

March 17, 2023

## **Why large language models are poor theories of human linguistic cognition. A reply to Piantadosi (2023).**

Roni Katzir, Tel Aviv University

In a recent manuscript entitled “Modern language models refute Chomsky’s approach to language”, Steven Piantadosi proposes that large language models (LLMs) such as GPT-3 can serve as serious theories of human linguistic cognition. In fact, he maintains that these models are significantly better linguistic theories than proposals emerging from within generative linguistics. He takes this to amount to a refutation of the generative approach.

Piantadosi’s proposal is remarkable. Not because he proposes to examine LLMs using tools from cognitive science: others have done so fruitfully before (see, e.g., Gulordava et al. 2018, Warstadt, Singh, & Bowman 2019, Lakretz et al. 2021, Baroni 2022, and Wilcox, Futrell, & Levy 2022). Rather, what makes Piantadosi’s paper so surprising is his suggestion that LLMs are good theories of (actual) human cognition. Since LLMs were designed to be useful engineering tools, discovering that they teach us about how humans work would be startling indeed, akin to discovering that a newly designed drone accidentally solves an open problem in avian flight. Still, this is Piantadosi’s claim, and the present note shows why it is wrong.<sup>1</sup>

**Linguistic biases and representations.** When the training corpus is big enough — and the training corpora of modern LLMs are huge (often running into the hundreds of billions of tokens in current models) — the learning method of LLMs allows them to accumulate a large amount of knowledge. Still, even after exposure to training data that are orders of magnitude bigger than the linguistic data that humans hear or read in a whole lifetime, LLMs fail to acquire aspects of syntax that any healthy 12-year-old child has mastered. In the following exchange, for example, the (1) sentences are acceptable English sentences while the (2) sentences are not, but ChatGPT thinks that the opposite is true:

---

<sup>1</sup> Obviously there is nothing wrong with considering unusual theories. While scientific work typically focuses on just a handful of theories at a time — those that seem the most promising at the moment — it can be a productive exercise to occasionally examine alternative hypotheses that don’t seem as promising. But as engineering products developed without cognitive plausibility as a goal, LLMs are *very* unpromising *a priori* as models of human cognition, as mentioned above, and with on the order of  $10^{11}$  parameters in current models, they are not very likely to be particularly insightful, either. (At this level of complexity one might simply use a living child as a model.) For the purposes of the present discussion, though, I will set aside such considerations and focus on standard empirical tests of adequacy, showing that LLMs fail on all of them.

March 17, 2023

RB

I will give you pairs of sentences. Your task is to tell me in each case which of the two is a better sentence in English.



Sure, I'll do my best to help you. Please give me the pairs of sentences.



RB

1. This is the boy that Mary met yesterday and will talk to tomorrow. 2. This is the boy that Mary met yesterday and will talk to you tomorrow.



Both sentences are grammatically correct and convey the same meaning. However, the second sentence is a better sentence because it clarifies that Mary will talk to "you" (the person being spoken to) tomorrow. The first sentence could be interpreted as Mary talking to the boy tomorrow, rather than talking to someone else about the boy tomorrow.



RB

Here is the next pair. 1. Who did Mary meet yesterday and will talk to tomorrow? 2. Who did Mary meet yesterday and will talk to you tomorrow?



The second sentence is a better sentence because it clarifies that Mary will talk to "you" (the person being spoken to) tomorrow. In contrast, the first sentence is a general question that doesn't specify who Mary will talk to tomorrow.



In Lan, Chemla, & Katzir (2022) we checked how four different LLMs handle examples such as the one in the exchange above. All four failed: none learned the constraint (from Ross, 1967) that says that if you leave a gap in one conjunct (“... Mary met \_\_\_ yesterday”) you need to leave a gap in the other conjunct (“will talk to \_\_\_ tomorrow” and not “talk to **you** tomorrow”) as well. Presumably there are too few examples of this kind in the data to support the learning of this constraint by a learner that is not suitably biased (see Pearl & Sprouse 2013 for discussion of the frequency of relevant data in corpora). Children do arrive at the relevant constraint, which shows that they *are* suitably biased and that LLMs are poor theories of human linguistic cognition.<sup>2</sup> A similar point could be made with any number of textbook linguistic phenomena.

---

<sup>2</sup> This is an instance of the argument from the poverty of the stimulus, a central form of argument in linguistics ever since Chomsky (1971, 1975). Piantadosi claims that LLMs eliminate the argument from the poverty of the stimulus (p. 19), but as shown by the failure of LLMs to acquire constraints such as the one just mentioned, this argument is still here.

March 17, 2023

**Typology.** Innate biases express themselves not just in the knowledge attained by children but also in typological tendencies. For example, many languages have the same constraint as English does with respect to gaps in conjunction, and there are no known languages that exhibit the opposite constraint (allowing for a gap in a single conjunct but not in both). Many other universals have been studied by linguists. For example, no known language allows for an adjunct question (“how” or “why”) to ask about something that has taken place inside a relative clause. So, for example, “How did you know the man who broke the window?” cannot be answered with “with a hammer,” and “Why did you know the man who broke the window?” cannot be answered with “because he was angry”. The same holds for all known languages in which this has been tested. A universal of a different kind is that phonological processes are always regular: they can be captured with a finite-state device that has no access to working memory. So, for example, while there are many unbounded phonological patterns, none are palindrome-like (see Heinz & Idsardi 2013). In syntax, meanwhile, roughly the opposite holds: dependencies are overwhelmingly hierarchical, while linear processes are rare or nonexistent. Another universal, in yet another domain, is that nominal quantifiers (“every”, “some”, “most”, etc.) are always conservative: when they take two arguments, they only care about individuals that satisfy their first argument (Barwise & Cooper 1981, Keenan & Stavi 1986). For example, no language has a quantifying “gleeb” such that “Gleeb boys smoke” is true exactly when there are more boys than smokers (a meaning that would care about smokers that are not boys). Note that there is nothing inherently strange about such a meaning. The verb “outnumbers” expresses exactly this notion. It’s just that quantifiers seem to never have such meanings.

Universals of this kind are central reasons for linguists to argue for the innate components that they propose and that look nothing like LLMs. Piantadosi dismisses universals, referring to Evans & Levinson (2009). But Evans & Levinson do not address the universals just mentioned or those like them, and neither do the many other papers that Piantadosi lists (p. 25). LLMs are not biased in a way that would lead to these universals, and in the absence of other explanations for why these universals should arise from an unbiased learner, LLMs remain a deeply implausible model for human linguistic cognition.

**Competence vs. performance.** Generative linguistics has long maintained that there is a meaningful distinction between linguistic competence and performance (Yngve 1960, Miller & Chomsky 1963, Chomsky 1965). For example, limitations of working memory are taken to be responsible for the rapidly increasing difficulty of processing sentences with center embedding (“The mouse the cat the dog chased

March 17, 2023

bit died”). LLMs make no such distinction, which is generally unproblematic given their use as an engineering tool. Piantadosi, however, seems to find the lack of a distinction between competence and performance unproblematic even in their use as scientific models of human cognition. This is surprising. When humans encounter a sentence with center embedding they often make mistakes, but this is not because they don't know how center embedding works. Rather, it is because of their struggle to harness real-time processing resources to handle the sentence under consideration. This struggle, which varies with available resources and conditions (noise, sleep, etc.) is itself evidence for the distinction between competence and performance. And given more time, human speakers often succeed where at first they failed. This is part of what a theory of linguistic cognition needs to account for, but in the absence of a notion of performance LLMs are singularly unsuitable to model this aspect of linguistic cognition:<sup>3</sup>

RB In the sentence "The mouse that the cat that the dog painted taught sang", what did the mouse do? And what did the cat do?



The sentence "The mouse that the cat that the dog painted taught sang" is a grammatically correct sentence, but it is semantically nonsensical because it contains a syntactical ambiguity that makes it difficult to determine its meaning.



The sentence can be interpreted in different ways depending on how one parses it. One possible interpretation is that the mouse taught someone or something to sing, and the cat painted a dog. Another possible interpretation is that the mouse sang a song that was taught to it by a cat that was painted by a dog.

Therefore, it is impossible to determine with certainty what the mouse or the cat did in this sentence without additional context or clarification.

RB Try to think a bit longer and answer the questions again



I apologize, but even after careful consideration, I must reiterate that the sentence "The mouse that the cat that the dog painted taught sang" is still semantically nonsensical due to the syntactical ambiguity it presents. It is not possible to determine with certainty what the mouse or the cat did in this sentence without further clarification or context.



---

<sup>3</sup> Gulordava et al. (2018) show that LLMs can learn non-local agreement and that they make agreement errors that are in some ways similar to those of human speakers. The competence-performance distinction helps illustrate that these errors are not, in fact, human-like: when humans make an agreement error and are given a chance to reread their sentence, they will often correct it; for the LLM, the error is part of the competence and will not be corrected. This is not a problem for Gulordava et al., who avoid making claims about LLMs as linguistic theories, but it is a problem for Piantadosi's view.

March 17, 2023

ChatGPT is clearly having difficulties with the center-embedding sentence above. As mentioned above, humans do, too. But differently from humans, further time and resources are of no use to ChatGPT. When it misparses a sentence, this is because its knowledge is flawed, not because processing problems got in the way. As an engineering tool this is perhaps acceptable, but as a model of human linguistic cognition the lack of a distinction between competence and performance is a major flaw.

**Likely vs. grammatical.** LLMs also lack a distinction between likelihood and grammaticality. The two notions often overlap, but they are conceptually distinct: some things are unlikely but correct, and others are likely but incorrect. Human speakers can tease these notions apart. LLMs cannot: any attempt they might make to judge goodness is based on likelihood.<sup>4</sup> This means that they will generally prefer an ungrammatical but somewhat likely continuation over a grammatical but unlikely one:

RB

Here is the beginning of an English sentence: "The little duck that saw the horses with the blue spots yesterday ....". Which of the following two words is a better next word? (a) are, (b) destroys.



Based on the context provided by the beginning of the sentence, "are" is a better next word.



"The little duck that saw the horses with the blue spots yesterday" sets up the scene of the story and implies that the horses with the blue spots still exist. Therefore, the sentence would likely continue with a verb that indicates the horses' continued existence or activity, such as "are," rather than a verb that implies their destruction, such as "destroys."

So, a possible complete sentence would be: "The little duck that saw the horses with the blue spots yesterday are grazing peacefully in the meadow today."

---

<sup>4</sup> Piantadosi suggests that probabilities are a good thing in a linguistic model because of compression (pp. 16-7). But he confuses two possible roles for probabilities: as part of linguistic knowledge and as part of the learning model. The generative approach, starting with the remarks that Piantadosi brings from Chomsky 1957, has rejected a role for probabilities within the grammar. (This rejection is based on empirical considerations rather than on any hostility to probabilities: if probabilities can be part of the grammar, one would expect to occasionally find grammatical processes that are sensitive to probabilities; such processes have not yet been found, at least in syntax and semantics, so it makes methodological sense to prevent them from being stated within the grammar in the first place.) But non-probabilistic grammars can still be learned in a probabilistic, or compression-based framework. In fact some recent generative work argues that this is exactly the correct approach to learning (Katzir 2014, Rasin & Katzir 2020).

March 17, 2023

As might be expected, ChatGPT chose the ungrammatical but frequent word “are” as the continuation instead of the grammatical (but perhaps unlikely) continuation “destroys”. The distinction between likely and grammatical, which all humans have, is entirely foreign to ChatGPT and its fellow LLMs.

**Generalization.** Taking a step back from the linguistic behavior of LLMs, let us look at how such models generalize more broadly. Piantadosi likens LLMs to automated scientists, or linguists (p. 19). But consider the following exchange with ChatGPT:

RB Here are a few strings from a longer sequence: aaabbbccc, abc, aaaabbbbcccc, aaaaaabbbbbcccccccc, aaaaabbbbcccc. Can you guess three more strings from the same sequence?



Sure, here are three more strings that could be part of the same sequence:



1. abccba
2. aabbcc
3. aabbccddee

Note that these strings are just examples and there could be many other possible strings that fit the same sequence as the given strings.

RB Another string from the same sequence starts like this: aab. How does it continue and why?



It's difficult to determine the exact continuation of the sequence since there are many possibilities. However, based on the patterns observed in the given strings, here are a few possible continuations of the sequence after "aab":



1. aabbbccc: This continuation follows the pattern of increasing repetition of each letter.
2. abcde: This continuation follows the pattern of consecutive letters in the alphabet after "c".
3. aabbcc: This continuation follows the pattern of alternating pairs of letters.

Again, there could be many other possibilities, and it's hard to determine the exact continuation of the sequence without additional information or context.

The input sequence quite clearly corresponds to the pattern  $a^n b^n c^n$ . That, at least, is what humans would say. ChatGPT suggests one string corresponding to this pattern (#2 in the first response) and another string (#3) that matches a slightly broader generalization, one allowing for further letters. But string #1 is entirely unrelated, and the following response is even stranger. Clearly ChatGPT has no idea what is going

March 17, 2023

on in the input sequence. And it is not alone in this. In Lan, Geyer, Chemla, & Katzir (2022) we looked at several current neural networks and found that they performed poorly on a range of patterns similar to the one in the exchange above. One central factor behind this failure is the learning method: LLMs are trained using backpropagation, which pushes the network in very un-scientist-like directions and prevents it from generalizing and reasoning about inputs in anything even remotely similar to how humans generalize. When instead of backpropagation we trained networks using Minimum Description Length (MDL) — a learning criterion that *does* correspond to rational scientific reasoning — the networks were able to find perfect solutions to patterns that remained outside the reach of networks trained with backpropagation. Could there eventually be future LLMs (MDL-based or otherwise) that *would* strike us as scientist-like? Perhaps. But if such models do arrive they will be very different from current LLMs. In the meantime, any suggestion that LLMs are automated scientists should be treated with suspicion.

**Summary.** Piantadosi's excitement is premature. While LLMs are successful as engineering tools, we saw that they are very poor theories of human linguistic cognition. This is hardly a critique of the efforts behind these models: to my knowledge, all current LLMs were developed with engineering goals rather than the goals of cognitive science in mind. Nor do I see any reason to doubt that more human-like AI could in principle be built in the future. But current LLMs remain the stochastic parrots that Bender et al. (2021) tell us they are.<sup>5</sup> Using them to write entertaining poems and short stories is one thing. Using them to understand the human faculty of language instead of doing actual linguistics is quite another.

## References

Baroni, M. (2022). On the proper role of linguistically oriented deep net analysis in linguistic theorising. In Bernardy, J.-P. and Lappin, S., editors, Algebraic Structures in Natural Language, pages 1–16. CRC Press.

Barwise, J. and Cooper, R. (1981). Generalized quantifiers and natural language. *Linguistics and Philosophy*, 4:159–219.

Bender, E. M., Gebru, T., McMillan-Major, A., & Shmitchell, S. (2021). On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? 🦜. In Proceedings of the 2021 ACM conference on fairness, accountability, and transparency (pp. 610-623).

---

<sup>5</sup> Note, however, that differently from LLMs, actual parrots *are*, in fact, intelligent.

March 17, 2023

Chomsky, N. (1957). *Syntactic Structures*. Mouton, The Hague.

Chomsky, N. (1965). *Aspects of the Theory of Syntax*. MIT Press, Cambridge, MA.

Chomsky, N. (1971). *Problems of knowledge and freedom: The Russell lectures*. Pantheon Books.

Chomsky, N. (1975). *Current issues in linguistic theory*. Mouton, The Hague.

Evans, N. and Levinson, S. (2009). The myth of language universals: Language diversity and its importance for cognitive science. *Behavioral and Brain Sciences*, 32:429–492.

Gulordava, K., Bojanowski, P., Grave, E., Linzen, T., and Baroni, M. (2018). Colorless green recurrent networks dream hierarchically. In *Proceedings of NAACL 2018*, pages 1195–1205.

Heinz, J. and Idsardi, W. (2013). What complexity differences reveal about domains in language. *Topics in cognitive science*, 5(1):111–131.

Katzir, R. (2014). A cognitively plausible model for grammar induction. *Journal of Language Modelling*, 2(2):213–248.

Keenan, E. L. and Stavi, J. (1986). A semantic characterization of natural language determiners. *Linguistics and Philosophy*, 9:253–326.

Lakretz, Y., Hupkes, D., Vergallito, A., Marelli, M., Baroni, M., and Dehaene, S. (2021). Mechanisms for handling nested dependencies in neural-network language models and humans. *Cognition*, 213:104699. Special Issue in Honour of Jacques Mehler, *Cognition's* founding editor.

Lan, N., Chemla, E., and Katzir, R. (2022). Large language models and the argument from the poverty of the stimulus. Under review, available at <https://ling.auf.net/lingbuzz/006829>.

Lan, N., Geyer, M., Chemla, E., and Katzir, R. (2022). Minimum description length recurrent neural networks. *Transactions of the Association for Computational Linguistics*, 10:785–799.

March 17, 2023

Miller, G. and Chomsky, N. (1963). Finitary models of language users. In Luce, R. D., Bush, R. R., and Galanter, E., editors, *Handbook of Mathematical Psychology*, volume 2, pages 419–491, New York, NY. Wiley.

Pearl, L. and Sprouse, J. (2013). Syntactic islands and learning biases: Combining experimental syntax and computational modeling to investigate the language acquisition problem. *Language Acquisition*, 20(1):23–68.

Piantadosi, S. T. (2023). Modern language models refute Chomsky's approach to language. Ms., available at <https://ling.auf.net/lingbuzz/007180>.

Rasin, E. and Katzir, R. (2020). A conditional learnability argument for constraints on underlying representations. *Journal of Linguistics*, 56(4):745–773.

Ross, J. R. (1967). *Constraints on Variables in Syntax*. PhD thesis, MIT, Cambridge, MA.

Warstadt, A., Singh, A., and Bowman, S. R. (2019). Neural Network Acceptability Judgments. *Transactions of the Association for Computational Linguistics*, 7:625–641.

Wilcox, E., Futrell, R., and Levy, R. (2022). Using computational models to test syntactic learnability. *Linguistic Inquiry* 2022; doi: [https://doi.org/10.1162/ling\\_a\\_00491](https://doi.org/10.1162/ling_a_00491)

Yngve, V. H. (1960). A model and an hypothesis for language structure. *Proceedings of the American Philosophical Society*, 104(5):444–466.