

Evidence of Generative Syntax in Large Language Models

Anonymous ACL submission

Abstract

The syntactic probing literature has been largely limited to shallow structures like dependency trees, which are unable to capture the subtle differences in sub-surface syntactic structures that yield semantic nuances. These structures are captured by theories of syntax like generative syntax, but have not been researched in the LLM literature due to the difficulties in probing these complex structures that have many silent, covert nodes. Our work presents a method for overcoming this limitation by deploying [Hewitt and Manning \(2019\)](#)’s dependency-trained probe on sentence constructions whose structural representation is identical in a dependency parse, but differs in theoretical syntax. If a pretrained language model has captured the theoretical syntax structure, then the probe’s predicted distances should vary in syntactically-predicted ways. Using this methodology and a novel dataset, we find evidence that LLMs have captured syntactic structures far richer than previously realized, indicating LLMs are able to capture the nuanced meanings that result from sub-surface differences in structural form.

1 Introduction

Large Language Models (LLMs) have demonstrated remarkable language capabilities that have been steadily increasing ever since BERT ([Devlin et al., 2019](#)). This impressive performance has prompted a body of research interested in investigating *why* these models are so successful. From this came a subset of research seeking to understand what, if any, linguistic features or knowledge these models have acquired ([Jawahar et al., 2019](#); [Belinkov and Glass, 2019](#); [He et al., 2024](#); [Kallini et al., 2024](#)) as means to better understand their language performance. The focus of linguistic inquiry can vary from the semantic ([Nikolaev and Padó, 2023](#); [Kamath et al., 2024](#)) to the morphological ([Coleman, 2020](#); [Anh et al., 2024](#)) to the syntactic

([Clark et al., 2019](#); [Chi et al., 2020](#); [Kulmizev et al., 2020](#); [Maudslay and Cotterell, 2021](#)), the latter of which our current research seeks to extend.

Much of the research into syntactic representation in LLMs have utilized dependency parses to represent a sentence’s syntactic structure ([Hewitt and Manning, 2019](#); [Maudslay and Cotterell, 2021](#); [Tucker et al., 2022](#); [Eisape et al., 2022](#); [Buder-Gröndahl, 2024](#)). However, these relatively shallow representations can fail to capture features of a sentence and the nuanced differences in meaning that result from different sub-surface syntactic structures. Adopting a dependency framework makes theoretical assumptions and imposes limitations on the richness of meaning that can be expressed, the consequences of which are often not addressed. These simple, compact representations stand in stark contrast to the deeper, hierarchically-complex structures that are posited in theoretical syntax, particularly the generative frameworks, which postulate these complexities in order to account for the difference in semantic meaning and syntactic patterning of certain syntactic constructions. Because the structures posed by theoretical syntax are far more complex with more tree nodes than words in the sentence, attempts to probe for generative syntactic structures have been stymied, and it remains unclear whether LLMs have captured these richer sub-surface structures.

Our work seeks to overcome this through our unprecedented application of [Hewitt and Manning \(2019\)](#)’s dependency-trained probe to test for theory-backed syntactic structures. To implement this, we identified two sentences structures—Subject Raising (SR) and Subject Control (SC)—whose surface and dependency representations are identical, but *whose generative structures differ* as their complement sizes differ (control predicates take larger complements than raising predicates). Using this method, we are able to circumvent the issue of handling empty nodes in the syntactic tree

as the probe only recovers dependency distances. Our hypothesis holds that if the probe predicts distances for the two sentence types that significantly differ in syntactically-predicted ways, then this is evidence that LLMs have captured the more complex structures of generative syntax, which can yield structures to explain the pattern variations and semantic nuances of certain constructions.

To test this, we developed a novel dataset of over 33,000 SR/SC sentences, which we fed into our dependency-trained probes using pretrained models of BERT, RoBERTa, and GPT2. When we probe the content words, our results find strong evidence of the SR/SC difference that aligns with SC taking a larger complement than SR as generative syntax proposes. However, probing the distances with infinitival "to" suggest that the functional word may encode syntactic structure in an aberrant manner. Together, our work suggests LLMs have encoded structure that is more complex than previously realized, and provides a novel method to probe for theoretical syntactic structure in LLMs.

2 Related Work

The impressive language abilities of recent LLMs have prompted researchers to ask whether this performance is due to some probabilistic modeling, or if these language models have managed to capture linguistic structures. To answer this question, a line of research known as **probing** was developed. This methodology feeds the model's contextualized vector representations into a neural network whose training objective is to predict a targeted linguistic structure from the representations alone (see [Alain and Bengio, 2017](#) or [Conneau et al., 2018](#) for example). The argument follows that if such a neural network probe is in fact able to predict the target pattern or structure, then it can be concluded that the language model has indeed implicitly learned that linguistic feature; otherwise, the probe task would have been doomed to failure.

This area of research has largely focused specifically on investigating whether models have learned to properly encode syntactic phenomenon ([Mueller et al., 2020](#); [Hu et al., 2020](#); [Warstadt et al., 2020](#); [Ravfogel et al., 2021](#); [Davis et al., 2022](#)). However, much of this structural syntactic research has relied on dependency parses as a means of representing syntactic structure ([Hewitt and Manning, 2019](#); [Chi et al., 2020](#); [Maudslay and Cotterell, 2021](#); [Tucker et al., 2022](#); [Eisape et al., 2022](#)), with one notable

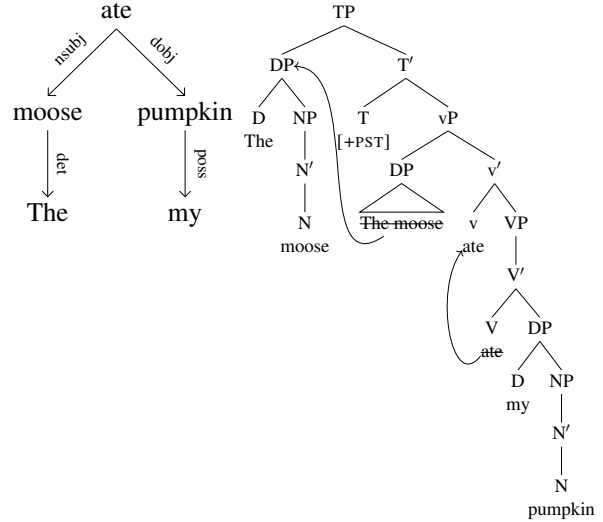


Figure 1: An example of the dependency tree (left) and generative syntax tree (right) for the sentence "The moose ate my pumpkin." Note how the dependency tree has a flatter structure with a one-to-one mapping of words to nodes in the tree. Compare this to the deeper generative tree where there are far more nodes in the tree than words in the sentence.

exception being [Arps et al. \(2022\)](#), which sought to (and largely succeeded) in training a probe to reconstruct a skeletal constituency tree. While this line of research is of value and great interest, there are theoretical assumptions made by using dependency parses, and there are limitations to using that particular syntactic framework.

2.1 Syntactic Theories

Dependency parses derive from French linguist Lucien Tesnière's (1959) theory of syntax known as Dependency Grammar (DG), which focuses on the head-dependent relationship between words (see Figure 1). In these trees, each word can have one and only one incoming arc that indicates it is the dependent of its head, excepting the root of the sentence (often the matrix verb), which has no head.

DG trees are relatively flat structures with one-to-one mappings between words in the sentence and nodes in the tree. The appeal of such trees are largely three-fold: (1) the representations are compact and efficient due to the one-to-one mapping, (2) learning to parse a dependency tree is relatively easy once one understands the head-dependent relationships that exist, and (3) the dependency tree does not need to capture the sentence's linear order of words. The last factor makes DG an appealing theory for researchers working on languages with freer word-order ([Müller, 2019](#)); however, this "fea-

ture" can become a bug when it loses nuance or creates ambiguous parses (see Figure 2).

An alternative, more structurally-rich approach to syntax has built off the theories of Chomsky (1957; 1981; 1986; 1995) and others who have refined this phrase-structure (also known as a constituency-based) framework to build up the syntactic framework known as generative syntax (GS). This family of syntactic theories are built on the X-bar theory, which proposes the operations *External Merge* and *Internal Merge* (formerly known as "Move"), and stipulates that nodes are binary-branching and that every phrase has a head (Chomsky, 1995). After all operations are applied in the course of derivation, the end result is the linearization of the sentence when read from left to right along the children nodes.¹ The generative framework is concerned with identifying the operations and rules that together generate licit sentences but do *not* generate illicit constructions.

Unlike DG, generative syntax and other phrase-structure grammars thus yield deeper, more complex trees with hierarchical structures and phonologically null nodes whose presence must be deduced through testing. While this complexity is well-warranted (in that it can generate sentences that are grammatical and explain what causes ungrammaticality), the tree's size and complexity creates a complicated and unwieldy structure that is difficult for non-linguists to implement.

The result of this has been a limitation in the scope of feasible GS research in Natural Language Processing (NLP). Even work that has sought to test for the deeper, more complex phrase-structures in NLP has largely focused either on only seeking to recover a phrase's boundaries (Tenney et al., 2019; Kallini et al., 2024) or has otherwise trained their probe on the overly-simplified, *n*-branching constituency trees of the English Penn Treebank (PTB) (Marcus et al., 1993), which are only annotated with "skeletal" syntactic structure schema that is relatively atheoretical. The PTB's annotation is often used for automatic conversion into a dependency parse with little issue since the simplicity

¹In earlier versions of GS, it was argued that when an element moves, it leaves behind a coindexed trace element (Chomsky, 1973; Fiengo, 1977). Because such an derivation introduces a *new* element, Chomsky (1993) revised the approach to the Copy Theory of Movement, claiming that movement leaves behind a copied element that is not phonologically realized. Thus, the sentence used above would be "Do you [do] [you] know [know] who [who] [who] stole [stole] the Crown Jewels?" For a discussion on which elements are phonologically realized and why, see Corver (2007).

- (1) Do you know who stole the Crown Jewels?
- (2) Who do you know stole the Crown Jewels?

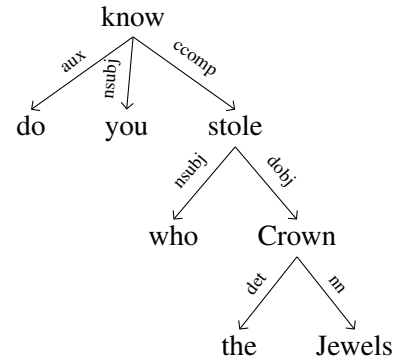


Figure 2: Above are two different sentences that yield *identical* dependency parses. While similar, the sentences have different meanings (imagine you are being questioned about the theft of the Crown Jewels: the first question merely seeks to inquire whether or not you know who the thief is, while the latter presumes you know who the thief is and seeks to learn the identity). This illustrates both the ways in which dependency trees do not capture linear order, and highlights some of the limitations of dependency trees.

of PTB syntax is non-problematic for the simpler structures and principles of an NLP dependency parse. However, the result of this is a corpora that is not theoretically-sound for many of the deeper linguistic inquiries into phrase-structure grammars.

Because of this, it has yet to be discovered whether LLMs have managed to capture the deeper, hierarchical structures of GS. However, there are three major barriers to testing whether models have captured the richer hierarchical structures as proposed by generative frameworks:

1. GS and similar frameworks often have "empty" nodes that are not overtly realized.² As such, they are not overtly present in the texts LLMs train on, and so probing at their presence is difficult because it raises the question: how can you probe at something that is not overtly represented?
2. Probes largely require a gold tree that indicates the correct structure or parse. Human annotation, while crucial when handling such fine-grained analysis, is laborious and costly

²This can be due to movement (see Footnote 1) or the feature not having an overt representation (e.g. there is no specific word or morpheme that indicates present tense for plural subjects as in "They *walk_{pres}* to the store"). See Figure 1 for demonstration.

in resources and time.

3. Even if one is able to secure the resources necessary to create such a gold standard, there are competing theories even within the generative framework that would change a sentence’s representation. As such, a gold parse would be subject to great theoretical scrutiny and likely face present or future dissenting opinions.

Our research develops a method that circumvents these obstacles while still addressing the fundamental research question of whether LLMs have captured something of the deeper, sub-surface syntactic representation theorized by many linguists.

3 Methodology

3.1 Probing Method

To combat the first issue of accounting for structures that are not phonologically realized, we have opted for the novel approach of re-purposing the original Hewitt and Manning (2019), which was trained to recover dependency trees, to investigate whether LLMs have encoded theoretically-motivated generative phrase-structure trees.

The structural probe developed by Hewitt and Manning (2019) proposes a model M that encodes a sequence of vector representations $\mathbf{h}_{1:n}^l$ from an input sequence of n words $\mathbf{w}_{1:n}^l$ where l identifies the sentence index. From there, they define a linear transformation matrix $\mathbf{B} \in \mathbb{R}^{k \times n}$ to parameterize the parse tree-encoding distances:

$$d_B(\mathbf{h}_i^l, \mathbf{h}_j^l)^2 = (\mathbf{B}(\mathbf{h}_i^l - \mathbf{h}_j^l))^T (\mathbf{B}(\mathbf{h}_i^l - \mathbf{h}_j^l))$$

where i and j are the words in the sentence and where the transformation matrix \mathbf{B} ’s objective is to reproduce the gold parse distances between each pair of words (w_i^l, w_j^l) for all sentences l in the parsed training corpus T^l . The training uses a gradient descent objective:

$$\min_{\mathbf{B}} \sum_l \frac{1}{|s^l|^2} \sum_{i,j} |d_{T^l}(w_i^l, w_j^l) - d_B(h_i^l, h_j^l)|^2$$

In this equation, $|s^l|$ is the length of sentences, which the function normalizes using the square of the sentence’s length as each sentence contains $|s^l|^2$ pairs of words. The probe’s objective thus seeks to approximate a matrix of distances that most closely resembles the gold-standard distances. Because Hewitt and Manning (2019) use a

dependency parse, gold-standard parses were converted into gold-standard distance matrices where distances are defined such that the distance between a parent node and its child nodes is 1, the distance between a child node and its grandparent node is 2, the distance between a child node and a so to speak “aunt” or “uncle” node is 3, and so on and so forth. Evaluation of the probe involved calculating the minimum spanning tree for each sentence’s predicted distances to derive the sentence’s predicted undirected, unlabeled attachment score (UAS) compared to the gold tree, and the average Spearman correlation of the predicted matrix of distances compared to the gold-standard matrix.

We chose this method specifically *because* it is a probe trained *only* to capture dependency parses with their one-to-one mappings between a sentence’s words and a tree’s nodes. Though there are critical limitations to dependency parses as discussed in Section 2.1, we argue that its simplicity and overgeneralization can in fact be converted into a benefit. It is because the probe is superficially only supposed to capture shallow-level, generalized syntactic structures that we can use the method to tease apart syntactic structures whose representations are **identical** in a *dependency* parse but **vary** in a *generative* framework.

3.2 Syntactic Structures of Interest

Our method hinges on testing syntactic structures whose representations are crucially different in generative accounts, but are invariant in a dependency parse. In doing so, we propose turning the limitations of a dependency probe to an asset. If the probe’s predicted dependency distances vary between the sentences in question in ways that align with generative theoretical predictions, then we have evidence that not only do LLMs’ contextualized vector representations capture generative syntactic structures, but that a probe trained only to recover dependency parses is additionally sensitive to hierarchical phrase-structure distances.

To test this, we have selected the well-researched Subject Control (SC) and Subject Raising (SR) constructions as our experimental condition. Observed first by Rosenbaum (1967), SR constructions are those that consist of two clauses: a matrix clause and an infinitival Tense Phrase (TP) complement. Since its initial observation by Rosenbaum (1967), it’s largely been accepted that the subject position of the embedded clause is occupied by a trace element (later revised to a copy element, see footnote

1) due to the subject being raised into the matrix clause by the EPP features³ in the matrix clause. SC constructions, meanwhile, are assumed to take a larger complement than a raising verb, with many typically assuming an SC complement to be a Complement Phrase (CP). Many theories follow Chomsky and Lasnik (1993) and posit a silent PRO element that is co-indexed with and controlled by the matrix’s subject. This PRO receives its theta-role from the embedded verb while the matrix subject receives its theta-role from the matrix verb, thus satisfying the Theta Criterion (Chomsky, 1957).⁴

For the purposes of our experiment, the crucial things to know are: Subject Raising takes a Tense Phrase (TP) as its complement, while Subject Control takes the larger Complement Phrase (CP) as its complement, which inherently contains a TP itself. Thus, the result are two structures whose surface forms and dependency parses are *identical*, but whose hierarchical syntactic representations are **different**. Thus, we would expect that if the LLM has not acquired any knowledge of deeper syntactic representations or if the dependency-trained probe is insensitive to phrase-structure representations, then the probe’s predicted distances between relevant word-pairs should *not differ* between the two structures. However, if such hierarchical representations are indeed captured and if the probe is sensitive to these structures, then we would anticipate that the distances between certain word-pairs in an SC construction are **longer** than the equivalent word-pairs in an SR construction due to SCs containing the larger CP complement as opposed to the smaller TP complement of SR predicates.

4 Experiments

4.1 Generating Data

For our experiment, we identified 6 subject-raising verbs and 6 subject-control verbs, which we permutationally paired with a set of 8 subject words, 61 embedded verbs, and a set of possible direct objects (either a single pronominal direct object or a two-word definite nominal object that was matched to a specific embedded verb). Thus, we yielded 33,120

³Chomsky (1995) proposed the Extended Projection Principle, which stipulates that Tense bears a strong D-feature that requires a subject in its Specifier. This can be satisfied by either moving the subject to Spec,TP or by inserting an expletive like "it."

⁴For further discussion on Subject Control and Subject Raising and their structural and semantic differences, see Appendix B.

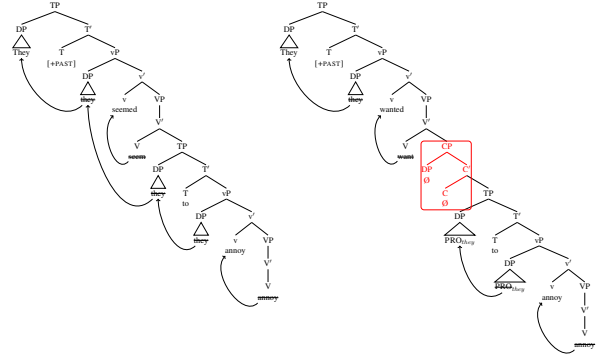


Figure 3: Syntactic trees of the SR sentence "They seemed to annoy him" (left) and the SC sentence "They wanted to annoy him" (right). The two structures are nearly identical, except SC contains a CP above the TP (red), which makes the hierarchical distances between the subject and the embedded clause’s elements (i.e., the infinitive "to", the embedded verb "annoy", and the direct object "him") longer in the SC sentence.

unique sentences, such as "They wanted/seemed to annoy him."

Metrics Should the LLMs *not* have any awareness of the deeper hierarchies or should the probe be insensitive to such differences, then should be *no difference* between SR’s and SC’s distances between words in the matrix clauses and words in the complement clauses. However, if such structures are captured and if the probe is sensitive to this, then we anticipate that the distance between a word in the matrix clause and a word in the complement clause will be **longer** in an SC construction compared to an SR construction since the CP complement is larger (see Figure 3).

For this reason, we opted to investigate the probe’s predicted distances between the following word-pairs: subject and the infinitive (subj-infin, e.g., "they" and "to"), subject and the embedded verb (subj-embed, e.g., "they" and "annoy"), subject and the direct object (subj-dobj, e.g., "they" and "him"), and lastly, embedded verb and the direct object (embed-dobj, e.g., "annoy" and "him"), which serves as our baseline. We should acknowledge at this point that excepting our baseline comparison, none of our word-pairs have any direct dependent or syntactic relationship to each other. This is not a problem. Recall that the probe was trained on a the gold parses for dependency trees in which the distance between two nodes can be counted as the number of edges between the two. Because of this design, we are able to probe the distances between the words in the matrix clause and

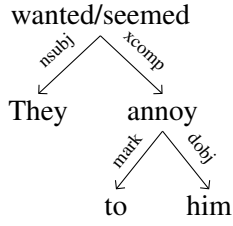


Figure 4: Dependency parse for the two sentences "They seemed/wanted to annoy him." The two trees are identical, and the distance between the subject and embedded verb is 2 while the distance between the subject and the infinitive or direct object is 3.

the words in the complement clause, despite these words not having a direct dependency or syntactic relationship.

As the dependency parses do not differ between the two structures, the gold-parse distances are also stable: subj-embed has a dependency distance of 2 while subj-infin and subj-dobj have dependency distances of 3 (see Figure 4). For this reason, if the LLMs do not capture generative syntactic hierarchies or if the probe is insensitive to such differences, then we should see *no difference in predicted distances between the two experimental conditions*. If, however, the models **do** capture this deep structural difference and if the probe is an adequate tool to measure this, then we should anticipate that the SC distances should be **longer** than their equivalent SR distances. To verify that our probe is working as anticipated, we included the baseline word-pair embed-dobj, which should *not* show any differences in distances as these words are *not* affected by the SC/SR distinction.

4.2 Experimental Setup

Models We probed three pre-trained Transformer (Vaswani et al., 2017) models: **BERT** (Devlin et al., 2019), **RoBERTa** (Liu et al., 2019), and **GPT-2** (Radford et al., 2019). We constrained our probing to models with hidden dimensions of 768 and 1024, which corresponded to the bert-base-cased, roberta-base, and gpt2 for the smaller models and bert-large-cased, roberta-large, and gpt2-medium for the larger models, all of which were accessed using the Huggingface Transformers library (Wolf et al., 2020). The probe was developed using the parsing train/dev/test splits of the Penn Treebank (Marcus et al., 1993).

Following Hewitt and Manning (2019), a probe was trained to convergence (maximum of 40 epochs) on each layer with a batch size of 20. Anal-

ysis was conducted on the best-performing layer.

Once the best-performing layer⁵ was selected, we fed our novel dataset to that probe and obtained the predicted distances for our word-pairs of interest. Analysis was conducted on the predicted distances for our word-pairs *if and only if* the probe properly established the necessary dependency relationships. That is to say, if the probe *misparsed* the tree in a relevant manner, that word-pair’s predicted distance was excluded from analysis. Using Figure 4 as a gold-parse, if the probe’s minimum spanning tree situated "him" as a dependent of "wanted," then we excluded the subj-dobj word-pair as the tree was misparsed in a critical way for that word-pair. We currently do not have strong reason to suspect that such a misparse would have consequent effects on non-affected word-pairs; therefore, the subj-embed and subj-infin’s predicted distances would still be used for statistical analysis since the probe would have still correctly parsed the subject and embedded verb as being dependents of "wanted" and still parsed the infinitive as the dependent of "annoy."⁶

5 Results

As mentioned, we generated 33,120 sentences for which we gathered a total of 704,761 distances across our four word-pairs and all six language models. Overall, this represents an 88.66% accuracy score. The accuracy for our four word-pairs can be found in Table 1 where we may observe that while the accuracy for the SC condition is slightly higher, both showed high accuracy with the lowest being attributable to the subj-dobj word-pair, which was due to the direct object not being tied to the embedded verb, hence the equivalent scores with embed-dobj.

Due to the large size of the data, we split the data by word-pair for statistical analyses. Mixed effect models were developed with the lmer function from lme4 (v. 1.1-31) (Bates et al., 2015) and lmerTest (v. 3.1-3) (Kuznetsova et al., 2017). Fixed effects were identified as the condition (SC or SR) and the linear

⁵See Appendix C, Figures 5 and 6 for model performances.

⁶For our current study, the best-performing layer was selected as the probe with the highest UUAS. As such, we opted only to use word-pairs in which the minimum spanning tree established the correct necessary dependencies for the word-pair in question. However, future work may investigate selecting the probe based on the Spearman correlation, in which case, the motivation to reject data based on improper parses disappears as the Spearman metric does not utilize the minimum spanning tree and instead seeks to globally reduce the differences between the gold distances and the predicted distances.

WordPair	Condition	Acc	PredDist (avg)	Fixed Effects					
				Coefficient	$\hat{\beta}$	SE($\hat{\beta}$)	t	df	p
Subj-Embed	Cont	97.61%	1.93	(Intercept)	1.82691	0.02809	21.75698	65.032	<2.2e-16
	Raise	93.44%	1.72	Condition	-0.21346	0.04026	15.14671	-5.302	8.58e-5
Subj-Infin	Cont	95.99%	2.81	(Intercept)	2.74962	0.06965	9.41633	39.475	8.83e-12
	Raise	89.50%	2.70	Condition	-0.11665	0.08133	14.66730	-10.434	0.172
Subj-Dobj	Cont	85.26%	2.87	(Intercept)	2.64070	0.26273	64.27496	10.051	8.24e-15
	Raise	81.12%	2.64	Condition	-0.56805	0.07609	47.30316	-7.465	1.56e-9
				LinDist	0.01660	0.04944	67.53294	0.336	0.738
				Interaction	0.08320	0.01324	35.36045	6.285	3.12e-7
Embed-Dobj (baseline)	Cont	85.26%	1.50	(Intercept)	1.54400	8.921e-2	58.310	17.304	<2e-16
	Raise	81.12%	1.51	Condition	-3.929e-2	2.702e-2	9.254	-1.454	0.179
				LinDist	-3.964e-2	4.203e-2	64.88	0.943	0.349
				Interaction	3.534e-2	3.108e-3	1.652e+5	11.370	<2e-16

Table 1: Results of the probes’ predicted squared Euclidean distances between the word-pairs of interest. Accuracy records what percentage of the sentences properly established the necessary dependency relationships for that particular word-pair. The right side of the table reports the fixed effects findings for the linear mixed-effect models that were built for each word-pair. See Appendix C, Figures 7 and 8 for visuals.

distance (the number of intervening words plus 1 to avoid issues of 0 multiplication) as well as their interaction. The latter two only applied to word-pairs with the direct object as the direct object could be a single pronominal like "it" (in which case the linear distance would be 1) or a full nominal phrase like "the car" (in which case the linear distance would be 2). For the other word-pairs, there was no variation in linear distance, hence its exclusion as a fixed effect. Condition was contrast-coded with SC being -0.5 and SR being 0.5.

Random effects were identified via model comparison and included by-MatrixVerb, by-SubjectWord, by-EmbedVerb, by-ObjectWord, and by-LanguageModel random slopes for our factor of interest (Condition), excepting by-MatrixVerb, which warranted only a random intercept.⁷

To recap, our hypothesis is that the probe’s predicted distances between the matrix subject and elements in the embedded clause (i.e., the infinitive, the embedded verb, and the direct object) should be **longer** in the SC condition compared to the SR condition. Should this be the case, this effect should appear in all of our word-pairs (excepting our baseline of embed-dobj). In this regard, our study uses conjunction testing in that we require all tests be significant in order to reject the null hypothesis (Weber, 2007). We thus follow Rubin (2021) and do not adjust our alpha level.

⁷Word-pairs with direct objects made for more complicated linear models due to the addition of a by-ObjectWord grouping factor for random effects. Because of this, the linear model for subj-dobj included random intercepts for all grouping factors mentioned, but only warranted random slopes for the grouping factor of language model.

Table 1 reports our results where we find a main effect for Condition in our subj-embed and subj-dobj data ($p = 8.58e^{-5}$ and $p = 1.56e^{-9}$), but *not* for subj-infin ($p = 0.172$). This paints a puzzling picture, which we discuss further in Section 6. For now, it is abundantly clear that the predicted distances from an SC construction are significantly longer than an SR construction when considering the distance between the subject and either the embedded verb or the direct object.

Crucially, we do **not** find Condition to be a significant predictor for our baseline, suggesting the probe is not spuriously attributing higher distances to SCs than SRs in ways that are not predicted by the syntax. However, interaction between Linear Distance and Condition *is* found to be a main effect for embed-dobj. To conduct follow-up models to investigate this result, we split the data by linear distance, meaning sentences were grouped into those that took a pronominal direct object such as "it" (linear distance of 1) and those that took a nominal phrase object such as "the car" (linear distance of 2). In doing so, we do **not** find Condition to be a main effect in *either* group. Analysis of the data further reveals that linear distance decreases the predicted distance for subject raising verbs only. When we control for linear distance by splitting up the data by direct object type, though, our follow-up analyses find that the predicted distance between the embedded verb and the direct object do not significantly vary between the two conditions ($p = 0.875$ for pronominal direct objects and $p = 0.346$ for nominal phrase objects).

The significantly longer predicted distances of

the SC condition in subj-embed and subj-dobj, paired with Condition *not* being a significant predictor for our baseline comparison of embed-dobj (even when accounting for interaction effects), together show strong evidence to reject our null hypothesis. However, the SC/SR condition anomalously does *not* significantly predict the probe distance. Though this result does not negate the significant findings in our other word-pairs, it is an outcome which is not predicted by the syntax, nor is it one that is readily explained when combined with the significant findings in subj-embed and subj-dobj.

6 Discussion & Conclusion

If the structures proposed by generative syntax to account for Subject Raising and Subject Control are indeed captured by LLMs and subsequently by the probe, then we would anticipate that *all* three word-pairs of interest should show significantly longer predicted distances in the SCs compared to SRs while the baseline comparison of embed-dobj (which is *not* affected by an SC or SR construction) should not. While these predictions are evidently borne out by subj-embed and subj-dobj along with our baseline of embed-dobj, it does not hold true for subj-infin.

This finding is particularly puzzling. Should it be that the language models do not capture the SC/SR distinction, then *none* of the word-pairs should have significant differences in distances. Additionally, there is no theory in any school of syntax (generative or otherwise) we are aware of that suggests SC verbs take larger complements *below* the TP head of "to." We might then posit that the infinitive's seemingly imperviousness to the SC/SR distinction may arise from the language models somehow building a novel and alien structure in which the infinitival "to" sits in the matrix clause while the complement size distinctions are displayed beneath it. Again, however, we are resistant to such a notion as we know of no theory postulating such an arbitrary and alien structure.

It is evident this matter requires further investigation, but it is possible the aberrant behavior of the infinitive is due to the nature of infinitives themselves. Infinitival "to" is semantically vacuous: there is little to any rich semantic meaning to the word, which is entirely functional in nature—denoting either non-finiteness as an infinitive⁸ or

directionality or telicity as a preposition. For this reason, we suspect the lack of semantic-richness of purely functional words may impact the ways in which structure is encoded by embedding vectors.⁹

While further work is needed to investigate the outlier behavior of the infinitive, we overall find strong support that LLMs have captured the deep syntactic hierarchies proposed by generative syntax. We furthermore demonstrate that a dependency-trained probe is sensitive to such structures and can provide a means by which to probe for more complex syntactic representations that do not enjoy the benefits of a one-to-one mapping.

The implications of this work have impact on both the field of NLP and the field of linguistics. Our work suggests that LLMs have learned to capture deeper and more complex syntactic structures within their embeddings than previously realized and thus have the ability to capture the semantic nuances that result from sub-surface structural differences. Our findings therefore further the interpretability research of language models to discover what these models have actually learned regarding the features and structures of language. We also find evidence that neural networks trained using the dependency framework can still capture deeper syntactic structures, suggesting these simpler representations may be adequate for downstream tasks as they appear to be capable of reaping the benefits of deep structure without needing to explicitly train on deep structure. As for linguists, the findings of our work warrant further investigation into the viability of using language models as a means to test syntactic structures. There are many competing theories to explain different syntactic phenomenon, which researchers have spent decades gathering data and judgments to help explain. Our work begins to open up the possibility of utilizing LLMs as another source of data to help augment, build, and perhaps even test syntactic theories.

Taken together, we situate our work as a realization of Linzen (2019) and Futrell and Mahowald (2025)'s claim that the skillsets and knowledge of the fields of NLP and linguistics complement each other, and that the two stand primed to advance each other's respective fields through collaboration.

differences between different types of infinitives.

⁹We exempt pronouns from this hypothesis. Our dataset subjects were pronominal and our single-word direct objects were also pronouns. Unlike infinitival "to," pronouns pick out referents in the real world, and can furthermore carry information such as Case, Gender, and Number as opposed to infinitival "to," which indicates non-finiteness only.

⁸See (Satik, 2022) for discussion on the subtle semantic

References

- Guillaume Alain and Yoshua Bengio. 2017. Understanding intermediate layers using linear classifier probes. In *The 5th International Conference on Learning Representations*.
- Dang Anh, Limor Raviv, and Lukas Galke. 2024. [Morphology matters: Probing the cross-linguistic morphological generalization abilities of large language models through a wug test](#). In *Proceedings of the Workshop on Cognitive Modeling and Computational Linguistics*, pages 177–188, Bangkok, Thailand. Association for Computational Linguistics.
- David Arps, Younes Samih, Laura Kallmeyer, and Hassan Sajjad. 2022. [Probing for constituency structure in neural language models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 6738–6757, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Douglas Bates, Martin Mächler, Ben Bolker, and Steve Walker. 2015. [Fitting linear mixed-effects models using lme4](#). *Journal of Statistical Software*, 67(1):1–48.
- Yonatan Belinkov and James Glass. 2019. [Analysis methods in neural language processing: A survey](#). *Transactions of the Association for Computational Linguistics*, 7:49–72.
- Tommi Buder-Gröndahl. 2024. [What does parameter-free probing really uncover?](#) In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 327–336, Bangkok, Thailand. Association for Computational Linguistics.
- Ethan A. Chi, John Hewitt, and Christopher D. Manning. 2020. [Finding universal grammatical relations in multilingual BERT](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5564–5577, Online. Association for Computational Linguistics.
- Noam Chomsky. 1957. *Syntactic Structures*. Mouton.
- Noam Chomsky. 1973. Conditions on transformations. In *A Festschrift for Morris Halle*. Holt, Rinehart Winston.
- Noam Chomsky. 1981. *Lectures on government and binding*. Foris Publications.
- Noam Chomsky. 1986. *Knowledge of Language: Its Nature, Origin, and Use*. Praeger.
- Noam Chomsky. 1993. A minimalist program for linguistic theory. In Kenneth Locke Hale and Samuel Jay Keyser, editors, *The View From Building 20: Essays in Linguistics in Honor of Sylvain Bromberger*. MIT Press.
- Noam Chomsky. 1995. *The Minimalist Program*. MIT Press, Cambridge, MA.
- Noam Chomsky and Howard Lasnik. 1993. *Syntax: An international handbook of contemporary research*. In Joachim Jacobs, Arnim von Stechow, Wolfgang Sternefeld, and Theo Vennemann, editors, *The theory of principles and parameters*. De Gruyter.
- Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. 2019. [What does BERT look at? an analysis of BERT’s attention](#). In *Proceedings of the 2019 ACL Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 276–286, Florence, Italy. Association for Computational Linguistics.
- Haley Coleman. 2020. [This is a BERT. now there are several of them. can they generalize to novel words?](#) In *Proceedings of the Third BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pages 333–341, Online. Association for Computational Linguistics.
- Alexis Conneau, German Kruszewski, Guillaume Lample, Loïc Barrault, and Marco Baroni. 2018. [What you can cram into a single \\$&!#* vector: Probing sentence embeddings for linguistic properties](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2126–2136, Melbourne, Australia. Association for Computational Linguistics.
- N.F.M. Corver. 2007. From trace theory to copy theory. In *The Copy theory of movement*, pages 1–10, Netherlands. John Benjamins.
- Mark Davies. 2008–. [The corpus of contemporary american english \(coca\)](#).
- Christopher Davis, Christopher Bryant, Andrew Caines, Marek Rei, and Paula Buttery. 2022. [Probing for targeted syntactic knowledge through grammatical error detection](#). In *Proceedings of the 26th Conference on Computational Natural Language Learning (CoNLL)*, pages 360–373, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Tiwalayo Eisape, Vineet Gangireddy, Roger Levy, and Yoon Kim. 2022. [Probing for incremental parse states in autoregressive language models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 2801–2813, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Robert Fiengo. 1977. On trace theory. *Linguistic Inquiry*, 8:35–61.

746	Richard Futrell and Kyle Mahowald. 2025. How linguistics learned to stop worrying and love the language models . <i>Preprint</i> , arXiv:2501.17047.	802
747		803
748		804
749	Linyang He, Peili Chen, Ercong Nie, Yuanning Li, and Jonathan R. Brennan. 2024. Decoding probing: Revealing internal linguistic structures in neural language models using minimal pairs . In <i>Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)</i> , pages 4488–4497, Torino, Italia. ELRA and ICCL.	805
750		806
751		807
752		808
753		809
754		810
755		811
756		
757	John Hewitt and Christopher D. Manning. 2019. A structural probe for finding syntax in word representations . In <i>Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)</i> , pages 4129–4138, Minneapolis, Minnesota. Association for Computational Linguistics.	812
758		813
759		814
760		815
761		
762		816
763		817
764		818
765	Jennifer Hu, Jon Gauthier, Peng Qian, Ethan Wilcox, and Roger Levy. 2020. A systematic assessment of syntactic generalization in neural language models . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 1725–1744, Online. Association for Computational Linguistics.	819
766		820
767		821
768		822
769		
770		823
771		824
772	Ganesh Jawahar, Benoît Sagot, and Djamé Seddah. 2019. What does BERT learn about the structure of language? In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 3651–3657, Florence, Italy. Association for Computational Linguistics.	825
773		826
774		827
775		828
776		829
777		
778	Julie Kallini, Isabel Papadimitriou, Richard Futrell, Kyle Mahowald, and Christopher Potts. 2024. Mission: Impossible language models . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 14691–14714, Bangkok, Thailand. Association for Computational Linguistics.	830
779		831
780		832
781		
782		833
783		834
784		835
785	Gaurav Kamath, Sebastian Schuster, Sowmya Vajjala, and Siva Reddy. 2024. Scope ambiguities in large language models . <i>Transactions of the Association for Computational Linguistics</i> , 12:738–754.	836
786		837
787		838
788		839
789	Artur Kulmizev, Vinit Ravishankar, Mostafa Abdou, and Joakim Nivre. 2020. Do neural language models show preferences for syntactic formalisms? In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 4077–4091, Online. Association for Computational Linguistics.	840
790		841
791		842
792		
793		843
794		844
795		
796	Alexandra Kuznetsova, Per B. Brockhoff, and Rune H. B. Christensen. 2017. lmerTest package: Tests in linear mixed effects models . <i>Journal of Statistical Software</i> , 82(13):1–26.	845
797		846
798		847
799		
800	Idan Landau. 2024. <i>Elements in Generative Syntax</i> , chapter Control. Cambridge University.	848
801		849
		850
		851
		852
		853
		854
		855
		856

- Mark Rubin. 2021. When to adjust alpha during multiple testing: A consideration of disjunction, conjunction, and individual testing. *Synthese*, 199:10969–11000.
- Deniz Satik. 2022. [The semantics of infinitival tense](#). Under review.
- M.E. Sánchez, Y. Sevilla, and A. Bachrach. 2016. [Agreement processing in control and raising structures. evidence from sentence production in spanish](#). *Lingua*, 177:60–77.
- Ian Tenney, Patrick Xia, Berlin Chen, Alex Wang, Adam Poliak, R Thomas McCoy, Najoung Kim, Benjamin Van Durme, Samuel R. Bowman, Dipanjan Das, and Ellie Pavlick. 2019. [What do you learn from context? probing for sentence structure in contextualized word representations](#). *Preprint*, arXiv:1905.06316.
- Lucien Tesnière. 1959. *Eléments de Syntaxe Structurale*. Klincksieck, Paris.
- Mycal Tucker, Tiwalayo Eisape, Peng Qian, Roger Levy, and Julie Shah. 2022. [When does syntax mediate neural language model performance? evidence from dropout probes](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5393–5408, Seattle, United States. Association for Computational Linguistics.
- Ashish Vaswani, 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*, volume 30. Curran Associates, Inc.
- Alex Warstadt, Alicia Parrish, Haokun Liu, Anhad Mohananey, Wei Peng, Sheng-Fu Wang, and Samuel R. Bowman. 2020. [BLiMP: The benchmark of linguistic minimal pairs for English](#). *Transactions of the Association for Computational Linguistics*, 8:377–392.
- Rene Weber. 2007. Responses to matsunaga: To adjust or not to adjust alpha in multiple testing: That is the question. guidelines for alpha adjustment as response to o’keefe’s and matsunaga’s critiques. *Communication Methods and Measures*, 1:281–289.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. [Transformers: State-of-the-art natural language processing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.