## Standard Language Ideology in Al-Generated Language

Genevieve Smith\* genevieve.smith@berkeley.edu UC Berkeley Berkeley, CA, USA

Madeline Bossi UC Berkeley Berkeley, CA, USA

Ishita Rustagi UC Berkeley Berkeley, CA, USA Xavier Yin UC Berkeley Berkeley, CA, USA

### **KEYWORDS**

language models, generative ai, language ideology, dialects, natural language processing

#### **ACM Reference Format:**

Genevieve Smith\*, Eve Fleisig\*, Madeline Bossi, Ishita Rustagi, and Xavier Yin. 2018. Standard Language Ideology in AI-Generated Language. *ACM/IMS J. Data Sci.* 37, 4, Article 111 (August 2018), 7 pages.

## **ABSTRACT**

In this position paper, we explore standard language ideology in language generated by large language models (LLMs). First, we outline how standard language ideology is reflected and reinforced in LLMs. We then present a taxonomy of open problems regarding standard language ideology in AI-generated language with implications for minoritized language communities. We introduce the concept of standard AI-generated language ideology, the process by which AI-generated language regards Standard American English (SAE) as a linguistic default and reinforces a linguistic bias that SAE is the most "appropriate" language. Finally, we discuss tensions that remain, including reflecting on what desirable system behavior looks like, as well as advantages and drawbacks of generative AI tools imitating-or often not-different English language varieties. Throughout, we discuss standard language ideology as a manifestation of existing global power structures in and through AI-generated language before ending with questions to move towards alternative, more emancipatory digital futures.

## 1 INTRODUCTION

Since its public release in November 2022, ChatGPT has drawn over 100 million weekly active users globally. While the tool brings immense benefits, those benefits are not distributed equally across speakers of different language varieties. Rather, ChatGPT—and other generative AI language technologies—reflect standard language ideology, which reinforces a hierarchy between language varieties. In this paper, we illustrate how and in what ways standard language

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 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 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 servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

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

ACM 2831-3194/2018/8-ART111 https://doi.org/XXXXXXXXXXXXXXX ideology is reinforced in AI-generated language and present a taxonomy of open problems regarding standard language ideology in AI-generated language.

## 2 BACKGROUND

### 2.1 Standard language ideology

Standard language ideology is a construct that reinforces a hierarchy between language varieties. As defined by [42], it is "a bias toward an abstracted, idealized, homogenous spoken language which is imposed from above" whose goal is the "suppression of variation of all kinds." In other words, it is a common and false belief that some language varieties—usually those used by communities with more social prestige—are "better" or "more complex" than others. Linguistically, however, all language varieties are equally valid. The idea that certain languages are "better" than others ignores the fact that all language varieties are equally capable of expression [31]. There is no "correct" or "incorrect" way of using the English language or any language; in fact, "standard" language is not spoken by any real community, but is an abstracted variety that can only be defined in contrast to the speech of marginalized communities.

Despite this, certain language varieties have been institutionally privileged as more "standard" and viewed as more "appropriate" or "professional" than others. This privileged status is linked to the association with people in power that the language variety holds. In particular, economic globalization has cast English as a lingua franca (ELF), or the common language adopted by speakers of different languages [48]. The English language has been granted this dominant position in international business as well as other domains, giving it and its speakers a privileged position in international business communications [14, 54, 61]. As [65] notes, this "leads to portrayal of English as a neutral solution for overcoming linguistic diversity in these relationships, as exemplified by the concept of Business English as a lingua franca." The privileged position of English is a result of the social and historical power that English speakers have had throughout history, especially as it relates to colonization across many parts of the world. The spread of English as the default has been framed as "linguistic imperialism" that threatens other languages and language varieties [35, 53].

Importantly, even within the English language, there is a range of language varieties, with certain varieties being granted more privileged "standard" positions. In the context of the United States

 $<sup>^{1}</sup>https://www.demandsage.com/chatgpt-statistics/\\$ 

(US), "Standard" American English (SAE) is considered the dominant language variety and reflects an abstracted collection of linguistic norms of middle-class, white men who have held disproportionate levels of power within the country [3, 31, 56]. Other language varieties have been devalued through their institutional subordination [3, 31, 42]. African American English (AAE) is one such language variety, among many others (e.g. Irish English, [49]; Indian English, [39]; Chicano English, [29]). For example, using AAE has been linked to being denied housing due to "sounding Black" [45, 66]. Even when using "standard" varieties such as SAE, marginalized people can be subjected to linguistic bias [1]. Ultimately, language, identity and power are linked and affect people's lives in various ways [13].

By misleading people to believe that some languages and language varieties are better than others, standard language ideology can perpetuate harmful patterns of linguistic discrimination and the oppression of speakers of "non-standard" varieties. Linguistic discrimination can often serve as a proxy for other forms of discrimination along lines of race, gender, nationality, class, and more [19]. This discrimination resulting from the promotion of language hierarchy can be subtle under perceptions of benevolence, such as encouraging others to speak more "appropriately," or more obvious bigotry, such as associating certain ways of speaking with lacking intelligence [5]. Relatedly, standard language expectations dictate access to social capital through means such as education, employment, or public office [64]. Given this, those who speak more closely to the standard varieties benefit from better access to such resources and opportunities.

# 2.2 Standard language ideology in large language models

We outline how language models reinforce standard language ideology and explore why LLMs perpetuate hierarchies between language varieties. In particular, we discuss how training data for LLMs overrepresents English, particularly SAE, alongside certain voices and perspectives. We then explore how SAE is treated as the default at a higher level, which is further reinforced by the demographics of those leading generative AI research and tech companies.

Research on language models preceding the release of ChatGPT has highlighted ways in which these models perform worse for certain speakers. In particular, language models perform worse for AAE on tasks including text generation, sentiment analysis, and parsing [7, 24, 28]. [60] discuss how content moderation tools fail to capture the semantic richness of AAE, such as by making blanket assumptions about the complex semantics of reclaimed slurs. Beyond worse performance for minoritized language communities, language models can also advance stereotypes regarding speakers of a particular language variety [32, 50].

Performance discrepancies in LLMs are linked to the language data that underpins these technologies. English is the default language for data powering large LLMs, linked to LLMs relying on language data from the Internet. Internet data overwhelmingly represents English, and Internet use varies due to social factors. An estimated 60% of all language content on the Internet is in English,

despite only 17% of people speaking English globally [50]. Meanwhile, 88% of languages have "exceptionally limited resources" in digital spaces [37].

Smith et al.

Even if languages are well-represented in digital corpora, certain perspectives are over- or under- represented. For example, on Reddit, users are 67% male and 70% White, resulting in potential reinforcement of White, male perspectives ([6], citing [4]). Other perspectives may be actively targeted, harassed, cyberbullied, and otherwise marginalized online. Research exposes growing concerns of online gender-based harassment with clear links between online spaces and misogyny, which can be further amplified for Black and brown women [20, 22, 26]. This harassment and resulting safety and mental health implications can result in decreased use of online spaces. Meanwhile, although Black people have historically been overrepresented on Twitter compared to other demographics in the general US population [8], Black Tweets are still often considered "inappropriate" and are more often inaccurately flagged as hateful by automatic hate speech detection tools [15]. Beyond being inaccurate, this disproportionately censors Black speakers.

Research on language models treats English-particularly SAEas the status quo, with work on other languages "often considered 'language specific' and thus reviewed as less important than equivalent work on English" [6]. This default of SAE partly reflects the people who hold disproportionate power in these spaces. White men are overrepresented in AI research both industry and academia [63]. Tech companies leading the generative AI charge, such as Open AI, Google, Microsoft, and Meta, are largely headquartered in the United States with employees and corporate leadership that skew White and male [52]. The demographics of NLP research are also skewed, with researchers disproportionately affiliated with North American and European institutions; researchers from Latin America, Africa, the Middle East, and Southeast Asia are particularly underrepresented [57]. While the demographics of those developing and leading NLP research and companies do not directly result in the reinforcement of SAE at the level of language models, this lack of diversity in the technologies' developers and managers can lead to similar thinking in ways that miss potential biases the technology can have [11] with a default towards an (inequitable) status quo.

## 2.3 Standard language ideology extends to and manifests in AI-generated language

Standard language ideology extends to AI-generated language and technologies, in which hierarchies of language are reinforced. While still a relatively nascent space, there is some research that explores how generative AI technologies—including ChatGPT—perform for and exhibit biases of different language varieties. [16] discuss how language models including GPT-3, ChatGPT, and GPT-4 have higher perplexity for AAE, interpreted as greater difficulty in understanding AAE, and have trouble producing natural-sounding and semantically accurate AAE text. [67] find inconsistent performance for ChatGPT's production of Singlish² and code-switched text involving Southeast Asian languages, finding that even when generated text was perceived as natural (as though produced by a native speaker), it sometimes contained "semantic inaccuracies... discernible

 $<sup>^2 {\</sup>rm Singlish}$  is an English-based creole spoken in Singapore.

by native speakers." [67] illustrate that language models produce text with harmful stereotypes about speakers of AAE through "dialect prejudice." Another study examining performance of several widely-used GPT detectors found that non-native English writing samples were consistently misclassified as AI-generated compared to native writing samples [41], though the strength of this effect is uncertain.<sup>3</sup>

Examining ChatGPT performance across ten English varieties globally, [25] found that ChatGPT responses to inputs in the different English language varieties significantly reduced the occurrence of linguistic features of all varieties except SAE and British English. SAE had the least reduction in linguistic features (22.1%) followed by British English (27.8%). Meanwhile, the other eight varieties had reductions in linguistic features of over 84%, with four languages experiencing over 96% reduction (Jamaican English, Singaporean English, Scottish English and AAE). The same study also found that, compared to Standard American and British English, ChatGPT responses to minoritized varieties of English expressed more stereotyping, demeaning, and condescending content, as well as lack of comprehension. These results illustrate that "standard" English varieties, particularly SAE, are the default output in ChatGPT, and that ChatGPT perpetuates linguistic discrimination against speakers of other varieties. Taken together, this amplifies the hierarchical position of SAE in digital spaces and technologies.

## 3 TAXONOMY OF OPEN PROBLEMS REGARDING STANDARD LANGUAGE IDEOLOGY IN AI-GENERATED LANGUAGE

We introduce a taxonomy of problems that standard language ideology in AI-generated language technologies present with implications for different language communities globally (see Table 1). These problems of AI-generated language technologies include: the default production of "standard" language varieties reinforces "correct" ways of communicating; outputs have lower quality of service for minoritized varieties; producing minoritized varieties can result in stereotyping of languages; producing minoritized varieties can result in appropriation and/or manipulation; and preventing outputs in minoritized varieties can result in limited quality of service and erasure.

The default production of "standard" language varieties reinforces "correct" ways of communicating. At a high level, generative AI technologies default to "standard" language varieties in their responses, particularly SAE, as illustrated by [25]. This subtly reinforces the idea that the "correct" or "more appropriate" way of communicating uses "standard" languages, and particularly SAE. This can have a cascading effect that impacts people's own linguistic biases and perceptions of language hierarchies. Lessons can be learned from the dominant position that SAE has held in educational settings in the United States. Despite the prioritization of

SAE in classrooms often being seen by teachers as a means of unifying students through a common language [18, 58], it has negative implications for minoritized students. For example, Black students who speak AAE feel compelled to adopt SAE and dominant ideologies, which then impact their identity expression, academic achievement, and self-perception [43].

Additionally, due to automation bias—a psychological tendency to over-rely on automation that can result in complacency regarding automation outputs [27]—humans can have greater trust and overconfidence in AI [34]. It is not clear how defaulting to SAE in generative AI language outputs impacts the mental perceptions of minoritized language speakers. However, when combined with automation bias, the default use of SAE in LLM outputs has the potential to exacerbate the idea that "standard" languages are more "correct" ways of speaking.

Outputs have lower quality of service for minoritized varieties. Generative AI technologies can produce lower-quality responses to speakers of minoritized language varieties. At a high level, these tools may not comprehend the user's prompt as well, and result in an incorrect or unhelpful response (see Appendix A for an example). Relatedly, generative AI tools can more often inaccurately determine that inputs of certain minoritized language varieties are hateful or offensive speech, resulting in refusal to respond [15].

These behaviors result in speakers of minoritized varieties experiencing more difficulty in using language models than speakers of "standard" varieties, which is a form of allocational harm [12]. Furthermore, if users of generative AI tools are aware that using a minoritized variety results in poorer performance, they may be incentivized to use a "standard" variety instead, which can reinforce the stigmatization of non-standard varieties and result in "digital code switching." Code switching, which results from a pressure minoritized speakers can feel to conform to more standard (more "appropriate") language, causes psychological tolls [47]. While additional research is needed, these psychological tolls could extend to the digital sphere.

Producing minoritized varieties can result in stereotyping of language use. Generative AI technologies may be prompted by the user to produce minoritized language varieties, or may be designed to respond to the particular language variety of the input. This may not always be a source of harm, particularly if users of that particular variety are asking for responses back in their language variety. However, responses in minoritized varieties can disproportionately carry or convey stereotypes or demeaning content.

In response to inputs of minoritized languages, models prompted to imitate the input variety may produce a stereotyped version of that language variety that does not accurately reflect the range of linguistic features used by that community. This behavior perpetuates the association between speaking a non-standard variety and stereotypical traits of that speaker community, which often serves as a covert form of racism, xenophobia, or other widespread harms ([55]; see also Section 2.1). The outputs could also convey demeaning content related to speakers of that language community. Examining ChatGPT outputs that imitated the language variety of the inputs, [25] found that native speakers of minoritized varieties

<sup>&</sup>lt;sup>3</sup>[62] suggested that the non-native writing samples being shorter had a potential confounding effect and found comparatively strong performance across models on other non-native speaker writing samples.

<sup>&</sup>lt;sup>4</sup>Language varieties included: "Standard" American English (SAE), African American English (AAE), "Standard" British English (SBE), Indian English, Irish English, Jamaican English, Kenyan English, Nigerian English, Scottish English, and Singaporean English.

### Open problem Societal implications • Subtly reaffirms language hierarchies and can reinforce or amplify people's internal linguis-The default production of "standard" tic biases language varieties reinforces "correct" ways of communicating. • Speakers of minoritized varieties experience more difficulty in using language models than Outputs have lower quality of service speakers of "standard" varieties for minoritized varieties. Speakers of minoritized varieties practice "digital code switching" to achieve better quality of service, with potential psychological tolls • Inaccurate reflection of the linguistic features of the particular variety Producing minoritized varieties can re- Reinforcement of limiting or harmful stereotypes about minoritized language communities sult in stereotyping of languages. • Appropriating the speech of minoritized groups decontextualize community norms and Producing minoritized varieties can reerases the history of the group's struggles sult in appropriation and/or manipula- Non-native speakers of minoritized varieties prompting and using minoritized language tion. varieties can result in "linguistic minstrelsy" or "figurative blackface" • Facilitating production of minoritized varieties lowers the barrier for bad actors' infiltration or manipulation of minoritized communities • Lower quality of service from being unable to fully interact with the model in one's native Preventing outputs in minoritized varivarietv eties can result in limited quality of ser-• Erasure of minoritized language varieties and reinforcement of language hierarchies vice and erasure.

Table 1: Taxonomy of harms regarding standard language ideology in AI-generated language.

rated outputs as more often conveying demeaning content than responses to standard English varieties.

Producing minoritized varieties can result in appropriation and/or manipulation. Even if a language model correctly produces text that reflects the grammar of that minoritized variety, there may be concerns over appropriation (particularly if nonnative speakers of that variety are asking for and using this text). Because speaking a non-standard variety can carry covert prestige (i.e., social value associated with use of canonically minoritized language varieties; [40]), use of non-standard varieties by non-native speakers, particularly in the case of White speakers appropriating AAE, has been discussed as "linguistic minstrelsy" or "figurative blackface" [9, 21]. Parallel concerns have been raised regarding the use of speech, music, and image generation tools to imitate people of color, sometimes called "digital blackface" [23, 33, 46, 59]. As LLMs improve at mimicking minoritized varieties, these models may increasingly contribute to similar harms.

In addition, use of uncredited language from marginalized groups by language models can further cultural appropriation and erase the linguistic history of speaker communities. A popular "Gen Z Translator" extension of ChatGPT<sup>5</sup> that claims to "transform any short text to be Gen Z slang filled" frequently produces AAE. However, rebranding expressions that emerged in Black communities as generational vernacular, as opposed to vernacular linked to ethnic and racial communities, erases the history and contexts of these

expressions. It risks "appropriating Black culture and perpetuating racism as [speakers] take on Black speech without assuming Black Americans' struggle" [10].

Producing minoritized varieties could also result in manipulation of those language communities by people outside of them. Readily available ways of producing text that imitates a non-standard variety could help agents to feign in-group membership for malicious purposes, such as the case of Russian misinformation bots posing as Black people online [36].

Preventing outputs in minoritized varieties can result in limited quality of service and erasure. Preventing language models from producing anything but prestige or "standard" varieties of a language avoids the potential harms discussed regarding appropriation and manipulation. However, if speakers of a marginalized variety want the model to reply in that variety, then it may constitute a quality of service harm because, unlike speakers of the "standard" variety, they are unable to fully interact with the model in their native variety. This can also contribute to erasure of minoritized languages. Thus, preventing production of AI-generated text in minoritized languages could further reinforce language hierarchies and the false notion of "correct" ways of speaking.

### 4 DISCUSSION

## 4.1 Standard AI-generated language ideology

Our taxonomy outlines issues of standard language ideology in Algenerated language, illustrating how popular language models and

 $<sup>^5</sup> https://chat.openai.com/g/g-AbhjZGbhY-gen-z-translator \\$ 

associated technologies grant more power to "standard" language varieties while opening opportunities for harm to speakers of minoritized language varieties. At a high level, default production of certain varieties reinforces "correct" ways of communicating that can have ensuing psychological implications for speakers of minoritized varieties. It reinforces an idea of standard AI-generated language ideology that holds SAE as a linguistic default and reinforces a linguistic bias that using SAE is the most "appropriate" way of speaking. AI-generated language thus subtly reinforces the belief that some language varieties (i.e. SAE) are better than others.

This is further reinforced by AI-generated language outputs being of lower quality for inputs in marginalized varieties compared to inputs in "standard" varieties. Given this, native speakers of minoritized languages may be incentivized to provide inputs in "standard" varieties even if they are less fluent in them, a form of digital code switching. If users prompt outputs to be in the particular minoritized language variety of the input, outputs can disproportionately carry stereotypes, produce demeaning content, facilitate appropriation, or support manipulation. Taken together, these harms can impact access to resources or opportunities when generative AI tools are used as gateways or checkpoints (e.g., chatbots for scheduling healthcare appointments or educational tutors) and as they become increasingly integrated into daily life. It illustrates the powerful-yet subtle-ways in which Western hegemony, and particularly American hegemony, manifests in and through emerging AI technologies.

## 4.2 No clear way for popular language models to "win"

These issues prompt a question regarding what constitutes desirable model behavior when different potential behaviors may result in different harms. Our taxonomy illustrates the challenges that technologists face in attempts to identify appropriate behavior regarding LLMs and different language varieties. What is the "appropriate" way that language models—and ensuing tools producing AI-generated text or voices—should behave given the current technological context?

Central to these dilemmas is the fact that people use language to reflect aspects of their identity. For example, speakers may switch to a non-standard variety when speaking to someone to signal their shared membership in an in-group [2, 30]. However, the "identity" of a language model is unclear: is it assumed to reflect that of its creators (often a White, Western, and male-dominated group)? Is it, instead, a neutral tool with a "view from nowhere," if such a view is even possible? Does the range of appropriate language use depend on who created the model, whose data was used to train it, who profits from it, or other factors? Furthermore, these decisions are typically made by model creators and tech industry members, but they fundamentally depend on the nuanced norms of speaker communities that rarely have a say in how these models are designed.

For popular generative AI chatbots and voice tools developed and owned by large corporations (e.g., ChatGPT, Gemini, Voice Engine) there is no clear "correct" behavior. These technologies draw on data from the Internet (which is largely comprised of "standard" language varieties, including SAE), resulting in inherent issues and

inequities, including reinforcement of standard language ideology. If a model is designed to automatically mimic the input variety unprompted, as opposed to defaulting to the language variety it is largely trained on (i.e. SAE), this could lead to issues outlined in the taxonomy, such as appropriation and manipulation. Instead, some may argue that the default production of a "standard" language is acceptable or even desirable. However, this line of thought has been documented to cause harm to minoritized language speakers in other contexts. For example, in education systems, prioritization of SAE results in psychological and academic impacts for Black students [43], while digital code switching may carry psychological tolls. The prioritization of SAE in ChatGPT and other popular generative AI technologies could result in harmful impacts at an even greater, global scale. Therefore, it is not clear what immediate approach should be taken for behavior of generative AI language technologies. This lack of clarity is concerning, particularly as those who hold decision-making power regarding the largest generative AI language technologies are their for-profit corporate owners with incentive structures tied to delivering value to shareholders.

## 4.3 Moving towards emancipatory outcomes

Instead of asking what the "appropriate" way is for generative AI models to behave, perhaps a more important question is: how might generative AI models be developed to support more emancipatory outcomes, and what do emancipatory outcomes from generative AI look like?

Perhaps most simply, the fact that language models perpetuate discrimination on the basis of language variety means that evaluation of language models for potential harms against minoritized groups (e.g. "toxicity" evaluations), which often focus on common demographic attributes such as race or gender, should expand to capture discrimination based on the user's language variety. Measurement of harms due to standard language ideology is particularly important because language use often serves as a proxy for other harms—a way that racism, sexism, or xenophobia may fly under the radar.

Issues surrounding appropriate responses to minoritized varieties, and when/if responding in the variety is appropriate, stem in part from the fact that these models are developed by organizations that are not composed of members of these speech communities. Instead, the development of language models by speaker communities, for speaker communities, helps to avoid this top-down imposition of linguistic norms (see, for example, [44, 51], and the broader discussion of indigenous data sovereignty by [38]). However, given the dominance of a small set of LLMs, such as ChatGPT, there is also an onus for developers of these widely used models to improve model quality for speakers of different language varieties. Despite the difficulty of implementing participatory practices in model design, setting "ambitious, yet reasonable goals" that help to move from one of the "extreme poles of transactional consulting and transformative empowerment" to the other [17] could help to improve power dynamics involving speakers of minoritized varieties. This includes innovating around participatory processes that

truly center the agency of minoritized people in the design, development, and ownership necessary to move towards more emancipatory digital futures.

In answering these open questions and advancing opportunities, it is important to consider what language data these models are learning from, how they can better learn from minoritized language varieties in non-extractive ways, and how marginalized communities can be enabled as central actors in the development and ownership of generative AI tools.

## 5 LIMITATIONS AND ETHICAL CONSIDERATIONS

This paper highlights the ways in which standard language ideology is present in AI-generated language, with implications for minoritized communities and society more broadly. While we believe that our taxonomy on open problems is comprehensive and we consider a range of issues, the list and associated implications are likely not exhaustive. A potential consequence of overlooking some work is narrowing the range of perspectives considered in future research based on the avenues we mention and the structure of our taxonomy. Nonetheless, it represents a first step in the development of a complete taxonomy and introduces critical questions to inform forward paths. Our discussion also centers largely on standard language ideology in English; potential differences in the manifestation of standard language ideology in other languages is an important area for future work.

## **REFERENCES**

- ALIM, H. S. Critical hip-hop language pedagogies: Combat, consciousness, and the cultural politics of communication. *Journal of language, identity, and educa*tion 6, 2 (2007), 161–176.
- [2] AUER, P., Ed. Code-Switching in Conversation. Routledge, July 2013.
- [3] BAKER-BELL, A. Linguistic justice. NCTE-Routledge Research Series. Routledge, London, England, May 2020.
- [4] BARTHEL, M., STOCKING, G., HOLCOMB, J., AND MITCHELL, A. Reddit news users more likely to be male, young and digital in their news preferences. *Pew Research Center 25* (2016).
- [5] BAUGH, J. Linguistic profiling. In *Black linguistics*. Routledge, 2005, pp. 167–180.
- [6] BENDER, E. The# benderrule: On naming the languages we study and why it matters. The Gradient 14 (2019), 34.
- [7] BLODGETT, S. L., WEI, J., AND O'CONNOR, B. Twitter universal dependency parsing for african-american and mainstream american english. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers) (2018), Association for Computational Linguistics.
- [8] BROCK, A. From the blackhand side: Twitter as a cultural conversation. Journal of Broadcasting & Electronic Media 56, 4 (2012), 529–549.
- [9] BUCHOLTZ, M., AND LOPEZ, Q. Performing blackness, forming whiteness: Linguistic minstrelsy in hollywood film1. *Journal of Sociolinguistics* 15, 5 (Nov. 2011), 680–706.
- [10] CHERY, S. Black english is being misidentified as gen z lingo, speakers say. Washington Post 17 (2022).
- [11] COWGILL, B., DELL'ACQUA, F., DENG, S., HSU, D., VERMA, N., AND CHAINTREAU, A. Biased programmers? or biased data? a field experiment in operationalizing ai ethics. In Proceedings of the 21st ACM Conference on Economics and Computation (2020), pp. 679–681.
- [12] CRAWFORD, K. The trouble with bias. keynote talk at Neural Information Processing Systems (NIPS '17), 2017.
- [13] CRENSHAW, K., GOTANDA, N., PELLER, G., AND THOMAS, K. Critical race theory: The key writings that formed the movement. The New Press, 1995.
- [14] CRYSTAL, D. English as a Global Language. Cambridge University Press, July
- [15] DAVIDSON, T., BHATTACHARYA, D., AND WEBER, I. Racial bias in hate speech and abusive language detection datasets. In *Proceedings of the Third Workshop on Abusive Language Online* (Florence, Italy, Aug. 2019), S. T. Roberts, J. Tetreault, V. Prabhakaran, and Z. Waseem, Eds., Association for Computational Linguistics, pp. 25–35.

- [16] DEAS, N., GRIESER, J., KLEINER, S., PATTON, D., TURCAN, E., AND MCKEOWN, K. Evaluation of African American language bias in natural language generation. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing (Singapore, Dec. 2023), H. Bouamor, J. Pino, and K. Bali, Eds., Association for Computational Linguistics, pp. 6805–6824.
- [17] DELGADO, F., YANG, S., MADAIO, M., AND YANG, Q. The participatory turn in ai design: Theoretical foundations and the current state of practice. In Proceedings of the 3rd ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization (New York, NY, USA, 2023), EAAMO '23, Association for Computing Machinery.
- [18] DELPIT, L. D. Education in a multicultural society: Our future's greatest challenge. The Journal of Negro Education 61, 3 (1992), 237–249.
- [19] DROŻDŻOWICZ, A., AND PELED, Y. The complexities of linguistic discrimination. Philosophical Psychology (Mar. 2024), 1–24.
- [20] DUGGAN, M. Online harassment 2017, 2017.
- [21] EBERHARDT, M., AND FREEMAN, K. 'first things first, i'm the realest': Linguistic appropriation, white privilege, and the hip-hop persona of iggy azalea. *Journal* of Sociolinguistics 19, 3 (2015), 303–327.
- [22] EUROPEAN UNION AGENCY FOR FUNDAMENTAL RIGHTS. Violence against women: An eu-wide survey. Luxembourg: Publications Office of the European Union (2014).
- [23] FEFFER, M., LIPTON, Z. C., AND DONAHUE, C. Deepdrake ft. bts-gan and taylorve: an exploratory analysis of musical deepfakes and hosting platforms. In Proceedings of the 2nd Workshop on Human-Centric Music Information Retrieval (2023).
- [24] FIELD, A., BLODGETT, S. L., WASEEM, Z., AND TSVETKOV, Y. A survey of race, racism, and anti-racism in NLP. 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) (Online, Aug. 2021), C. Zong, F. Xia, W. Li, and R. Navigli, Eds., Association for Computational Linguistics, pp. 1905–1925.
- [25] FLEISIG, E., SMITH, G., BOSSI, M., RUSTAGI, I., YIN, X., AND KLEIN, D. Linguistic bias in ChatGPT: Performance and stereotype discrepancies among english language varieties, 2024, forthcoming.
- [26] GING, D., AND SIAPERA, E. Special issue on online misogyny. Feminist Media Studies 18, 4 (2018), 515–524.
- [27] GODDARD, K., ROUDSARI, A., AND WYATT, J. C. Automation bias: a systematic review of frequency, effect mediators, and mitigators. *Journal of the American Medical Informatics Association* 19, 1 (2012), 121–127.
- [28] GROENWOLD, S., OU, L., PAREKH, A., HONNAVALLI, S., LEVY, S., MIRZA, D., AND WANG, W. Y. Investigating african-american vernacular english in transformerbased text generation. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP) (2020), Association for Computational Linguistics.
- [29] GUERRERO JR, A. 'you speak good english for being mexican'east los angeles chicano/a english: Language & identity. Voices 2, 1 (2014).
- [30] GUMPERZ, J. J. Discourse strategies. Cambridge University Press, 1982.
- [31] HILL, J. H. The Everyday Language of White Racism. Wiley, Oct. 2008.
- [32] HOFMANN, V., KALLURI, P. R., JURAFSKY, D., AND KING, S. Dialect prejudice predicts ai decisions about people's character, employability, and criminality, 2024.
- [33] HOLLIDAY, N. Siri, you've changed! acoustic properties and racialized judgments of voice assistants. Frontiers in Communication 8 (Apr. 2023).
- [34] HOROWITZ, M. C., AND KAHN, L. Bending the Automation Bias Curve: A Study of Human and AI-Based Decision Making in National Security Contexts. *Inter*national Studies Quarterly 68, 2 (04 2024), sqae020.
- [35] JENKINS, J. English as a lingua Franca in the international university. Routledge, London, England, Aug. 2013.
- [36] JERKINS, M. Black or bot? the long, sordid history of co-opting blackness online. Mother Jones (2022).
- [37] JOSHI, P., SANTY, S., BUDHIRAJA, A., BALI, K., AND CHOUDHURY, M. The state and fate of linguistic diversity and inclusion in the NLP world. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (Online, July 2020), D. Jurafsky, J. Chai, N. Schluter, and J. Tetreault, Eds., Association for Computational Linguistics, pp. 6282–6293.
- [38] KUKUTAI, T., AND TAYLOR, J. Indigenous data sovereignty: Toward an agenda. ANU press, 2016.
- [39] KUTLU, E., AND WILTSHIRE, C. Where do negative stereotypes come from? the case of indian english in the usa. Proceedings of the Linguistic Society of America 5, 1 (2020), 74–82.
- [40] LABOV, W. The social stratification of English in New York city. Cambridge University Press, 2006.
- [41] LIANG, W., YUKSEKGONUL, M., MAO, Y., WU, E., AND ZOU, J. Gpt detectors are biased against non-native english writers. *Patterns* 4, 7 (July 2023), 100779.
- [42] LIPPI-GREEN, R. Accent, standard language ideology, and discriminatory pretext in the courts. *Language in Society 23*, 2 (1994), 163–198.
- [43] LYN, K. O. Negotiating african american language, identity, and culture in the urban classroom. *Journal of Black Studies* 53, 8 (2022), 780–795.
- [44] MAHELONA, K., LEONI, G., DUNCAN, S., AND THOMPSON, M. Openai's whisper is another case study in colonisation. *Papa Reo* (2023).

- [45] MAKONI, S., SMITHERMAN, G., BALL, A. F., AND SPEARS, A. K. Black linguistics. Routledge, London, England, Dec. 2002.
- [46] MATAMOROS-FERNÁNDEZ, A., RODRIGUEZ, A., AND WIKSTRÖM, P. Humor that harms? examining racist audio-visual memetic media on tiktok during covid-19. Media and Communication 10, 2 (May 2022), 180–191.
- [47] McCluney, C. L., Robotham, K., Lee, S., Smiths, R., and Durkee, M. The costs of code-switching. Harvard Business Review (2019).
- [48] MICHAUD, M., COLPITTS, B. D. F., MATTHEW, M., AND COLPITTSBRADLEYD., F. English as a lingua franca: Globalization, ownership, and the diversification of english. Kwansei Gakuin University humanities review 20 (2015), 125–131.
- [49] MOORE, R. "if i actually talked like that, i'd pull a gun on myself": Accent, avoidance, and moral panic in irish english. Anthropological Quarterly 84, 1 (2011), 41–64.
- [50] NEE, J., SMITH, G. M., SHEARES, A., AND RUSTAGI, I. Linguistic justice as a framework for designing, developing, and managing natural language processing tools. *Big Data & Society 9*, 1 (2022), 20539517221090930.
- [51] NEKOTO, W., MARIVATE, V., MATSILA, T., FASUBAA, T., FAGBOHUNGBE, T., AKINOLA, S. O., MUHAMMAD, S., KABONGG KABENAMUALU, S., OSEI, S., SACKEY, F., NIYONGABO, R. A., MACHARM, R., OGAYO, P., AHIA, O., BERHE, M. M., ADEYEMI, M., MOKGESI-SELINGA, M., OKEGBEMI, L., MARTINUS, L., TAJUDEEN, K., DEGILA, K., OGUEJI, K., SIMINYU, K., KREUTZER, J., WEBSTER, J., ALI, J. T., ABBOTT, J., ORIFE, I., EZEANI, I., DANGANA, I. A., KAMPER, H., ELSAHAR, H., DURU, G., KIOKO, G., ESPOIR, M., VAN BILJON, E., WHITENACK, D., ONYEFULUCHI, C., EMEZUE, C. C., DOSSOU, B. F. P., SIBANDA, B., BASSEY, B., OLABIYI, A., RAMKILOWAN, A., ÖKTEM, A., AKINFADERIN, A., AND BASHIR, A. Participatory research for low-resourced machine translation: A case study in African languages. In Findings of the Association for Computational Linguistics: EMNLP 2020 (Online, Nov. 2020), T. Cohn, Y. He, and Y. Liu, Eds., Association for Computational Linguistics, pp. 2144–2160.
- [52] O'MARA, M. Why can't tech fix its gender problem? MIT Technology Review (2022).
- [53] PENNYCOOK, A. The myth of english as an international language. English in Australia, 139 (2004), 26–32.
- [54] PHILLIPSON, R. Linguistic imperialism continued. Routledge, London, England, Mar. 2010.
- [55] RICKFORD, J. R., AND KING, S. Language and linguistics on trial: Hearing rachel jeantel (and other vernacular speakers) in the courtroom and beyond. *Language* (2016), 948–988.
- [56] ROSA, J., AND FLORES, N. Unsettling race and language: Toward a raciolinguistic perspective. *Language in Society* 46, 5 (2017), 621–647.
- [57] RUNGTA, M., SINGH, J., MOHAMMAD, S. M., AND YANG, D. Geographic citation gaps in NLP research. In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing (Abu Dhabi, United Arab Emirates, Dec. 2022), Y. Goldberg, Z. Kozareva, and Y. Zhang, Eds., Association for Computational Linguistics, pp. 1371–1383.
- [58] SMITHERMAN, G. Talkin' and testifyin'. Wayne State University Press, Detroit, MI, June 1986.
- [59] SOBANDE, F. Spectacularized and branded digital (re)presentations of black people and blackness. *Television & New Media 22*, 2 (2021), 131–146.
- [60] THYLSTRUP, N., AND TALAT, Z. Detecting 'Dirt' and 'Toxicity': Rethinking Content Moderation as Pollution Behaviour. SSRN Electronic Journal (2020).
- [61] TIETZE, S. Spreading the management gospel-in english. Language and Intercultural Communication 4, 3 (2004), 175–189.
- [62] Verma, V., Fleisig, E., Tomlin, N., and Klein, D. Ghostbuster: Detecting text
- ghostwritten by large language models, 2024.

  [63] West, S. M., Whittaker, M., and Crawford, K. Discriminating systems: Gender, race and power in ai. AI Now Institute (2019), 1–33.
- [64] WILEY, T. G., AND LUKES, M. English-only and standard english ideologies in the u.s. TESOL Quarterly 30, 3 (1996), 511.
- [65] WILMOT, N. V. Language and the faces of power: A theoretical approach. International Journal of Cross Cultural Management 17, 1 (Jan. 2017), 85–100.
- [66] WRIGHT, K. E. Housing policy and linguistic profiling: An audit study of three american dialects. *Language 99*, 2 (June 2023), e58–e85.
- [67] YONG, Z. X., ZHANG, R., FORDE, J., WANG, S., SUBRAMONIAN, A., LOVENIA, H., CAHYAWIJAYA, S., WINATA, G., SUTAWIKA, L., CRUZ, J. C. B., TAN, Y. L., PHAN, L., PHAN, L., GARCIA, R., SOLORIO, T., AND AJI, A. Prompting multilingual large language models to generate code-mixed texts: The case of south East Asian languages. In Proceedings of the 6th Workshop on Computational Approaches to Linguistic Code-Switching (Singapore, Dec. 2023), G. Winata, S. Kar, M. Zhukova, T. Solorio, M. Diab, S. Sitaram, M. Choudhury, and K. Bali, Eds., Association for Computational Linguistics, pp. 43–63.

## A EXAMPLE RESPONSE

Here, we provide an example of an incorrect and unhelpful response to a prompt in Indian English. Part of an input in Indian English stated: "I hope you have received my earlier letter. I would

have liked to take up your kind invitation to act as resource person for the ELTC seminar, but for reasons already mentioned, I won't be able to make it this time."

When prompted to respond, Chat GPT's response simply rephrased part of the input: "Thank you for your letter dated Feb. 16th. I appreciate your kind invitation to act as a resource person for the ELTC seminar. However, due to reasons mentioned in my previous letter, I regret to inform you that I won't be able to attend this time."

