

# LARGE LINGUISTIC MODELS: ANALYZING THEORETICAL LINGUISTIC ABILITIES OF LLMs

GAŠPER BEGUŠ<sup>a</sup>, MAKSYMILIAN DĄBKOWSKI<sup>a</sup>, RYAN RHODES<sup>b</sup>

<sup>a</sup>UNIVERSITY OF CALIFORNIA, BERKELEY, <sup>b</sup>RUTGERS UNIVERSITY

MANUSCRIPT *as of* MAY 1, 2023

**ABSTRACT** The performance of large language models (LLMs) has recently improved to the point where the models can generate valid and coherent meta-linguistic analyses of data. This paper illustrates a vast potential for analyses of the meta-linguistic abilities of large language models. LLMs are primarily trained on language data in the form of text; analyzing their meta-linguistic abilities is informative both for our understanding of the general capabilities of LLMs as well as for models of linguistics. In this paper, we propose several types of experiments and prompt designs that allow us to analyze the ability of GPT-4 to generate meta-linguistic analyses. We focus on three linguistics subfields with formalisms that allow for a detailed analysis of GPT-4’s theoretical capabilities: theoretical syntax, phonology, and semantics. We identify types of experiments, provide general guidelines, discuss limitations, and offer future directions for this research program.

## 1 INTRODUCTION

GPT-4 (OpenAI, 2023) may be the first large language model (LLM) that can generate coherent syntactic, phonological, and semantic analyses of linguistic data using some of the popular theoretical approaches in these subfields (Chomsky, 2014; Heim and Kratzer, 1998; Prince and Smolensky, 1993/2004).<sup>1</sup> This paper illustrates a vast potential to explore not only language competence, but also metalinguistic abilities of LLMs. We explore GPT-4’s ability to analyze linguistic constructions at syntactic, phonological, prosodic, and semantic levels.

LLMs have seen enormous success in recent years in a variety of fields. The architecture that underlies most LLMs are *transformers* (Vaswani et al., 2017)—a deep neural network architecture for sequential data processing. While the exact details about the architecture of GPT-4 is not available, we know that the underlying model is a transformer trained to predict the next token in a sequence (OpenAI, 2023). GPT-4’s predecessor, GPT-3 has 175 billion parameters (Brown et al., 2020); the exact number of parameters in GPT-4 is unknown.

There exists a large body of work testing the linguistic abilities of neural networks trained on text, from recurrent neural networks (RNNs) (Matuskevych and Culbertson, 2022), long-short term memory networks (LSTMs) to, most recently, transformers (Wilcox et al., 2018; Yedetore et al., 2023). LLMs have also ignited a debate on what their outputs mean for linguistic theory in general (Katzir, 2023; Piantadosi, 2023). Most of these studies, however, test the models’

<sup>1</sup> This discussion is inspired by Beguš’s (2023) tweet posted on the day of GPT-4’s release (May 14) as well as subsequent discussions that emerged from the tweet.

performance with behavioral tasks and do not elicit metalinguistic analysis—the ability to generate analyses of data within a given theoretical linguistic framework. Until recently, LLMs were unable to generate coherent complex linguistic analyses.

The purpose of this paper is not to provide evidence in favor of or against individual theoretical proposals in linguistics or even make claims in favor of or against functional or generative approaches to language. Theoretical linguistics is an approach to modeling language and LLMs are trained primarily on language data. This paper argues that testing the metalinguistic abilities of these models is now a possible line of inquiry and that the results of such tests can provide valuable insights for understanding the general metacognitive abilities of LLMs.

The paper outlines a methodology for such studies, identifies pitfalls, and offers future directions in this line of inquiry. We present three preliminary case studies: on syntax, phonology, and semantics; and evaluate how well GPT-4 performs on each task. We show that the model's performance can be evaluated from a developmental perspective: we can compare the performance of different versions (e. g. GPT-3.5 vs. GPT-4) or the performance of different models (e. g. GPT-4 vs. Bard).

## 2 RESEARCH PROGRAM OUTLINE

Why is this line of work important? The majority of studies thus far perform *behavioral* tests of LLMs. Tasks like asking whether a sentence is (un)grammatical or testing LLMs on whether they perform correctly various syntactic operations (such as agreement, movement, embedding) are behavioral tasks that test *language performance*. Here, we outline a research program where LLMs like GPT-4 are tested on higher level *metalinguistic abilities*. We use the term *metalinguistic ability* to refer to the ability to generate formal, theoretical analyses of linguistic phenomena—in other words, to refer to the work that linguists do. Metalinguistic ability is cognitively more complex compared to behavioral use of language (Tunmer and Herriman, 1984). Metalinguistic ability is acquired later, and linguistic competence is its precondition. Applying theoretical linguistic analyses and formalisms from the training data to the model's metalinguistic abilities is a process that is quite complex. This line of work thus tests the next level of cognitive abilities of LLMs.

Testing LLMs on more complex cognitive tasks can inform us about LLMs' general capabilities as well as provide a metric to compare these abilities across versions and models. We argue that metalinguistic tasks are a very useful testing ground for these kinds of evidence.

Many studies attempt to test whether linguistic structures are learnable from surface statistics (the relative poverty of a human learner's input notwithstanding, given a sufficient quantity of input data, can the target grammar be acquired?). LLMs acquire their linguistic competence from the surface statistics of their training data. We want to know whether this is a sufficient basis to perform the cognitive task of analyzing language itself.

Prompting the networks to analyze structure in linguistic terms might give us clues about whether the model has, at some level at least, access to linguistic structure that could enable/inform behavioral outputs. In other words, if the model is performing well on behavioral

tasks, are its outputs felicitous because of distributional knowledge, or does the network have access to linguistic structure from which it generates its output? Although the internal representational status of LLMs is currently unknown (and hotly contested), we hope that this line of inquiry may provide a useful direction for this overarching question.

Human linguists have innovated the range of analytical tools that comprise the science of linguistics by reasoning from their own knowledge of their native languages (that is, from their mental grammar). It remains to be seen whether LLMs can ever achieve a similar capacity to reason from their “knowledge” about language to analyze their own grammar.

The proposed research program will also test whether LLMs are capable of innovative solutions that were not hypothesized by humans thus far. Deep neural networks have been instrumental in shrinking the hypothesis spaces and offer new solutions to problems in fields as diverse as protein design (Jumper et al., 2021), geometry (Davies et al., 2021), and understanding unknown communication systems (Beguš, Leban, and Gero, 2023). It remains to be seen whether any new insights can be gained by metalinguistic prompting of LLMs.

Finally, this line of work can be useful for educational purposes. Evaluating the performance of LLMs’ metalinguistic abilities is an effective pedagogic task for both learning about linguistic theory and for learning to evaluate the abilities of deep neural models.

### 2.1 General guidelines

Any metalinguistic ability can be tested with GPT-4. We will outline several experimental design choices to elicit metalinguistic behavior from LLMs.

First, the vast majority of training data is likely in English, but GPT-4 returns coherent analyses of other languages as well. Testing the model’s performance on metalinguistic tasks in non-English languages can be very informative. Additionally, the training data for GPT in particular is so expansive that it almost certainly includes linguistic analyses—particularly simple analyses and problem sets that one would find in linguistics educational materials. While we can use these examples as a starting point, it will become important to vary the tasks and data sets used to prompt the LLM to probe whether it has the capacity to generalize beyond the exemplars included in its training. Crucially, we are interested in whether the LLM has simply memorized a linguistic analysis or a particular problem or data set, or whether it can, in some sense, think *creatively* as well as analytically.

The best test case for this purpose may be to test LLMs on artificial languages. Some work has recently begun testing LLMs using artificial grammar learning (AGL) paradigms that have been used in human research for decades (Matussevych and Culbertson, 2022; Yedetore et al., 2023). These studies are structured like human behavioral AGL studies, but they can also potentially be adapted to test metalinguistic abilities as well. This is a crucial test case,

because it avoids the problem of memorization—we can tailor data sets explicitly to be as *unlike* the exemplars in the training data as possible.<sup>2</sup>

Reproducibility is a major concern for this type of research. Unfortunately, GPT is still a black box, and there is a lot we don't know about its architecture, training data, memory, or how it was optimized for human conversation. GPT's outputs can be noisy and unpredictable—subsequent prompts, even with the same wording, can generate substantially different responses. Given the opacity of the model, it will be difficult to make inferences about its abilities and internal representational states. Data transparency will be extremely important for reproducibility and cross-study comparison. Prompts and responses need to be disclosed in full. Our own data will be available for download, comparison, and reanalysis.

It can also be informative to compare models, and to compare current models with their past versions. In this paper, we will make some comparisons between GPT-4 and GPT-3.5 (the previous model inevitably performs much worse on the same prompts). Again, transparency and prompt consistency is key for making these kinds of comparisons.

### 3 SYNTACTIC THEORY

Syntax is a subfield of linguistics which studies the structure of sentences and the relations or rules which derive complex sentences from simpler ones. While there are many competing theoretical frameworks in syntax (e. g. *construction grammar*: Fillmore, Kay, and O'Connor, 1988; *categorial grammar*: Ajdukiewicz, 1935; *relational grammar*: Perlmutter, 1980), one frequently adopted set of approaches evolved from the *minimalist program* (Chomsky, 2014), which in turn traces its intellectual lineage back to *X-bar theory* (Chomsky, Jacobs, and Rosenbaum, 1970) and *government and binding* (Chomsky, 1993). X-bar theory was used for testing the GPT-4's linguistic abilities because its formalism is relatively well-established and straightforward to both prompt and evaluate. Since this research tradition is prominently represented in the literature, it is most likely present in the data set used to train large language models.

We tested GPT on a set of syntactic problems, involving syntactic ambiguity and center embedding in English and movement in German. We find that GPT-4 is capable of providing mostly correct analyses of relatively simple structures. GPT-3.5 fares considerably worse on the same tasks.

#### 3.1 Ambiguity and center-embedding

To evaluate GPT-4's ability to perform syntactic analyses couched in these frameworks, we prompted the model to draw syntactic tree diagrams for a number of sentences in various LaTeX packages such as the forest package (Živanović, 2017). Other packages such as tikz (Tantau, 2007) or tikz-tree can also be used.

---

<sup>2</sup> There are many other reasons to pursue this approach, including the possibility that LLMs have acquired human-like cognitive *biases* alongside their distributional knowledge of human language. However, this is beyond the scope of the current paper.

In [Prompt 1](#), we asked GPT-4 to draw a syntactic analysis of the famous sentence “I saw an elephant with binoculars” to test if GPT-4 can recognize two possible structures of this sentence based on exactly the same surface form.

**Prompt 1:** Draw a theoretical syntactic analysis of the sentence “I saw an elephant with binoculars” in LaTeX using the forest package. Use the X-bar theory. Draw all possible structures of this sentence.

Broadly speaking, the structures generated by GPT-4 successfully capture the difference between the two readings of the sentence ([Reply 1](#)).<sup>3</sup> In the first structure, the prepositional phrase “with binoculars” attaches to the verb phrase. This corresponds to the reading where seeing the elephant is done with binoculars, i. e. the speaker is wearing the binoculars. In the second structure, the PP “with binoculars” attaches to the noun “elephant.” This corresponds to the reading where the elephant is wearing the binoculars. GPT-4’s explanation correctly tracks the relationship between syntactic structure and meaning.

GPT-4 consistently places tensed verbs under T (or I) and leaves the V-head position empty. This seems to reflect V-to-T raising, a type of head movement that has been argued for languages like French. This is a mistake since the [+PAST] tense feature in English is commonly analyzed as lowering to V (as opposed to the verb raising to T).

Finally, while GPT-4 correctly identified the general relationship between syntactic structure and meaning, the response skips many bar levels and does so inconsistently. According to the standard X-bar schema ([Figure 1](#)), each phrasal category has a three-level structure: X, X', and XP (Jackendoff, 1977). With this elaborate structure, X-bar trees may get very unwieldy. Consequently, skipping vacuous levels of structure as shorthand is common in X-bar theory, but it is worth noting that GPT-4 is inconsistent. For example, adjuncts are supposed to combine with an X'-level node without going up a level (i. e. adjunction yields another X'-level node). GPT-4 deviates from this structure in two different ways. In the tree on the left, the PP adjoins to V' (correct), but the next node is labeled as VP (as opposed to V'). In the tree on the right, the PP adjoins to N (not N') and goes up two bar levels, yielding an NP (as opposed to staying at N').

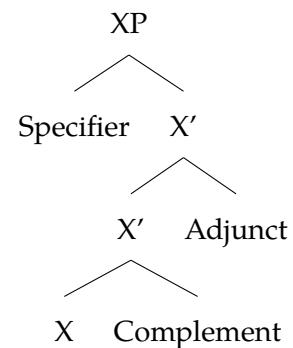


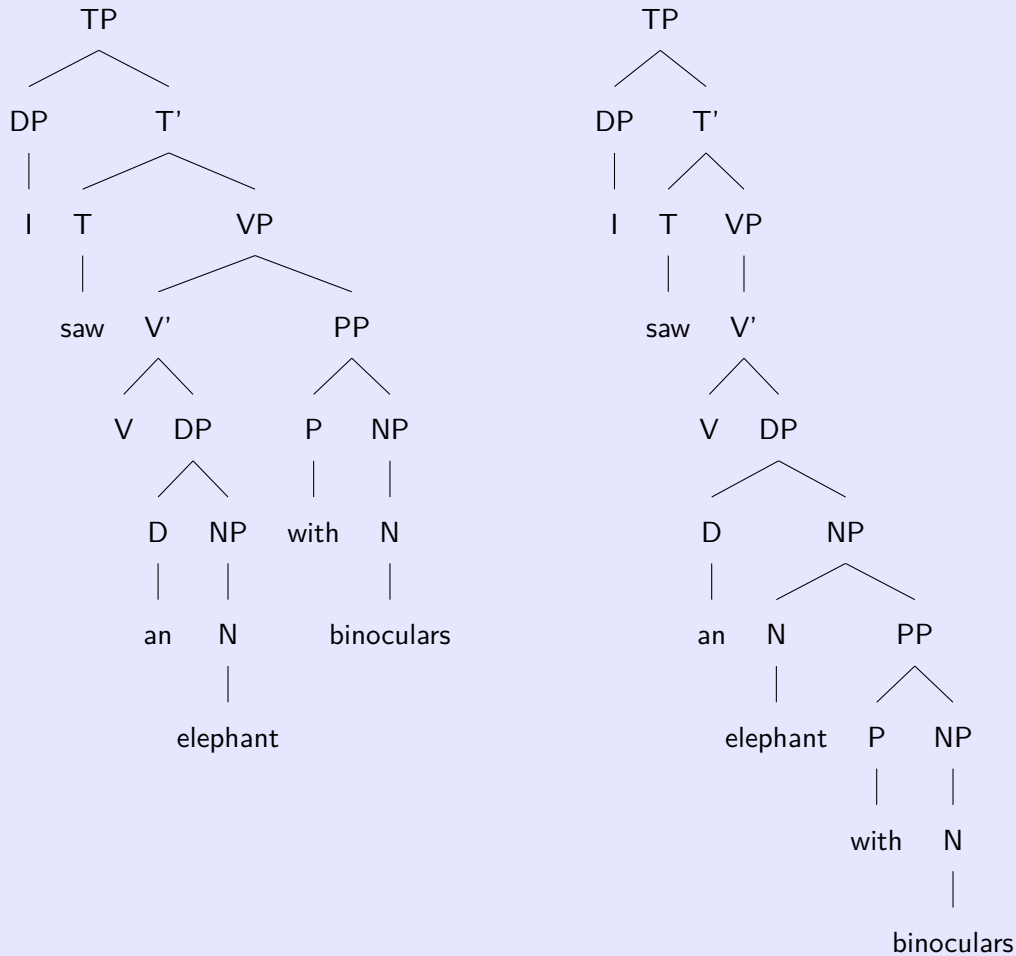
Figure 1: X-bar schema.

**Reply 1:** As a text-based AI, I am unable to directly draw diagrams. However, I can provide you with the LaTeX code using the forest package to create the syntactic trees based on the X-bar theory. You can copy and paste the code into your LaTeX editor to visualize the tree structures. Below are two possible structures for the sentence “I saw an elephant with binoculars.”

<sup>3</sup> Some of GPT’s replies have been formatted and edited for conciseness. For complete transcripts of full exchanges, see Supplementary Materials.

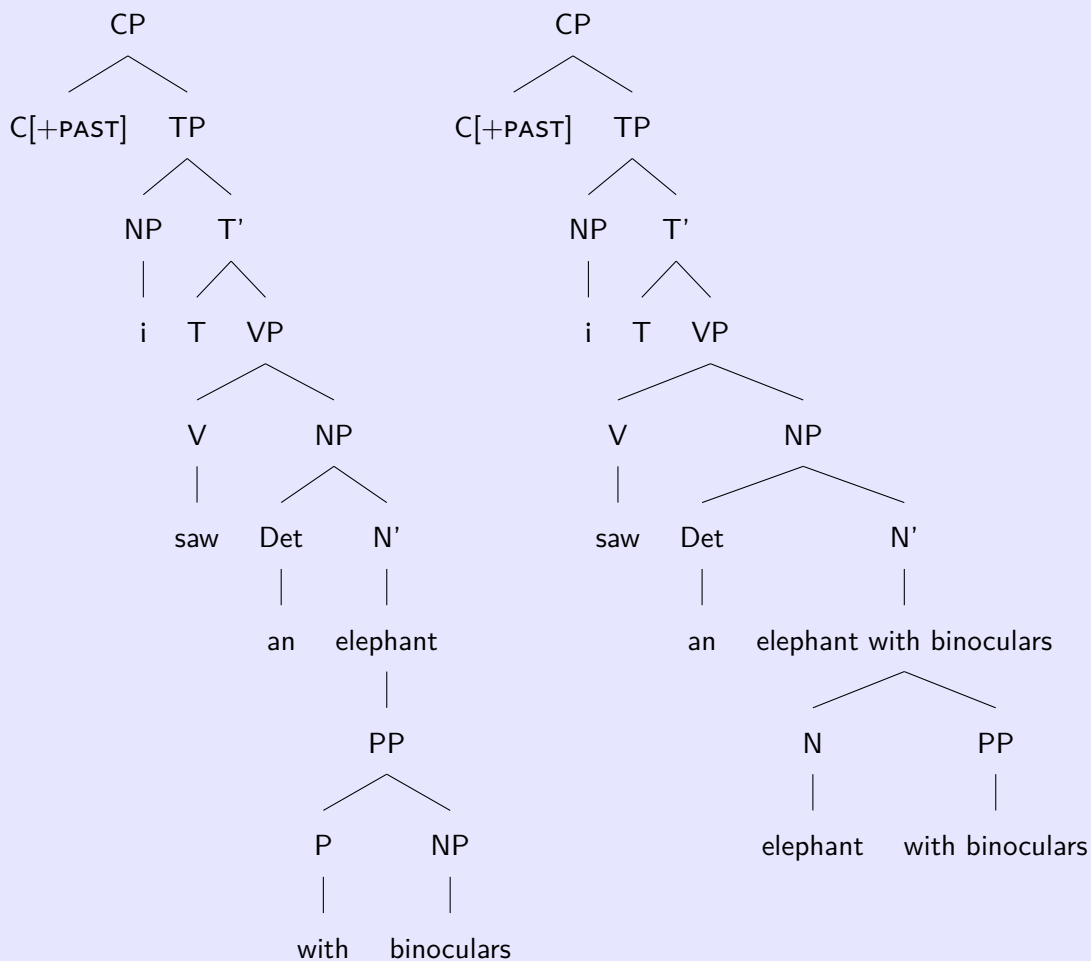
The first structure assumes that the “with binoculars” phrase is modifying the verb “saw,” indicating that the person saw the elephant using binoculars.

The second structure assumes that the “with binoculars” phrase is modifying the noun “elephant,” indicating that the person saw an elephant that had binoculars.



As a comparison, GPT-3.5 performs considerably worse on the same prompt ([Reply 2](#)). It generates an analysis for only one of the two possible interpretations (left-hand tree). The tense feature appears in C, and the modifying PP branches directly underneath *elephant* (which should be a terminal node).

When asked to diagram the second structure (right-hand tree), it generates an even less plausible tree. Again, the tense feature is in C, and again the PP is inappropriately placed under a terminal node. But here, the entire phrase *elephant with binoculars* is placed under N' as a single node, with N and PP nodes consisting of a repetition of *elephant* and *with binoculars* underneath it.

**Reply 2:**

The explanation produced by GPT-3.5 ([Reply 3](#)) is also substantially worse. The first paragraph describes the structure of the first tree, but does not discuss the relationship between meaning and syntactic structure. The second paragraph alludes to the possibility of providing different analyses, but fails to relate them to semantic ambiguity. While GPT-3.5 observes that the PP node may be “attached to the VP node” or be “modifier with N’ node,” it points to individual linguists’ analytical preferences and fails to explain the differences in meaning.

**Reply 3:** This tree structure represents one possible interpretation of the sentence. The CP (complementizer phrase) node contains the C (complementizer) node with the feature [+PAST], indicating that the sentence is in the past tense. The TP (tense phrase) node contains the subject NP (noun phrase) “I” and the T’ (tense phrase) node, which contains the T (tense) node and the VP (verb phrase) node. The VP node contains the verb “saw” and the object NP, which consists

of the determiner “an” and the N’ (noun phrase) node “elephant”, which in turn contains the PP (prepositional phrase) node “with binoculars”.

Note that there are other possible structures for this sentence, depending on the analysis chosen. For example, some linguists might treat “with binoculars” as a separate sentence constituent and place it in a separate PP node attached to the VP node, while others might analyze it as a modifier within the N’ node. The X-bar theory allows for flexibility in representing sentence structure, and different analyses can capture different aspects of the sentence’s meaning and syntax.

It’s crucially important to also test LLMs on constructions that are unlikely to have been part of their training data. The *elephant* sentence is a very common example sentence that appears in many introductory linguistics texts and problem sets. It’s likely that this sentence in particular appeared (at some level of frequency) in GPT’s training data. For this reason, in [Prompt 4](#), we asked GPT-4 to analyze a similarly ambiguous sentence (adapted from a joke in a Cyanide & Happiness comic; McElfatrick, 2008).

**Prompt 4:** Do a theoretical syntactic analysis of the sentence “I fed my baby turtles” in LaTeX using the forest package. Use the X-bar theory. Draw all possible structures of this sentence.

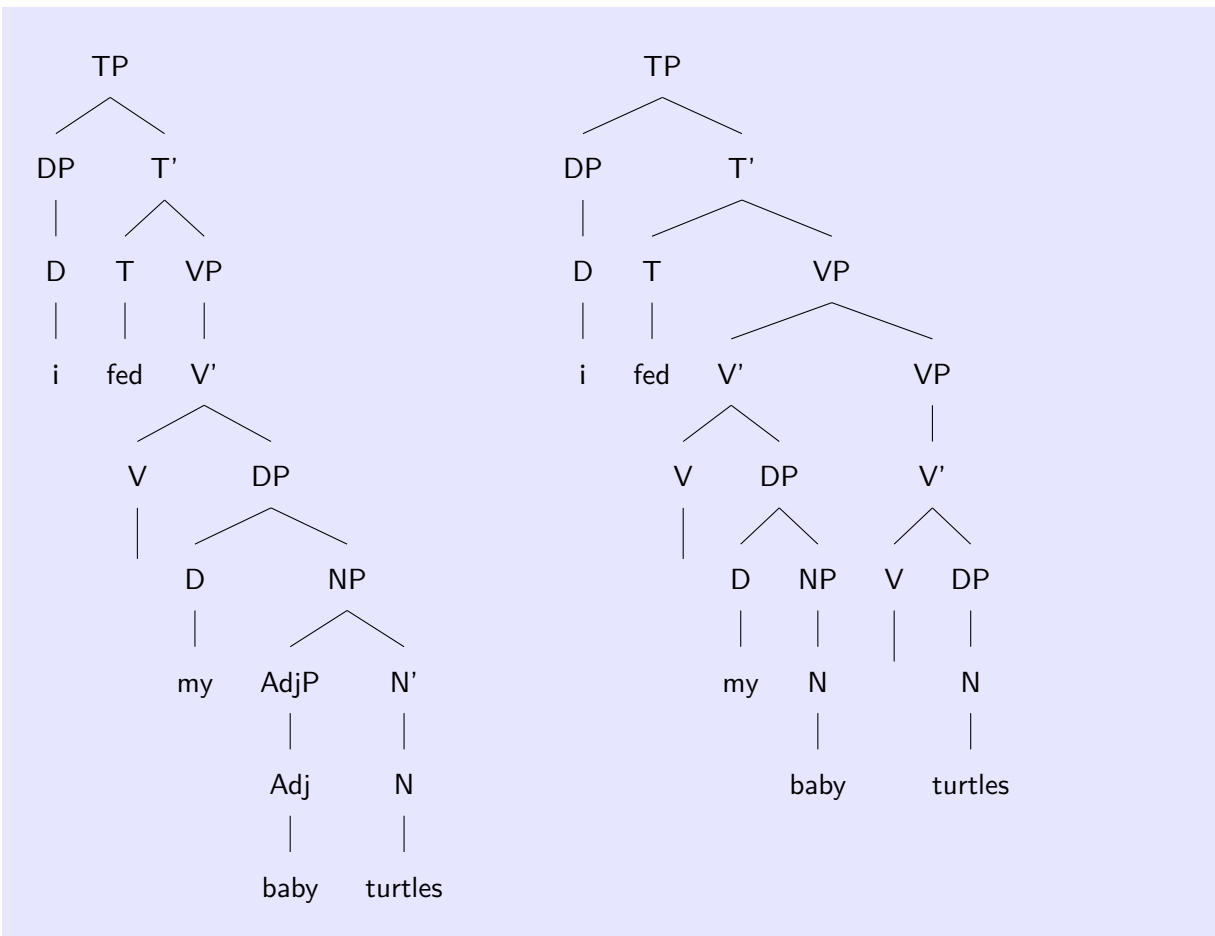
GPT-4 performs surprisingly well in this task, generating two tree structures from the initial prompt ([Reply 4](#)). Each tree corresponds to an interpretation of the sentence, which GPT accurately describes. The trees (particularly the second tree) are not perfect—GPT gets some details of the X-Bar analysis wrong.<sup>4</sup> But importantly, it largely gets the constituent structure right in a way that neatly delineates the two possible interpretations.

The same structure obtains for other variations of this basic structure (using different object nouns like *mice*), as well as other ambiguous sentences like *Joan baked cookies for the baby in the oven*.

**Reply 4:** The sentence “I fed my baby turtles” can have two possible interpretations, depending on whether “baby” is modifying “turtles” or not. If “baby” is modifying “turtles,” the sentence means “I fed the baby turtles that I have.” If “baby” is not modifying “turtles,” the sentence means “I fed my baby some turtles (to eat).” Here is the LaTeX code for the two possible structures using the X-bar theory and the forest package:

<sup>4</sup> GPT is still placing the verb in T, as though through V-to-T raising, and although *turtles* adjoins at an approximately correct level to modify the VP, it for some reason projects its own VP.



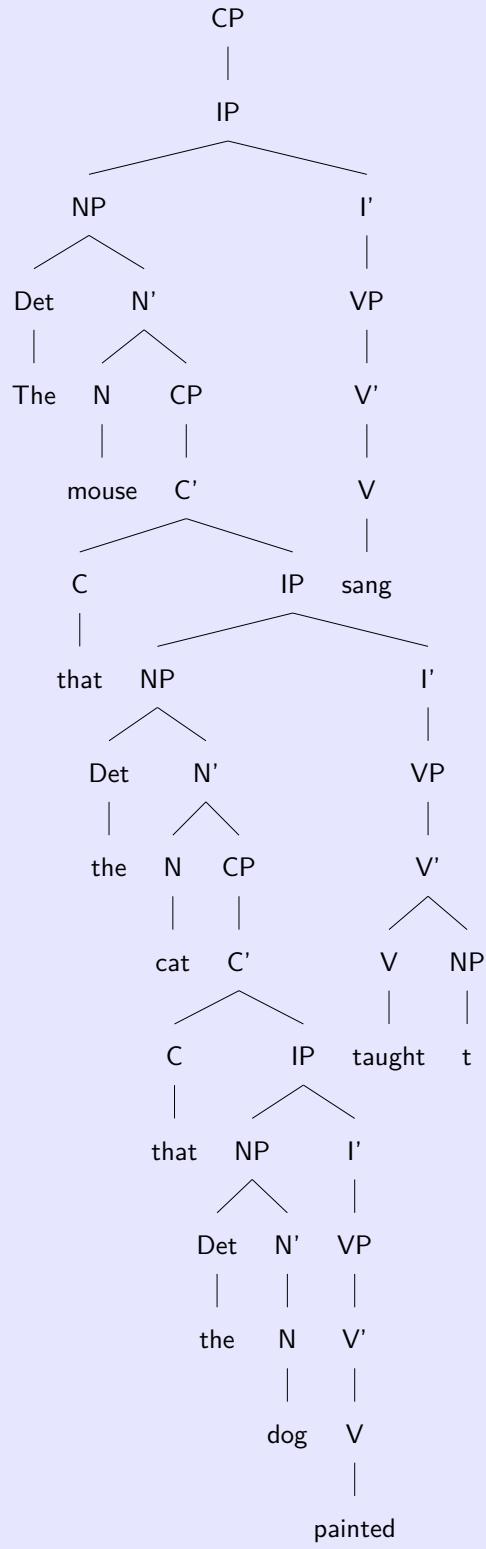


Finally, we tested GPT-4's performance on a center-embedding sentence ([Prompt 5](#); the example sentence adopted from Katzir, 2023). Center embedding refers to a specific type of structure where a phrase is nested in the middle of another phrase of the same type. Center-embedding sentences are difficult to process, resulting in an impression of ill-formedness. Thus, although grammatical, center-embedding sentences are virtually unattested in natural discourse. This means that LLMs have almost certainly not been exposed to center-embedding structures. As such, they are a useful test for the models' generative syntactic ability.

**Prompt 5:** Draw a theoretical syntactic analysis of the sentence "The mouse that the cat that the dog painted taught sang" in LaTeX using the forest package. Use the X-bar theory. Draw all possible structures of this sentence.

GPT-4 performs remarkably well on this task ([Reply 5](#)). The CP "that the dog painted" modifies "cat" and the CP "that the cat that the dog painted taught" modifies "mouse." I. e., the model's analysis is correct. [Reply 5](#) furthermore improves on the analyses above in some minor respects. For example, it does not show incorrect V-to-T (or I) raising.

Reply 5:



### 3.2 Movement in German

In [Prompt 6](#), we asked GPT-4 to analyze a German sentence with multicausal wh-movement. This construction is analyzed under an X-Bar theoretic approach to involve successive cyclic movement, with traces of the wh-element at each phase edge (under the minimalist framework, these are rendered as unpronounced copies, but they are functionally equivalent).

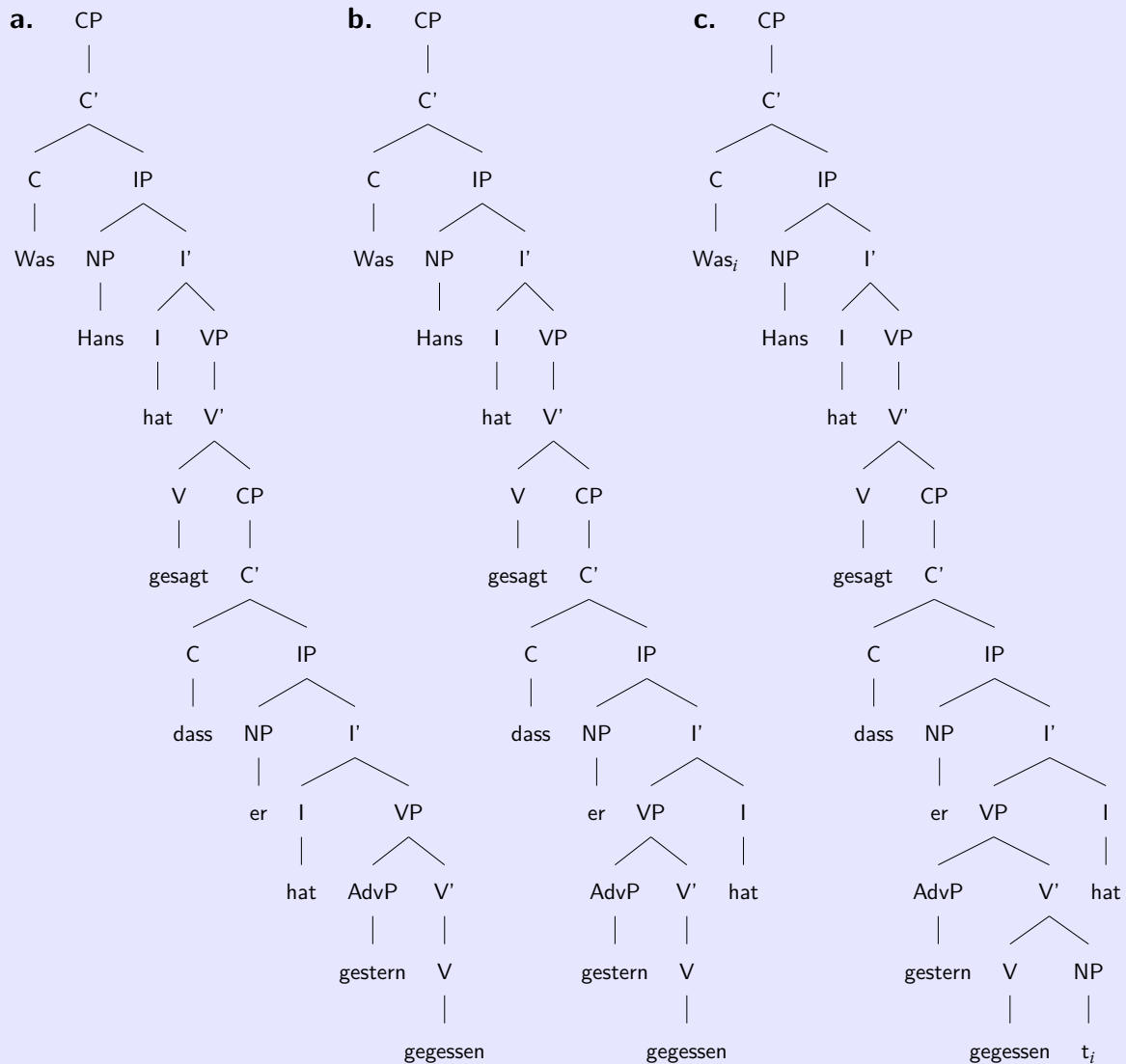
GPT-4 consistently gets some details wrong (wh-elements are typically argued to occupy a Spec-CP position; rather, the verb/auxiliary should occupy Head-CP). With the initial [Prompt 6a](#), although GPT (mostly) gets the hierarchical structure right, it fails to produce a branch ordering that reflects the surface word order (e. g. the daughters of I' and V' have the wrong branch ordering to reflect verb-final word order). And there is no reference to either traces or copies of the wh-element ([Reply 6a](#)).

When specifically prompted to account for the correct surface word order ([Prompt 6b](#)), GPT revises its tree. It correctly orders the branches under I' in the embedded clause, but not in the matrix clause ([Reply 6b](#)). When GPT is prompted to specify the base-generated position of the wh-element ([Prompt 6c](#)), it introduces a trace in the verb phrase of the embedded clause ([Reply 6c](#)). Still, there is no explicit reference to cyclic movement (no intermediate trace at Spec-CP of the embedded clause). This may reflect a bias in GPT's training data for older versions of X-Bar theory, or pedagogical materials aimed at introductory linguistics students (materials that presumably lack these theoretical details).

#### Prompt 6:

- a. Draw a theoretical syntactic analysis of the German sentence "Was hat Hans gesagt dass er gestern gegessen hat?" in LaTeX using the forest package. Use the X-bar theory. Draw all possible structures of this sentence.
- b. Draw the same sentence, but make sure you account for the correct word order.
- c. Where is "Was" base-generated?

## Reply 6:



In summary, although GPT-4 makes some unconventional assumptions about syntactic structure, the model is capable of providing analyses of syntactically ambiguous sentences, center embedding structures, and complex structures which involve both  $A'$ -bar and head movement. Notably, GPT-4 far outperforms GPT-3.5, which often fails to generate well-formed syntactic trees altogether.

## 4 PHONOLOGICAL THEORY

Phonology is a subfield of linguistics which studies the organization of sounds within a language. In phonological theory, facts of sound distribution or sound alternations between

sets of morphologically related words can be modeled as rules transforming one sound into another.

We tested GPT-4 on two phonological problems by prompting the model with small data sets: a palatalization process in Korean and a spirantization process in an artificial language. These were the type of questions that a linguistics student could encounter in a homework problem set. GPT-4 performs well on tasks involving phonological theory and successfully analyses familiar data sets as well as novel artificial languages.

#### 4.1 Korean palatalization

In [Prompt 7](#), we asked GPT-4 to solve a simple phonological neutralization problem in Korean. In Korean, /s/ palatalizes before /i/ and is realized as [ʃ]. This is a simple problem that is often given to introductory linguistics students. The model did very well at this problem ([Reply 7](#)), correctly recognizing the conditioning environment and writing a phonological rule to describe the pattern in the data set.

**Prompt 7:** Consider the following (simplified) words from Korean and answer the questions about the data below – focus on the [s] and [ʃ] sounds.

**KOREAN** ʃi kasu sal kasəl miso ʃinmun kaʃi paŋʃik oʃip so  
**ENGLISH** poem singer flesh hypothesis smile newspaper thorn method fifty cow

Write out the environments for [s] and [ʃ] in a LaTeX table. Include the phonemes that appear on either side of [s] and [ʃ] in the data set.

Write a rule describing the transformation that's occurring in Korean. Example:  $A \rightarrow B / C \_ D$

**Reply 7:** Here is a LaTeX table showing the environments for [s] and [ʃ] in the given Korean words:

Word	Phoneme	Preceding Phoneme	Following Phoneme
[ʃi]	[ʃ]	# (word boundary)	[i]
[kasu]	[s]	[a]	[u]
[sal]	[s]	[a]	[l]
[kasəl]	[s]	[a]	[ə]
[miso]	[s]	[i]	[o]
[ʃinmun]	[ʃ]	# (word boundary)	[i]
[kaʃi]	[ʃ]	[a]	[i]
[paŋʃik]	[ʃ]	[ŋ]	[i]
[oʃip]	[ʃ]	[o]	[i]
[so]	[s]	# (word boundary)	[o]

Based on the table above, we can make the following generalization about the distribution of [s] and [ʃ] in Korean:

[s] → [ʃ] / # \_ [i]

This rule states that the [s] sound changes to [ʃ] when it appears at the beginning of a word (word boundary, represented by "#") and is followed by [i].

#### 4.2 Spirantization in an artificial language

In [Prompt 8](#), we provided GPT-4 with a small data set demonstrating a phonological alternation between the voiceless stops (p t k) and voiceless fricatives (ɸ θ x) in morphologically related forms. We asked the model to identify one phonological process that operates in the language and analyze it.

**Prompt 8:** Here are some singular nouns and their corresponding plural form in a language you're encountering for the first time:

**SINGULAR** rap rat vik sik lum zen jot kit gut kip zep vanup roθep milik hep  
**PLURAL** raɸa raθa vixa sixa luma zena joθa kiθa guθa kiɸa zeɸa vanuɸa roθeɸa milixa heɸa  
 suril mirik riθit delop raðat vop  
 surila mirixa riθiθa deloɸa raðaða voɸa

The plural is formed by adding /-a/ to the singular form. Identify one phonological process that operates in the language and analyze it.

GPT-4 first analyzed the alternation as assimilation to the following vowel with respect to the place of articulation ([Reply 8](#)). Since the spirantization of a stop in the proximity of a vowel can be seen as a type of assimilation (the increase in sonority makes the sound more similar to the vowel), the answer is partly correct. However, the assimilated feature is the manner, not the place, of articulation.<sup>5</sup>

**Reply 8:** Based on the provided data, one phonological process that seems to operate in this language is assimilation. Assimilation is a process where a sound becomes more like a neighboring sound in terms of some phonetic features. In this case, it appears that the final consonant of the singular form assimilates to the /-a/ suffix in the plural form, with regard to the place of articulation.

Here are some examples:

<sup>5</sup> Additionally, there are various notational issues which we do not focus on in evaluating the models' answers. For example, GPT-4 uses slashes / / to represent both underlying and surface forms. (The latter are conventionally given in brackets [ ].) GPT-4 uses the right-pointing arrow (→) to represent (i) a morphological operation followed by (ii) a phonological process, whereas conventionally the arrow is used to represent only the latter.

1. rap → raφa: /p/ becomes /φ/ (bilabial voiceless fricative) before /a/
2. rat → raθa: /t/ becomes /θ/ (dental voiceless fricative) before /a/
3. vik → vixa: /k/ becomes /x/ (velar voiceless fricative) before /a/
4. sik → sixa: /k/ becomes /x/ (velar voiceless fricative) before /a/
5. lum → luma: /m/ remains the same (bilabial nasal) before /a/
6. vanup → vanuφa: /p/ becomes /φ/ (bilabial voiceless fricative) before /a/
7. roθep → roθeφa: /p/ becomes /φ/ (bilabial voiceless fricative) before /a/

The assimilation process in this language is regressive, meaning that the features of the following sound (the /-a/ suffix) affect the preceding consonant, causing it to change its place of articulation to match the following vowel. This results in the final consonant of the singular form becoming a fricative or remaining a nasal in the plural form, depending on the original consonant.

When pressed about the identity of the assimilated feature ([Prompt 9](#)), GPT-4 completely changed its answer and correctly identified the phonological process as “a form of consonant mutation called spirantization” ([Reply 9](#)). The model correctly illustrated the process with examples from the data set and generalized over the natural class of alternating segments by observing that “[t]he stop consonants at the end of the singular nouns become fricatives when the plural suffix /-a/ is added.”

**Prompt 9:** If it's assimilation, what is the feature that the consonant assimilates with respect to?

**Reply 9:** Apologies for the confusion in my previous response. Upon closer examination, the phonological process that operates in this language is not assimilation. Instead, it appears to be a form of consonant mutation called spirantization. Spirantization is a process where a stop consonant turns into a fricative or approximant in certain phonological environments.

In this language, spirantization occurs when a stop consonant (/p/, /t/, /k/) appears at the end of a singular noun form and is followed by the plural suffix /-a/. The stop consonant turns into a fricative with the same place of articulation as the original stop:

1. /p/ → /φ/: bilabial voiceless stop becomes bilabial voiceless fricative
  - rap → raφa
  - vanup → vanuφa
  - roθep → roθeφa
2. /t/ → /θ/: alveolar voiceless stop becomes dental voiceless fricative
  - rat → raθa

- rađat → rađaθa
3. /k/ → /x/: velar voiceless stop becomes velar voiceless fricative
- vik → vixa
  - sik → sixa

In this case, the feature that changes during the spirantization process is the manner of articulation. The stop consonants at the end of the singular nouns become fricatives when the plural suffix /-a/ is added. The place of articulation remains the same for the affected consonants.

Next, we performed a wug test (Berko, 1958) to establish if GPT-4 is able to correctly apply the generalizations it identified to previously unencountered forms (Prompt 10).

**Prompt 10:** Given the singular nouns “zefop,” “trewen,” “nenek,” and “gesel,” how would you pluralize them?

GPT-4 fully succeeds at applying the spirantization rule to new forms (Reply 10) and explaining its reasoning. (Curiously, the model represented the velar fricative with “χ,” which typically stands for the uvular fricative, and “motivated” its unconventional notation.)

**Reply 10:** Based on the observed pattern of spirantization in the provided data, we can predict the plural forms of the given singular nouns as follows:

1. zefop → zefoφa: The voiceless stop /p/ in the final position of the singular noun “zefop” undergoes spirantization, turning into the bilabial voiceless fricative /φ/ before the plural suffix /-a/.
2. trewen → trewena: In the singular noun “trewen,” the final consonant is a nasal /n/ rather than a voiceless stop. Since spirantization applies only to voiceless stops, we would simply add the plural suffix /-a/ without any changes to the final consonant.
3. nenek → nenexa: The voiceless stop /k/ in the final position of the singular noun “nenek” undergoes spirantization, turning into the velar voiceless fricative /x/ (represented here as χ for clarity) before the plural suffix /-a/.
4. gesel → gesela: In the singular noun “gesel,” the final consonant is a liquid /l/ rather than a voiceless stop. Since spirantization applies only to voiceless stops, we would simply add the plural suffix /-a/ without any changes to the final consonant.

So, the plural forms of “zefop,” “trewen,” “nenek,” and “gesel” would likely be “zefoφa,” “trewena,” “nenexa,” and “gesela,” respectively.

Then, we tested GPT-4’s ability to predict further alternations by reasoning about natural classes of sounds (Prompt 11). To do so, we asked the model to generate plural forms of nouns



ending in voiced stops (b d g). Note that such words do not appear in the data set. Thus, the task is more difficult than the standard wug test, since the model has to guess what class of sounds the spirantization rule it identified actually applies to. It is also open-ended, in that more than one answer is justifiable.

**Prompt 11:** Given the singular nouns “lob,” “weg,” “nemed,” how would you pluralize them?

Given the lack of voiced stop-final words in the data set, GPT-4’s [Reply 11](#) erred on the side of rule non-application for /d/ and /g/. For /b/, the model suggested three possible outputs: [ϕ], [β], and [b]. The first option ([ϕ]) assumes that the output of spirantization for the voiced /b/ is the same as for the voiceless /p/. While phonetically implausible (since it would involve intervocalic devoicing), this guess is not unreasonable given that only spirantization to [ϕ] was observed in the data set. The second option ([β]) assumes that the spirantization applies to voiceless and voiced stops alike and that voicing is preserved. Given phonetic plausibility and phonological typology, this is arguably the most likely outcome. The third option ([b]) assumes that spirantization applies only to voiceless stops. This again is a reasonable guess given no counter-indications in the original data set. Overall, the model’s performance on the task is impressively nuanced.

**Reply 11:** Given the pattern of spirantization observed in the provided data, we can predict the plural forms of the given singular nouns as follows:

1. lob → loϕa: The voiceless stop /b/ in the final position of the singular noun “lob” is not one of the voiceless stops we’ve observed undergoing spirantization in the data provided (/p/, /t/, /k/). However, since /b/ is the voiced counterpart of /p/, we could assume that it might undergo a similar process, turning into the bilabial voiced fricative /β/ before the plural suffix /-a/. Although this is a reasonable assumption based on the data provided, it is also possible that “lob” simply takes the plural suffix without any change: loba.
2. weg → wega: In the singular noun “weg,” the final consonant is a voiced velar stop /g/, which is not one of the voiceless stops we’ve observed undergoing spirantization. Since we do not have enough data to predict how voiced stops behave in this language, we would simply add the plural suffix /-a/ without any changes to the final consonant.
3. nemed → nameda: In the singular noun “nemed,” the final consonant is a voiced alveolar stop /d/, which is not one of the voiceless stops we’ve observed undergoing spirantization. Since we do not have enough data to predict how voiced stops behave in this language, we would simply add the plural suffix /-a/ without any changes to the final consonant.

So, the plural forms of “lob,” “weg,” and “nemed” would likely be “loϕa” (or “loba” if /b/ does not undergo spirantization), “wega,” and “nameda,” respectively. It is important to note that without more data on how voiced stops behave in this language, these predictions may not be entirely accurate.

When explicitly instructed to “assume that voiced stops behave in a way analogous to voiceless stops” (Prompt 12), GPT-4 flawlessly extends the spirantization rule to /b d g/, predicting [β ð ɣ] in the plural forms (Reply 12).

**Prompt 12:** Assume that voiced stops behave in a way analogous to voiceless stops. Given the singular nouns “lob”, “weg”, “nemed”, how would you pluralize them?

**Reply 12:** If we assume that voiced stops behave in a way analogous to voiceless stops, we can predict the plural forms of the given singular nouns as follows:

1. lob → loβa: The voiced stop /b/ in the final position of the singular noun “lob” undergoes a process analogous to spirantization, turning into the bilabial voiced fricative /β/ before the plural suffix /-a/.
2. weg → weɣa: The voiced stop /g/ in the final position of the singular noun “weg” undergoes a process analogous to spirantization, turning into the velar voiced fricative /ɣ/ before the plural suffix /-a/.
3. nemed → nemeða: The voiced stop /d/ in the final position of the singular noun “nemed” undergoes a process analogous to spirantization, turning into the dental voiced fricative /ð/ before the plural suffix /-a/.

So, under the assumption that voiced stops behave in a way analogous to voiceless stops, the plural forms of “lob,” “weg,” and “nemed” would be “loβa,” “weɣa,” and “nemeða,” respectively.

In summary, GPT-4 performs remarkably well on the spirantization problem. Although the task was relatively simple, the language of the task was constructed. Consequently, we can be certain that the model’s success was not due to memorization of the original data set. Moreover, the model correctly extends the spirantization rule to a class of sounds missing from the data set. This demonstrates a considerable ability for abstraction and generalization in the domain of phonological analysis.

GPT-4 outclasses GPT-3.5. Given the initial Prompt 8, the latter model hallucinates a vowel harmony process, completely failing the task (Reply 13). It is not the case that GPT-3.5 unexpectedly picks up on subtle vowel distribution asymmetries in the data set. Rather, the model’s answer is completely nonsensical. For example, it identifies the last vowel in [raɸa] as back and in [vixa] as front despite the fact that they are the same vowel.

**Reply 13:** One phonological process that may be operating in this language is vowel harmony, where vowels in a word are influenced by the vowels in neighboring syllables. Based on the plural forms provided, it seems that there may be a front-back vowel harmony system, where the suffix /-a/ has a front or back variant depending on the vowel in the previous syllable.

For example, the singular noun “rap” has a back vowel /a/, and the plural form “raϕa” also has a back vowel in the suffix. Similarly, the singular noun “vik” has a front vowel /i/, and the plural form “vixa” also has a front vowel in the suffix.

We can see this pattern consistently across the plural forms provided, suggesting that vowel harmony is a regular feature of the language. Further analysis would be needed to determine the specific rules and constraints of the vowel harmony system, such as whether it is restricted to certain syllables or environments, or whether certain vowels are dominant over others in determining harmony.

## 5 LAMBDA CALCULUS

Formal semantics is a subfield of both linguistics and logic which studies how the meaning of complex expressions is derived from the meaning of their constituent parts and the structure of those expressions. Semanticists employ mathematical tools to analyze and manipulate linguistic expressions with precision. One such tool is lambda calculus, which is a formal system for expressing computation using function abstraction and application. Lambda calculus plays a crucial role in modeling the compositional structure of meaning in natural language.

In this section, we tested GPT-4’s ability to construct lambda calculus analyses of English sentences, including cases of scopal ambiguity. Overall, the model performs well on simpler sentences and understands scopal ambiguity, but makes some substantive mistakes in using the lambda calculus formalism.

In [Prompt 14](#), we asked GPT-4 to provide an analysis of a simple quantified English sentence. The model’s [Reply 14](#) is largely correct, although there are notational and presentational issues (to be discussed at the end of this section).

**Prompt 14:** Write the lambda calculus analysis of the sentence “Every plant is yellow” in LaTeX.

**Reply 14:** Every plant is yellow  
 $\equiv (\forall x)((\text{Plant}(x) \rightarrow \text{Yellow}(x)))$   
 $\equiv \forall x((\lambda x.\text{Plant}(x))(x) \rightarrow (\lambda x.\text{Yellow}(x))(x))$

Next, we asked GPT-4 to analyze a sentence which involves scopal ambiguity ([Prompt 15](#)). The two readings of “Every student likes a classmate” are “For every student, there is a classmate that they like” (with universal quantification scoping over existential) and “There exists a classmate such that every student likes them” (with the opposite scope).

**Prompt 15:** Write the lambda calculus analysis of the sentence “Every student likes a classmate” in LaTeX.

The answer provided by the GPT-4 presents an analysis of the former reading, with the universal quantifier scoping high ([Reply 15](#)). Nevertheless, the answer contains two mistakes. The most important error is that the model existentially quantifies over “Student” in the consequent. Instead, there should only be an existential quantification over the relational noun “Classmate.” Additionally, if “Likes” and “Classmate” are interpreted as two-place functions, they must be applied to their arguments one by one, not to two arguments at the same time.

**Reply 15:** Every student likes a classmate

$$\begin{aligned} &\equiv (\forall x)(\text{Student}(x) \rightarrow (\exists y)(\text{Student}(y) \wedge \text{Classmate}(x, y) \wedge \text{Likes}(x, y))) \\ &\equiv \forall x((\lambda x.\text{Student}(x))(x) \rightarrow (\exists y)((\lambda y.\text{Student}(y))(y) \wedge \\ &\quad (\lambda x.\lambda y.\text{Classmate}(x, y))(x, y) \wedge (\lambda x.\lambda y.\text{Likes}(x, y))(x, y))) \end{aligned}$$

When explicitly asked about another meaning of the sentence ([Prompt 16](#)), GPT-4 recognizes the existence of scopal ambiguity and provides an analysis of the other reading. [Reply 16](#) doesn’t involve extraneous quantification but repeats the mistake of applying two-place functions to both of their arguments at the same time.

**Prompt 16:** Is there another reading for this sentence?

**Reply 16:** Yes, there is another reading for the sentence “Every student likes a classmate.” The initial interpretation assumes that each student likes a possibly different classmate. However, another reading could be that there exists a particular classmate whom every student likes. Here’s the lambda calculus analysis for this alternative interpretation:

Every student likes a classmate

$$\begin{aligned} &\equiv (\exists y)(\text{Student}(y) \wedge (\forall x)(\text{Student}(x) \rightarrow (\text{Classmate}(x, y) \wedge \text{Likes}(x, y)))) \\ &\equiv \exists y((\lambda y.\text{Student}(y))(y) \wedge (\forall x)((\lambda x.\text{Student}(x))(x) \rightarrow \\ &\quad \lambda x.\lambda y.\text{Classmate}(x, y)(x, y) \wedge (\lambda x.\lambda y.\text{Likes}(x, y))(x, y))) \end{aligned}$$

Finally, the style and notation used by the model leave a lot to be desired. First, there should be no parentheses around the quantifiers;  $(\forall x)$  and  $(\exists y)$  are not permitted by the syntax of predicate logic. Second, the English sentences should be enclosed in double brackets  $\llbracket \quad \rrbracket$ , which represent the interpretation function. Lastly, equals signs  $(=)$  should be used instead of the equivalence symbol  $(\equiv)$ .

In summary, GPT-4 can provide semantic analyses of English sentences and recognize scopal ambiguities. The model makes lambda calculus errors and uses unconventional notation.

These issues parallel the model’s problems with consistent presentation of syntactic trees and correct application of phonological notation.

## 6 PITFALLS AND LIMITATIONS

There are several potential pitfalls when analyzing the metalinguistic abilities of LLMs. The first to consider is the problem of *memorization*. It is likely that many of the examples we chose existed in the corpus that GPT was trained on. Simple problem sets that can be found in linguistics textbooks and Wikipedia articles also have readily available solutions printed online, which GPT has almost definitely seen. This is why it is crucial to also test LLMs on data that they are very *unlikely* to have seen—either by creating new versions of these problems with similar structures, or by wholly inventing new languages and new grammars.

It is also possible that the learning in LLMs is not fully unsupervised: it is possible that GPT-4 is trained on large amounts of syntactically labeled data. Due to OpenAI’s policy to not release any details about training data, it will be difficult to evaluate confounding factors based on training data. We can ask GPT-4 itself this question ([Prompt 17](#)), but of course its answer ([Reply 17](#)) cannot be trusted.

**Prompt 17:** Were you trained on any labeled syntactic data?

**Reply 17:** As an AI language model based on the GPT-4 architecture, I was trained on a large dataset of text from various sources. While the primary training data consists of raw, unlabeled text, the model does benefit indirectly from labeled syntactic data present in the training corpus.

The training data may include texts that contain syntactic annotations or examples, such as educational materials or linguistic research papers. By learning from these examples, the model can develop an understanding of syntactic structures and linguistic concepts. However, it’s important to note that the primary focus of GPT models is on learning patterns in natural language rather than being explicitly trained on labeled syntactic data. (This may be less of a problem in phonology, since phonetic or phonological annotations are much rarer in the written corpus.)

LLM output can also be highly variable and unpredictable. This is particularly true for GPT. Using the same prompts in new chats can result in substantially different outputs. It’s not clear whether this is due to an inherent instability in GPT’s next-token-prediction process or a programmed discourse function that GPT is emulating with its limited memory. GPT may treat separate, repeated prompts—even across different chat instances—as requests for new or different information. We speculate that it may, in effect, be performing a kind of pragmatic inference by inferring our intentions as we repeatedly ask the same prompts.

## 7 CONCLUSION

Large language models not only perform well on predicting the next sequence to generate coherent text, but they also are able to generate theoretical linguistic analyses with various degrees of accuracy. This paper proposes a new research program that tests the metalinguistic abilities of LLMs. We argue that GPT-4 is the first LLM that can generate coherent metalinguistic analyses. We present several prompt designs that elicit GPT-4's performance on syntactic, phonological, and semantic tasks.

We believe the metalinguistic abilities of LLMs will also be useful in the classroom. Many instructors have already begun incorporating ChatGPT into their coursework, typically by using it to generate text which students can then critique or improve. This can be especially useful in linguistics education. GPT is capable of generating theoretical analyses in several different linguistic domains. Here, we show that these analyses are often far from perfect, which makes them useful educational tools, by giving students the opportunity to take the role of the instructor and critically evaluate the model's mistakes. This task is not only useful for linguistic purposes, but for learning how to evaluate and interpret models' performance—a task that we believe will be increasingly important in the future.

## BIBLIOGRAPHY

- Ajdukiewicz, Kazimierz (1935). "Die syntaktische Konnexität." In: *Studia philosophica* 1. English translation in S. McCall (ed.), *Polish logic 192–1939*, 207–231. Oxford: Oxford University Press, pp. 1–27.
- Beguš, Gašper (May 2023). "GPT4 draws a syntactic tree for the sentence: ..." Accessed on April 17, 2023. URL: <https://twitter.com/begusgasper/status/1635719937275871233?s=20>.
- Beguš, Gašper, Andrej Leban, and Shane Gero (2023). "Approaching an unknown communication system by latent space exploration and causal inference." In: *arXiv* 2303.10931. eprint: [2303.10931](https://arxiv.org/abs/2303.10931) (stat.ML).
- Berko, Jean (1958). "The child's learning of English morphology." In: *Word* 14.2-3, pp. 150–177.
- Brown, Tom B. et al. (2020). *Language Models are Few-Shot Learners*. arXiv: [2005.14165](https://arxiv.org/abs/2005.14165) [cs.CL].
- Chomsky, Noam (1993). *Lectures on government and binding: The Pisa lectures*. Studies in Generative Grammar [SGG] 9. Walter de Gruyter. doi: <https://doi.org/10.1515/9783110884166>.
- Chomsky, Noam (2014). *The Minimalist Program*. MIT press.
- Chomsky, Noam, Roderick A Jacobs, and Peter S Rosenbaum (1970). "Remarks on nominalization." In: *1970* 184, p. 221.
- Davies, Alex et al. (2021). "Advancing mathematics by guiding human intuition with AI." In: *Nature* 600.7887, pp. 70–74. doi: [10.1038/s41586-021-04086-x](https://doi.org/10.1038/s41586-021-04086-x). URL: <https://doi.org/10.1038/s41586-021-04086-x>.
- Fillmore, Charles J, Paul Kay, and Mary Catherine O'connor (1988). "Regularity and idiomaticity in grammatical constructions: The case of let alone." In: *Language*, pp. 501–538.
- Heim, Irene and Angelika Kratzer (1998). *Semantics in generative grammar*. Wiley-Blackwell. ISBN: 978-0-631-19713-3.

- Jackendoff, Ray (1977). *X-bar Syntax: A study of phrase structure*. Cambridge, MA: MIT press.
- Jumper, John et al. (2021). “Highly accurate protein structure prediction with AlphaFold.” In: *Nature* 596.7873, pp. 583–589. DOI: [10.1038/s41586-021-03819-2](https://doi.org/10.1038/s41586-021-03819-2). URL: <https://doi.org/10.1038/s41586-021-03819-2>.
- Katzir, Roni (2023). “Why large language models are poor theories of human linguistic cognition. A reply to Piantadosi (2023).” In: *Ms., Tel Aviv University*. [lingbuzz/007190](https://lingbuzz.007190.com).
- Matushevych, Yevgen and Jennifer Culbertson (2022). “Trees neural those: RNNs can learn the hierarchical structure of noun phrases.” In: *Proceedings of the Annual Meeting of the Cognitive Science Society* 44.44, pp. 1848–1854.
- McElfatrick, Dave (Mar. 2008). *Cyanide & Happiness Comic*. URL: <https://files.explosm.net/comics/Dave/comicshthead3.png>.
- OpenAI (2023). *GPT-4 Technical Report*. arXiv: [2303.08774](https://arxiv.org/abs/2303.08774) [cs.CL].
- Perlmutter, David M. (1980). “Relational grammar.” In: *Current Approaches to Syntax*. Brill, pp. 195–229.
- Piantadosi, Steven (Mar. 2023). “Modern language models refute Chomsky’s approach to language.” In: [lingbuzz 007180](https://lingbuzz.007180.com).
- Prince, Alan and Paul Smolensky (1993/2004). *Optimality Theory: Constraint Interaction in Generative Grammar*. First published in Tech. Rep. 2, Rutgers University Center for Cognitive Science. Malden, MA: Blackwell.
- Tantau, Till (2007). *TikZ and pgf Manual for version 1.18*. URL: <https://www.bu.edu/math/files/2013/08/tikzpgfmanual.pdf>.
- Tunmer, William E. and Michael L. Herriman (1984). “The Development of Metalinguistic Awareness: A Conceptual Overview.” In: *Metalinguistic Awareness in Children: Theory, Research, and Implications*. Ed. by William E. Tunmer, Christopher Pratt, and Michael L. Herriman. Berlin, Heidelberg: Springer Berlin Heidelberg, pp. 12–35. ISBN: 978-3-642-69113-3. DOI: [10.1007/978-3-642-69113-3\\_2](https://doi.org/10.1007/978-3-642-69113-3_2). URL: [https://doi.org/10.1007/978-3-642-69113-3\\_2](https://doi.org/10.1007/978-3-642-69113-3_2).
- Vaswani, Ashish, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin (2017). “Attention is All you Need.” In: *Advances in Neural Information Processing Systems*. Ed. by I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett. Vol. 30. Curran Associates, Inc. URL: <https://proceedings.neurips.cc/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf>.
- Wilcox, Ethan, Roger Levy, Takashi Morita, and Richard Futrell (Nov. 2018). “What do RNN Language Models Learn about Filler–Gap Dependencies?” In: *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*. Brussels, Belgium: Association for Computational Linguistics, pp. 211–221. DOI: [10.18653/v1/W18-5423](https://doi.org/10.18653/v1/W18-5423). URL: <https://aclanthology.org/W18-5423>.
- Yedetore, Aditya, Tal Linzen, Robert Frank, and R. Thomas McCoy (2023). “How poor is the stimulus? Evaluating hierarchical generalization in neural networks trained on child-directed speech.” In: *arXiv preprint arXiv:2301.11462*.
- Živanović, Sašo (2017). *Forest*. <https://ctan.org/pkg/forest>. Version 2.1.5. LaTeX package for drawing (linguistic) trees.

ACKNOWLEDGEMENTS G.B. would like to thank students in his *Linguistics 111: Phonology* and *Linguistics 211A: Advanced Phonology I* classes at UC Berkeley. We also thank Zuzanna Fuchs and Yimei Xiang for their feedback on syntactic and semantic analyses. This research was partially funded by the Georgia Lee Fellowship in the Society of Hellman Fellows to G.B.