

How poor is the stimulus? Evaluating hierarchical generalization in neural networks trained on child-directed speech

Anonymous ACL submission

Abstract

When acquiring syntax, children consistently choose hierarchical rules over competing non-hierarchical possibilities. Is this preference due to a learning bias for hierarchical structure, or due to more general biases that interact with hierarchical cues in children’s linguistic input? We explore these possibilities by training LSTMs and Transformers—two types of neural networks without a hierarchical bias—on data similar in quantity and content to children’s linguistic input: text from the CHILDES corpus. We then evaluate what these models have learned about English yes/no questions, a phenomenon for which hierarchical structure is crucial. We find that, though they perform well at capturing the surface statistics of child-directed speech (as measured by perplexity), both model types generalize in a way more consistent with an incorrect linear rule than the correct hierarchical rule. These results suggest that human-like generalization from text alone requires stronger biases than the general sequence-processing biases of standard neural network architectures.

1 Introduction

Syntax is driven by hierarchical structure, yet we typically encounter sentences as linear sequences of words. How do children come to recognize the hierarchical nature of the languages they acquire? Some argue that humans must have a hierarchical inductive bias—an innate predisposition for hierarchical structure (Chomsky, 1965, 1980). An alternate view (e.g., Lewis and Elman, 2001) is that no such bias is necessary: there may be clear evidence for hierarchical structure in children’s input, so that children would choose hierarchical rules even without a hierarchical bias.

At first blush, recent work in natural language processing (NLP) may seem to indicate that no hierarchical bias is necessary. Neural networks trained on naturally-occurring text perform impressively

on syntactic evaluations even though they have no explicit syntactic structure built into them (e.g., Gulordava et al., 2018; Wilcox et al., 2018; Warstadt et al., 2020). However, these results do not provide strong evidence about the learning biases required to learn language from the data available to humans because these models receive very different training data than humans do. First, NLP models are typically trained on far more data than children receive, so models have more opportunities to encounter rare syntactic structures (Linzen, 2020). Second, most training sets in NLP are built from online text (e.g., Wikipedia), which differs qualitatively from the utterances that children typically hear; e.g., sentences in Wikipedia are on average 25 words long (Yasseri et al., 2012), compared to 5 words for sentences in the North American English subset of the CHILDES corpus of child-directed speech (MacWhinney, 2000).

In this work, to evaluate if neural networks without a hierarchical bias generalize like children do, we train models on text¹ comparable to the sentences in children’s linguistic input: English data from CHILDES. We then analyze what they have learned about the relationship between declarative sentences, such as (1a), and their corresponding yes/no questions, such as (1b):

- (1) a. Those **are** your checkers.
b. **Are** those your checkers?

Crucially, nearly all naturally-occurring yes/no questions are consistent with two rules: one based on hierarchical structure (2), and one based on linear order (3):^{2,3}

- (2) HIERARCHICALQ: The auxiliary at the start

¹Section 7.3 discusses other input types (e.g., visual input).

²In past work these rules have been framed as transformations named MOVE-FIRST and MOVE-MAIN (McCoy et al., 2020). We instead follow Berwick et al. (2011) and frame the child’s knowledge as a relationship between sentences.

³Though these two rules are the most prominent in prior literature, other rules are possible; see Section 5.2.

076 of a yes/no question corresponds to the **main**
077 auxiliary of the corresponding declarative.

- 078 (3) LINEARQ: The auxiliary at the start of a
079 yes/no question corresponds to the **first** auxi-
080 liary of the corresponding declarative.

081 Despite the scarcity of evidence disambiguating
082 these rules, children reliably favor HIERARCHI-
083 CALQ (Crain and Nakayama, 1987), albeit with
084 occasional errors consistent with LINEARQ (Am-
085 bridge et al., 2008). Yes/no questions thus are a
086 prime candidate for an aspect of English syntax
087 for which human-like generalization requires a hi-
088 erarchical bias. We evaluate yes/no question per-
089 formance in LSTMs and Transformers, two neural-
090 network architectures that have no inherent hierar-
091 chical inductive bias (McCoy et al., 2020; Petty and
092 Frank, 2021). These architectures employ different
093 computational mechanisms, so consistent results
094 across both would indicate that our results are not
095 due to idiosyncrasies of one particular architecture.

096 To investigate if models generalize more con-
097 sistently with the hierarchical or linear rule, we
098 evaluate them on cases where the rules make dif-
099 ferent predictions, such as (4): under HIERARCHI-
100 CALQ, the question that corresponds to (4a) is (4b),
101 whereas under LINEARQ it is (4c).

- 102 (4) a. The boy who **has** talked **can** read.
103 b. **Can** the boy who **has** talked ___ read?
104 c. ***Has** the boy who ___ talked **can** read?

105 We find that across several ways of framing the
106 learning task, models fail to learn HIERARCHI-
107 CALQ. Instead, they generalize in ways that de-
108 pend on linear order and on the identities of spe-
109 cific words. These results suggest that children’s
110 training data, if taken to be words alone, may not
111 contain enough hierarchical cues to encourage hier-
112 archical generalization in a learner without a hierar-
113 chical bias. Thus, explaining human acquisition of
114 syntax may require postulating that humans have
115 stronger inductive biases than those of LSTMs and
116 Transformers, or that information other than word
117 sequences plays a crucial role.⁴

118 2 Background

119 Though HIERARCHICALQ and LINEARQ often
120 make the same predictions, the evidence in chil-
121 dren’s input may still favor HIERARCHICALQ.

⁴To facilitate further research, we have uploaded our datasets and trained models at [LINK ANONYMIZED].

The most straightforward evidence would be ut-
terances that directly disambiguate the rules, such
as (4b). Pullum and Scholz (2002) show that disam-
biguating examples appear in the *Wall Street Jour-
nal*, in literature, and arguably in child-directed
speech, but direct evidence may still be too rare to
robustly support HIERARCHICALQ (Legate and
Yang, 2002). Nonetheless, children might con-
clude that yes/no questions obey HIERARCHI-
CALQ rather than LINEARQ based on *indirect*
evidence—evidence that *other* syntactic phenom-
ena are hierarchical (Mulligan et al., 2021).

To test if the cues favoring HIERARCHICALQ
render a hierarchical bias unnecessary, we study
how well non-hierarchically-biased models acquire
English yes/no questions. Several prior papers have
used this approach, but their training data differed
from children’s input in important ways: some used
synthetic datasets (Lewis and Elman, 2001; Frank
and Mathis, 2007; Clark and Eyraud, 2007; McCoy
et al., 2020), others used massive Internet corpora
(Lin et al., 2019; Warstadt and Bowman, 2020),
and those that used child-directed speech simpli-
fied the data by replacing each word with its part
of speech (Perfors et al., 2011; Bod et al., 2012).
We used training data closer to children’s input,
namely sentences from CHILDES with word iden-
tities preserved, rather than being converted to parts
of speech. Two other recent works have also trained
neural networks on CHILDES data (Pannitto and
Herbelot, 2020; Huebner et al., 2021), but neither
investigated yes/no questions.

154 3 Overview of Experimental Setup

We evaluated models on yes/no questions in two
ways. First, we used relative acceptability judg-
ments (Experiment 1): We trained neural networks
on the task of language modeling (predicting the
next word at every point in the sentence) and evalu-
ated whether they assigned a higher probability to
sentences consistent with LINEARQ or HIERAR-
CHICALQ. Our second approach was based on text
generation (Experiment 2): We trained networks
to take in a declarative sentence and output the
corresponding question, and tested whether they
generalized in a way more consistent with LIN-
EARQ or HIERARCHICALQ. Under both framings,
we trained models on data from CHILDES and
evaluated them on targeted datasets constructed to
differentiate LINEARQ and HIERARCHICALQ.

The size of the dataset that we extracted from

CHILDES was plausibly within the range from which children can acquire HIERARCHICALQ. Crain and Nakayama (1987) found that children between ages 3 and 5 behaved much more consistently with HIERARCHICALQ than LINEARQ. By age 3, American children from families with a lower socioeconomic status receive approximately 10 million words of input (Hart and Risley, 1995), similar to the 8.5 million words of our training set. Thus, it is reasonable to suppose that a learner that generalizes as children do would favor HIERARCHICALQ after being trained on our training set.

4 Experiment 1: Relative Acceptability

4.1 Dataset

To train models on data as similar as possible to the sentences children receive, we extracted data from CHILDES (MacWhinney, 2000). We used the North American English portion. We wished to replicate children’s *input*, so we excluded the children’s own utterances, leaving a 9.6-million-word corpus. We allocated 90% of the data to training, 5% to validation, and 5% to testing. We replaced words that appeared two or fewer times in the training set with <unk>, giving a replacement rate of 0.3%. See Appendix A for more details.

4.2 Task: Next-Word Prediction

We trained models on next-word prediction, also known as language modeling. We chose this task for two reasons. First, it is clear empirically that next-word prediction can teach neural networks a substantial amount about syntax (e.g., Hu et al., 2020). Second, it is plausible that humans perform some version of next-word prediction during sentence processing (Altmann and Kamide, 1999; Hale, 2001; Levy, 2008; Kutas et al., 2011) and that such prediction may play a role in acquisition (Elman, 1991). Thus, while next-word prediction is certainly not the only goal of human language learners, we view this task as a reasonable first step in emulating human language acquisition.

4.3 Architectures

We used LSTMs (Hochreiter and Schmidhuber, 1997) and Transformers (Vaswani et al., 2017). We chose these models for two reasons. First, they have been the most successful architectures in NLP. Thus, we have reason to believe that, of the types of low-bias models invented, these two are the ones most likely to discover linguistic regularities in

our CHILDES training data. Second, the two architectures process sequences very differently (via recurrence vs. via attention). Thus, if both generalize similarly, we would have evidence that what was learned is strongly evidenced in the data, rather than due to a quirk of one particular architecture.

For our LSTMs, we used 2 layers, a hidden and embedding size of 800, a batch size of 20, a dropout rate of 0.4, and a learning rate of 10. For our Transformers, the corresponding values were 4, 800, 10, 0.2, and 5, and we used 4 attention heads. We chose these values based on a hyperparameter search described in Appendix B. All following results are averaged across 10 runs with different random seeds.

4.4 Results: Language Model Quality

Before testing models on questions, we used perplexity to evaluate how well they captured the basic structure of their training domain. For a baseline, we used a 5-gram model with Kneser-Ney smoothing (Kneser and Ney, 1995) trained with KenLM (Heafield, 2011). The test set perplexity for the 5-gram baseline was 24.37, while the average test set perplexity for the LSTMs and Transformers was 20.05 and 19.69, respectively. For perplexity, lower is better. Thus, both neural network types outperformed the strong baseline of a smoothed 5-gram model, showing that they performed well at capturing the basic statistics of their training domain.⁵ We now test whether these models have also successfully learned yes/no questions.

4.5 Yes/No Questions

Evaluation Dataset: Forced-Choice Acceptability Judgments As a first way to test whether our models have learned HIERARCHICALQ, we evaluate whether they assign higher probabilities to sentences consistent with HIERARCHICALQ than to minimally different sentences that are ungrammatical. For this purpose, we create an evaluation dataset containing groups of 6 questions, each created by starting with a declarative sentence, such as (5), and then deleting the **first**, **main**, or neither auxiliary, and inserting the **first** or **main** auxiliary at the front of the sentence.⁶ For instance, in (6b), the **first** auxiliary has been preposed, and the **main** auxiliary has been deleted.

⁵For an intuitive illustration of our model quality, see the sample text generated by them in Appendix H.

⁶It would be possible to also use a ‘prepose other’ category, where an auxiliary not in the input is inserted (McCoy et al., 2018). We excluded this category because using it would raise complications about which ‘other’ auxiliary to choose.

- (5) The dog who **has** seen a boy **did** try.
- (6) a. **Has** the dog who seen a boy **did** try?
 b. **Has** the dog who **has** seen a boy try?
 c. **Has** the dog who **has** seen a boy **did** try ?
 d. **Did** the dog who seen a boy **did** try?
 e. **Did** the dog who **has** seen a boy try?
 f. **Did** the dog who **has** seen a boy **did** try?

Within each group, we evaluate which question the model assigned the highest probability to. If a model has correctly learned HIERARCHICALQ, it should assign the highest probability to the question consistent with this rule, such as (6e).

Several past papers about yes/no questions have used the same general approach (Lewis and Elman, 2001; Reali and Christiansen, 2005). However, these papers considered only pairs of sentences, whereas we consider groups of 6 to allow for a wider range of possible generalizations that a model might have learned.

To generate the declaratives from which we formed groups of 6 questions, we used the context-free grammar (CFG) in Appendix F, which has a vocabulary selected from the most common words in CHILDES. Each declarative generated by the CFG (e.g., (5)) contains two auxiliary verbs: one before the sentence’s main verb and one inside a relative clause modifying the subject. One potential problem is that some questions are consistent with both HIERARCHICALQ and LINEARQ. For instance, (7a) can be formed from (7b) with the HIERARCHICALQ-consistent steps PREPOSE-MAIN,DELETE-MAIN, or from (7c) with the LINEARQ-consistent steps PREPOSE-FIRST,DELETE-MAIN.

- (7) a. Did the boy who did see the person laugh?
 b. The boy who did see the person did laugh.
 c. The boy who did see the person can laugh.

To avoid this problem, we required that the auxiliary before the main verb must select for a different verb inflection than the one in the relative clause. For instance in (5), **did** selects for the verb’s bare form, while **has** selects for the past participle form. Thus, the auxiliary at the start of the question could only correspond to whichever auxiliary in the declarative has the same selectional properties.⁷

Results: Relative Question Acceptability For each sentence group, we used per-word perplex-

⁷A model could succeed on this dataset with a rule that relates the auxiliary at the start of a question with the *last* auxiliary in the declarative form. Since our models fail on this dataset, this consideration is not relevant here.

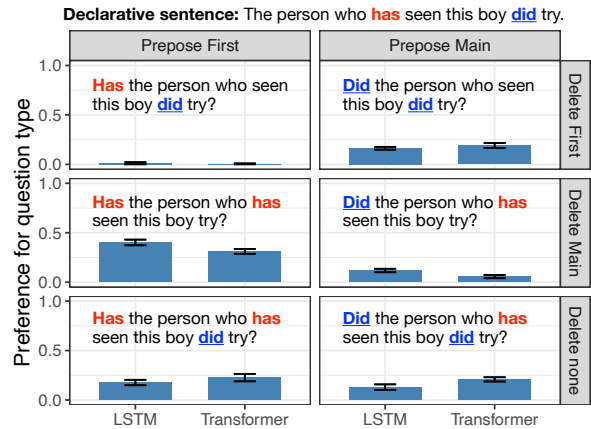


Figure 1: The question types that models prefer when offered a choice between 6 questions. These 6 questions are formed by modifying a declarative with a relative clause on the subject according to ‘prepose’ and ‘delete’ rules. Within each architecture, the proportions across all 6 question types necessarily sum to 1. Each bar shows the average across 10 model re-runs, with single-standard-deviation error bars.

ity to see which of the 6 candidates the models scored most highly.⁸ For both LSTMs and Transformers, the correct category (PREPOSE MAIN, DELETE MAIN) was the second-rarest choice, and the most frequent preference was for PREPOSE FIRST, DELETE MAIN, a category that is only partially correct because it references linear order in addition to hierarchical structure. (Figure 1). Thus, neither model displays preferences consistent with the correct, fully-hierarchical generalization. The two model types showed similar scores, which may mean that these results are largely driven by the statistics of the training data that both models share, rather than the models’ differing inductive biases.

One of the incorrect categories—PREPOSE MAIN, DELETE NONE, such as (6f)—only requires reference to hierarchical structure, so it could be said to capture the hierarchical nature of yes/no questions. Nonetheless, this category was also relatively rare: combining the two fully hierarchical possibilities (PREPOSE MAIN, DELETE MAIN and PREPOSE MAIN, DELETE NONE) accounts for only 26% of LSTM preferences and 27% of Transformer preferences, meaning that both models over 70% of the time favored a sentence generated at least partially based on linear order.

⁸We also explored evaluation of the models with a more complex measure called SLOR where we additionally normalized scores by word frequency (Pauls and Klein, 2012). Both metrics produced qualitatively similar results, so we only report the simpler metric here. See Appendix C.1.

5 Experiment 2: Question Formation

The previous experiment was designed to operate entirely in the next-word-prediction paradigm, motivated by arguments from past literature about the strength and relative ecological validity of next-word-prediction as a training objective (see Section 4.2). However, one of this setup’s shortcomings is that the conclusions are based on the relative acceptability of questions alone, whereas HIERARCHICALQ describes the acceptability of a correspondence between a declarative and a question.

In this second experiment, to better capture that HIERARCHICALQ is defined over sentence pairs, we trained models on a sentence-pair task: transforming a declarative into a question (McCoy et al., 2020). For instance, given *the child did learn* the model must produce *did the child learn ?*

We evaluated models in two ways. First, we checked if the models’ predictions fully matched the correct questions. This full-sentence evaluation is demanding, and models might fail this evaluation for reasons unrelated to our core hypotheses. For instance, given *the child did learn* the model might produce *did the baby learn*, which would be marked as incorrect, even though this lexical error is not relevant to HIERARCHICALQ.

As a metric that is less demanding and that also more directly targets HIERARCHICALQ, we measured if the first word of the output question corresponded to the first or main auxiliary of the input. Critically, LINEARQ and HIERARCHICALQ make different predictions for the first word of a question so long as the two auxiliaries are distinct: see (4). Because this framing lets the model freely generate its output (instead of choosing one option from a pre-specified set), we allow for the possibility that the rule learned by models may not be identical to any of our manually-generated hypotheses.

Solely training models to perform this transformation involves the implicit assumption that, when children acquire English yes/no questions, the only evidence they leverage is English yes/no questions. However, other types of sentences may also provide useful evidence (Pearl and Mis, 2016): e.g., *wh*-questions also illustrate subject-auxiliary inversion (Pullum and Scholz, 2002), while, more generally, many types of sentences could provide evidence that the syntax as a whole is hierarchical (Perfors et al., 2011). To explore this possibility, we compared a condition in which models were only trained to perform question formation (the

QUESTION FORMATION condition) to another in which models were first pre-trained on next-word prediction with the exact same setup as in Experiment 1 before being further trained to perform question formation (the NEXT-WORD PREDICTION + QUESTION FORMATION condition).

5.1 Dataset

Training Set Our question formation dataset consisted of the yes/no questions in the CHILDES Treebank (Pearl and Sprouse, 2013a,b), a parsed subset of CHILDES containing 189,359 sentences. We used these parses to extract all yes/no questions from the CHILDES Treebank and derive their corresponding declarative forms. The resulting declarative was concatenated with the question. An example declarative/question pair is:

(8) you can spell your name . can you
spell your name ?

The training set consisted of 10,870 declarative/question pairs, the validation set 1,360 pairs, and the test set 1,358 pairs (we will call this test set the *randomly-partitioned test set* to distinguish it from two other evaluation sets discussed below). We trained models to perform next-word prediction on such concatenated sentence pairs. The first-word accuracy of the trained model was then computed based on the model’s prediction for the word after the period in each test example, while the full-sentence accuracy was computed based on the model’s predictions for all tokens after the period. All the questions in the randomly-partitioned test set were withheld from both the question-formation training set and the next-word-prediction training set. Thus, models had not seen these test examples in their training, even in the NEXT-WORD PREDICTION + QUESTION FORMATION condition in which they were trained on both tasks.

Evaluation Sets In addition to the randomly-partitioned test set, we used CFGs to generate two targeted evaluation sets. As in Experiment 1, we selected the CFGs’ vocabulary from common words in our CHILDES data. In sentences generated from the first CFG, the sentence’s first auxiliary was also its main auxiliary, so LINEARQ and HIERARCHICALQ make the same predictions. (8) exemplifies the type of declarative-question pair in this dataset. We call this dataset FIRST-AUX = MAIN-AUX. For sentences generated by the second CFG, the main auxiliary was the *second* auxiliary in the sentence; thus, these examples disambiguate LINEARQ and

HIERARCHICALQ. Example (9) is a declarative-question pair from this evaluation set.

(9) a boy who is playing can try . can a boy who is playing try ?

We call this dataset $\text{FIRST-AUX} \neq \text{MAIN-AUX}$. See Appendix F for the CFGs used. We sampled 10,000 declarative sentences from these grammars and transformed them into questions according to HIERARCHICALQ to create our evaluation sets.

5.2 Results

Randomly-Partitioned Test Set The LSTMs and Transformers in the QUESTION FORMATION condition performed well on the randomly-partitioned test set, with a full-question accuracy of 0.68 ± 0.014 and 0.87 ± 0.005 (averaged across 10 reruns with margins indicating one standard deviation). The models in the NEXT-WORD PREDICTION + QUESTION FORMATION condition performed similarly well, with a full-question accuracy of 0.66 ± 0.008 for the LSTMs and 0.93 ± 0.004 for the Transformers. For both model types, the first-word accuracy for the question was nearly 1.00 across re-runs. We suspect that Transformers have a stronger full-question accuracy because producing the question requires copying all words from the declarative (but in a different order). Copying is likely easy for Transformers because they can attend to specific words in the prior context, while our LSTMs must compress the entire context into a fixed-size vector, which may degrade the individual word representations. Because both model types achieved near-perfect performance on the crucial first-word accuracy metric, we conclude that our models have successfully learned how to handle the types of declarative/question pairs that we extracted from the CHILDES Treebank.

Targeted Evaluation Sets On our two targeted evaluation sets, models almost never produced the complete question correctly. Turning to the more lenient measure of first-word accuracy, for examples on which LINEARQ and HIERARCHICALQ predict the same first output word ($\text{FIRST-AUX} = \text{MAIN-AUX}$), the Transformer trained only on question formation performed strongly, while the Transformer trained on both tasks, and both LSTMs, performed reasonably well (Figure 2; note models could choose any word in their vocabulary to begin the output, so chance performance is near 0.00). For the crucial cases that disambiguate the two rules ($\text{FIRST-AUX} \neq \text{MAIN-AUX}$), both mod-

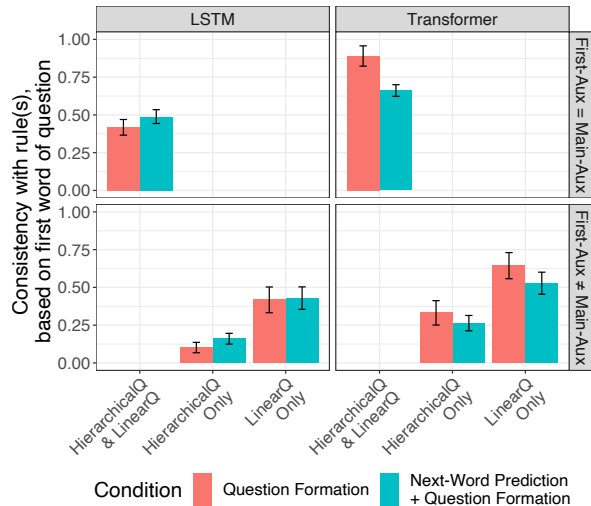


Figure 2: Proportion of model-produced questions that were consistent with the linear rule LINEARQ and/or the hierarchical rule HIERARCHICALQ. In the $\text{FIRST-AUX} = \text{MAIN-AUX}$ dataset, the first auxiliary is the main auxiliary, so both LINEARQ and HIERARCHICALQ produce the correct question string. The $\text{FIRST-AUX} \neq \text{MAIN-AUX}$ dataset disambiguates the two rules. Each bar shows the average across 10 model reruns, with error bars showing one standard deviation.

els in both conditions performed more consistently with LINEARQ than HIERARCHICALQ. Training on next-word prediction before question formation had inconsistent effects: it modestly increased the likelihood of hierarchical generalization in LSTMs, yet it decreased that likelihood in Transformers.

Lexical Specificity In Appendix G, we further break down the $\text{FIRST-AUX} \neq \text{MAIN-AUX}$ results based the auxiliaries' identity. The generalization pattern varied considerably across auxiliary pairs. For some auxiliary pairs, the auxiliary chosen to begin the question was usually neither auxiliary in the input (Figure 3, left facet). For other pairs, models usually chose the first auxiliary, regardless of lexical identity (Figure 3, middle facet). Finally, for some pairs, the auxiliary chosen was usually the same one, regardless of whether it was the first or main auxiliary (Figure 3, right facet).

Generalization based on lexical identity is rarely considered in past discussions of English yes/no question acquisition. Of the papers on this phenomenon (see Clark and Lappin (2010), Lasnik and Lidz (2017), and Pearl (2021) for overviews), the only one to our knowledge that discusses lexical specificity is Frank and Mathis (2007), which studied models trained on synthetic data. Our results highlight the importance of testing for a broad

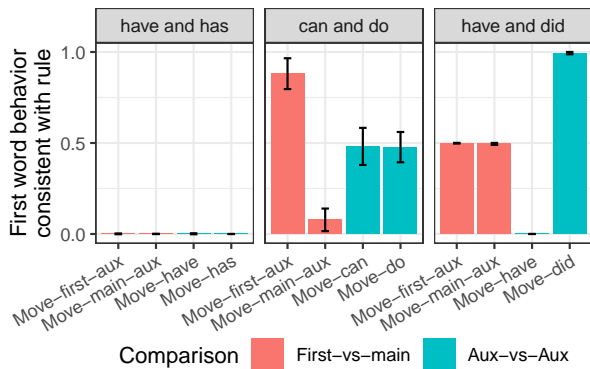


Figure 3: Lexical specificity in model behavior. Each facet considers only the evaluation examples containing the two auxiliaries in the facet heading; e.g., the *can and do* facet includes, for example, the inputs *the children who can play do learn* and *the children who do play can learn*. The bars show the proportion of model predictions for the first word of the output that are consistent with four potential movement rules, averaged across 10 model re-runs and with error bars showing one standard deviation above and below the mean. This plot only shows an illustrative subset of auxiliary pairs for one model type (Transformers in the NEXT-WORD PREDICTION + QUESTION FORMATION condition); see Appendix G for the full results.

range of generalizations: Lexically-specific hypotheses appear attractive for our low-bias learners, so an account of what biases can yield human-like learning should rule out these lexically-specific hypotheses along with linear ones.

6 Evaluating on Other Phenomena

We have found that our models consistently failed to learn HIERARCHICALQ. Does this failure result from a general failure to learn syntax? If so, this could indicate that our training setup is flawed, since Huebner et al. (2021) showed that Transformers can score well on certain syntactic evaluations after being trained on CHILDES data (though they did not evaluate models on yes/no questions).

To test this possibility, we evaluated our models on the Zorro dataset (Huebner et al., 2021), which is based on BLiMP (Warstadt et al., 2020). Zorro contains 24 evaluations, each of which targets one syntactic phenomenon (e.g., subject-verb agreement) and involves sentence pairs for which one sentence is grammatical, and the other is minimally different but ungrammatical (e.g., by violating subject verb agreement). We evaluated models as we did in Section 4.5: a model gets a sentence pair correct if it has a lower perplexity for the grammatical sentence than for the ungrammatical sentence.

See Appendix D for full results. For each syntactic phenomenon, most model re-runs scored above 0.9, though at least one scored near the chance level of 0.5. For each re-run of each architecture there is at least one phenomenon for which the model scores over 0.97, and many models score 1.00 on some phenomena. Thus, our models’ failure on yes/no questions cannot be explained by a general failure to learn syntax, since all models score well on at least some syntactic evaluations.

7 Discussion

We have found that, when trained on child-directed speech, two types of standard neural networks performed reasonably well at capturing the statistical properties of the dataset, yet their handling of English yes/no questions was more consistent with a linear rule LINEARQ than the correct hierarchical rule HIERARCHICALQ. These results support the hypothesis that a learner requires a hierarchical bias to consistently learn hierarchical rules when learning from the linguistic data children receive.

7.1 Takeaways for LSTMs and Transformers

When trained on massive corpora, LSTMs and Transformers perform impressively on some syntactic evaluations. Based on such results, it is tempting to conclude that the general-purpose biases of these architectures suffice to yield human-like syntax acquisition. Our results caution against this interpretation: When we trained the same architectures on data more similar to children’s input, they failed to learn the structure of English yes/no questions. Thus, at least when learning from text alone, LSTMs and Transformers do not display human-like language learning—they do not generalize as humans do from the data that humans receive.

7.2 Takeaways for the Poverty of the Stimulus Debate

Below we specify four possible positions in the poverty-of-the-stimulus debate about the adequacy of children’s input for inducing hierarchical rules in low-bias learners, arranged from assuming the most limited to the most expansive innate component:

- (10) **Any inductive biases:** Any learner trained on CHILDES will generalize like humans do.
- (11) **Any inductive biases that enable in-distribution learning:** Any learner that captures the statistical patterns of the training distribution will generalize to HIERARCHICALQ.

589 (12) **Some non-hierarchical inductive biases:**
590 Some general-purpose learners will generalize
591 as humans do, but others will not.

592 (13) **Only a hierarchical inductive bias:** No
593 general-purpose learners will generalize as
594 humans do: hierarchical biases are necessary.

595 Position (10) is clearly false: many learners can-
596 not learn certain aspects of syntax, no matter their
597 training data (e.g., bigram models cannot capture
598 long-distance dependencies). Our work shows that
599 position (11) is also false: Though our models per-
600 formed well on the in-distribution test sets of Exper-
601 iments 1 and 2, they did not generalize in human-
602 like ways. This leaves positions (12) and (13),
603 which our existing results cannot differentiate. It is
604 possible that only learners with hierarchical induc-
605 tive biases can demonstrate human-like language
606 learning (position (13)), but also that some learners
607 without this bias can succeed (position (12))—just
608 not the learners we tested. For further discussion
609 of how computational modeling can bear on learn-
610 ability arguments, see [Wilcox et al. \(2021\)](#).

611 One potential solution supporting position (12)
612 would be that learners leverage the hierarchical
613 structure of some syntactic phenomenon to help
614 conclude that other, impoverished phenomena are
615 hierarchical ([Perfors et al., 2011](#); [Mulligan et al., 2021](#)).
616 However, our results from Experiment 2
617 show that giving learners access to a wider range
618 of phenomena does not automatically improve hi-
619 erarchical generalization: Models’ performance on
620 question formation was not substantially improved
621 (and in some cases was even harmed) when they
622 were trained not just on question formation but also
623 on next-word prediction on the entire CHILDES
624 corpus. Thus, although training on text that con-
625 tains many linguistic phenomena can give mod-
626 els a hierarchical inductive bias when the training
627 is done over large Internet corpora ([Warstadt and
628 Bowman, 2020](#); [Mueller et al., 2022](#)), our results
629 provide evidence that this conclusion does not ex-
630 tend to models trained on child-directed speech.

631 7.3 Comparison of Our Training Data to 632 Children’s Input

633 Our training set was both qualitatively and quanti-
634 tatively closer to children’s input than the massive
635 Internet corpora standardly used to train models in
636 NLP ([Linzen, 2020](#)). This difference is important:
637 [Lin et al. \(2019\)](#), [Warstadt and Bowman \(2020\)](#),
638 and [Mueller et al. \(2022\)](#) all found evidence that

639 models trained on large Internet corpora performed
640 well on yes/no questions evaluations, whereas our
641 models trained on CHILDES performed poorly—
642 though we cannot be certain the differences in re-
643 sults are solely due to differences in the training
644 data, since these prior papers used different model
645 architectures, training tasks, and evaluation setups.

646 Though our training data are more similar to
647 children’s input than massive Internet corpora are,
648 differences remain. Our experiments omit several
649 aspects of a child’s experience that might help them
650 acquire syntax, such as prosody ([Morgan and De-
651 muth, 1996](#)), visual information ([Shi et al., 2019](#)),
652 and meaning ([Fitz and Chang, 2017](#); [Abend et al.,
653 2017](#)), all of which might correlate with syntactic
654 structure and thus provide additional cues to the
655 correct hierarchical generalization. On the other
656 hand, our dataset might present an easier learning
657 scenario than children are faced with, because chil-
658 dren must learn to segment the speech stream into
659 words ([Lakhotia et al., 2021](#)), while our models do
660 not need to learn this. Further, although real-world
661 grounding could provide access to syntactically-
662 relevant information, learners might struggle to
663 leverage this information because of difficulties in
664 determining what is being discussed in the physical
665 world ([Gleitman et al., 2005](#)).

666 8 Conclusion

667 In this work, we trained two types of neural net-
668 works (LSTMs and Transformers) on sentences of
669 the types available to children and then analyzed
670 what they had learned about English yes/no ques-
671 tions. Across several evaluation paradigms, these
672 models failed to generalize in human-like ways:
673 Humans display hierarchical generalization, while
674 the models’ generalization was instead based on
675 linear order and individual words’ identities. Our
676 results support the hypothesis that human-like lin-
677 guistic generalization requires biases stronger than
678 those of LSTMs and Transformers. Future work
679 should investigate what inductive biases enable suc-
680 cessful generalization. One approach would be to
681 test architectures with built-in hierarchical struc-
682 ture; past work has shown that such architectures
683 have a hierarchical bias ([McCoy et al., 2020](#)) and
684 generalize better on the hierarchical phenomenon
685 of subject-verb agreement ([Kuncoro et al., 2018](#);
686 [Lepori et al., 2020](#)), so they may also generalize
687 better on English yes/no questions.

Ethics Statement

Use of human data: While we did not collect any new human data ourselves, many of our analyses involved the use of prior datasets within the CHILDES database. All of these datasets were collected in accordance with IRB policies at the institutions of the data collectors, and all followed standard practices in obtaining informed consent and deidentifying data.⁹

Risks and limitations: The main risk of our proposed analyses is that future work using the same analyses might draw overly strong conclusions based on increased model performance, leading to overestimates of model strength. Such overestimates are an issue because they can lead users to place more trust in a model than is warranted.

To clarify, we view strong performance on our evaluation datasets as necessary but not sufficient to demonstrate human-like learning. Thus, if models perform poorly on our datasets (as the models we evaluated did), then we have strong reason to conclude that models are not learning in human-like ways. If future models perform better, such results would be consistent with human-like learning but would not conclusively establish that models learn as humans do, as they might instead be using some shallow heuristic that is not controlled for in our datasets. In other words, a criterion that is necessary but not sufficient facilitates strong conclusions about failure but does not facilitate strong conclusions about success. If future papers are faced with models that are more successful, such papers would ideally supplement results based on our datasets with analyses of models' internal strategies in order to more conclusively establish that what they have learned is not a spurious heuristic.

References

- Omri Abend, Tom Kwiatkowski, Nathaniel J Smith, Sharon Goldwater, and Mark Steedman. 2017. Bootstrapping language acquisition. *Cognition*, 164:116–143.
- Gerry TM Altmann and Yuki Kamide. 1999. Incremental interpretation at verbs: Restricting the domain of subsequent reference. *Cognition*, 73(3):247–264.
- Ben Ambridge, Caroline F Rowland, and Julian M Pine. 2008. Is structure dependence an innate constraint? New experimental evidence from children's

complex-question production. *Cognitive Science*, 32(1):222–255.

Robert Berwick, Paul Pietroski, Beracah Yankama, and Noam Chomsky. 2011. [Poverty of the stimulus revisited](#). *Cognitive science*, 35:1207–42.

Rens Bod, Margaux Smets, et al. 2012. Empiricist solutions to nativist problems using tree-substitution grammars. *Workshop on Computational Models of Language Acquisition and Loss: EACL*.

Noam Chomsky. 1965. *Aspects of the Theory of Syntax*, 50 edition. The MIT Press.

Noam Chomsky. 1980. *Rules and representations*. Columbia University Press.

Alexander Clark and Rémi Eyraud. 2007. Polynomial identification in the limit of substitutable context-free languages. *Journal of Machine Learning Research*, 8(8).

Alexander Clark and Shalom Lappin. 2010. *Linguistic Nativism and the Poverty of the Stimulus*. John Wiley & Sons.

Stephen Crain and Mineharu Nakayama. 1987. Structure dependence in grammar formation. *Language*, pages 522–543.

Jeffrey L Elman. 1991. Distributed representations, simple recurrent networks, and grammatical structure. *Machine learning*, 7(2):195–225.

Hartmut Fitz and Franklin Chang. 2017. Meaningful questions: The acquisition of auxiliary inversion in a connectionist model of sentence production. *Cognition*, 166:225–250.

Robert Frank and Donald Mathis. 2007. Transformational networks. *Models of Human Language Acquisition*, 22.

Lila R Gleitman, Kimberly Cassidy, Rebecca Nappa, Anna Papafragou, and John C Trueswell. 2005. Hard words. *Language learning and development*, 1(1):23–64.

Kristina Gulordava, Piotr Bojanowski, Edouard Grave, Tal Linzen, and Marco Baroni. 2018. [Colorless green recurrent networks dream hierarchically](#).

John Hale. 2001. [A probabilistic Earley parser as a psycholinguistic model](#). In *Second Meeting of the North American Chapter of the Association for Computational Linguistics*.

Betty Hart and Todd R Risley. 1995. *Meaningful differences in the everyday experience of young American children*. Paul H Brookes Publishing.

Kenneth Heafield. 2011. [KenLM: Faster and smaller language model queries](#). In *Proceedings of the Sixth Workshop on Statistical Machine Translation*, pages 187–197, Edinburgh, Scotland. Association for Computational Linguistics.

⁹<https://talkbank.org/share/irb/>

787	Sepp Hochreiter and Jürgen Schmidhuber. 1997.	Roger Levy. 2008. Expectation-based syntactic comprehension. <i>Cognition</i> , 106(3):1126–1177.	842
788	Long short-term memory. <i>Neural computation</i> ,		843
789	9(8):1735–1780.		
790	Jennifer Hu, Jon Gauthier, Peng Qian, Ethan Wilcox,	John Lewis and Jeffrey Elman. 2001. Learnability and	844
791	and Roger Levy. 2020. A systematic assessment	the statistical structure of language: Poverty of stim-	845
792	of syntactic generalization in neural language mod-	ulus arguments revisited. <i>Proceedings of the 26th</i>	846
793	els. In <i>Proceedings of the 58th Annual Meeting</i>	<i>Annual Boston University Conference on Language</i>	847
794	<i>of the Association for Computational Linguistics</i> ,	<i>Development</i> , 1.	848
795	pages 1725–1744, Online. Association for Compu-		
796	tational Linguistics.	Yongjie Lin, Yi Chern Tan, and Robert Frank. 2019.	849
797	Philip A. Huebner, Elior Sulem, Cynthia Fisher, and	Open sesame: Getting inside BERT’s linguistic	850
798	Dan Roth. 2021. BabyBERTa: Learning more gram-	knowledge. In <i>Proceedings of the 2019 ACL Work-</i>	851
799	mar with small-scale child-directed language. In	<i>shop BlackboxNLP: Analyzing and Interpreting Neu-</i>	852
800	<i>Proceedings of CoNLL</i> .	<i>ral Networks for NLP</i> , pages 241–253, Florence,	853
		Italy. Association for Computational Linguistics.	854
801	Xuân-Nga Cao Kam, Iglia Stoynevska, Lidiya Torny-	Tal Linzen. 2020. How can we accelerate progress to-	855
802	ova, Janet D Fodor, and William G Sakas. 2008. Bi-	wards human-like linguistic generalization? In <i>Pro-</i>	856
803	grams and the richness of the stimulus. <i>Cognitive</i>	<i>ceedings of the 58th Annual Meeting of the Asso-</i>	857
804	<i>Science</i> , 32(4):771–787.	<i>ciation for Computational Linguistics</i> , pages 5210–	858
		5217, Online. Association for Computational Lin-	859
805	Reinhard Kneser and Hermann Ney. 1995. Improved	guistics.	860
806	backing-off for m-gram language modeling. <i>1995</i>	Brian MacWhinney. 2000. <i>The CHILDES project:</i>	861
807	<i>International Conference on Acoustics, Speech, and</i>	<i>Tools for analyzing talk</i> . Lawrence Erlbaum Asso-	862
808	<i>Signal Processing</i> , 1:181–184 vol.1.	ciates.	863
809	Adhiguna Kuncoro, Chris Dyer, John Hale, Dani Yo-	R. Thomas McCoy, Robert Frank, and Tal Linzen.	864
810	gatama, Stephen Clark, and Phil Blunsom. 2018.	2018. Revisiting the poverty of the stimulus: hier-	865
811	LSTMs can learn syntax-sensitive dependencies	archical generalization without a hierarchical bias in	866
812	well, but modeling structure makes them better. In	recurrent neural networks.	867
813	<i>Proceedings of the 56th Annual Meeting of the As-</i>	R. Thomas McCoy, Robert Frank, and Tal Linzen.	868
814	<i>sociation for Computational Linguistics (Volume 1:</i>	2020. Does syntax need to grow on trees? sources of	869
815	<i>Long Papers)</i> , pages 1426–1436, Melbourne, Aus-	hierarchical inductive bias in sequence-to-sequence	870
816	tralia. Association for Computational Linguistics.	networks.	871
817	Marta Kutas, Katherine A DeLong, and Nathaniel J	James L. Morgan and Katherine Demuth. 1996. <i>Signal</i>	872
818	Smith. 2011. A look around at what lies ahead: Pre-	<i>to syntax: Bootstrapping from speech to grammar in</i>	873
819	dition and predictability in language processing. In	<i>early acquisition</i> . Psychology Press.	874
820	<i>Predictions in the brain: Using our past to generate</i>	Aaron Mueller, Robert Frank, Tal Linzen, Luheng	875
821	<i>a future</i> .	Wang, and Sebastian Schuster. 2022. Coloring the	876
822	Kushal Lakhota, Eugene Kharitonov, Wei-Ning Hsu,	blank slate: Pre-training imparts a hierarchical in-	877
823	Yossi Adi, Adam Polyak, Benjamin Bolte, Tu-Anh	ductive bias to sequence-to-sequence models. In	878
824	Nguyen, Jade Copet, Alexei Baevski, Abdelrahman	<i>Findings of the Association for Computational Lin-</i>	879
825	Mohamed, et al. 2021. On generative spoken lan-	<i>guistics: ACL 2022</i> , pages 1352–1368, Dublin, Ire-	880
826	guage modeling from raw audio. <i>Transactions of the</i>	land. Association for Computational Linguistics.	881
827	<i>Association for Computational Linguistics</i> , 9:1336–	Karl Mulligan, Robert Frank, and Tal Linzen. 2021.	882
828	1354.	Structure here, bias there: Hierarchical generaliza-	883
829	Howard Lasnik and Jeffrey L Lidz. 2017. The argu-	tion by jointly learning syntactic transformations. In	884
830	ment from the poverty of the stimulus. <i>The Oxford</i>	<i>Proceedings of the Society for Computation in Lin-</i>	885
831	<i>handbook of universal grammar</i> , pages 221–248.	<i>guistics 2021</i> , pages 125–135, Online. Association	886
832	Julie Anne Legate and Charles D Yang. 2002. Em-	for Computational Linguistics.	887
833	pirical re-assessment of stimulus poverty arguments.	Ludovica Pannitto and Aurélie Herbelot. 2020. Recur-	888
834	<i>The Linguistic Review</i> , 19(1-2):151–162.	rent babbling: evaluating the acquisition of gram-	889
835	Michael Lepori, Tal Linzen, and R. Thomas McCoy.	mar from limited input data. In <i>Proceedings of</i>	890
836	2020. Representations of syntax [MASK] useful:	<i>the 24th Conference on Computational Natural Lan-</i>	891
837	Effects of constituency and dependency structure in	<i>guage Learning</i> , pages 165–176, Online. Associa-	892
838	recursive LSTMs. In <i>Proceedings of the 58th An-</i>	tion for Computational Linguistics.	893
839	<i>annual Meeting of the Association for Computational</i>	Adam Pauls and Dan Klein. 2012. Large-scale syntac-	894
840	<i>Linguistics</i> , pages 3306–3316, Online. Association	tic language modeling with treelets. In <i>Proceedings</i>	895
841	for Computational Linguistics.	<i>of the 50th Annual Meeting of the Association for</i>	896

897		<i>Computational Linguistics (Volume 1: Long Papers)</i> , pages 959–968, Jeju Island, Korea. Association for Computational Linguistics.	
898			
899			
900	Lisa Pearl. 2021. Poverty of the stimulus without tears.		
901		<i>Language Learning and Development</i> , pages 1–40.	
902	Lisa Pearl and Benjamin Mis. 2016. The role of indirect positive evidence in syntactic acquisition: A look at anaphoric one.		
903		<i>Language</i> , 92:1–30.	
904			
905	Lisa Pearl and Jon Sprouse. 2013a. Computational models of acquisition for islands.		
906		<i>Experimental syntax and islands effects</i> , pages 109–131.	
907			
908	Lisa Pearl and Jon Sprouse. 2013b. Syntactic islands and learning biases: Combining experimental syntax and computational modeling to investigate the language acquisition problem.		
909		<i>Language Acquisition</i> , 20(1):23–68.	
910			
911			
912			
913	Andrew Perfors, Josh Tenenbaum, and Terry Regier. 2011. The learnability of abstract syntactic principles.		
914		<i>Cognition</i> , 118:306–338.	
915			
916	Jackson Petty and Robert Frank. 2021. Transformers generalize linearly.		
917		<i>arXiv preprint arXiv:2109.12036</i> .	
918			
919	Geoffrey K. Pullum and Barbara C. Scholz. 2002. Empirical assessment of stimulus poverty arguments.		
920		<i>The Linguistic Review</i> , 18(1-2):9–50.	
921			
922	Florencia Reali and Morten H. Christiansen. 2005. Uncovering the richness of the stimulus: Structure dependence and indirect statistical evidence.		
923		<i>Cognitive Science</i> , 29(6):1007–1028.	
924			
925			
926	Haoyue Shi, Jiayuan Mao, Kevin Gimpel, and Karen Livescu. 2019. Visually grounded neural syntax acquisition.		
927		In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 1842–1861, Florence, Italy. Association for Computational Linguistics.	
928			
929			
930			
931			
932	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need.		
933		In <i>Advances in neural information processing systems</i> , pages 5998–6008.	
934			
935			
936			
937	Alex Warstadt and Samuel R Bowman. 2020. Can neural networks acquire a structural bias from raw linguistic data?		
938		<i>Proceedings of the 42nd Annual Conference of the Cognitive Science Society</i> .	
939			
940			
941	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.		
942		<i>Transactions of the Association for Computational Linguistics</i> , 8:377–392.	
943			
944			
945			
946	Ethan Wilcox, Richard Futrell, and Roger Levy. 2021. Using computational models to test syntactic learnability.		
947		<i>lingbuzz preprint lingbuzz/006327</i> .	
948			
	Ethan Wilcox, Roger Levy, Takashi Morita, and Richard Futrell. 2018. What do RNN language models learn about filler–gap dependencies?		
		In <i>Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP</i> , pages 211–221, Brussels, Belgium. Association for Computational Linguistics.	
	Taha Yasseri, András Kornai, and János Kertész. 2012. A practical approach to language complexity: A Wikipedia case study.		
		<i>PLoS ONE</i> , 7(11):e48386.	
	A CHILDES preprocessing details		
	The train, test, and validation split kept each document in the corpora intact to allow for learning of context. Since a document roughly correspond to a single recording session, and the sentence order within each document was not randomized, the networks could utilize cross sentence context while predicting the next word.		
	Generally, we kept the data as close to the actual input that the child receives as possible. However, in some cases we modified tokenization to match CHILDES Treebank, a syntactically parsed subset of the CHILDES corpora. For instance, contractions were split, e.g. we replaced <i>don't</i> with <i>do n't</i> ,		
	The ages of the children vary by corpus, ranging from six months to twelve years. Almost 95% (49/52) of the corpora consist of transcriptions with children between one and six years of age.		
	Note that for Experiment 2, we used the same vocabulary as we used in Experiment 1, which means that the words that were not present in the Experiment 1's vocabulary were replaced with <unk> tokens.		
	The unprocessed CHILDES datasets were downloaded in XML format from the online XML version of the CHILDES database (MacWhinney, 2000). A modified NLTK CHILDESCorpusReader was used to parse the XML into plain text for training.		
	The CHILDES dataset is licensed for use under a CC BY-NC-SA 3.0 license (https://talkbank.org/share/rules.html). Under the terms of this license, the data can be freely used and adapted, as long as it is not used for commercial purposes and as long as attribution (https://creativecommons.org/licenses/by-nc-sa/3.0/) is provided. Our usage fits these criteria.		
	Though CHILDES contains many corpora of many languages, we use only corpora from the		

North American English subset of CHILDES, which contains child-directed speech with many different North American children. See the [CHILDES database](#) for more details.

By the [CHILDES rules for data citation](#), research that relies on more than 6 of the corpora need only cite the overall database, not each individual corpus.

All the data on CHILDES must [adhere to IRB guidelines](#), including a requirement for anonymity.

The final dataset may be downloaded from [LINK ANONYMIZED]. This dataset is not intended for commercial use.

CHILDES corpora included Bates, Bernstein, Bliss, Bloom70, Bloom73, Bohannon, Braunwald, Brent, Brown, Carterette, Clark, Cornell, Demetras1, Demetras2, EllisWeismer, Evans, Feldman, Garvey, Gathercole, Gelman, Gillam, Gleason, HSLLD, Haggerty, Hall, Higginson, Kuczaj, MacWhinney, McCune, McMillan, Morisset, NH, Nelson, NewEngland, NewmanRatner, Normal, POLER, Peters, Post, Rollins, Sachs, Sawyer, Snow, Soderstrom, Sprott, Suppes, Tardif, Valian, VanHouten, VanKleeck, Warren, Weist.

B Hyperparameter Search and Model Implementation

B.1 Hyperparameter search

LSTMs For LSTMs we explored the following hyper-parameters via a grid search for a total of 144 models.

1. layers: 2
2. hidden and embedding size: 200, 800
3. batch size: 20, 80
4. dropout rate: 0.0, 0.2, 0.4, 0.6
5. learning rate: 5.0, 10.0, 20.0
6. random seed: 3 per parameter combination, unique for each LSTM

The LSTM model with the lowest perplexity on the validation set after training had 2 layers, a hidden and embedding size of 800, a batch size of 20, a dropout rate of 0.4, and a learning rate of 10.¹⁰ A LSTM model with these hyperparameters has 37,620,294 parameters.

¹⁰The hyperparameters we explored for the LSTMs were those of [Gulordava et al. \(2018\)](#), the code for which can be found at <https://github.com/facebookresearch/colorlessgreenRNNs>

LSTMs	prepose first	prepose main
delete first	0.01072	0.14408
delete main	0.38672	0.11982
delete none	0.20099	0.13767

Table 1: Analysis of models’ preference for questions consistent with combinations of ‘prepose’ and ‘delete’ rules. Within each architecture, the proportion preferences across all 6 question types necessarily sum to 1.

Transformers For the Transformers we performed a hyperparameter sweep over the following hyper-parameters for a total of 84 models.

1. layers: 2, 4, 8, 16
2. context size: 50, 100, 500
3. hidden and embedding size: 200, 800, 1600
4. heads: 2, 4, 8, 16
5. batch size: 20, 80, 160
6. dropout rate: 0.0, 0.2, 0.4, 0.6
7. learning rate: 0.5, 1.0, 5.0, 10.0, 20.0
8. random seed: 3 per parameter combination

The Transformer model with the lowest perplexities after training had 4 layers, a context size of 500, a hidden size of 800, a batch size of 10, 4 heads, a dropout rate of 0.2, and a learning rate of 5.0. A Transformer model with these parameters has 42,759,494 parameters.

B.2 Implementation

All models were implemented in PyTorch by building on code from [here](#) and [here](#), and trained using Nvidia k80 GPUs. The final models may be downloaded from [LINK ANONYMIZED]. These models are not intended for commercial use.

C PREPOSE-ONE&DELETE-ONE Full Results

See Table 1 and Table 2 for these results.

C.1 Results using SLOR

See Table 3 and Table 4 for these results.

D BabyBERTa dataset evaluation

For an illustrative subset of the results on the Zorro evaluation dataset (discussed in Section 6), see Figure 4. For the full results, see Figure 5.

Transformers	prepose first	prepose main
delete first	0.00662	0.15964
delete main	0.31436	0.06482
delete none	0.24538	0.20918

Table 2: Analysis of models’ preference for questions consistent with combinations of ‘prepose’ and ‘delete’ rules. Within each architecture, the proportion preferences across all 6 question types necessarily sum to 1.

LSTMs	Prepose First	Prepose Main
Delete First	0%	14%
Delete Main	33%	8%
Delete None	26%	18%

Table 3: Analysis of LSTMs’ preference for questions consistent with combinations of ‘prepose’ and ‘delete’ rules, evaluated using SLOR. Within each architecture, the proportion preferences across all 6 question types necessarily sum to 1.

Transformers	Prepose First	Prepose Main
Delete First	0%	15%
Delete Main	27%	4%
Delete None	29%	24%

Table 4: Analysis of Transformers’ preference for questions consistent with combinations of ‘prepose’ and ‘delete’ rules, evaluated using SLOR. Within each architecture, the proportion preferences across all 6 question types necessarily sum to 1.

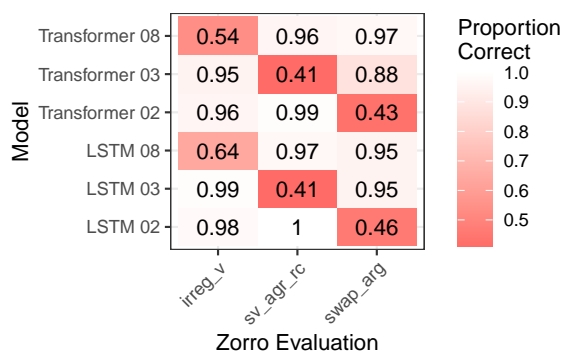


Figure 4: The performance of a selected subset of model re-runs on a selected subset of the Zorro evaluations. Each Zorro evaluation targets a specific syntactic phenomenon—in the cases shown here, irregular verbs, subject-verb agreement across relative clauses, and correct argument ordering.

E Move-One Dataset Results

One approach used in several past papers (e.g., Lewis and Elman (2001) and Real and Christiansen (2005)) evaluate models using pairs of sen-

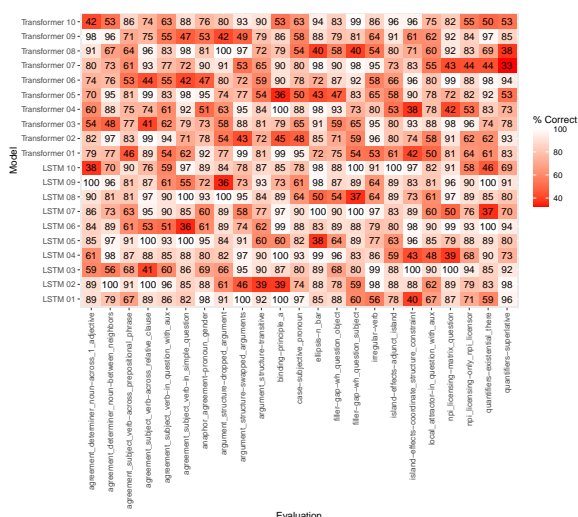


Figure 5: Results on the targeted syntactic evaluations in Huebner et al. (2021) in percent accuracy. Evaluation names in Figure 4 were shortened.

tences that can be formed by starting with a declarative sentence (e.g., (14)) and moving one of its auxiliaries to the front of the sentence. The first sentence in each pair (e.g., (15a)) follows HIERARCHICALQ, because the *main* auxiliary is moved, while the second (e.g., (15b)), follows LINEARQ because the *first* auxiliary is moved.

- (14) The children who **are** talking **are** sleeping. 1079
- (15) a. **Are** the children who **are** talking sleeping? 1080
- b. **Are** the children who talking **are** sleeping? 1081

If a model assigns a higher probability to (15a) than (15b), that is evidence that the models favors HIERARCHICALQ over LINEARQ. While this preference is a necessary component of correctly learning HIERARCHICALQ, it is by no means sufficient: indeed, Kam et al. (2008) showed that models can prefer sentences consistent with HIERARCHICALQ over sentences consistent with LINEARQ due to shallow *n*-gram statistics rather than due to knowledge of hierarchical structure. More generally, there are infinitely many other incorrect hypotheses besides LINEARQ, and demonstrating successful learning of HIERARCHICALQ would require ruling out all of them. Investigating all possibilities is intractable, but we can at least investigate a few additional plausible ones. Thus, in the main paper we depart from prior work by considering a greater number of candidate sentences than just the pairs of sentences used in prior work.

To create the MOVE-ONE dataset, we randomly sampled 10,000 declarative sentences from

1110 our CFGs for which the first and main auxiliary
1111 were identical and then modified them to give
1112 10,000 sentence pairs. To create the PREPOSE-
1113 ONE&DELETE-ONE dataset, we randomly sam-
1114 pled a different 10,000 declarative sentences from
1115 our CFGs for which the first and main auxiliary
1116 were different and then we modified them to give
1117 10,000 6-tuples of sentences. See Appendix F for
1118 more details about the CFGs.

1119 **F Context Free Grammars**

1120 The context free grammars used to generate the
1121 evaluation datasets appear in Figure 7, Figure 6 ,
1122 Figure 8, and Figure 9.

1123 **G Breakdown by lexical identity**

1124 Here we further break down models' predictions
1125 for the $\text{FIRST-AUX} \neq \text{MAIN-AUX}$ evaluation set
1126 based on the identities of the two auxiliaries in the
1127 input sentence. Figure 11 gives the results for the
1128 LSTM in the NEXT-WORD PREDICTION + QUES-
1129 TION FORMATION condition; Figure 10 for the
1130 LSTM in the QUESTION FORMATION condition;
1131 Figure 13 for the Transformer in the NEXT-WORD
1132 PREDICTION + QUESTION FORMATION condi-
1133 tion; and Figure 12 for the for the Transformer in
1134 the QUESTION FORMATION condition.

1135 **H Example generated text**

1136 Figure 14 gives some example text generated by our
1137 models. Models trained on next-word prediction
1138 produce their predictions as a probability distribu-
1139 tion over the vocabulary. To use such models to
1140 generate text, we sample a word from this distribu-
1141 tion then use that word as the model's input for the
1142 next time step.

Det_S	→ {the some this }
Det_P	→ {the some those}
N_S	→ {baby girl boy animal child person horse }
N_P	→ {babies girls boys animals children people horses }
IV	→ {play read draw sit fall talk sleep try work walk}
IV_IS	→ {playing reading drawing sitting falling talking sleeping trying working walking}
IV_HAS	→ {played read drawn sat fallen talked slept tried worked walked}
TV	→ {call see find help feed know pick visit watch reach}
TV_IS	→ {calling seeing finding helping feeding knowing picking visiting watching reaching}
TV_HAS	→ {called seen found helped fed known picked visited watched reached}
Aux_P	→ {do did can would shall}
Aux_S	→ {does did can would shall}
Aux_S_BE	→ {is was}
Aux_P_BE	→ {are were}
Aux_S_HAS	→ {has}
Aux_P_HAS	→ {have}
Prep	→ {by behind }
Rel	→ {who that }

Figure 6: Vocabulary used for the PREPOSE-ONE-AND-DELETE-ONE, FIRST-AUX \neq MAIN-AUX, and FIRST-AUX = MAIN-AUX evaluation datasets

S → {NP_M_S VP_M_S | NP_M_P VP_M_P}
NP_M_S → {Det_S N_S | Det_S N_S Prep Det_S N_S | Det_S N_S Prep Det_P N_P}
NP_M_P → {Det_P N_P | Det_P N_P Prep Det_S N_S | Det_P N_P Prep Det_P N_P}
NP_O → {Det_S N_S | Det_P N_P | Det_S N_S Prep Det_S N_S | Det_S N_S Prep
Det_P N_P | Det_P N_P Prep Det_S N_S | Det_P N_P Prep Det_P N_P | Det_S
N_S RC_S | Det_P N_P RC_P }
VP_M_S → {Aux_S IV }
VP_M_S → {Aux_S TV NP_O}
VP_M_S → {Aux_S_BE IV_IS}
VP_M_S → {Aux_S_BE TV_IS NP_O}
VP_M_S → {Aux_S_HAS IV_HAS}
VP_M_S → {Aux_S_HAS TV_HAS NP_O}
VP_M_P → {Aux_P IV}
VP_M_P → {Aux_P TV NP_O}
VP_M_P → {Aux_P_BE IV_IS}
VP_M_P → {Aux_P_BE TV_IS NP_O}
VP_M_P → {Aux_P_HAS IV_HAS}
VP_M_P → {Aux_P_HAS TV_HAS NP_O}
RC_S → {Rel Aux_S IV | Rel Det_S N_S Aux_S TV | Rel Det_P N_P Aux_P TV |
Rel Aux_S TV Det_S N_S | Rel Aux_S TV Det_P N_P}
RC_S → {Rel Aux_S_BE IV_IS | Rel Det_S N_S Aux_S_BE TV_IS | Rel Det_P
N_P Aux_P_BE TV_IS | Rel Aux_S_BE TV_IS Det_S N_S | Rel Aux_S_BE
TV_IS Det_P N_P}
RC_S → {Rel Aux_S_HAS IV_HAS | Rel Det_S N_S Aux_S_HAS TV_HAS | Rel
Det_P N_P Aux_P_HAS TV_HAS | Rel Aux_S_HAS TV_HAS Det_S N_S |
Rel Aux_S_HAS TV_HAS Det_P N_P}
RC_P → {Rel Aux_P IV | Rel Det_S N_S Aux_S TV | Rel Det_P N_P Aux_P TV |
Rel Aux_P TV Det_S N_S | Rel Aux_P TV Det_P N_P}
RC_P → {Rel Aux_P_BE IV_IS | Rel Det_S N_S Aux_S_BE TV_IS | Rel Det_P
N_P Aux_P_BE TV_IS | Rel Aux_P_BE TV_IS Det_S N_S | Rel Aux_P_BE
TV_IS Det_P N_P}
RC_P → {Rel Aux_P_HAS IV_HAS | Rel Det_S N_S Aux_S_HAS TV_HAS | Rel
Det_P N_P Aux_P_HAS TV_HAS | Rel Aux_P_HAS TV_HAS Det_S N_S |
Rel Aux_P_HAS TV_HAS Det_P N_P}

Figure 7: CFG used to generate FIRST-AUX = MAIN-AUX evaluation dataset

S → {NP_M_S VP_M_S | NP_M_P VP_M_P}
NP_M_S → {Det_S N_S | Det_S N_S Prep Det_S N_S | Det_S N_S Prep Det_P N_P}
NP_M_P → {Det_P N_P | Det_P N_P Prep Det_S N_S | Det_P N_P Prep Det_P N_P}
NP_O → {Det_S N_S | Det_P N_P | Det_S N_S Prep Det_S N_S | Det_S N_S Prep
Det_P N_P | Det_P N_P Prep Det_S N_S | Det_P N_P Prep Det_P N_P | Det_S
N_S RC_S | Det_P N_P RC_P }
VP_M_S → {Aux_S IV }
VP_M_S → {Aux_S TV NP_O}
VP_M_S → {Aux_S_BE IV_IS}
VP_M_S → {Aux_S_BE TV_IS NP_O}
VP_M_S → {Aux_S_HAS IV_HAS}
VP_M_S → {Aux_S_HAS TV_HAS NP_O}
VP_M_P → {Aux_P IV}
VP_M_P → {Aux_P TV NP_O}
VP_M_P → {Aux_P_BE IV_IS}
VP_M_P → {Aux_P_BE TV_IS NP_O}
VP_M_P → {Aux_P_HAS IV_HAS}
VP_M_P → {Aux_P_HAS TV_HAS NP_O}
RC_S → {Rel Aux_S IV | Rel Det_S N_S Aux_S TV | Rel Det_P N_P Aux_P TV |
Rel Aux_S TV Det_S N_S | Rel Aux_S TV Det_P N_P}
RC_S → {Rel Aux_S_BE IV_IS | Rel Det_S N_S Aux_S_BE TV_IS | Rel Det_P
N_P Aux_P_BE TV_IS | Rel Aux_S_BE TV_IS Det_S N_S | Rel Aux_S_BE
TV_IS Det_P N_P}
RC_S → {Rel Aux_S_HAS IV_HAS | Rel Det_S N_S Aux_S_HAS TV_HAS | Rel
Det_P N_P Aux_P_HAS TV_HAS | Rel Aux_S_HAS TV_HAS Det_S N_S |
Rel Aux_S_HAS TV_HAS Det_P N_P}
RC_P → {Rel Aux_P IV | Rel Det_S N_S Aux_S TV | Rel Det_P N_P Aux_P TV |
Rel Aux_P TV Det_S N_S | Rel Aux_P TV Det_P N_P}
RC_P → {Rel Aux_P_BE IV_IS | Rel Det_S N_S Aux_S_BE TV_IS | Rel Det_P
N_P Aux_P_BE TV_IS | Rel Aux_P_BE TV_IS Det_S N_S | Rel Aux_P_BE
TV_IS Det_P N_P}
RC_P → {Rel Aux_P_HAS IV_HAS | Rel Det_S N_S Aux_S_HAS TV_HAS | Rel
Det_P N_P Aux_P_HAS TV_HAS | Rel Aux_P_HAS TV_HAS Det_S N_S |
Rel Aux_P_HAS TV_HAS Det_P N_P}

Figure 8: CFG used to generate FIRST-AUX \neq MAIN-AUX evaluation dataset

S	→ {NP_S RC_S_BARE MAIN-AUX VP_S_PAST}
S	→ {NP_S RC_S_PAST MAIN-AUX VP_S_BARE}
S	→ {NP_S RC_S_BARE MAIN-AUX VP_S_PROG}
S	→ {NP_S RC_S_PROG MAIN-AUX VP_S_BARE}
S	→ {NP_S RC_S_PAST MAIN-AUX VP_S_PROG}
S	→ {NP_S RC_S_PROG MAIN-AUX VP_S_PAST}
S	→ {NP_P RC_P_BARE MAIN-AUX VP_P_PAST}
S	→ {NP_P RC_P_PAST MAIN-AUX VP_P_BARE}
S	→ {NP_P RC_P_BARE MAIN-AUX VP_P_PROG}
S	→ {NP_P RC_P_PROG MAIN-AUX VP_P_BARE}
S	→ {NP_P RC_P_PAST MAIN-AUX VP_P_PROG}
S	→ {NP_P RC_P_PROG MAIN-AUX VP_P_PAST}
NP_S	→ {Det_S N_S}
NP_P	→ {Det_P N_P}
NP_O	→ {Det_S N_S Det_P N_P Det_S N_S Prep Det_S N_S Det_S N_S Prep Det_P N_P Det_P N_P Prep Det_S N_S Det_P N_P Prep Det_P N_P}
VP_S_BARE	→ {Aux_S IV }
VP_S_BARE	→ {Aux_S TV NP_O}
VP_S_PROG	→ {Aux_S_BE IV_IS}
VP_S_PROG	→ {Aux_S_BE TV_IS NP_O}
VP_S_PAST	→ {Aux_S_HAS IV_HAS}
VP_S_PAST	→ {Aux_S_HAS TV_HAS NP_O}
VP_P_BARE	→ {Aux_P IV }
VP_P_BARE	→ {Aux_P TV NP_O}
VP_P_PROG	→ {Aux_P_BE IV_IS}
VP_P_PROG	→ {Aux_P_BE TV_IS NP_O}
VP_P_PAST	→ {Aux_P_HAS IV_HAS}
VP_P_PAST	→ {Aux_P_HAS TV_HAS NP_O}
RC_S_BARE	→ {Rel Aux_S IV Rel Det_S N_S Aux_S TV Rel Det_P N_P Aux_P TV Rel Aux_S TV Det_S N_S Rel Aux_S TV Det_P N_P}
RC_S_PROG	→ {Rel Aux_S_BE IV_IS Rel Det_S N_S Aux_S_BE TV_IS Rel Det_P N_P Aux_P_BE TV_IS Rel Aux_S_BE TV_IS Det_S N_S Rel Aux_S_BE TV_IS Det_P N_P}
RC_S_PAST	→ {Rel Aux_S_HAS IV_HAS Rel Det_S N_S Aux_S_HAS TV_HAS Rel Det_P N_P Aux_P_HAS TV_HAS Rel Aux_S_HAS TV_HAS Det_S N_S Rel Aux_S_HAS TV_HAS Det_P N_P}
RC_P_BARE	→ {Rel Aux_P IV Rel Det_S N_S Aux_S TV Rel Det_P N_P Aux_P TV Rel Aux_P TV Det_S N_S Rel Aux_P TV Det_P N_P}
RC_P_PROG	→ {Rel Aux_P_BE IV_IS Rel Det_S N_S Aux_S_BE TV_IS Rel Det_P N_P Aux_P_BE TV_IS Rel Aux_P_BE TV_IS Det_S N_S Rel Aux_P_BE TV_IS Det_P N_P}
RC_P_PAST	→ {Rel Aux_P_HAS IV_HAS Rel Det_S N_S Aux_S_HAS TV_HAS Rel Det_P N_P Aux_P_HAS TV_HAS Rel Aux_P_HAS TV_HAS Det_S N_S Rel Aux_P_HAS TV_HAS Det_P N_P}

Figure 9: CFG used to generate PREPOSE-ONE-AND-MOVE-ONE evaluation dataset



Figure 10: Breakdown by the identities of the two auxiliaries for outputs in the $\text{FIRST-AUX} \neq \text{MAIN-AUX}$ evaluation set for LSTMs first trained on next-word prediction and then question formation.

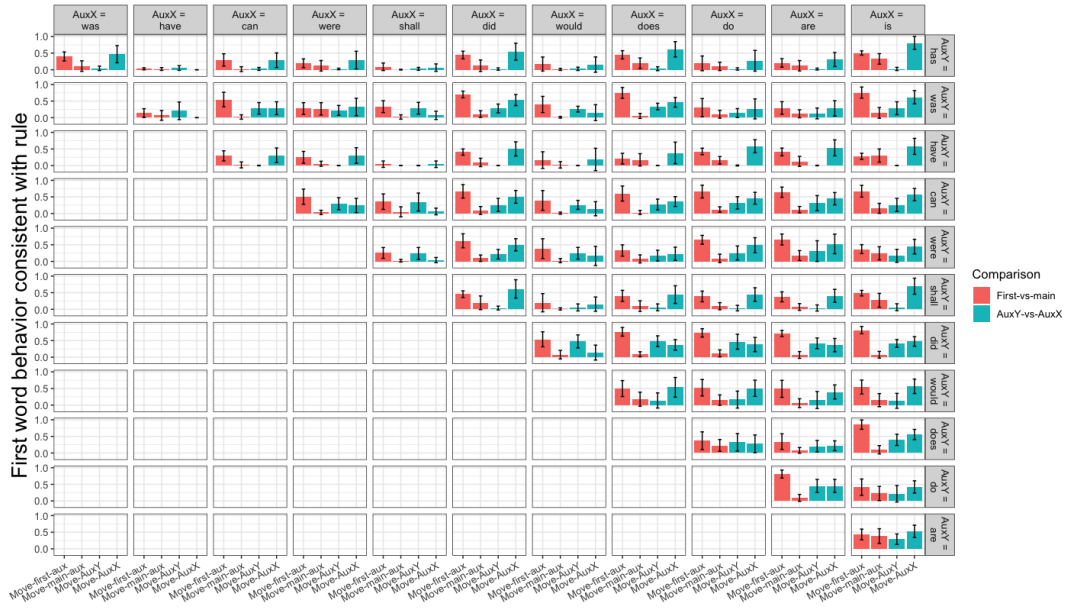


Figure 11: Breakdown by the identities of the two auxiliaries for outputs in the $\text{FIRST-AUX} \neq \text{MAIN-AUX}$ evaluation set for LSTMs trained only on question formation.

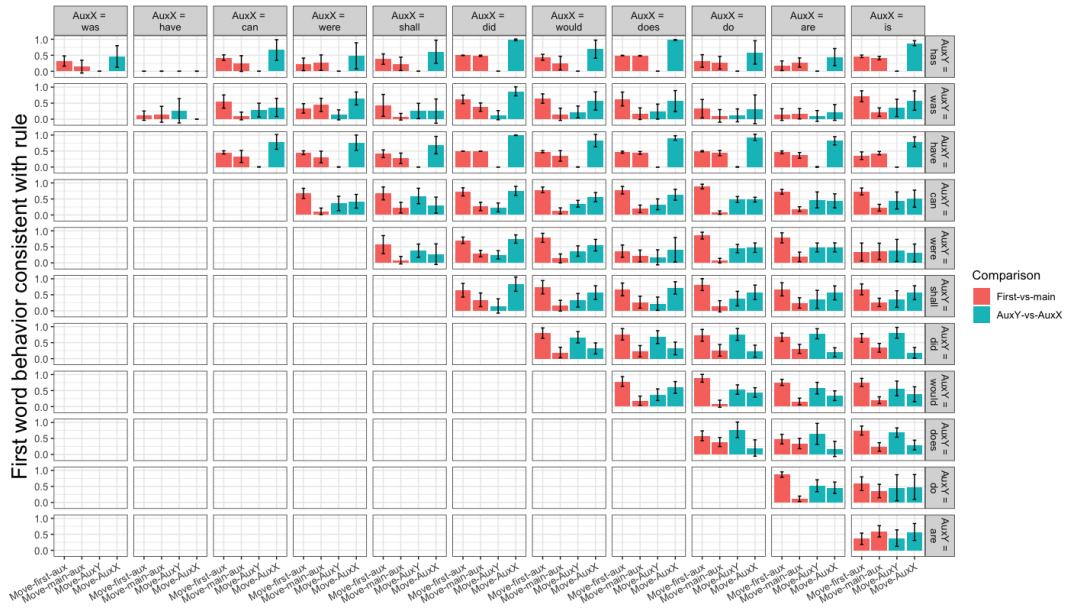


Figure 12: Breakdown by the identities of the two auxiliaries for outputs in the $\text{FIRST-AUX} \neq \text{MAIN-AUX}$ evaluation set for Transformers first trained on next-word prediction and then question formation.

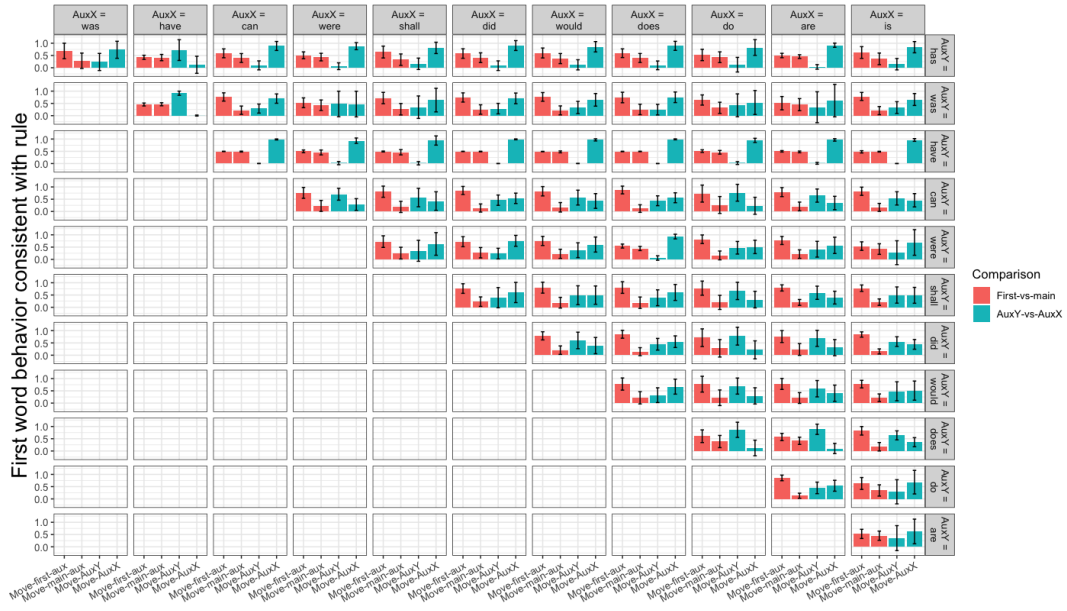


Figure 13: Breakdown by the identities of the two auxiliaries for outputs in the $\text{FIRST-AUX} \neq \text{MAIN-AUX}$ evaluation set for Transformers trained only on question formation.

Text from CHILDES	do you eat in nursery school too ? oh . it 's nursery school . xxx . what do you get to eat ? what else do you get to eat besides a book ? do n't play with those books . what else do you have ? what do you have in nursery school ? alright - you can take it . why do n't you put the turtle back . well - you 're not going to change his water now . tonight . tonight - you change it about ten times a day . no - you can do it later . forget it now . too bad . too bad . that poor turtle . you do n't let him sit there five minutes . that was five minutes ?
LSTM generated text	<u>do you eat in nursery school too ? oh . it 's nursery school .</u> no ! tomorrow 's Friday . right . yeah - yeah . yeah I know what . I told you that he liked snipsnap that way . well this is address the grayish kings . ok . what kind of bird do you want to play with ? what else do you want to do ? your crayons . got a circle . I 'm going to use these . there . remember ? ok do you want to put all of these there ? ok let 's put everything away . what do you want to do ? you want to play in the box ? it 's just like a ball because it has a ball on it .
Transformer generated text	<u>do you eat in nursery school too ? oh . it 's nursery school .</u> mm-hm . and it looks like a little girl . mm-hm . like a little girl . hm , mm-hm . looks like she 's concerned about baseball this is a picture of a double football player . mm-hm . um , it has a uniform on it . mm-hm . anything more ? no . ok great . got a hat on . and a flower and a little flower . remember that ? do you remember that word ? what do you remember that word ? do you remember what it word says ? it says seven - eight . what 's that word ? it says eight . look at that word . three - four - six ... eight - nine ...

Figure 14: Comparison of text generated by the LSTM and Transformer models with a block of text chosen randomly from the training data. The LSTMs and Transformers were both seeded with the first three sentences of the text taken from CHILDES, which is the underlined in the two model generated texts. Note that neither of the model generated text were cherry picked either for quality or to be representative of the models' usual output: rather they were the first things they generated when seeded with the above underlined portion.