jamie Hall,
Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze
Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du,
YaGuang Li, Hongrae Lee, Huaixiu Steven Zheng,
Amin Ghafouri, Marcelo Menegali, Yanping Huang,
Maxim Krikun, Dmitry Lepikhin, James Qin, De-
hao Chen, Yuanzhong Xu, Zhifeng Chen, Adam
Roberts, Maarten Bosma, Vincent Zhao, Yanqi
Zhou, Chung-Ching Chang, Igor Krivokon, Will
Rusch, Marc Pickett, Pranesh Srinivasan, Laichee
Man, Kathleen Meier-Hellstern, Meredith Ringel
Morris, Tulsee Doshi, Renelito Delos Santos, Toju
Duke, Johnny Soraker, Ben Zevenbergen, Vinod-
kumar Prabhakaran, Mark Diaz, Ben Hutchinson,

Kristen Olson, Alejandra Molina, Erin Hoffman-John, Josh Lee, Lora Aroyo, Ravi Rajakumar, Alena Butryna, Matthew Lamm, Viktoriya Kuzmina, Joe Fenton, Aaron Cohen, Rachel Bernstein, Ray Kurzweil, Blaise Aguera-Arcas, Claire Cui, Marian Croak, Ed Chi, and Quoc Le. 2022. [Lamda: Language models for dialog applications](#). 1

Ann Yuan, Andy Coenen, Emily Reif, and Daphne Ippolito. 2022. [Wordcraft: Story writing with large language models](#). In *27th International Conference on Intelligent User Interfaces, IUI '22*, page 841–852, New York, NY, USA. Association for Computing Machinery. 1

Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. 2022. [Opt: Open pre-trained transformer language models](#). 1

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 the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 15–20, New Orleans, Louisiana. Association for Computational Linguistics. 2