

# Do LLMs Understand Syntactic Center Embedding? Should They?

Anonymous ACL submission

## Abstract

We consider the case of syntactic center embedding, where an embedding phrase contains material on both sides of the embedded phrase. While a single center embedding is easily understandable for human language users, multiple center embeddings are generally uninterpretable. Despite this, it has been claimed that multiple embeddings are in fact grammatically acceptable. We construct sentences with center embeddings of varying levels, ranging from 1-4, and we find that GPT-3.5, like humans, interprets level 1 sentences correctly, but fails with higher levels. On the other hand, GPT-4 achieves superhuman accuracy levels, with nearly perfect results even with 3 or 4 levels of embeddings. We suggest that this raises relevant questions about the relation of LLMs to the human language faculty.

## 1 Introduction

Recursive syntactic structures are fundamental to natural language. A propositional verb like “believe” can take a sentence as its complement to its right, and that sentential complement might itself involve such a structure, as in (1):

- (1) a. [John believes [Harry likes fish]]
- b. [John believes [Tom said [everyone knows ... [Harry likes fish]]]]

An adverbial phrase like “in the library” can modify a verb phrase to its left; the modified verb phrase might itself contain such a modifier, as shown by (2):

- (2) a. Col. Mustard [[[killed Mr Boddy] in the library]
- b. Col. Mustard [[[[killed Mr Boddy] with the candlestick] in the library] ... without remorse.]

The above cases illustrate the potential for unbounded levels of embedding. In example (1), the

embedding clause contains material to the left of the embedded clause, and in (2), the embedding clause contains material to the right. A third possibility is center embedding, where the embedding clause contains material both to the left and right of the embedded clause. This is illustrated by (3). Here, a nominal expression, “teacher”, is modified by a relative clause, “the student saw”.<sup>1</sup>

- (3) [The teacher [the student saw  $t$ ] is happy.]  
**Level 1**

Multiple levels of center embedding are readily constructed. Examples (4) - (6) represent levels 2-4 of center embedding.

- (4) [The teacher [the student [the driver hit  $s$ ] saw  $t$ ] is happy.] **Level 2**
- (5) [The teacher [the student [the driver [the girl likes  $d$ ] hit  $s$ ] saw  $t$ ] is happy.] **Level 3**
- (6) [The teacher [the student [the driver [the girl [the man hates  $g$ ] likes  $d$ ] hit  $s$ ] saw  $t$ ] is happy.] **Level 4**

Such multiple center embeddings, while easy to construct, are generally uninterpretable for human language users, and are virtually non-existent in normal texts. This is strikingly different from multiple left and right embeddings, which are generally easy to interpret, and not at all unusual.

Although syntactic center embedding has received little recent attention in the NLP literature, it has special significance in theoretical linguistics. Despite the evident inability of human language users to interpret multiple center embeddings, it has been widely claimed that they are in fact fully grammatical. Famously, Chomsky has explained

<sup>1</sup>The relative clause “the student saw” includes a trace or variable, which we indicate with  $t$  to show that it in this case is bound by “the teacher”, and similarly with the variables  $s$ ,  $d$ , and  $g$  in examples (4) - (6), standing for “student”, “driver” and “girl”, respectively.

072 this apparent paradox by arguing that center em- 121  
073 beddings are completely acceptable according to 122  
074 human linguistic *competence*, attributing their evi- 123  
075 dent difficulty to limitations in *performance*. These 124  
076 claims are central to the very founding of modern 125  
077 linguistics (Chomsky, 1957; Chomsky et al., 1963).

078 In this paper, we explore whether large language 126  
079 models (LLMs) can interpret such structures. We 127  
080 find that GPT3-5 is rather similar to humans, per- 128  
081 forming very well with level 1 center embeddings, 129  
082 but very poorly with any higher levels. On the other 130  
083 hand, GPT-4 performs extremely well at all levels, 131  
084 from 1 to 4. We consider two possible explana- 132  
085 tions for this; the first is simply that GPT-4 has 133  
086 achieved superhuman linguistic abilities. The sec- 134  
087 ond explanation is that GPT-4 has exactly captured 135  
088 human linguistic *competence*, but is not subject to 136  
089 the same *performance* limitations as humans. 137

## 090 2 Related Work 138

### 091 2.1 Syntactic Center Embedding 139

092 Karlsson (2007, p. 365) notes that “A common 140  
093 view in theoretical syntax and computational lin- 141  
094 guistics holds that there are no grammatical re- 142  
095 strictions on multiple center-embedding of clauses.” 143  
096 Indeed, Karlsson (p. 368) sees this as “the main- 144  
097 stream view...voiced by many linguists from dif- 145  
098 ferent camps”. This view derives from the ear- 146  
099 liest work in modern linguistics; most famously, 147  
100 Chomsky (1957) argues that the grammar of En- 148  
101 glish permits unbounded center-embedding. This 149  
102 claim plays a central role in Chomsky’s argument 150  
103 that English is a context-free rather than a finite- 151  
104 state language. For example, Chomsky et al. (1963) 152  
105 present sentence (7), which is an example of level 153  
106 2 center embedding: 154

107 (7) The rat the cat the dog chased killed ate the 155  
108 malt. 156

109 In the view of Chomsky et al., example (7) “is 157  
110 surely confusing and improbable but it is perfectly 158  
111 grammatical and has a clear and unambiguous 159  
112 meaning.” This argument relies on the Chom- 160  
113 skyan distinction between competence and perfor- 161  
114 mance, where competence is an idealized theory 162  
115 of the “mental reality underlying actual behavior”. 163  
116 (Chomsky, 2014)[p. 4] Performance factors, such 164  
117 as memory limitations, might make the underlying 165  
118 linguistic competence difficult to observe, much as 166  
119 friction makes it difficult to observe the underly-  
120 ing nature of Newton’s law of gravity. The theory

of linguistic competence, on this view, correctly 121  
permits unbounded center embedding. The fact 122  
that humans nevertheless encounter difficulty, is 123  
ascribed to performance factors. 124

### 125 2.2 Linguistic Probing of LLMs 126

127 There is a large literature describing the probing 128  
129 of LLMs for specific linguistic capabilities or char- 130  
131 acteristics. Mahowald et al. (2023) has suggested 132  
133 that current LLMs have largely mastered what they 134  
135 call “formal linguistic competence”. However, sev- 136  
137 eral recent works have shown that there remain 138  
139 specific capabilities that pose difficulties for some 140  
141 of the most powerful current models. For example 142  
143 Hardt (2023) probes LLMs in their understanding 144  
145 of elliptical sentences by posing a Yes-No question 146  
147 that relies on a correct understanding of an ellipti- 148  
149 cal construction. Hardt concludes that LLMs still 150  
151 struggle with the phenomenon of ellipsis. Simi- 152  
153 larly, Cui et al. (2023) probe LLMs with construc- 154  
155 tions involving “respectively”; testing models on 156  
157 their ability to draw correct inferences based on 158  
159 the logic of respectively. They find that the models 160  
161 they tested have substantial difficulties in this tasks. 162

## 163 3 Data 164

165 We construct a synthetic dataset, consisting of a 166  
167 context, a prompt and a question. <sup>2</sup> 168

### 169 3.1 Context 170

171 The context consists of synthetic examples of cen- 172  
173 ter embedding of levels 1-4, as illustrated above by 174  
175 examples (3) - (6). The form of these examples is 176  
177 as follows, where N is noun, TV is transitive verb 177  
178 and IV is intransitive verb: 178

179 **Level 1:** The N the N TV IV. 179

180 **Level 2:** The N the N the N TV TV IV. 180

181 **Level 3:** The N the N the N the N TV TV TV 181  
182 IV. 182

183 **Level 4:** The N the N the N the N the N TV TV 183  
184 TV TV IV. 184

185 See A.2 for instantiations of N, TV, and IV. 185

### 186 3.2 Prompt 187

188 We define the prompt shown in figure 1, which 189  
190 we designate P1. The prompt includes a single 190  
191 example, exhibiting level 1 center embedding. This 191  
192 can be seen as 1-shot learning. 192

<sup>2</sup>Data and associated code will be made available on 193  
Github upon acceptance. 194

You will be given an example consisting of a context and a question to answer. The answer should always be of this form "The N V the N", where N stands for a single word that is a noun, and V stands for a single word that is a verb. Here is a sample:

Context: The student the man saw is happy  
 Question: Who saw who?  
 Answer: The man saw the student.

Context: {context}  
 Question: {question}  
 Now answer the question:

Figure 1: Prompt P1, containing one sample case

We wish to investigate whether the provision of the sample has an effect on model performance. Thus we define a second prompt that lacks a sample. This prompt is designated as P0 (figure 2).

You will be given an example consisting of a context and a question to answer. The answer should always be of this form "The N V the N", where N stands for a single word that is a noun, and V stands for a single word that is a verb.

Context: {context}  
 Question: {question}  
 Now answer the question:

Figure 2: Prompt P0, with no sample

### 3.3 Question

For all our examples, we formulate a question of the form "Who TV'ed who", where the verb TV is taken from the most deeply embedded clause. We designate this question as Q0 (figure 3). We define an alternative question that targets the next most deeply embedded clause, which we designate Q1 (figure 4). Note that Q1 is not applicable for level 1.

## 4 Test

For each embedding level (1-4), we construct 500 synthetic examples, and we test both GPT-3.5 and GPT-4 (GPT). Our initial test uses prompt P1 and question Q0. We also report on tests with alternative versions of both the prompt and question in different combinations.

**Level 1**  
 Context: The teacher the student saw is happy  
 Q: Who saw who?  
 A: the student saw the teacher.

**Level 2**  
 Context: The teacher the student the driver saw hit is happy  
 Q: Who saw who?  
 A: the driver saw the student.

**Level 3**  
 Context: The teacher the student the driver the girl saw hit likes is happy  
 Q: Who saw who?,  
 A: the girl saw the driver.

**Level 4**  
 Context: The teacher the student the driver the girl the man saw hit likes hates is happy  
 Q: Who saw who?  
 A: the man saw the girl.

Figure 3: Four Embedding Levels with Question Q0, targeting most deeply embedded structure

**Level 2**  
 Context: The teacher the student the driver saw hit is happy  
 Q: Who hit who?  
 A: the student hit the teacher.

**Level 3**  
 Context: The teacher the student the driver the girl saw hit likes is happy  
 Q: Who hit who?  
 A: the driver hit the student.

**Level 4**  
 Context: The teacher the student the driver the girl the man saw hit likes hates is happy  
 Q: Who hit who?  
 A: the girl hit the driver.

Figure 4: Embedding Levels 2-4 with Question Q1, targeting the next most deeply embedded structure

In figure 5 we present results for GPT-4 and GPT-3.5 for the four levels of embedding, with prompt P1 and question Q0. Both models are perfectly accurate for level 1 examples. Such examples tend to be very easy for humans. For GPT-3.5, accuracy falls sharply for levels 2 and 3, and is even lower for level 4. GPT-4 is far more accurate with higher levels of embedding – nearly perfect for levels 2 and 3, and still highly accurate (0.85) for level 4. This is striking, as these levels of embedding are not interpretable by human language users. Furthermore, multiple embeddings are almost certainly vanishingly rare in the training data for these models. In an extensive corpus study, Karlsson (2007)[p. 378] found that “in ordinary language use, written C3s [level 3] and spoken C2s [level 2] are almost non-existent”.

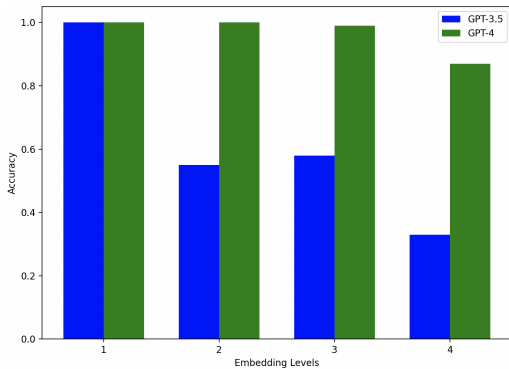


Figure 5: Accuracy of Center Embedding at levels 1-4, with Prompt P1 and Question Q0. GPT-4 is highly accurate even up to level 4, while GPT-3.5 is degraded at all levels above 1. (500 examples for each model, for each level)

#### 4.1 Alternative Prompts and Questions

Model	P	Q	L1	L2	L3	L4
GPT-3.5	P0	Q0	0.86	0.49	0.34	0.14
GPT-3.5	P0	Q1	-	0.03	0.03	0.03
GPT-3.5	P1	Q0	1.00	0.55	0.58	0.33
GPT-3.5	P1	Q1	-	0.16	0.03	0.06
GPT-4	P0	Q0	1.00	0.55	0.28	0.11
GPT-4	P0	Q1	-	0.17	0.02	0.00
GPT-4	P1	Q0	1.00	1.00	0.99	0.87
GPT-4	P1	Q1	-	0.74	0.05	0.00

Table 1: Accuracy by Model and Embedding Level. (500 examples for each model, for each level)

In table 1 we present the accuracy of the two models with alternative prompt and question forms.

In general, it is clear that both models are quite sensitive to these variations, in ways we are not in a position to explain. We would, however, like to draw attention to one specific observation: while the GPT-4 model achieves extremely high levels of accuracy with prompt P1 and question Q0, these levels drop precipitously with prompt P0 and question Q0, for all except level 1. We find this rather astonishing, since the only difference is that the model in the former case is provided with a single level 1 example, which is absent in the latter case. Somehow a single level 1 example has enabled GPT-4 to master higher levels of center embedding.

## 5 Conclusions

While multiple embedding structures are ubiquitous in human language, multiple center embeddings are different: they almost never occur, and are almost always uninterpretable for human language users. It has nonetheless been steadfastly maintained that they are grammatical, according to mainstream theories of human linguistic competence. In this paper, we have shown that GPT-3.5 struggles with center embeddings of any level greater than 1, much like humans, while GPT-4 performs very well with all four levels of center embeddings, thus apparently far exceeding human abilities. Why should this be?

One straightforward response is that GPT-4 is simply too big – at least with respect to its linguistic competence, the size of training data and number of system parameters is simply larger than needed, since it can now process linguistic structures that are far too complicated for humans.

There is another way to look at this, however. Chomsky famously argued that center embeddings are completely grammatical according to the theory of human linguistic competence. Humans, on this view, have a grammar that allows deeply embedded center embeddings, but this fact is obscured by performance factors – limitations on the general computational system in which the human language faculty is implemented. If the same linguistic competence could be implemented in a more powerful system, it would be easier to observe its true nature, since some of the performance limitations would be removed. Perhaps GPT-4 is just such a system: it has largely duplicated human linguistic competence, but is not subject to the same performance limitations as humans.

## 6 Limitations

The paper seeks to determine whether LLMs understand syntactic center embedding, but this general question is explored in only a few particular ways. First, only two LLMs are considered, and we suspect that other models might give quite different results. However, GPT-4 is the most powerful model we had access to, and we suspect that other less powerful models would, like GPT.3-5, have great difficulty with the tests reported on here. There are also several important limitations with respect to the data. First, the data is solely English. Second, it is synthetic data, constructed according to a template that reflects one specific form of center embedding, in which a noun phrase is modified by a relative clause. We believe this is the form of center embedding that is most familiar from the linguistics literature. However, there are other forms of center embedding that could also be considered. Furthermore, while we explored various combinations of different prompt and question forms, there are other forms and combinations that would be well worth exploring. Finally, we have made claims about the general uninterpretability of multiple center embeddings for humans; while these generally echo claims made in the literature, they are claims that would benefit from rigorous empirical examination.

## References

- [GPT Models Overview](#). Accessed on 2023-12-10.
- Noam Chomsky. 1957. *Syntactic structures*. The Hague: Mouton.
- Noam Chomsky. 2014. *Aspects of the Theory of Syntax*. 11. MIT press.
- Noam Chomsky, George Armitage Miller, R Luce, R Bush, and E Galanter. 1963. Introduction to the formal analysis of natural languages. *1963*, pages 269–321.
- Ruixiang Cui, Seolhwa Lee, Daniel Hershcovich, and Anders Søgaard. 2023. [What does the failure to reason with “respectively” in zero/few-shot settings tell us about language models?](#) In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8786–8800, Toronto, Canada. Association for Computational Linguistics.
- Daniel Hardt. 2023. Ellipsis-dependent reasoning: a new challenge for large language models. In *Proceedings of the 61st Annual Meeting of the Association for*

*Computational Linguistics*, pages 39–47. Association for Computational Linguistics.

Fred Karlsson. 2007. Constraints on multiple center-embedding of clauses. *Journal of Linguistics*, 43(2):365–392.

Kyle Mahowald, Anna A Ivanova, Idan A Blank, Nancy Kanwisher, Joshua B Tenenbaum, and Evelina Fedorenko. 2023. Dissociating language and thought in large language models: a cognitive perspective. *arXiv preprint arXiv:2301.06627*.

## A Appendix

### A.1 Error Analysis

In all cases, the system is expected to produce answers of the form N1 V N2. We define three types of errors:

- Type 1: N1 is incorrect, N2 is correct
- Type 2: N1 is correct, N2 is incorrect
- Type 3: N1 is incorrect, N2 is incorrect

We consider selected settings based on a manual evaluation of the first 10 examples. Table 2 shows the percentage of errors of each type.

Model	P	Q	L	T1	T2	T3
GPT-3.5	P0	Q0	L1	0.00	0.00	1.00
GPT-3.5	P0	Q0	L2	0.10	0.90	0.00
GPT-3.5	P1	Q0	L2	0.00	0.90	0.10
GPT-3.5	P1	Q0	L3	0.00	0.90	0.10
GPT-3.5	P1	Q0	L4	0.00	0.90	0.10
GPT-4	P0	Q0	L2	0.00	0.90	0.10
GPT-4	P0	Q0	L3	0.00	0.90	0.10
GPT-4	P0	Q0	L4	0.00	0.80	0.20
GPT-4	P1	Q0	L4	0.00	0.90	0.10
GPT-4	P1	Q1	L2	0.40	0.00	0.60

Table 2: Error Types, T1, T2, T3 for selected settings of model, prompt type, question type and level of embedding (based on manual analysis of first 10 errors for each setting)

For all but two of the settings in table 2, nearly all the errors are of type T2, as in the following example:

Context: The man the girl the driver knows hates is glad. Question: Who knows who? Model Answer: The driver knows the man. Correct Answer: The driver knows the girl.
--

329 Since the verb “knows” is explicit in the question,  
330 the model could simply assume that N1 is the noun  
331 phrase preceding “knows” in the context. This as-  
332 sumption ensures that a model avoids T1 errors,  
333 for question Q0. A T2 error arises in the above ex-  
334 ample, because the model selects “the man” rather  
335 than “the girl” as the second NP. Interestingly, GPT-  
336 3.5 has *only* T3 errors in the setting, P0, Q0, L1. In  
337 each case, it simply reverses N1 and N2, as in the  
338 following example:

339 Context: The woman the man hates left.  
Question: Who knows who?  
Model Answer: The woman hates the man.  
Correct Answer: the man hates the woman.

340 Finally, GPT-4 has *only* T1 or T3 error types  
341 on the setting P1, Q1, L2. The following example  
342 illustrates a T3 error for this setting:

343 Context: The student the man the driver hates  
saw is glad.  
Question: Who saw who?  
Model Answer: The student saw the man.  
Correct Answer: the man saw the student.

344 We have, of course, no direct insight into the  
345 strategies employed by these large language mod-  
346 els in any of these settings. It seems intuitively  
347 plausible that models employ a strategy would nor-  
348 mally get N1 right and N2 wrong, and this is indeed  
349 the pattern that arises with this limited error analy-  
350 sis. At this point we will offer no speculation about  
351 the two settings for which we observe different  
352 error patterns.

## 353 A.2 Sample Instantiations

354 We have the following substitutions for N and TV.  
355 N: (teacher, student, driver, girl, man), and TV:  
356 (saw, hit, likes, hates, knows). IV is always substi-  
357 tuted with the phrase, "is happy".