

Proceeding Paper

Towards an Android Linguistics: Pragmatics, Reflection, and Creativity in Machine Language [†]

Evan Donahue

Tokyo College, University of Tokyo, Tokyo 113-8654, Japan; evan.donahue@tc.u-tokyo.ac.jp

[†] Presented at the Conference on Theoretical and Foundational Problems in Information Studies, IS4SI Summit 2021, online, 12–19 September 2021.

Abstract: Contemporary natural language processing (NLP) emphasizes comparing machine language performances to standards defined by static corpora of human text. However, despite some successes, current models remain weak in areas such as pragmatics. Using scholarship on neosentience as a foundation, this essay proposes an alternative view of machine language that emphasizes generativity rather than stasis and draws on historical work on computational reflection in artificial intelligence to outline an alternative architecture for conversational systems. It concludes by proposing an “android linguistics” that takes human-machine linguistic communication as its object of study.

Keywords: artificial intelligence; natural language processing; cybernetics; neosentience; pragmatics; speech acts; reflection; self-reference

1. Introduction

Is language a noun or a verb? Contemporary NLP views language as a static thing—a standard against which machine performance can be measured; language as an encyclopedia that always contains the appropriate response if only it can be made large enough. The cybernetically-inflected view of language offered by neosentience and recombinant poetics, by contrast, views language use as an act of creation, and so offers an alternative perspective on designing conversational machines [1].

Expanding on the neosentient view of language as an act of creation, I suggest that machine language researchers’ perennial difficulties with pragmatics—with accounting for the influence of context on interpretation—cannot be solved with scale, but instead require a different, self-reflexive architecture. Pragmatics, I contend, is inseparable from self-reflection. In this essay, I illustrate the need for such an architecture, suggest some preliminary requirements for its design, and call for a deeper consideration of language as a phenomenon that is not limited to human speakers, with implications for how we define linguistic performance for machines.

2. Pragmatics as Self-Reflection

AI researcher Douglas Lenat, although supportive of scaling up language systems, tempers his support with a note of caution. In observing how the slightest change in the placement of a comma can radically alter how a sentence is interpreted by inviting in a host of assumptions about the context and purpose for which it was written, he writes despairingly that, “There’s always this annoying residue of pragmatics, which ends up being the lower 99% of the iceberg . . . lurking in the empty spaces around the letters, words, and sentences.” [2] (p. 2). Like dark matter, pragmatics constitutes, for Lenat, the vast and unseen majority of the reality of language.

That such a reality should so trouble Lenat speaks to the seeming hopelessness of fitting the inexhaustible totality of language into a finite computer system. Through the lens



Citation: Donahue, E. Towards an Android Linguistics: Pragmatics, Reflection, and Creativity in Machine Language. *Proceedings* **2022**, *81*, 156. <https://doi.org/10.3390/proceedings2022081156>

Academic Editor: Mark Burgin

Published: 19 July 2022

Publisher’s Note: MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Copyright: © 2022 by the author. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

of recombinant poetics, however, the infinite creativity of language is precisely what makes such a project possible. Shifting the focus of machine language research from attempting to build machines that know in advance what words mean to negotiating that meaning within the context of their creation, however, requires a shift in our understanding of language itself in the machinic context.

If generativity, rather than stasis, is taken to be the nature of language, then language systems must be designed less to know vast quantities of language and more to be reflexive in its performance. As Seaman [3] argues, agents in conversation generate meaning not by dredging it whole cloth out of a database, but rather by attending to their own and others' efforts to bring it into being with the words, gestures, and objects available to them. By existing and conversing, a conversational agent creates new utterances with new meanings not contained in any dataset. The one thing a dataset cannot contain, no matter how vast, is its own model. It is only by understanding how it itself and its own speech are implicated in the process of making meaning that the agent can make sense of these new meanings made.

Self-reflexivity has a long history in AI. The idea that a computer could program itself and thereby exceed the limitations imposed by its programmers has attained, at times, an almost mystical quality in the field. Seaman [4] focuses on a more contained form of self-reference in language through a discussion of Givón's work on pragmatics. As Givón [5] argues, it is language's ability to refer to itself that underlies much of its function. Correspondingly, I contend that the central architectural principle that must underwrite any attempt at building a pragmatically competent conversational system is that its objects of discourse—the system's linguistic inputs and outputs—must themselves be first-class objects within that discourse. It must be possible, using reflection, to move fluidly between considering an utterance as a meaning-bearing element of conversation and viewing it as an object about which other meaning-bearing utterances can be made. From this principle, many key points of computational pragmatics follow.

3. Pragmatics in AI

AI inherited from philosophy a Russellian epistemology in which sentences map neatly onto reality. Words refer to objects in reality, statements are true or false depending on whether they correspond to this reality, and questions can be answered correctly or incorrectly depending on whether their correspondence with that reality aligns with those of their answers. Consequently, self-reflexive statements that frustrate efforts to parse them as either true or false, such as "this statement is false," have haunted the field's efforts to conceptualize meaning in machine language.

A pragmatic explanation for this phenomenon is that the human escapes to another level of interpretation. While capable of playing the language game of true and false, the human is also capable of recognizing such a self-referential statement as an utterance made by a logician attempting to prove a point about a particular language game, as in this very essay, even though the logician making the point appears nowhere in the statement. AI researchers have long intuited this change of levels, as when Marvin Minsky observes that the usual reaction to such "liar's paradoxes" is not to become trapped in endless deduction, but rather to laugh and reject the language game of true and false altogether [6] (p. 139). In an instant, the paradoxical statement is reduced to a sequence of sounds that have no meaning and therefore present no paradox.

Richard Weyhrauch, whose work with Carolyn Talcott and others on reflective computer architectures drew inspiration from self-reference in language, remarks that much of what we would like to talk to machines about is language itself [7] (p. 155). Although he does not invoke pragmatics explicitly, his illustrations of the contextual nature of meaning align closely with its intuitions. As he observes, the truth value of even the most seemingly inescapable universal proposition, such as " $2 + 2 = 4$ ", becomes suspect when growled angrily across a barroom counter by the head of a rowdy biker gang [8] (p. 14). Even the

staunchest defender of universal meaning would have to think twice about whether they were in fact caught up in some larger language game.

Weyhrauch's example highlights the limitations of any attempt to evaluate machine language by sorting the world into clear categories against which to test the performance of AI systems. Even the most intuitively unambiguous utterances can always have their meanings utterly displaced in the context of the right language game. As in a spy novel, it is always possible for a seemingly innocuous utterance to signify the transmission of a secret code. Whether such language games represent exceptional circumstances that can be safely ignored until the basics of machine language have been figured out, or whether they point to more fundamental mechanisms, the omission of which will doom the whole project, is the question that must be addressed and the point on which poetic and encyclopedic theories of machine language most differ.

This difference is apparent even in simple sentences. A recent project by researchers at Facebook offers one perspective on the nature of machine language by attempting to enumerate a set of tests of fundamental reasoning abilities a machine must possess to answer simple questions [9]. These tests include an understanding of true and false, of elements and sets, of logical conjunction and disjunction, of negation, and of numbers. Taken together, these tests assess an artificial agent's ability to emulate a formal reasoning system. While not necessarily an unreasonable ability for an artificial agent to possess, from the point of view of pragmatics it omits a more fundamental level of analysis. Namely, it rests on assumption of an unproblematic mapping from sentences of English to sentences of logic. The sentences themselves are not objects of discourse within the proposed test environment to be reasoned about.

The importance of being able to speak at this meta-level becomes evident when considering Searle's famous question, "can you pass the salt?" [10]. A system trained to respond to questions may go wrong if it assumes that the question is a request for information about the machine's capabilities. The purpose of this question is of course to illustrate that what appears to be a question about passing salt may, in fact, be a request to effect the passing. Then again, in another context, the question may not even be asked in good faith, but rather as part of a test of the system's understanding of pragmatics, perhaps designed in response to an essay such as this one, in which case the "correct" response would depend on the level at which one interprets the question.

Any test of an AI's linguistic competence inevitably rests on a set of assumptions on the part of the researcher about the context in which the language is to be interpreted and what correctness or incorrectness looks like in such a context. From the system's point of view, however, it is never told it is being evaluated. The parameters of the evaluation are never explained to the system in language, in part because, in most cases, it lacks the representational machinery to even recognize the discursive elements of a language test as entities with which it shares a reality. It is in the position of the chess-playing automaton that does not know it is playing chess; it may play a fair game, but a human operator must carry it in, set it up next to the board, and face it in the right direction. As long as this remains the default experimental paradigm in AI, it is likely that special-purpose systems that exhibit virtuoso linguistic performances without actually being communicatively competent will continue to be the norm.

Searle's question underscores that the sign is not just arbitrary in the Saussurean sense that there is no necessary relationship between the form of the signifier and the signified it has historically come to represent, but that it is radically arbitrary. No matter the history of the signifier or the conventions of its interpretation in other contexts, it is always possible to subvert that history with the appropriate language game in the appropriate context. The neosentient view suggests that this subversion happens rapidly and continuously as part of the creative flux of language in practice. Starting from this premise, any attempt to learn the correspondence between the form of the utterance and its meaning is to build castles on sand. What is needed is an approach that underlies and precedes the performances the

Facebook researchers sought to measure—one that situates the act of interpretation as prior to the consideration of form.

4. Towards an Android Linguistics

Contemporary NLP depends heavily on quantifying the correctness of language. However, such quantification is a compromise with which perhaps no one in the field has ever been truly happy. Moreover, the field's history offers a wealth of examples of alternative conceptions of the work of studying machine language that may offer inspiration. In particular, with respect to the question of pragmatics, a body of work responding to the speech act theory of Searle, Austin, and others, that emerged in the late 1970s offers a distinctive approach to conceptualizing machine language [11].

While a full review of this literature is beyond the scope of this essay, one key insight that offers concrete guidance on the design and evaluation of contemporary systems is that form must be held entirely apart from meaning. "Can you pass the salt" should not immediately resolve into either a request for information or a request for action. Rather, the words must be evaluated based on what is known about the speaker and the environment and what they reveal the beliefs, intentions, and desires of that speaker. Once it is determined that the speaker desires the salt and believes the system can obtain it for them, then the system can exercise its agency by passing the salt or withholding it, by speaking or remaining silent. This action is undertaken not on the basis of the force of the utterance itself but on the basis of what the utterance has revealed about its speaker. Moreover, that revelation is a product not only of interpretation but of conversational interaction—of the poiesis of meaning—enabled by the explicit self-representation of the field of discourse.

Several important conversational behaviors are enabled by representing discursive objects explicitly alongside other domain objects in a reflective manner. These might stand alongside the behaviors outlined by the Facebook researchers as heuristics with which to probe whether a conversational system is representationally sufficient to capture, even in principle, the important pragmatic dimensions of communication.

The most basic requirement for a pragmatic conversational system is the ability to refer to the words of the conversation itself. Many current systems, if they had never encountered the word "salt", would simply fail to process Searle's question, rather than being able to formulate a question of which the word itself was the subject. Even asking for clarification of how a known word or construction is being used creatively in a new context requires the ability to refer to the word in question.

The second property is the ability to refer to the system's own interpretations of prior utterances as first-class discursive objects. Inevitably, in conversation, misunderstandings will arise. Repairing them requires the ability to discuss what was previously understood. Explaining to the system that Searle's question was intended as a request for salt rather than for information necessitates the ability to ground references to the erroneous interpretation, even if such an interpretation is only a discursive fiction rather than a technical reality in the underlying system.

Repairing the conversation, in turn, is not a simple matter of correcting a previous misunderstanding to what it originally should have been, but rather of determining how that meaning has been changed in light of statements made and actions taken on the basis of that misunderstanding. In order to ask how to proceed, and whether the speaker still wants the salt, the system must be able to discuss its plans and adjust them based on the conversation that follows. Discussion of such plans, in turn, requires the ability to refer to the future worlds that such plans might bring about. Although the mental or physical reality of plans and possible worlds has been a longstanding point of debate within AI, their reality as objects of the discursive universe requires no underpinning outside of language to validate their utility.

Finally, the ability to project hypothetical possible worlds invites consideration of the ability to project fictional ones—worlds that do not exist, and moreover that could not in principle exist, or that may not even make complete sense. If an AI system were to read

a work of fiction, it should be able to do so while neither confusing the fictional world with the real one nor keeping the two so wholly separate that the linguistic and cultural materials with which fictions are constructed become unintelligible. Moreover, as work on narrative theory and story worlds seems to hint, it may be worth viewing reality itself as a collection of fictions we tell and retell, inventing in the process that which can never be captured by a single totalizing view of language [12].

5. Conclusions

To treat language as a static whole, rather than a dynamic process in which the researchers themselves are implicated, is, to paraphrase Givón, to rescue the study of machine language by abandoning its purpose [5] (p. 4). Indeed, any artificially intelligent system not intelligent enough to know it is being tested is unworthy of the name. Centering context in communication promises more conversationally capable machines, and this essay has offered, as a starting point, the narrow technical requirement that objects of discourse should have first-class representations in the system. More speculatively, because consideration of context calls attention to speakers and their standpoints, it is perhaps worth contemplating whether interrogating how it normalizes certain language as “correct” might force the field to confront the assumptions about race, gender, and disability inscribed on its datasets and artifacts that scholars have documented for decades. This attention to the language of machines and its place in broader human language communities as a basis for the design of socio-technical systems is what I refer to as an android linguistics.

Funding: This research received no external funding.

Institutional Review Board Statement: Not applicable.

Informed Consent Statement: Not applicable.

Data Availability Statement: Not applicable.

Conflicts of Interest: The author declares no conflict of interest.

References

1. Seaman, B. Generative Works: From Recombinant Poetics to Recombinant Informatics. In Proceedings of the 2014 International Conference on Cyberworlds, Santander, Spain, 6–8 October 2014; IEEE: Piscataway, NJ, USA; 2014. [CrossRef]
2. Lenat, D. Sometimes the Veneer of Intelligence Is Not Enough. *Cognitive World*. December 2018. Available online: <https://cognitiveworld.com/articles/sometimes-veneer-intelligence-not-enough> (accessed on 12 May 2022).
3. Seaman, B. Towards a Dynamic Heterarchical Ecology of Conversations. Heinz von Foerster Lecture. Heinz von Foerster Society in Cooperation with the Institute for Contemporary History. November 2017. Available online: https://billseaman.com/Papers/Towards%20A%20Dynamic%20Heterarchical%20Ecology%20Of%20Conversations_Final_2%20copy.pdf (accessed on 12 May 2022).
4. Seaman, B. Neosentience and the Abstraction of Abstraction. *Systems. Connecting Matter, Life, Culture and Technology*. 2013. Volume 1, p. 51. Available online: <https://billseaman.com/Papers/Abstraction%20of%20AbstractionFinal%20Edit.pdf> (accessed on 12 May 2022).
5. Givón, T. *Mind, Code and Context: Essays in Pragmatics*; Psychology Press: Hove, UK, 2014.
6. Minsky, M. A Framework for Representing Knowledge. In *Mind Design 2: Philosophy, Psychology, Artificial Intelligence*; John, H., Ed.; MIT Press: Cambridge, UK, 1997; pp. 111–142.
7. Weyhrauch, R.W. Prolegomena to a Theory of Mechanized Formal Reasoning. *Artif. Intell.* **1980**, *13*, 133–170. [CrossRef]
8. Weyhrauch, R. Ideas on Building Conscious Artifacts. The FOL Project. 1991. Available online: <http://www-formal.stanford.edu/FOL/rww-91con-work> (accessed on 12 May 2022).
9. Weston, J.; Bordes, A.; Chopra, S.; Rush, A.M.; Van Merriënboer, B.; Joulin, A.; Mikolov, T. Towards AI-Complete Question Answering: A Set of Prerequisite Toy Tasks. *arXiv* **2015**, arXiv:1502.05698.
10. Searle, John R. Indirect Speech Acts. In *Speech Acts*; Brill: Leiden, The Netherlands, 1975; pp. 59–82.
11. Cohen, P.R.; Perrault, C.R. Elements of a Plan-Based Theory of Speech Acts. *Cogn. Sci.* **1979**, *3*, 177–212. [CrossRef]
12. Turner, M. *The Origin of Ideas: Blending, Creativity, and the Human Spark*; Oxford University Press: Oxford, UK, 2014.