

Multilingualism and AI: The Regimentation of Language in the Age of Digital Capitalism

Britta Schneider, *Europa-Universität Viadrina, Germany*

ABSTRACT

This article examines the effects of commercial digital language technologies on the regimentation of language. Language technologies based on the exploitation of large data sets—from machine translation and automatic text generation to digital voice assistants—are a particular form of human-made sign practice in which traditional language norms interact with the affordances of digital devices and the capitalist interests of those who design them. Such sociotechnological practices construct language hierarchies within the realm of commercially based language technology and can shape both dominant discourses about language in society and epistemologies of language in linguistics. The article focuses on interrelationships between digital language technology and metasemiotic interpretations of language that pertain to multilingualism, language variation, and language prestige. It examines languages as discursive constructs and reviews the role of media technology in shaping language ideology, showing that writing and print have had a crucial impact on modern language concepts. It draws on expert discourse and qualitative interviews with programmers and users and examines ideological effects of digital language technology and the potential epistemological reconfigurations of concepts of language that may emerge as a result.

The study of signs in today's societies cannot ignore the role of commercial language technologies. From machine translation and automatic spell checkers to the use of voice assistants, language technologies are omnipresent in many people's lives, and particularly in the lives of those who hold

Contact Britta Schneider at Kulturwissenschaftliche Fakultät, Große Scharrnstrasse 59, Frankfurt Oder, Brandenburg 15230, Germany (bschneider@europa-uni.de).

This article has been inspired by talks and discussions realized in the Ideologies, Beliefs, Attitudes working group within the European Cooperation in Science and Technology (COST) Language in the Human-Machine Era project. I want to thank all members and participants for sharing their thoughts and ideas. Also, I want to thank the editor as well as an anonymous reviewer for constructive feedback. All remaining errors are my own.

Signs and Society, volume 10, number 3, fall 2022.

© 2022 Semiosis Research Center at Hankuk University of Foreign Studies. All rights reserved. Published by The University of Chicago Press for the Semiosis Research Center, Hankuk University of Foreign Studies. <https://doi.org/10.1086/721757>

dominant social positions and whose job it is to create contents, values, and knowledge. As scholars who study the social functions of linguistic and other signs, we aim for an understanding of “sociolinguistic economies” (Blommaert et al. 2005), that is, the socially meaningful hierarchies that exist with regard to different sets of signs in different societies. It has to be assumed that the pervasiveness and convenience of many digital language technologies will have an impact on how individuals and groups understand the value and the form of linguistic signs.

The research field that studies discourses on linguistic signs refers to conceptualizations regarding form and value as *language ideologies* (e.g., Woolard 1998; Kroskrity 2000; Gal and Irvine 2019), and it is from the perspective of language ideology research that I scrutinize discursive constructions of language and sociolinguistic hierarchies in the context of commercial language technologies. I focus particularly on the question of what happens to ideas regarding language variation and multilingual practices in contexts of language technology design and use. Language variation is here understood as both variation in the form of using more or less established “varieties” (conceptualized as geographically or socially based sets of linguistic form) and variation *within* specific “varieties” or “languages,”¹ that is, lexical, pragmatic, or syntactic variation of individual speakers in, for example, contexts of language contact or diversity. What are the effects of the design, affordances, and materialities of language technological tools on normative concepts of language? Do digital language technologies reproduce “languages” as normed and bounded entities, and are societies understood as “normally” being monolingual? And what are the social indexical hierarchies relating to language variation that digital technologies reproduce or create? Finally, what are the epistemological effects of language technologies on what we believe language to be?

In the following, I first establish the theoretical foundations for understanding language ideologies and for understanding *languages* as discursively constructed language ideological phenomena. In this context, I also discuss the relevance of material-technological practices in shaping language ideologies, where the cultural practices of writing and print have been particularly influential in coining Western normative concepts of language (see the section “Theoretical Background: Languages as Material-Technological Discourse and Nonessentialist Understandings of Human Meaning-Making Practices” below). The language ideologies developed in the age of print capitalism (e.g., Anderson 1985)—such as the ideal of

1. Note that I use the terms *language* and *variety* to refer to discursive concepts that imagine but also construct particular sets of linguistic signs to index particular social persona (for discussion, see, e.g., Schneider, forthcoming).

nationally normed pronunciation, orthography, and grammar (Agha 2003)—continue to be highly influential in digital contexts, and, yet, these are also reconfigured.

The third and main section of the article is devoted to understanding discourses and technological practices that construct language in the age of digital capitalism. Machine translation, text generation, and speech recognition are some popular examples. Large, machine-readable data sets form the foundation for the functioning of these tools. I discuss the potentials of reification and dominance of particular *languages* (sets of linguistic signs) that are associated with such data sets and with their exploitation, in which machine-learning algorithms often come into play. Both data sets and algorithms are framed by language ideological discourses of programmers, to which I give insight based on qualitative interviews and publications by the language technology industry. These discourses hint at the strong dominance of English; have a tendency to construct digital technologies as agentive beings, equipped with cognitive abilities; and display the underlying commercial logics of the tools where efficiency, data availability, and capitalist interest have an impact on how digital language technologies are designed. To round up this explorative study on the possible effects of language technologies on language ideologies, I give access to the user's side. Interviews with multilingual users of voice assistants confirm the dominance of English in private, everyday human-machine interaction and bring to the fore surprising indexical associations between language and culture that overcome traditional nationalist language paradigms. The article ends with concluding thoughts on benefits and potential problems of commercial language technologies and discusses their epistemological effects on concepts of language in more general terms.

Theoretical Background: *Languages* as Material-Technological Discourse and Nonessentialist Understandings of Human Meaning-Making Practices

On the basis of contemporary language ideology research and sociolinguistic theory, I understand dominant concepts of *languages*, that is, constructions of language as appearing in systemic and bounded entities, to be an effect of specific sociocultural discourses (e.g., Irvine 2001; Agha 2003; Pennycook 2004). Drawing on poststructuralist concepts of performativity, language practice may be described as a series of performative acts that bring into being *language*, rather than being based on a priori and immaterial grammatical systems (Pennycook 2004, 15). Grammatical, lexical, and phonetic systematicity, in this perspective, “is an illusion produced by the partial settling or *sedimentation* of frequently used forms into temporary subsystems” (Hopper 1998, 157–58, quoted in Pennycook 2004, 4).

In linguistic anthropological as well as in cognitive-oriented traditions, language is similarly considered as interactive practice, sometimes referred to as *linguaging* (e.g., Madsen et al. 2016; Love 2017). Thus, the human ability to communicate via sound and other means is understood not as based on internal grammatical systems but rather as an activity in which humans coordinate their behavior with their human and nonhuman environment: “Language can be traced to how living bodies co-ordinate with the world. On this perspective, far from being a synchronic “system,” language is a mode of organization that functions by linking people with each other, external resources and cultural traditions” (Cowley 2011, 2).

In the context of linguistic anthropology and various related research strands, the idea of languages as rule-based, orderly, and clearly delimitable cognitive systems that have primarily referential functions is understood as culturally contingent concepts with roots in European modernity (e.g., Joseph and Taylor 1990; Errington 2008; Irvine and Gal 2009). A cognitive and systemic understanding of language is, at the same time, intertwined with traditions of taking the communities that “use” these languages for granted. In this sense, it is related to what has been referred to as *methodological nationalism* in the social sciences (Wimmer and Schiller 2002)—the idea that the world is “naturally” divided into territorial entities in which reside culturally homogenous groups that use one (and ideally only one) language (for a discussion of the effects of methodological nationalism on linguistics, see Schneider [2019]). And yet, a monolingual nation has never existed, and the promotion and establishment of the idea has cost a lot in terms of money, material, and human resources—for example, in national educational systems or in the design and production of grammars and dictionaries.

One reason for the institutionalization of monolingual ideologies can be found in the development of national publics. Gal and Woolard, two pioneers in language ideology research, argue that national publics and national standard languages are in a dialectical relationship to each other and thus depend on each other. Both of these—national languages and national publics—bring about the idea that there are “voices from nowhere” (2001, 7), that is, social positions and voices of those in power in a particular national setting, which are understood as socially unmarked and as not conveying a specific perspective. Hegemonic social and linguistic dominance is here constructed as “neutral.”

When discussing languages as socially constructed entities that have come into being under particular historical, discursive, and political conditions, material and technological practices play an important role. In the context of producing dominant national voices from nowhere, the material practices of nationally normed

writing in printed documents are crucial but often have been taken for granted, also in linguistic research, where a “written language bias” has been attested (Linell 2005). In contemporary, globalized, and postnational types of publics, in which digital practices imply a destabilization of supposedly “neutral” dominant voices (Heyd and Schneider 2019), the fact that national language standards depend on specific materialities and material-technological practices has come to the fore. This ties in with a recent interest in so-called posthumanist approaches, which deconstruct humanist and human-centered perspectives that typically focus on the rational, cognitive abilities of humans and fail to take into account the material, nonhuman, and environmental conditions in which human communicative practice takes place (Pennycook 2018; Schneider 2021). It is important to note that the signs we use to produce and convey meaning are themselves of a material kind, as Gal and Woolard (2019, 89) have pointed out:

Unlike the usual Cartesian view, in which thought is rooted in radical doubt and introspection, our view is that thinking requires some sort of expressive form—signs—to convey the objects of thought. For Cartesians, communication is secondary, other people’s minds remain a mystery, and minds are separate from the materiality of bodies. For us, thinking starts *not* with doubt but with previous knowledge, with matters that at any historical moment are familiar to some knowers, to some extent. Signs are the products and tools of such knowers in social relations. Instead of a Cartesian split between mind and bodily matter, between individual thinkers and social groups, we are interested in how such realms—once separated in one major philosophical tradition—are connected, and how signs mediate the connection.

In the context of the study of linguistic signs, this means that we should consider the material and bodily aspects that come into play in language practice and in the conceptualization of language (consider also Gershon’s [2010] related discussion on *media ideologies*). Against this background, we can attest that the construction of national language standards as well as the construction of national publics are not conceivable without the technologies of writing and print. The nationwide distribution of the idea of national belonging and of nationally shared language standards would not have been possible without the practices of writing and the printing press (Eisenstein 1979; Anderson 1985; Giesecke 1991). Also, understanding languages as coherent, linear, and systemic entities is in itself an epistemological effect of writing. It is unlikely that purely oral cultures would have developed concepts that treat verbal signs produced by

human bodies in interaction (“speaking”) in an object-like, stable, and linear fashion (for a discussion, see also, e.g., Ong [1982]; Linell [2005]).

Overall, “the structure of a technology helps to shape the participant structure brought into being through its use, simultaneously enabling and limiting how communication can take place through that medium” (Gershon 2010, 285). The technologies of writing and print have co-constructed national language standards, which are often taken as a priori “given” in social contexts but also in many linguistic approaches to the study of language (for a critique, see, e.g., Bourdieu [1982, 27]). Such perspectives tend to coproduce an understanding of multilingualism, language change, and non-normative language as special cases with typically negative connotations (e.g., Lippi-Green 2012; Gramling 2016). And even though we can observe reconfigurations of monolingual national ideologies in the contexts of late modern capitalist culture due to the commodification of multilingualism as an asset in the global job market (Heller and Duchêne 2012), these discussions rarely have taken into consideration the material and technological bases of communicative practices.

Observing that, first, the material-technological practices of communication have influenced our concepts of language and that, second, these have been rather drastically reconfigured in the last 30 years due to digital media (for an early study, see, e.g., Herring [1996]), it is important to ask how our ideas about language—the social spaces, ideals, values, and hierarchies attached to specific linguistic practices—have been reconfigured, too. In the context of this article, I therefore study the language ideological practices found in the realm of the digital language technology industry and develop thoughts on how these may have an impact on language ideologies and on epistemologies of language. What happens to the monolingual ideologies associated with the public regimes and the materialities of the modern nation-state, and what happens to ideologies relating to multilingual practice and language hierarchies in the context of digital language technologies?

Commercial Language Technologies and Their Language Ideological Effects

Digital language technologies involve devices and programs that exploit the computing abilities of digital hardware and software with the aim of generating or decoding human languages, referred to as “natural languages” in this context (contrasting with computer code, also called “programming languages”). Some of the general aims of creating technologies that produce/decode human language are to translate across human languages or to provide other language-based services,

such as informing about something, controlling electronic devices, or classifying written contents. Similar to the effects of print technologies, human language is typically reified as consisting of linear, separable signs that occur one after the other. With the exception of digital tools for sign language, most language technologies render the multimodal and embodied nature of verbal interaction among humans invisible.

Earlier attempts of making computers “understand” or generate human language were often based on programming lexical elements and grammatical rules of human language into computer code.² Such “rule-based” systems are also referred to as “symbolic natural language processing,” or Natural Language Understanding (NLU; see, e.g., McShane and Nirenburg 2021, 33–36). Most contemporary digital language technologies do not, however, aim at programming algorithms to generate or decode the rules of human grammar. Rather, they are based on statistical analyses of large language corpora where “word embeddings”—the statistical probabilities of a word occurring in relation to other words—are calculated (see Jurafsky and Martin [2021, chap. 6] on statistical analyses and an introduction into how such methods works). The idea is to produce “models that process words in relation to all the other words in a sentence, rather than one-by-one in order” (Nayak 2019). In computational linguistics, such statistical word-embedding approaches are highly popular and are today often associated with Natural Language Processing (NLP; see, e.g., Jurafsky and Martin 2021).

Since 2012, more and more of such purely statistical calculations of word correlations have been realized by machine-learning techniques, that is, artificial neural networks (interrelated algorithms), that identify word embeddings on grounds of statistical analyses and do this in “unsupervised” or “self-supervised” learning (Bommasani et al. 2021, 4). This entails that human programmers do not necessarily decide which features are deemed relevant and that the data sets do not have to be labeled by humans (e.g., according to parts of speech) before computer algorithms start to calculate word relationships. In other words, multiple interconnected algorithms (“neural networks”) are used to statistically analyze very large language data sets, with little human intervention, with the aim of calculating the likelihood of words appearing next to each other. These so-called language models very often produce texts that are coherent according to human standards and are so far the most successful type of language technology. Generally, approaches that make use of neural networks are also referred to “deep

2. Verbs that refer to active cognitive abilities are used frequently in relation to digital tools. This brings about metatheoretical considerations concerning the question of what these abilities are, which has been problematized, studied, and discussed elsewhere (see, e.g., Bender and Koller 2020).

learning” techniques and are associated with “artificial intelligence” in public discourse (see Jurafsky and Martin [2021, chap. 7] on neural network techniques in language technology and Bommasani et al. [2021] and Crawford [2021] for critical discussions of these technologies).³

Two of the currently popular language models that apply machine learning to calculate word embeddings are Google’s Bidirectional Encoder Representations from Transformers, or BERT, published in 2018 (Devlin and Chang 2018) and the GPT series (based on an open initiative, today licensed by Microsoft; Hao 2020a). Given that neither syntactical rules nor semantic contents are part of the technology design in these approaches, it needs to be noted that the technologies are not meant to “understand” language, nor are they programmed to actually produce grammatical structures (see, e.g., Bender and Koller 2020). Rather, they generate strings of lexical elements that are put together on grounds of statistical calculations. Therefore, they have been criticized as “statistical parrots” (Bender et al. 2021).

Even though such digital computing technologies do not actually allow for decoding or producing syntactical rules and semantic meanings of human languages, they have shown surprising functions that can be very useful in everyday life and for commercial exploitation. Popular examples include

- Machine translation (decoding and generation of text);
- Voice assistant systems (decoding and generation of spoken language);
- Automatic text generation (e.g., in SMS, e-mails, weather forecasts, etc.);
- Localization (adaptation of texts, digital content, or software to specific cultural contexts);
- Automatic categorization of text (e.g., product descriptions for online market platforms); and
- Sentiment analysis (recognition and categorization of words associated with positive or negative emotions).

Many of these tools have positive, socially integrative effects and can support successful communication, particularly in multilingual contexts. Automatic translation provided by web tools or smartphone apps, such as Google Translate or DeepL, has become important in cross-language interaction (e.g., Randhawa et al. 2013; Asscher and Glikson 2021) and in, for example, trajectories of migration

3. For an early tutorial by Richard Socher, a central figure in applying deep learning to language data, see <https://youtu.be/IF5tGEgRCTQ>.

(Zijlstra and Liempt 2017). Such tools can also be used to translate governmental forms so that speakers who do not speak an official language can fill in documents and get access to resources provided by the state (e.g., Alatuno). Digital tools to support second language acquisition (e.g., Babbel, DuoLingo) are also on the rise and seem to be helpful for some language learners (even though they have also been critically commented upon due to a lack of social interaction and context by, for example, Piller, quoted in Hepworth [2021]).

At the same time, there are a number of problems attached to these technologies. Machine-learning technologies require a large amount of computing resources and a mass of data and thus have been developed mostly by large commercial companies that have access to such means (on environmental and human resources exploited in AI, see Crawford [2021]; see also Bommasani et al. [2021] on potential social problems and cultural homogenization). This also implies that the objective of developing language technologies is of a capitalist nature—things have to be profitable and work efficiently, and users should be satisfied. Following this logic, digital language technologies are based not on educational, aesthetic, or linguistic norms but rather on (machine-readable) language practices of customers. Thus, in order to make available huge data sets, language data are extracted from the web or from voice assistants. This implies that language data that have been digitized (e.g., newspaper archives, linguistic corpora, Google Books), produced on the internet (e.g., Wikipedia or social media data), or recorded by voice assistants are considered to represent *language* in general. Machine-readable language data are in this sense norm-providing, as they are understood as “normal use.” In addition, on grounds of the affordances of computing, which typically work according to quantitative logics, what is frequent in the data is more likely to be reproduced (see also the section “Machine Learning and the Amplification of Biases”).

This does not mean, however, that there are no moral-normative discourses in the context of language technology industries. There is a growing awareness that there are biases in machine-readable data, such as gendered or racial biases, that are reproduced and even amplified by statistical calculations.⁴ In the discourses that problematize these biases, they are typically referred to as “harms” and “risks” (e.g., Bommasani et al. 2021). These word choices have been influenced by “battlefield-oriented notions,” as the US military has been a foundational actor in the development of AI (Crawford 2021, 185), for example, in funding research on digital technologies as well as research on language (the interest of military

4. See also Crawford (2017); Zhao et al. (2017); Sun et al. (2019); Hewett and Nenno (2021).

actors in automatic decoding and translation of language is not new; see Heller and McElhinny [2017, chap. 6] on the role of the US military in funding syntax research during the Cold War).

Overall, the language ideological framing and the constructions of language norms in capitalist digital technologies differ from those of the traditional nation-state. This does not mean, however, that traditional language norms have become irrelevant. Language technologies reproduce concepts of language as appearing in separable and countable entities. Norms and ideals as they have been developed in the age of national print literacy, as well as imperialistic and postcolonial sociolinguistic hierarchies, are programmed into digital tools and interact with machine logics. How discourses and practices within the digital language industry reproduce languages as entities and how they have an impact on sociolinguistic hierarchies—that is, which linguistic resources are tacitly constructed as valuable and which are not—is related to the role of Big Data corpora, the mechanisms of algorithms, programmers' language ideologies, and resulting user practices. I will elaborate on each of these points in the following.

The Role of Language Corpora in the Regimentation of Language via Digital Tools

As mentioned above, most digital language technologies are based on large, machine-readable language data corpora. These are very predominantly monolingual, which means that only data that are considered to represent “one language” are used to program a specific tool, and so far all tools have been designed to produce monolingual forms of language. English plays a very central role in the development of language technologies given the dominance of the US technology industry. Thus, most digital language tools are first and foremost produced to work in English. This also has to do with the fact that English is the language for which by far the largest amount of data are available (more than 60 percent of data on the web were in English in 2021; see Hootsuite 2021).

In order to generate large language corpora, it is common to use so-called web crawlers (<https://www.oncrawl.com/technical-seo/introduction-web-crawler/>) that automatically search specific parts of the web for language data. In some cases, these are data that are assumed to comply with traditional language norms. Thus, Wikipedia pages, texts from EU institutions or sites from Google Patents have been reported to be commonly used to “train” machine-learning language technologies (Pakalski 2009; Dodge et al. 2021). Note that the notion of “training” here refers to the procedures by which algorithms “learn” typical correlations on the basis of a limited (but still very large) set of data on grounds of which

they produce a statistical model. The model is then used, for example, to translate, produce, or categorize other kinds of language data (see, e.g., Zweig [2019, chap. 5] on details of machine learning). As the aforementioned web pages are freely available, their exploitation represents an efficient and cheap way of “training” the algorithm. Given that data in English dominate the web and that the companies leading technology development are often based in English-speaking environments (e.g., Apple, Google, Amazon—note the exception of the very successful Germany-based machine translation company DeepL), it is unsurprising that most language tools are first “trained” with English.

Existing language ideological hierarchies that construct English as “normal,” more useful, and more valuable are here tacitly reproduced and amplified. As has been observed in research on voice-controlled digital assistants, for example, global hierarchies that derive from times of colonialism and imperialism continue to have an impact. The languages with the highest prestige are those for which most machine-readable data exist, which is why it is more likely that digital language tools are designed for these languages. Digital assistants like Siri and Alexa exist for working in “the standard variants of European/Western languages, with some exceptions,” and there are currently “no African language packages available from these leading companies in voice technology” (Lelebici 2021, 8; this is likely to change over time). In addition, nonstandard and hybrid forms of language use are not supported by technologies that construct languages as machine-readable monolingual patterns. As a result, “voice control technologies do not allow [for] bilingual practices or [nonstandard] accents and offer languages in a packaged form” (9).

A major problem attached to contemporary language technologies, therefore, is that linguistic practices that so far do not or only rarely appear in machine-readable data are clearly disadvantaged. This is so despite the fact that the variety of languages for which digital technologies are produced is constantly increasing (e.g., NLLB 2022). Some widely spoken nonstandard languages can now be decoded by digital voice assistants, as was reported to me, for example in interviews with users who speak South Tyrolian (a nonstandard variety of German spoken in the north of Italy) and told me that Alexa can “understand” them. However, the data sets for nonstandard varieties are obviously much smaller than for standard, written varieties, which can create problems. Kreuzer et al. (2021) discuss in this context that the data sets used for training tools to produce languages other than English, and particularly minority languages, are often of substandard quality so that the outputs may be inappropriate and sometimes simply wrong. For languages that do not conform to ideals of normed standardization and that are

intrinsically variable, it is even more difficult, culturally problematic, and maybe even impossible to construct language output on specific data sets (even if they existed). For example, it has been shown for some Caribbean creole languages that creativity, verbal poetry, and individual appropriation are important cultural values attached to language practice (see, e.g., French and Kernan 1981; Schneider 2021a). For such language cultures, the idea of producing normed language on the basis of limited data sets does not make sense and may even pose a threat to language cultural traditions.

Dominant, prestige-loaded, and standard forms (mostly from European languages), on the other hand, are further pushed in status as popular gadgets like machine translation and digital voice assistants are available and work best in these. As a consequence, “over 90% of the world’s languages used by more than a billion people currently have little to no support in terms of language technology” (Bender et al. 2021, 612), and insiders to the industry therefore argue that “most language technology is built to serve the needs of those who already have the most privilege in society” (613). Note that one of the authors of this 2021 article previously worked for Google and apparently lost her job due to its publication (for further discussion, see Hao [2020b]).

The Problem of Biases within Corpora

Besides the problem of unequal availability of machine-readable data sets, there are problems with biases within the respective data sets. In contrast to language corpora as they have been developed in linguistic research, there is often no intention in the language industry to create what linguists have called “balanced” or “representative” corpora—linguists typically design language corpora with the idea of representing a speech community, for example, by selecting texts on the basis of criteria like orality, literacy, register, genre age, educational status, or formality (see, e.g., Stefanowitsch 2020, chap. 2). As has been elaborated above, data sets typically used by the industry are based on data that are available on the web. Yet, the assumption that these data represent “neutral language” is faulty. Let me illustrate this with the case of Wikipedia, which has been problematized in several contexts (see, e.g., Barera 2020).⁵ Texts from English-language Wikipedia pages have been used to “train” algorithms that create models of word embeddings. However, the language used in Wikipedia is not “neutral” and does not represent the language use of all speakers of English. Rather, it is mostly male academics ages 18–30 who have no partner or family who produce the largest share of texts

5. For an internal discussion on Wikipedia (including statistics on authors), see https://meta.wikimedia.org/wiki/Community_Insights/2018_Report/Contributors.

on Wikipedia (Mitchell 2021). The fact that contents and topics on Wikipedia are therefore biased has been discussed—for example, that only 18 percent of biographies on Wikipedia are of women (Hewett and Nenno 2021). It is now well-known that “uncurated, Internet-based datasets encode the dominant/hegemonic view, which further harms people at the margins” (Bender et al. 2021, 613; see also UNESCO [2019] for a discussion of gender bias in the entire industry). In relation to language, this means that the written language practices of young, male, academically trained populations define what is understood and reproduced as “normal” language in digital language technologies that have been trained with such data sets. Culture- and class-specific monolingual, standard norms of writing are therefore reified by such technologies.

Some tools have been based on even more problematic data—GPT3, for example, has been partially trained with data from Reddit (Bender et al. 2021), a social media platform that is known for a male dominance (Hewett and Nenno 2021), as well as sexist and racist contents (on the Manosphere on Reddit, see, e.g., Ging [2019]). It could be shown 63,000 texts used for training GPT3 were from “subreddits” that had been banned due to their problematic content (Bender et al. 2021, 613). Language productions that entail structural racism and sexism are common in such data and “when the foundation is biased, there is a good chance that it spreads to the entire system” (Hewett and Nenno 2021; see also Gehman et al. 2020). Gendered biases in data, for example, lead to the automatic generation of text that reproduces highly stereotypical images of men and women (see, e.g., Savoldi et al. 2021). Thus, it has been demonstrated for some word-embedding tools that male subjects are associated not only with generally more words but also with more high-value words. In an online tool that illustrates the problem, the five most commonly associated adjectives with the pronoun *he* are for example “guy/guys, certainly, obviously, cocky, decent.” For *she*, these are “sassy, sexy, gorgeous, cute, lovely” (see <https://www.are.na/block/4754465> for documentation of the tool, which is inactive as of late 2022; see <https://jyzhao.net/files/naacl19.pdf> for another example showing highly stereotypical job associations). Another relatively famous example is that “computer programmer” is much more likely to be associated with “man,” while “home-maker” is more likely to be associated with “woman” (Bolukbasi et al. 2016).

In the industry, the bias problem is currently addressed mainly by trying to find technical, algorithmic solutions to avoid, for example, sexist or racist language productions of technological devices (e.g., Bolukbasi et al. 2016; Xu et al. 2021). This has been criticized, as it overlooks the fact that cultural categorizations are deeply engrained in language (or one could argue that these are actually the essence of language and culture; see Crawford 2017), so that

discriminative categorizations and common biases found in training data will not disappear (Crawford 2021, chap. 4). Categorizations of humans that are in one way or another problematic cannot be “solved” by so-called debiasing projects, in which the likelihood of a word appearing close to another word is statistically diminished. Some experts therefore suggest to generally forbid to use machine-learning techniques to categorize humans (Zweig 2019; Crawford 2021).

A tendency of marginalizing what is already marginalized is thus found due to the unequal availability of web corpora as well as because of biases found in the specific types of data. This tendency is enforced through the mechanisms of machine-learning algorithms.

Machine Learning and the Amplification of Biases

As has been explained above, the notion of “machine learning” is applied where correlations between machine-readable items (e.g., words but also images and other types of data) are detected automatically by algorithms. It is an often mentioned problem that humans do not always understand on what grounds or why the algorithms build their models (as discussed in, e.g., Bommasani et al. [2021]; Crawford [2021]). This is also true of today’s dominant Big Data models, like Google’s BERT, which are not only used for one specific purpose but where one algorithmic model is applied in many contexts, making use of so-called transfer learning techniques.⁶ This is why these influential and powerful models are also referred to as “foundation models” (Bommasani et al. 2021)—they are the foundation for many different types of applications. In these, there is the problem that humans generally do not clearly understand on which grounds the algorithms detect patterns, that is, declare items as being related to each other. This is discussed in a 2021 publication, written in a joint effort by a large number of specialists from the field, under the guidance of the Stanford University Center for Research on Foundation Models. The authors here argue that “at present, we emphasize that we do not fully understand the nature or quality of the foundation that foundation models provide; we cannot characterize whether the foundation is trustworthy or not.” (Bommasani et al. 2021, 7).

Something that is clear, however, is that the general logic of machine-learning algorithms is to detect correlations that are frequent. In other words, the aim of these tools is to overgeneralize what appears in the data often. As the training data tend to display problematic biases (see the previous section), this mechanism of overgeneralizing obviously may lead to even more problematic outputs. Additionally,

6. “The idea of transfer learning is to take the ‘knowledge’ learned from one task (e.g., object recognition in images) and apply it to another task (e.g., activity recognition in videos)” (Bommasani et al. 2021, 4).

the general desire of finding frequent patterns in data, and to produce language on the basis of what is frequent, has an intrinsic leaning toward homogenization. A study on the effects of machine translation on lexical variation, tellingly called “Lost in Translation: Loss and Decay of Linguistic Richness in Machine Translation” (Vanmassenhove et al. 2019), demonstrates this trend. The authors show, using the example of English-Spanish translations, that translations produced by neural network techniques are much less likely to contain variable translations of the same lexical item than translations produced by humans. Thus, while humans translate the English term *picture* in about 80 percent of cases by the Spanish term *imagen* (or its plural form), machine translation generated this translation nearly 100 percent of the time. The authors conclude that their “analysis shows that MT [machine translation] paradigms indeed increase/decrease the frequencies of more/less frequent words to such extent that a very large amount of words are completely ‘lost in translation’” and that neural network techniques are “among the worst performing” in terms of lexical diversity (230). Besides the issue of lexical homogenization, the authors assume that their example may explain the amplification of dominant forms on other levels of language as well, for example in the lexicon and in lexical associations: “The inability of MT systems to generate diverse outputs and its tendency to exacerbate already frequent patterns while ignoring less frequent ones, might be the underlying cause for, among others, the currently heavily debated issues related to gender biased output” (222).

The tendency of homogenization can also be related to other language-related questions. It has been shown that digital voice assistants confirm the pattern of reproducing and supporting dominant forms. For example, it could be demonstrated that they work better with male than with female voices—“Google’s widely used speech-recognition software is 70 percent more likely to accurately recognize male speech than female speech” (UNESCO 2019, 34, quoting Tatman 2016) because the devices were first trained with male speakers. Unsurprisingly, nonstandard varieties are also among the less well decoded forms of language in digital assistants. It has been empirically shown that these “systems exhibited substantial racial disparities, with an average word error rate (WER) of 0.35 for black speakers compared with 0.19 for white speakers” (Koencke et al. 2020, 7684). This means that it is much more likely for black speakers to not be understood by their digital assistants.

The problem applies also to the syntactical level, where the fact that most algorithms have been trained with the relatively analytic language English comes to the fore. It seems that it is more difficult to apply algorithmic models to languages that are typologically distinct from English. This is reinforced by a monolingual

habitus (Gogolin 1994) of the programming community, here elaborated on in an interview with a German language technology programmer:

Transcript 1

Das Problem ist, dass viele Systeme in Englisch konzipiert werden und initial realisiert werden	The problem is that many systems are designed and initially implemented in English
Und dann in anderen Sprachen aufgrund von sprachlicher Komplexität zusätzliche Ebenen auftreten, die teilweise in der Implementierung nachgerüstet werden können, teilweise eben aber auch nicht,	And then in other languages, due to linguistic complexity, additional levels occur, which can partly be retrofitted in the implementation, but partly not
Und da musst Du dir große Mühe geben diese Shortcomings irgendwie zu kaschieren	And you have to make a great effort to conceal these shortcomings somehow
Und insofern gibt es auch glaube ich an einigen Stellen schon Limitationen, dadurch, dass Systeme von englischen Muttersprachlern entwickelt werden, die keine anderen Sprachen können	And in this respect, I think there are indeed limitations in some places, due to the fact that systems are developed by native speakers of English who do not know any other languages

The interviewee hints at the fact that the grammatical structures of English display little complexity on certain levels (simple verb paradigms, lack of case system, etc.), so that it is difficult to implement such phenomena if the systems have not been programmed or trained with these features in mind. This is not, however, based on the fact that it would not be technically possible but the programmer holds monolingual ideologies and lack of competence in languages other than English on the side of programmers responsible for the problem.

We can thus postulate that the interaction of monolingual ideologies and the homogenizing effects of algorithmic techniques results in a lack of functioning of digital tools in languages typologically distinct from English and, at the same time, contribute to a constantly increasing performance of the tools in English. Indeed, language professionals report that, for example, machine translation is meagre for languages with a high number of grammatical cases like Finnish.⁷ In addition, it has been shown that the automatic generation of languages other than English may result in sentences whose grammatical structures are influenced by those of English, particularly where transfer learning is applied (see, e.g., Lauscher et al. 2020; Virtanen et al. 2019, quoted in Bommasani et al. 2021, 25). All in all, “multilingual models [models for machine translation that entail data sets from two or more languages—not models that would support multilingual practices]

7. Maarit Koponen, professional translator, personal communication with the author (at a talk in our working group on December 17, 2021).

show better performance in languages that are similar to the highest-resource languages in their training data” and researchers from the academic realm warn that there is a danger in “models that erase language variation and mostly conform to the linguistic majority in their training data” (Bommasani et al. 2021, 25).

The tendency to further support what already is dominant becomes even more pronounced if it is considered that at least currently, it is only very few models (among them the aforementioned foundation models BERT and GPT3) that are used in a large number of applications and that define how language technologies are developed, “Foundation models have led to an unprecedented level of *homogenization*, Almost all state-of-the-art NLP models are now adapted from one of a few foundation models . . . all AI systems might inherit the same problematic biases of a few foundation models” (Bommasani et al. 2021, 5). This tendency toward homogenization may be enforced through user practices.

Feedback Loops through User Practice

With regard to digital language technologies and the social biases they may produce, it is furthermore relevant to understand and explore user practices. There are many different approaches to studying user experiences of human-machine interaction, and the field is growing (e.g., Jones et al. 2015; Porcheron et al. 2018; Lind 2021). Let me here quote from a qualitative interview study (Leblebici 2021) in which six multilingual users of voice-controlled digital assistants in diasporic settings were invited to report on their experiences with the tools and on their language choices. The interviewees in question were of Turkish-speaking origin and had moved to Germany within the last 10 years. As an agglutinative language, Turkish is structurally different from English, and this is probably one reason why Amazon’s Alexa, for example, is at the moment of writing this text not available in Turkish (in contrast to Apple’s Siri). Those interviewed for the study thus frequently reported on using voice assistants in languages other than Turkish and above all in English. Some of their reasons can be found in transcripts 2 and 3:

Transcript 2

I also think that (.) original (.)
 English is more comfortable for me
 because e (.) it is a world language
 and when an update comes
 English gives the best opportunities
 and then German
 The last one would be Turkish, it feels like
 Because when I first listened to Siri in Turkish, I was like what is this?
 I mean, it is not (.) natural at all

The dominance of English in the field of language technologies implies that many tools function better if used in English. Therefore, the user reports that English is “more comfortable,” certainly referring to the fact that the tools work better if used in English and also that updates appear first for the English language and earlier for the German language than for Turkish (German is structurally similar to English, and the speaker has access to it since Germany is his place of residence). It is furthermore interesting to observe that the user perceives English to be the “original” language of these digital devices and that using Siri in Turkish is not “natural.” The argumentation is thus not only of an instrumental kind but is also related to cultural constructions of belonging. Overall, we here observe the co-construction of a transnational sociolinguistic economy (see “Theoretical Background” above; Blommaert et al. 2005) in which English has the highest rank.

Similarly, in the following transcript, the interviewee shows to the interviewer how Apple’s Siri sounds in different languages and associates its use with a cultural space that apparently does not fit with the language Turkish:

Transcript 3

Siri [in Turkish, male voice]	I am Siri, your virtual assistant
Siri [in Turkish, female voice]	I am Siri, your virtual assistant [. . .]
	I, it sounds very very very (.) perverted, really Turkish @ it doesn’t sound natural at all. Or maybe it’s because it’s my mother tongue, so that I can criticize it like this, but
Siri [in British English, male voice]	I am Siri, your virtual assistant
	I, this sounds a lot more to me okay we are in a different world, you know we are on Star Trek I can talk with it but when I speak Turkish I say you are Apple, why are you speaking to me in Turkish?

The user demonstrates how Siri sounds in Turkish and then in British English. Again, we here find the argument that Siri in Turkish is not “natural,” and it is even described as “perverted.” The interviewee also reflects on the fact that it could be the case that he is more critical of language forms in Turkish, as this is his first language, so that he is more sensitive to potential deviations. However, it is unlikely that there are incorrect grammatical forms or intonation patterns in

the very first sentence that the device has been produced to generate when switching it on. The perception that Siri in Turkish sounds “perverted” is most likely more attributable to the cultural discourses with which the speaker associates the respective languages. As a matter of fact, during the passage, he switches the device into its British English version and explains that for him this is culturally appropriate. Using a digital voice assistant is for him part of “a different world,” which he associates with particular productions of popular sci-fi culture (*Star Trek*) and with the company Apple. It seems to be irrelevant that neither *Star Trek* nor Apple derive from the UK—the interviewee associates British English with the general anglophone world. The cultural space that is here understood to be indexically associated with a prestige-loaded and standard form of English is not of a national and territorial kind; rather, it is a deterritorial cultural context that is framed in anglophone popular culture and the activities of a US-based globally active company. The traditional indexical links between language and national public space may still be applicable to the Turkish language but not, in this case, to English.

Even though the data corpus of this study is small, it is likely that the experiences reported on are not unusual and may be found in larger sections of the user population of digital voice assistants. The tools, first, work better, are more convenient to use, and are always up to date in English, more so than in any other language. At the same time, the entire cultural complex of AI and advanced digital technology is culturally associated with the anglophone world. Both aspects enforce the use of English on sides of users of language technologies. As user practices of this type of human-machine interaction are recorded and fed into the database of the companies’ servers, the data corpus for English is constantly growing. At the same time, the data base for a language like Turkish—which is structurally different from English and not associated with digital popular culture—is growing to a lesser extent. In this way, the tools eventually will work even better with English. We can assume that a kind of “feedback loop” develops, those language practices that are more frequent are more easily decoded by digital devices, which is why users adapt their language in order to use the tools efficiently. Forms that were frequent in the first place become even more frequent. The trends toward homogenization observed above are thus enforced also through user practices.

Discussion and Conclusion

We have learned from the above elaborations that language practices in the realm of commercial digital language technologies may have specific effects on language ideological discourses and practices. First of all, we can say that the monolingual

and standard ideologies associated with the public regimes and the material technologies of the modern nation-state do not disappear but are coded into digital devices, as data corpora are produced on grounds of data that are understood to represent one standardized language. In this sense, digital language technologies reproduce “languages” as normed and bounded entities. At the same time, depending on the availability of machine-readable data, discursive constructions of languages are hierarchically ordered. Because of digital language technology culture, there is thus a danger of reproducing language hierarchies that have been developed in colonial culture, with European-derived standard languages at the top. In addition, social indexical hierarchies relating to multilingual or nonstandard practice that stigmatize these as unusual or “wrong” are reproduced. This does not necessarily happen because of a desire of actors in the digital language industry to enforce social and sociolinguistic hierarchies but is coproduced by the affordances and materialities of digital language technology, as machine-learning tools require a predefined data set to be trained and amplify what is dominant. And yet, these tools have not been developed in a social void. Their functioning depends on the commercial intentions of companies—digital language technologies are typically designed to make money and stabilize the position of power of those who produce them. The observation that resources and users of resources that already enjoy a high prestige become even more privileged through capitalist desires in AI machine learning is not unique to language. It is important to note that the entire industry of digital AI systems is entangled in global socio-political power hierarchies: “AI systems are built to see and intervene in the world in ways that primarily benefit the states, institutions, and corporations that they serve. In this sense, AI systems are expressions of power that emerge from wider economic and political forces, created to increase profits and centralize control for those who wield them” (Crawford 2021, 211). One conclusion that thus can be drawn from the above elaborations is that commercial language technologies entail a threat for language practices that are not dominant on a global level while those linguistic resources that are dominant on a global scale may become even more dominant in the future—above all, this pertains to resources that are currently classified as English.

At the same time, one can attest that there is no reproduction of an idea of the social world as consisting of nations that are “normally” monolingual. In contrast to earlier language technologies of the nation-state (e.g., print, mass schooling, governmental language policing), the aim is not to homogenize language in order to create a linguistically homogenous national population. It is difficult to estimate what this may mean for newly forming types of public space and for

community formation. Societies are still often classified as “normally” being territorially based and monolingual (think of the fact that digital tools are designed to automatically detect the national-territorial space they are used in and adapt their settings accordingly), and yet the digital language industry produces a global digital public with a kind of Matthew’s Effect principle.⁸

There is no intention of educating speakers—typically referred to as “users” in these contexts—to speak “correctly” or making them conform to specific community- or territory-based language ideals. Rather, the tools are developed as a service product and are supposed to function efficiently. This means that it is not deemed necessary that they strictly and meticulously follow particular language rules. As long as the tools are practical, users are satisfied, and the technologies are profitable, the performance of such tools is regarded as adequate. Generally, there is a strong orientation toward actual language use—where “use” is understood as the production of specific types of language data. As they are currently being produced for a wider array of languages, language technologies therefore in a certain sense support multilingualism. This is, however, hindered by the data problem and the capitalist logics where this is realized only if it contributes to economic gains. Again, this has the overall effect that “the voices of people most likely to hew to a hegemonic viewpoint are also more likely to be retained” (Bender et al. 2021, 613). To what extent the orientation toward user data entails emancipatory potentials or trends toward linguistic diversification will depend on the aims of those who design (and pay for) the development of tools. Currently, it seems that trends toward homogenization are dominant.

I conclude with thoughts on the potential epistemological effects of language technologies on what we believe language to be. In summary, besides the amplification of sociolinguistic hierarchies, there is a trend toward understanding languages not as standardized voices from a national nowhere but as globally available and machine-readable data sets that can be classified statistically. Speakers are understood as “users,” and the patterns that they produce frequently and in ways that can be collected digitally on servers of a few globally acting companies are eventually understood as “correct.” It is the task of future research to scrutinize what this will imply for the enregisterment of language norms and for the discursive-technological construction of language boundaries, as well as what this entails for community formation and for the reproduction of power structures in society.

8. “For to him who has will more be given; and from him who has not, even what he has will be taken away” (Mark 2:25 [RSV]).

References

- Agha, Asif. 2003. "The Social Life of Cultural Value." *Language and Communication* 23:231–73.
- Anderson, Benedict. 1985. *Imagined Communities*. London: Verso.
- Asscher, Omri, and Ella Glikson. 2021. "Human Evaluations of Machine Translation in an Ethically Charged Situation." *New Media & Society*. <https://doi.org/10.1177/14614448211018833>.
- Barera, Michael. 2020. "Mind the Gap: Addressing Structural Equity and Inclusion on Wikipedia." University of Texas Arlington, Libraries Research Commons. <http://hdl.handle.net/10106/29572>.
- Bender, Emily, and Alexander Koller. 2020. "Climbing towards NLU: On Meaning, Form, and Understanding in the Age of Data." *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 5185–98. <https://aclanthology.org/2020.acl-main.463.pdf>.
- Bender, Emily, Angelina McMillan-Major, Timnit Gebru, and Shmargaret Shmitchell. 2021. "On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?" FAccT '21, March 3–10 (virtual), Canada. <https://doi.org/10.1145/3442188.3445922>.
- Blommaert, Jan, James Collins, and Stef Slembrouck. 2005. "Spaces of Multilingualism." *Language and Communication* 25:197–216.
- Bolukbasi, Tolga, Kai-Wei Chang, James Zou, Venkatesh Saligrama, and Adam Kalai. 2016. "Man Is to Computer Programmer as Woman Is to Homemaker? Debiasing Word Embeddings." 30th Conference on Neural Information Processing Systems (NIPS 2016), Barcelona. <https://papers.nips.cc/paper/2016/file/a486cd07e4ac3d270571622f4f316ec5-Paper.pdf>.
- Bommasani, Rishi, Percy Liang, et al. 2021. "On the Opportunities and Risks of Foundation Models." Center for Research on Foundation Models, Stanford University. <https://arxiv.org/abs/2108.07258>.
- Bourdieu, Pierre. 1982. *Ce que parler veut dire: L'économie des échanges linguistiques*. Paris: Fayard.
- Cowley, Stephen. 2011. "Distributed Language." In *Distributed Language*, edited by Stephen Cowley, 1–14. Amsterdam: John Benjamins.
- Crawford, Kate. 2017. "The Trouble with Bias." NIPS 2017 keynote address (video). Artificial Intelligence Channel. https://youtu.be/fMym_BKWQzk.
- . 2021. *Atlas of AI*. New Haven, CT: Yale University Press.
- Devlin, Jacob, and Ming-Wei Chang. 2018. "Open Sourcing BERT: State-of-the-Art Pre-training for Natural Language Processing." Google AI Blog. <https://ai.googleblog.com/2018/11/open-sourcing-bert-state-of-art-pre.html>.
- Dodge, Jesse, 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 Colossal Clean Crawled Corpus." <https://arxiv.org/pdf/2104.08758.pdf>.
- Eisenstein, Elizabeth L. 1979. *The Printing Press as an Agent of Change: Communications and Cultural Transformations in Early-Modern Europe*. Cambridge: Cambridge University Press.
- Errington, Joseph. 2008. *Linguistics in a Colonial World: A Story of Language, Meaning and Power*. Malden, MA: Blackwell.
- French, R., and K. T. Kernan. 1981. "Art and Artifice in Belizean Creole." *American Ethnologist* 8:238–58.

- Gal, Susan, and Judith T. Irvine. 2019. *Signs of Difference: Language and Ideology in Social Life*. Cambridge: Cambridge University Press.
- Gal, Susan, and Kathryn A. Woolard. 2001. *Languages and Publics: The Making of Authority*. Manchester: St. Jerome.
- Gehman, Samuel, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A. Smith. 2020. "REALTOXICITYPROMPTS: Evaluating Neural Toxic Degeneration in Language Models." *Findings of the Association for Computational Linguistics, EMNLP 2020*, 3356–69. <https://arxiv.org/pdf/2009.11462.pdf>.
- Gershon, Ilana. 2010. "Media Ideologies: An Introduction." *Journal of Linguistic Anthropology* 20:283–93.
- Giesecke, Michael. 1991. *Der Buchdruck in der frühen Neuzeit. Eine historische Fallstudie über die Durchsetzung neuer Informations- und Kommunikationstechnologien*. Frankfurt: Suhrkamp.
- Ging, Debbie. 2019. "Alphas, Betas, and Incels: Theorizing the Masculinities of the Manosphere." *Men and Masculinities* 22:638–57.
- Gogolin, Ingrid. 1994. *Der monolinguale Habitus der multilingualen Schule*. Münster: Waxmann.
- Gramling, David. 2016. *The Invention of Monolingualism*. New York: Bloomsbury.
- Hao, Karen. 2020a. "OpenAI Is Giving Microsoft Exclusive Access to Its GPT-3 Language Model." *MIT Technology Review*. <https://www.technologyreview.com/2020/09/23/1008729/openai-is-giving-microsoft-exclusive-access-to-its-gpt-3-language-model/>.
- . 2020b. "We Read the Paper That Forced Timnit Gebru out of Google. Here's What It Says. The Company's Star Ethics Researcher Highlighted the Risks of Large Language Models, Which Are Key to Google's Business." *MIT Technology Review*. <https://www.technologyreview.com/2020/12/04/1013294/google-ai-ethics-research-paper-forced-out-timnit-gebru/>.
- Heller, Monica, and Alexandre Duchêne. 2012. "Pride and Profit: Changing Discourses of Language, Capital and Nation-State." In *Language and Late Capitalism: Pride and Profit*, edited by Monica Heller and Alexandre Duchêne, 1–20. New York: Routledge.
- Heller, Monica, and Bonnie McElhinny. 2017. *Language, Capitalism, Colonialism*. Toronto: University of Toronto Press.
- Hepworth, Shelley. 2021. "People Flocked to Language Apps during the Pandemic—but How Much Can They Actually Teach You?" *The Guardian*. <https://www.theguardian.com/technology/commentisfree/2022/jan/01/people-flocked-to-language-apps-during-the-pandemic-but-how-much-can-they-actually-teach-you>.
- Herring, Susan C. 1996. *Computer-Mediated Communication: Linguistic, Social, and Cross-Cultural Perspectives*. Amsterdam: John Benjamins.
- Hewett, Freya, and Sami Nenko. 2021. "How to Identify Bias in Natural Language Processing." *HIIG – Digital Society Blog*. https://www.hiig.de/en/bias-in-natural-language-processing/?utm_source=mailpoet&utm_medium=email&utm_campaign=Monthly+Digest+November.
- Heyd, Theresa, and Britta Schneider. 2019. "The Sociolinguistics of Late Modern Publics." *Journal of Sociolinguistics* 23:435–49.
- Hootsuite. 2021. "Digital 2021: Global Overview Report." https://hootsuite.widen.net/s/zcdrtxwczn/digital2021_globalreport_en.
- Hopper, Paul. 1998. "Emergent Grammar." In *The New Psychology of Language*, edited by Michael Tomasello, 155–75. Mahwah, NJ: Lawrence Erlbaum.

- Irvine, Judith T. 2001. "Style' as Distinctiveness: The Culture and Ideology of Linguistic Differentiation." In *Style and Sociolinguistic Variation*, edited by Penelope Eckert and John R. Rickford, 21–43. Cambridge: Cambridge University Press.
- Irvine, Judith T., and Susan Gal. 2009. "Language Ideology and Linguistic Differentiation." In *Linguistic Anthropology: A Reader*, edited by Alessandro Duranti, 402–34. Oxford: Wiley Blackwell.
- Jones, Rodney H., Alice Chik, and Christoph A. Hafner. 2015. *Discourse and Digital Practices*. Milton Park: Routledge.
- Joseph, John E., and Talbot J. Taylor, eds. 1990. *Ideologies of Language*. London: Routledge.
- Jurafsky, Daniel, and James H. Martin. 2021. *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*. Unpublished manuscript, last modified December 29, 2021. https://web.stanford.edu/~jurafsky/slp3/ed3book_dec292021.pdf.
- Koenecke, Allison, Andrew Namb, Emily Lake, Joe Nudell, Minnie Quartey, Zion Mengesha, Connor Toupsc, John R. Rickford, Dan Jurafsky, and Sharad Goeld. 2020. "Racial Disparities in Automated Speech Recognition." *Proceedings of the National Academy of Sciences* 117:7684–89. <https://www.pnas.org/doi/epdf/10.1073/pnas.1915768117>.
- Kreutzer, Julia, Isaac Caswell, and Lisa Wang. 2021. "Quality at a Glance: An Audit of Web-Crawled Multilingual Datasets." *Transactions of the Association for Computational Linguistics*. <https://arxiv.org/pdf/2103.12028.pdf>.
- Kroskrity, Paul V. 2000. "Regimenting Languages: Language Ideological Perspectives." In *Regimes of language. Ideologies, politics, and identities*, edited by Paul V. Kroskrity, 1–34. Santa Fe: School of American Research Press.
- Lauscher, Anne, Vinit Ravishankar, Ivan Vulić, and Goran Glavaš. 2020. "From Zero to Hero, on the Limitations of Zero-Shot Language Transfer with Multilingual Transformers." *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 4483–99. <https://aclanthology.org/2020.emnlp-main.363.pdf>.
- Leblebici, Didem. 2021. "Language Ideologies in Human-Machine Interaction. A Qualitative Study with Voice Assistant Users." MA thesis, Europa-Universität Viadrina.
- Lind, Miriam. 2021. "Alexa, 3, Sprachassistentin, hat die Religion für sich entdeckt'—Eine framesemantische Korpusstudie zur Anthropomorphisierung von Sprachassistenten." In *Mensch – Tier – Maschine. Sprachliche Praktiken an und jenseits der Außengrenze des Humanen*, edited by Miriam Lind, 347–70. Bielefeld: transcript.
- Linell, Per. 2005. *The Written Language Bias in Linguistics: Its Nature, Origins and Transformations*. London: Routledge.
- Lippi-Green, Rosina. 2012. *English with an Accent*. New York: Routledge.
- Love, Nigel. 2017. "On Linguaging and Languages." *Language Sciences* 61:113–47.
- Madsen, Lian Malai, Martha Sif Karrebæk, and Janus Spindler Møller, eds. 2016. *Everyday Linguaging: Collaborative Research on the Language Use of Children and Youth*. Berlin: de Gruyter.
- McShane, Marjorie, and Sergei Nirenburg. 2021. *Linguistics for the Age of AI*. Cambridge, MA: MIT Press.
- Mitchell, Margaret. 2021. "Cementing a Foundation of Inequality in AI." Presentation, Workshop on Foundation Models, Stanford University. <https://crfm.stanford.edu/workshop.html>.

- Nayak, Pandu. 2019. "Understanding Searches Better than Ever Before." Google Blog. <https://blog.google/products/search/search-language-understanding-bert/>.
- NLLB Team et al. 2022. "No Language Left Behind: Scaling Human-Centered Machine Translation." <https://arxiv.org/abs/2207.04672>.
- Ong, Walter J. 1982. *Orality and Literacy: The Technologizing of the Word*. London: Routledge.
- Pakalski, Ingo. 2009. "Linguee: Suchmaschine für Übersetzungen." *Golem*. <https://www.golem.de/0904/66396.html>.
- Pennycook, Alastair. 2004. "Performativity and Language Studies." *Critical Inquiry in Language Studies* 1:1–19.
- . 2018. *Posthumanist Applied Linguistics*. London: Routledge.
- Porcheron, Martin, Joel E. Fischer, Stuart Reeves, and Sarah Sharples. 2018. "Voice Interfaces in Everyday Life." *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems, April 2018*. <https://www.semanticscholar.org/paper/Voice-Interfaces-in-Everyday-Life-Porcheron-Fischer/4ace92c22a895d5e23e58de8d738df8e500d8d79>.
- Randhawa, Gurdeeshpal, Mariella Ferreyra, Rukhsana Ahmed, Omar Ezzat, and Kevin Pottie. 2013. "Using Machine Translation in Clinical Practice." *Canadian Family Physician* 59:382–83.
- Savoldi, Beatrice, Marco Gaido, Luisa Bentivogli, Matteo Negri, and Marco Turchi. 2021. "Gender Bias in Machine Translation." *Transactions of the Association for Computational Linguistics* 9:845–74. <https://aclanthology.org/2021.tacl-1.51.pdf>.
- Schneider, Britta. 2019. "Methodological Nationalism in Linguistics." *Language Sciences* 76:101–69.
- . 2021a. "Creole Prestige beyond Modernism and Methodological Nationalism: Multiplex Patterns, Simultaneity and Non-closure in the Sociolinguistic Economy of Belize." *Journal of Pidgin and Creole Languages* 36:12–45.
- . 2021b. "Von Gutenberg zu Alexa—Posthumanistische Perspektiven auf Sprachideologie." In *Mensch – Tier – Maschine*, edited by Miriam Schmidt-Jüngst, 327–346. Bielefeld: Transcript.
- . Forthcoming. *Liquid Languages—Constructing Language in Late Modern Cultures of Diffusion*. Cambridge: Cambridge University Press.
- Stefanowitsch, Anatol. 2020. *Corpus Linguistics: A Guide to the Methodology*. Berlin: Language Science Press.
- Sun, Tony, Andrew Gaut, Shirlyn Tang, Yuxin Huang, Mai El Sherief, Jieyu Zhao, Diba Mirza, Elizabeth Belding, Kai-Wei Chang, and William Yang Wang. 2019. "Mitigating Gender Bias in Natural Language Processing: Literature Review." *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 1630–40. <https://www.aclweb.org/anthology/P19-1159.pdf>
- Tatman, Rachael. 2016. "Google's Speech Recognition Has a Gender Bias." Making Noise and Hearing Things. <https://makingnoiseandhearingthings.com/2016/07/12/googles-speech-recognition-has-a-gender-bias/>.
- UNESCO. 2019. "I'd Blush If I Could: Closing Gender Divides in Digital Skills through Education." <https://unesdoc.unesco.org/ark:/48223/pf0000367416.page=1>.

- Vanmassenhove, Eva, Dimitar Shterionov, and Andy Way. 2019. "Lost in Translation: Loss and Decay of Linguistic Richness in Machine Translation." *Proceedings of MT Summit XVII* 1:222–32. <https://arxiv.org/pdf/1906.12068.pdf>.
- Virtanen, Antti, Jenna Kanerva, Rami Ilo, Jouni Luoma, Juhani Luotolahti, Tapio Salakoski, Filip Ginter, and Sampo Pyysalo. 2019. "Multilingual Is Not Enough: BERT for Finnish." <https://arxiv.org/pdf/1912.07076.pdf>.
- Wimmer, Andreas, and Nina Glick Schiller. 2002. "Methodological Nationalism and Beyond: Nation-State Building, Migration and the Social Sciences." *Global Networks* 2 (4): 301–34.
- Woolard, Kathryn A. 1998. "Introduction: Language Ideology as Field of Inquiry." In *Language Ideologies: Practice and Theory*, edited by Bambi B. Schieffelin, Kathryn A. Woolard, and Paul V. Kroskrity, 3–47. Oxford: Oxford University Press.
- Xu, Albert, Eshaan Pathak, Eric Wallace, Suchin Gururangan, Maarten Sap, and Dan Klein. 2021. "Detoxifying Language Models Risks Marginalizing Minority Voices." *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics, Human Language Technologies*, 2390–97. <https://doi.org/10.18653/v1/2021.naacl-main.190>.
- Zhao, Jieyu, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2017. "Men Also Like Shopping: Reducing Gender Bias Amplification Using Corpus-Level Constraints." *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, 2979–89. <https://arxiv.org/pdf/1707.09457.pdf>.
- Zijlstra, Judith, and Ilse van Liempt. 2017. "Smart(phone) Travelling: Understanding the Use and Impact of Mobile Technology on Irregular Migration Journeys." *International Journal of Migration and Border Studies* 3:174–91.
- Zweig, Katharina. 2019. *Ein Algorithmus hat kein Taktgefühl*. Munich: Heyne.